# EXERCISES FOR STATISTICAL METHODS FOR CLIMATE SCIENTISTS

Timothy DelSole and Michael K. Tippett

## CONTENTS

1	Exercises: Basic Concepts in Probability and Statistics	1
2	Exercises for Statistical Inference	5
3	Exercises for Confidence Intervals	9
4	Exercises for Statistical Tests Based on Ranks	13
5	Exercises for Stochastic Processes	17
6	Exercises on the Power Spectrum	21
7	Exercises for Vectors, Matrices, and Geometry	25
8	Exercises for Linear Regression: Least Squares Estimation	29
9	Exercises for Linear Regression: Inference	35
		iii

#### iv contents

10	Exercises for Model Selection	39
11	Pitfalls of Statistical Inference	45
12	Exercises for Principal Component Analysis	47
13	Computational Exercises on Field Significance	53
14	Exercises for Multivariate Linear Regression	59
15	Exercises for Canonical Correlation Analysis	63
16	Exercises: Covariance Discriminant Analysis	67
17	Analysis of Variance	71
18	Predictable Component Analysis	77
19	Data Assimilation	81
20	Ensemble Square Root Filters	87
21	Exercises: Extreme Value Analysis	89
22	Answers to Exercises	91

#### Foreward

This document contains exercises and solutions to accompany *Statistical Methods for Climate Scientists* by Timothy DelSole and Michael K. Tippett. The usefulness of these solutions would be greatly diminished if they were posted on an insecure website. Therefore, we hope that instructors are mindful to not post these solutions on an insecure website. We will greatly appreciate users of this document notifying us about errors or ambiguities in these exercises and solutions.

v

# EXERCISES: BASIC CONCEPTS IN PROBABILITY AND STATISTICS

 Table 1.1
 Sample Outcome of 100 pairs of coin tosses

	$H_2$	$T_2$
$H_1$	30	21
$T_1$	22	27

**Exercise 1.1** (Empirical Probabilities). The *relative frequency* of an event is the fraction of times that the event occurs. Suppose a coin is tossed 200 times. The outcome of 100 consecutive pairs of tosses is summarized in the contingency table 1.1, where  $H_1$  and  $T_1$  indicate the number of heads and tails on the first toss, respectively, and  $H_2$  and  $T_2$  indicate the number of heads and tails on the second toss. For instance, the "21" listed in the top right entry of the table means that 21 out of 100 pairs was "heads-tails." The empirical probability of "heads-tails" is therefore 21/100.

- a) What is the empirical probability of heads on the first toss of a pair, denoted  $P(H_1)$ ?
- b) What is the empirical probability of a heads in the second toss of a pair, denoted  $P(H_2)$ ?
- c) What is the empirical probability of the joint event  $H_1$  and  $T_2$ , denoted  $P(H_1, T_2)$ ?

- 2 EXERCISES: BASIC CONCEPTS IN PROBABILITY AND STATISTICS
  - d) What is the empirical probability of  $T_2$  given  $H_1$ , denoted  $P(T_2|H_1)$ ? Compute this by examining only those cases in which  $H_1$  is true.
  - e) Verify that the empirical conditional probability computed in (d) is the ratio of the joint probability over the marginal:  $P(T_2|H_1) = P(T_2, H_1)/P(H_1)$ .

(Incidentally, this exercise suggests why unconditional probabilities are called marginal probabilities, namely because they can be calculated by summing values in a table along rows and columns and writing the results along the margins of the table.)

**Exercise 1.2.** If X and Y are independent, then show that p(x,y) = p(x)p(y) implies E[XY] = E[X]E[Y].

**Exercise 1.3.** Let X and Y be independent random variables with respective means  $\mu_X$  and  $\mu_Y$ , and respective variances  $\sigma_X^2$  and  $\sigma_Y^2$ . In terms of the population means and variances:

- (a) Compute E[XY]?
- (b) Compute cov[X, Y]?
- (c) Compute var[X Y]?
- (d) Compute var[XY]?

**Exercise 1.4** (Properties of Variance). Let k be a constant and X be a random variable with expectation  $\mu_X$  and variance  $\sigma_X^2$ . Using the definition of variance and the properties of expectation, prove the following:

- (a) var[k] = 0
- (b)  $var[kX] = k^2 var[X]?$
- (c)  $\operatorname{var}[k + X] = \operatorname{var}[X]$ ?

Exercise 1.5. Using the properties of expectation, show that

$$cov[X, Y] = E[XY] - E[X]E[Y].$$
 (1.1)

**Exercise 1.6.** If X is a random variable and k is a constant (i.e., independent of X), then show that  $E[(X - k)^2]$  is minimized when k = E[X]. (Hint: there are at least two ways to prove this: one way is to use calculus, the other way is to add and subtract an expectation. Can you find both proofs?) Be sure to prove that the solution is actually a *minimum*.

**Exercise 1.7** (Bounds on the Correlation Coefficient). Prove that the correlation coefficient between any two random variables X and Y is always between -1 and 1. Explain when the correlation is exactly equal to one, and when the correlation is exactly equal to -1, in terms of the relation between X and Y. Hint: use the fact that

$$E[(t(X - E[X]) + (Y - E[Y]))^2] = t^2 \operatorname{var}[X] + 2t \operatorname{cov}[X, Y] + \operatorname{var}[Y] \ge 0, \quad (1.2)$$

**Exercise 1.8** (Distribution of a Difference in Means). Suppose  $X_1, X_2, \ldots, X_N$  are independent random random variables drawn from a Gaussian distribution with mean  $\mu_X$  and variance  $\sigma^2$ . Similarly, let  $Y_1, Y_2, \ldots, Y_N$  be identically distributed random variables drawn from a Gaussian distribution with different mean  $\mu_Y$  but common variance  $\sigma^2$ . Also, assume that the X's and Y's are independent.

- What are the distributions of the sample means μ̂<sub>X</sub> and μ̂<sub>Y</sub>? (Do not just say "normal." Tell me all the parameters– e.g., the mean and variance.)
- What is the distribution of the difference in sample means  $\hat{\mu}_X \hat{\mu}_Y$ ?
- What is the distribution of the sample variances  $\hat{\sigma}_X^2$  and  $\hat{\sigma}_Y^2$ ? (Tell me all parameters.)
- What is the distribution of the sum of sample variances  $\hat{\sigma}_X^2 + \hat{\sigma}_Y^2$ ?

**Exercise 1.9** (Expectation of a Roll of a Die). If X is the outcome of the roll of a fair die, what is E[X]? What is var[X]?

**Exercise 1.10** (Expectations on Dice). Let X be the *sum* of the outcome of the roll of three fair dice. What is E[X]? What is var[X]? Assume the outcome of any individual die is independent of that of the others.

**Exercise 1.11** (Prove the Sample Variance is Unbiased). Let  $X_1, \ldots, X_N$  be independent random samples drawn from a population with expectation  $\mu$  and variance  $\sigma^2$ . Show that the expectation of the sample variance  $\hat{\sigma}_X^2$  is

$$E[\hat{\sigma}_X^2] = \sigma^2 \tag{1.3}$$

For this reason,  $\hat{\sigma}_X^2$  is called an *unbiased* estimate of  $\sigma_X^2$ . (This question is hard)

**Exercise 1.12.** Suppose you are predicting whether Y will be above or below normal. If "normal" is defined as the median, then, by definition, Y will be above normal 50% of the time (like a coin flip). Suppose, however, you are predicting Y given knowledge of another variable X. If X and Y are independent, then knowledge of X tells you nothing new about Y; i.e., the probability that Y is above normal is still 50%. However, if X and Y have a non-zero correlation, then knowledge of one tells you something new about the other. Suppose that the correlation between X and Y is 0.5. An example of such a sample is shown in fig. 1.3c. If X is above normal, what is the probability that Y is above normal? To answer this question, use example 1.7 to generate random numbers with a population correlation of  $\rho = 0.5$ , and then from these numbers compute the probability that Y > 0 given that X > 0. You should generate enough random samples to obtain *robust* estimates of the probability (i.e., ensure your answer does not change if you re-run your code with different random numbers).

**Exercise 1.13** (Random Walk in Two Dimensions). A particle takes random steps in the xy-plane. The step increment in the x-direction is determined by drawing independent random numbers from a normal distribution with zero mean and variance  $\sigma^2$ . Positive

#### 4 EXERCISES: BASIC CONCEPTS IN PROBABILITY AND STATISTICS

values correspond to steps in the positive x-direction and negative values correspond to steps in the negative x-direction. The step increment in the y-direction is determined in the analogous way. Find the equation for the radius of the circle such that there is a 95% probability that the particle lies within the circle after 10 steps. Evaluate this equation for the case  $\sigma = 2$ .

**Exercise 1.14.** In exercise 1.6, you proved that  $k = \mathbb{E}[X]$  minimized  $\mathbb{E}[(X - k)^2]$ . This says that the best prediction of X (in a mean square sense) is the mean. However, the population mean is unknown. Accordingly, consider a prediction of X based on the sample mean  $\hat{\mu}_X$ , derived from  $X_1, \ldots, X_N$ . It turns out that the sample mean is *not* the best prediction of X', when X' is independent of the sample used to construct  $\hat{\mu}_X$ . Introduce a scaling factor  $\alpha$ , and show that the  $\alpha$  that minimizes

$$\mathbb{E}[(X' - \alpha \hat{\mu})^2],$$

is

$$\label{eq:alpha} \alpha = \frac{1}{\left(\frac{\sigma_X^2}{N\mu_X^2}\right) + 1},$$

where  $\mu_X$  and  $\sigma_X^2$  are the population means and variances of X. Show that  $\alpha < 1$ , which explains why this predictor is sometimes called a *shrinkage* estimator.

#### EXERCISES FOR STATISTICAL INFERENCE

**Exercise 2.1.** A widely used measure of the state of the Pacific Ocean is the Pacific Decadal Oscillation (PDO) index. This index is a certain linear combination of sea surface temperatures in the Pacific Ocean poleward of 20N. More details and background information can be found at http://jisao.washington.edu/pdo/, but this information is not necessary in order to complete this homework assignment. The January-March mean value of this index is plotted in fig. 2.1. A glance at the figure reveals that the PDO index was predominantly negative during the period 1950-1976 and predominantly positive 1977-2017. In this homework set, you will address the following questions about the PDO:

- 1. Are the samples independent?
- 2. Has the mean PDO changed in recent decades?

To address the above questions, you need to download the data file PDO.latest.txt *from the class website.* The R tutorial discussed how to download this data set, but the data file from the PDO webpage changed slightly, so I have updated the data set and included R code for reading it called pdo.student.R. Both files can be downloaded from the class website<sup>1</sup> (click "data" and "repository of R programs"). This code should run without generating any errors, so please contact me if you have trouble running this code.

<sup>1</sup>http://cola.gmu.edu/delsole/clim762/webpage/clim762\_frontpage.html



JFM Average PDO Index

**Figure 2.1** The January-March mean Pacific Decadal Oscillation (PDO) index. Positive values are indicated by thick bars. Statistics of the January-March mean PDO Index are given in the table.

For this homework set, you may use the built-in functions mean and var, but do not use the built-in functions cov and cor; rather, you should compute these quantities explicitly using the sum command (as shown in the tutorial). The purpose of requiring you to calculate these directly is to ensure you understand how to do these calculations yourself. You can, however, use these functions to check your results. You may also use these built-in functions after this homework set.

You will write a *function* that performs a hypothesis test. You should read the appropriate sections of the tutorial to become familiar with functions. In particular, the tutorial illustrates a complete example for one particular hypothesis test (i.e., equality-of-variance test), which in turn can serve as a template for the other hypothesis tests in this homework.

In each problem, you should explicitly state your null hypothesis, test statistic, the observed value of the statistic, rejection region, and final decision about the hypothesis. You should also state all the assumptions that were made in performing the hypothesis test.

(a) **Are the Samples Independent?** Recall that one assumption in both the F-test and t-test is that the samples are independent. Climate time series have the property that the degree of dependence decays with time separation. For instance, 1-day forecasts are more accurate than 5-day forecasts. One approach to checking the independence assumption is by testing if the correlation *between two consecutive time steps* vanishes. If no correlation can be detected between consecutive time steps, then it is unlikely it can be detected for larger time separations. Write an all-purpose function called cor.equal.test that tests the

hypothesis that the correlation between two given data sets vanishes. Turn in a copy of your code. The call to this function should start as follows:

```
cor.equal.test = function(data1, data2, alpha=0.05) {
   ## THIS FUNCTION TESTS VANISHING CORRELATION BETWEEN
2
   ## TWO BI-VARIATE, NORMALLY DISTRIBUTED RANDOM VARIABLES
3
   ##
4
   # INPUT:
5
      DATA1: [N]-DIMENSIONAL VECTOR OF DATA
   #
6
       DATA2: [N]-DIMENSIONAL VECTOR OF DATA
   #
       ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
   #
8
   # OUTPUT LIST:
9
       RHO: SAMPLE CORRELATION BETWEEN DATA1 AND DATA2
   #
10
   #
       RHO.CRIT: 100*ALPHA% CRITICAL VALUE FOR THE CORRELATION
11
   #
       PVAL: P-VALUE OF THE STATISTIC RHO
12
```

(b) Using this function, test the hypothesis that during 1950-1977 the correlation between consecutive years is zero. This means that whatever you use for the time series in data1, data2 is the *same time series shifted by one year*. This also means you cannot input the entire time series, but instead the length of the time series is short by one year. Do the same test for 1978-2017. What is your answer to these questions? Turn in a copy of the output of your function, which should give numerical values for all quantities like  $\hat{\rho}$ ,  $\rho_{crit}$ , p-value.

Incidentally, an R function called cor.test also performs this test. Check that your function agrees with cor.test.

(c) Repeat the above test, but this time for the whole period 1950-2017. You should find that  $\rho = 0$  is rejected more strongly than for the two separate periods individually. Explain why. Hint: draw a scatter diagram for the three separate cases.

(d) Answer the other questions posed at the beginning (i.e., define your test statistic, state the value of the statistic, specify the rejection region, etc.).

(e) **Has the mean PDO changed in recent decades?** To address this question, we test the hypothesis of no difference in population mean between the periods 1950-1977 and 1978-2017? Write a function called mean.equal.test that performs the hypothesis test. Run your function on the PDO data. Turn in a copy of your code and its output. The function should begin as follows:

8 EXERCISES FOR STATISTICAL INFERENCE

```
mean.equal.test = function(data1, data2, alpha=0.05) {
1
  ## THIS FUNCTION TESTS EQUALITY OF MEANS OF TWO IID
2
   ## NORMALLY DISTRIBUTED RANDOM VARIABLES
3
   ##
4
   # INPUT:
5
   #
       DATA1: [N1]-DIMENSIONAL VECTOR OF DATA
6
   #
       DATA2: [N2]-DIMENSIONAL VECTOR OF DATA
7
   #
       ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
8
   # OUTPUT LIST:
0
   #
      DIFF.MEAN: DIFFERENCE IN MEANS (MEAN1 - MEAN2)
10
  #
      DIFF.MEAN.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE DIFFERENCE IN MEANS
11
12
  #
      PVAL: P-VALUE OF THE T-STATISTIC
13
   #
      T: T-STATISTIC FOR DIFFERENCE IN MEANS
  #
      T.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE T STATISTIC
14
   #
     MEAN1: ESTIMATE OF THE MEAN OF DATA1
15
       MEAN2: ESTIMATE OF THE MEAN OF DATA2
  #
16
       SPOOL: POOLED ESTIMATE OF THE STANDARD DEVIATION
   #
17
```

(f) An output quantity not discussed in class (or the notes) is diff.mean.crit: this is the threshold value of the absolute difference in means for deciding whether to reject the null hypothesis. Write an equation that gives this threshold value. State the numerical value for this problem.

(g) State your test statistic, the value of the statistic, and specify the rejection region. What is your decision? Be sure to explicitly state all assumptions made in your hypothesis test.

(h) Explain whether the  $\rho = 0$  test and equality-of-variance tests, which examined above and in the R tutorial, are relevant to testing equality of means.

#### EXERCISES FOR CONFIDENCE INTERVALS

**Exercise 3.1.** Suppose  $X_1, \ldots, X_N$  are independent and identically distributed as  $\mathcal{N}(10, 3)$ . Derive a formula for the 95% confidence interval for the mean. (The answer is given in (3.17), but I want to see you derive it.) Write an R code that generates N = 20 samples  $\mathcal{N}(10, 3)$ , computes the confidence interval, and repeats this many times (1000s). Compute the fraction of intervals that include the true mean  $\mu = 10$ . Verify that this fraction is close to 95%.

**Exercise 3.2.** The tutorial contains the function var.equal.test for testing equality of variance. Augment this function to also calculate a confidence interval for the ratio of variances. The preamble of the revised function should like the following:

**10** EXERCISES FOR CONFIDENCE INTERVALS

```
var.equal.test = function(data1, data2, alpha=0.05) {
  ### THIS FUNCTION TESTS EQUALITY OF VARIANCE OF TWO
2
  ### IID NORMALLY DISTRIBUTED RANDOM VARIABLES
3
  ###
   # INPUT:
      DATA1: [N1]-DIMENSIONAL VECTOR OF DATA
   #
6
   #
      DATA2: [N2]-DIMENSIONAL VECTOR OF DATA
7
   #
      ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
8
   # OUTPUT LIST:
   #
      F.MAX: RATIO OF VARIANCE, CONSTRUCTED TO BE GREATER THAN 1
10
  #
      F.CRIT: THE UPPER CRITICAL THRESHOLD OF SIGNIFICANCE
11
12
  #
      PVAL: P-VALUE OF THE F.MAX RATIO
13
      VAR1: UNBIASED ESTIMATE OF THE VARIANCE OF DATA1
      VAR2: UNBIASED ESTIMATE OF THE VARIANCE OF DATA2
14
  #
  #
     RATIO: VAR1/VAR2
15
  #
     RATIO.LOWER: LOWER LIMIT OF ALPHA% CONFIDENCE INTERVAL FOR VARIANCE RATIO
16
   # RATIO.UPPER: UPPER LIMIT OF ALPHA% CONFIDENCE INTERVAL FOR VARIANCE RATIO
17
```

The above is identical to the function in the tutorial, but with three additional lines at the end. Be sure the resulting confidence interval pertains to var1/var2 and not its reciprocal. Run this function on exactly the same PDO data you used in exercise 2.1 to test if the variance before and after 1977 differ. Report the 95% confidence interval. Comment on whether the result is consistent with the hypothesis test you performed in exercise 2.1.  $\Box$ 

**Exercise 3.3.** In the last exercise, you wrote the function cor.equal.test to test the correlation coefficient. Modify this function to also calculate a confidence interval for the correlation coefficient. The preamble of this function should look like the following:

```
cor.equal.test = function(data1, data2, alpha=0.05) {
1
   ## THIS FUNCTION TESTS VANISHING CORRELATION BETWEEN
   ## TWO BI-VARIATE, NORMALLY DISTRIBUTED RANDOM VARIABLES
3
   ##
4
  # INPUT:
5
   #
      DATA1: [N]-DIMENSIONAL VECTOR OF DATA
6
   #
      DATA2: [N]-DIMENSIONAL VECTOR OF DATA
7
   #
      ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
  # OUTPUT LIST:
     RHO: SAMPLE CORRELATION BETWEEN DATA1 AND DATA2
10
  #
  #
      RHO.CRIT: 100*ALPHA% CRITICAL VALUE FOR THE CORRELATION
11
  # PVAL: P-VALUE OF THE STATISTIC RHO
12
     RHO.LIMITS: [2] (1-ALPHA) *100% CONFIDENCE LIMITS OF THE CORRELATION
  #
13
```

Run this function on the same PDO data to test if the lag-1 correlation is significant before 1977, after 1977, and the full period. Report the 95% confidence intervals. Comment on whether the result is consistent with the hypothesis tests you performed in exercise 2.1.  $\Box$ 

**Exercise 3.4.** In the last exercise you wrote the function mean.equal.test to test equality of means. Modify this function to also compute a confidence interval for the difference in mean. The preamble of this function should look like the following:

```
mean.equal.test = function(data1, data2, alpha=0.05) {
1
   ## THIS FUNCTION TESTS EQUALITY OF MEANS OF TWO IID
2
   ## NORMALLY DISTRIBUTED RANDOM VARIABLES
3
   ##
   # INPUT:
5
       DATA1: [N1]-DIMENSIONAL VECTOR OF DATA
   #
6
   #
       DATA2: [N2]-DIMENSIONAL VECTOR OF DATA
7
   #
       ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
8
   # OUTPUT LIST:
0
   #
       DIFF.MEAN: DIFFERENCE IN MEANS (MEAN1 - MEAN2)
10
   #
       DIFF.MEAN.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE DIFFERENCE IN MEANS
11
12
   #
       PVAL: P-VALUE OF THE T-STATISTIC
13
   #
      T: T-STATISTIC FOR DIFFERENCE IN MEANS
   #
      T.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE T STATISTIC
14
   #
       MEAN1: ESTIMATE OF THE MEAN OF DATA1
15
       MEAN2: ESTIMATE OF THE MEAN OF DATA2
   #
16
   #
       SPOOL: POOLED ESTIMATE OF THE STANDARD DEVIATION
17
   #
       DIFF.MEAN.LIMITS: [2] (1-ALPHA)100% CONFIDENCE
18
   #
           LIMITS FOR THE DIFFERENCE IN MEANS
19
```

Apply this function to the same PDO data to test if the mean PDO index changed before and after 1977. State the 95% confidence limits. Comment on whether the result is consistent with the hypothesis tests you performed in the last exercise.

**Exercise 3.5** (Bootstrap Confidence Interval for a Correlation Coefficient). Write an R function to compute a confidence interval for a correlation coefficient using the bootstrap method. The preamble of this function should look like the following:

```
cor.equal.boot = function(x,y,alpha=0.05,nboot=100000) {
1
   ## ESTIMATES (1-ALPHA)100% CONFIDENCE INTERVAL
2
   ## FOR THE CORRELATION COEFFICIENT USING BOOTSTRAP METHODS
   ## INPUT:
   ##
      X[NSAMP]: FIRST DATA SET
   ##
       Y[NSAMP]: SECOND DATA SET
       ALPHA: EQUIVALENT SIGNIFICANCE LEVEL (DEFAULT = 5%)
   ##
7
      NBOOT: NUMBER OF BOOTSTRAP SAMPLES (DEFAULT = 100000)
   ##
   ## OUTPUT:
   ##
      CONFIDENCE LIMITS[2]
10
```

To test your code, generate N = 20 random samples of (X,Y) from a population with correlation  $\rho = 0.5$  (do you remember how to do this? The answer is in homework 1.). Show both the bootstrap interval and the interval computed from cor.test. Discuss.

The following are some coding hints. R has a function called sample that will randomly sample a vector with replacement. To compute a confidence interval for a correlation coefficient, you must sample the two data sets X and Y in such a way as to preserve their pairing. This can be done by sampling the *indices* rather than the values themselves. For instance, if the sample size is N = 10, then the indices  $1, \ldots, 10$  can be sampled using the following commands: 12 EXERCISES FOR CONFIDENCE INTERVALS

```
1 > nsamp = 10
2 > npic = sample(nsamp,nsamp,replace=TRUE)
3 > npic
4 [1] 4 6 6 4 6 1 1 2 2 8
```

The bootstrap sample and its correlation is then

```
x = rnorm(nsamp); y = rnorm(nsamp)
xy.cor.obs = cor(x,y)
xy.cor.boot = cor(x[npic],y[npic])
```

After you generate lots of bootstrap samples, the  $(1 - \alpha)100\%$  confidence limits can be obtained using the command

confidence.limits = quantile(xy.cor.boot,probs=c(alpha/2,1-alpha/2))

Because the bootstrap requires 1000s of iterations, you want a code that is fast. Unfortunately, R is slow whenever for loops are involved. Can you write a function that performs the bootstrap without for loops? If you cannot figure out how to do this, then just use loops-I much prefer that you turn in a slow code that is correct, rather than a fast code that is incorrect! Hint: use rowMeans or colMeans. Advice: set nboot to a small number (e.g., 100 or 1000) while you debug your code, so it will run fast, and then when you believe the code is working, increase nboot to 100,000 as requested in the above preamble.

## EXERCISES FOR STATISTICAL TESTS BASED ON RANKS

**Exercise 4.1.** Write a function to perform the Wilcoxon Rank-Sum test for a difference in medians. The preamble of this function should be as follows:

```
diff.mean.nonparametric = function(data1, data2, alpha=0.05) {
   ### PERFORMS MANN-WHITNEY TEST FOR A DIFFERENCE IN MEDIANS
2
   ### INPUT:
3
         DATA1[N1]: SAMPLE 1
   ###
4
   ###
         DATA2[N2]: SAMPLE 2
   ###
        ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 5%)
   ### OUTPUT: LIST$
   ###
         Z: Z-VALUE FOR THE NORMAL APPROXIMATION TO MANN-WHITNEY TEST
   ###
         Z.CRIT: CRITICAL VALUE FOR ALPHA*100% SIGNIFICANCE
         PVAL: P-VALUE OF THE TEST
   ###
10
```

The following R functions are useful for non-parametric testing: rank(), sort(), order(). You can read the documentation for these functions and play with them by generating the random numbers seeing what happens:

14 EXERCISES FOR STATISTICAL TESTS BASED ON RANKS

```
> x = rnorm(5)
  > x
2
       0.5726722 -0.6338714 0.3260635 -1.6519467 0.9540662
   [1]
  > sort(x)
   [1] -1.6519467 -0.6338714 0.3260635 0.5726722 0.9540662
  > order(x)
6
   [1] 4 2 3 1 5
7
  > x[order(x)]
8
   [1] -1.6519467 -0.6338714 0.3260635 0.5726722 0.9540662
  > rank(x)
10
   [1] 4 2 3 1 5
11
```

Rank is a very useful function for this homework set because it automatically deals with ties correctly.

```
1 > x[1] = x[2]
2 > x
3 [1] -0.6338714 -0.6338714 0.3260635 -1.6519467 0.9540662
4 > rank(x)
5 [1] 2.5 2.5 4.0 1.0 5.0
```

**Exercise 4.2.** Run this function on exactly the same PDO data you used in the previous 2 exercise sets to test if the median before and after 1977 differ. Clearly state the z-value, critical z value, and the p-value of the test. Clearly state whether you reject the null hypothesis or not, and whether this conclusion is consistent with your conclusions on previous homeworks. You can check your answer using wilcox.test(pdo1,pdo2)

**Exercise 4.3.** In exercise 2.1, you wrote the function cor.equal.test to test the correlation coefficient. Write a new function that performs a correlation test based on Spearman's rank correlation. The preamble of this function should look like the following:

```
cor.test.nonparametric = function(data1, data2, alpha=0.05) {
   ### PERFORMS RANK CORRELATION TEST
2
   ### INPUT:
         DATA1[N1]: SAMPLE 1
   ###
   ###
         DATA2[N2]: SAMPLE 2
         ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 5%)
   ###
   ### OUTPUT: LIST$
   ###
         RHO.SPEARMAN = SPEARMAN'S RANK CORRELATION
   ###
         ZVAL: Z-VALUE FOR THE NORMAL APPROXIMATION TO MANN-WHITNEY TEST
         Z.CRIT: CRITICAL VALUE FOR ALPHA*100% SIGNIFICANCE
10
   ###
   ###
         PVAL: P-VALUE OF THE TEST
11
```

**Exercise 4.4.** Run this function on the same PDO data you used in previous exercises to test if the lag-1 correlation is significant before 1977, after 1977, and the full period. State whether the procedure rejects the null hypothesis or not. Comment on whether the result

is consistent with the hypothesis tests you performed in exercise 2.1. You can check your answer to some extent using

cor.test(pdo1[10],pdo1[10+1],method='spearman',continuity=FALSE,exact=FALSE)

However, the p-value will not be exactly the same (but they are close). The reason they are not the same is because R computes p-values using

$$t = \frac{\hat{\rho}_{Sp}\sqrt{N-2}}{\sqrt{1-\hat{\rho}_{Sp}^2}},$$
(4.1)

and then it refers to a t-distribution with N-2 degrees of freedom.

**Exercise 4.5.** Write a function called diff.dispersion for testing a difference in dispersion. The preamble of this function should be as follows:

```
diff.dispersion = function(data1, data2, alpha=0.05) {
1
   ## PERFORMS WILCOXON SOUARED-RANK TEST
2
   ## ON ABSOLUTE DEVIATIONS |X - MEDIAN(X)|
3
   ### INPUT:
4
   ###
        DATA1[N1]: SAMPLE 1
   ###
         DATA2[N2]: SAMPLE 2
   ###
        ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 5%)
7
   ### OUTPUT: LIST$
   ###
         ZVAL: Z-VALUE FOR THE NORMAL APPROXIMATION TO MANN-WHITNEY TEST
         Z.CRIT: CRITICAL VALUE FOR ALPHA*100% SIGNIFICANCE
   ###
10
         PVAL: P-VALUE OF THE TEST
11
   ###
```

Run your function on exactly the same PDO data you used in previous exercises to test if the variance before and after 1977 differ. Clearly state the z-value, critical z value, and the p-value of the test. Clearly state whether you reject the null hypothesis or not, and whether this conclusion is consistent with your conclusions on previous homeworks. (Hint: most of this function is similar to your difference in medians test.)

#### EXERCISES FOR STOCHASTIC PROCESSES

**Exercise 5.1.** Write an R function to compute the autocorrelation function of a time series. The preamble of the function should be as follows:

```
acf.brute = function(x,lag.max=NULL) {
1
  ## COMPUTES THE AUTOCORRELATION FUNCTION OF A TIME SERIES
  ## INPUT:
       X: [NTOT]-LENGTH NUMERICAL VECTOR OF THE TIME SERIES
  ##
       LAG.MAX: MAXUMIM LAG TO COMPUTE
  ##
  ##
           (DEFAULT = MAX(1, FLOOR(10 * log10(LENGTH(X)))))
6
  ## OUTPUT:
7
       X.ACF: [1:(LAG.MAX+1)] AUTOCORRELATION FUNCTION OF X
  ##
8
           FOR LAGS 0 TO LAG.MAX
  ##
9
```

Test your function on the time series {set.seed(1); x = rnorm(20)} and state the correlation values for all lag.max lags. Show that these values equal those obtained from the built-in R function acf(x), by printing them out using the commands

```
acf.R = acf(x)
as.numeric(acf.R$acf)
```

18 EXERCISES FOR STOCHASTIC PROCESSES

Exercise 5.2. Write a function to generate a time series from the AR1 model

$$x_{t+1} = \phi_1 x_t + w_t + k \tag{5.1}$$

where  $w_t \sim GWN(0, 1)$ . The beginning of this function should be as follows:

```
ar1.ts = function(phil,ntot,stdv=1,constant=0,iseed=1) {
   ## GENERATE A REALIZATION FROM THE AR1 MODEL
   ## X(T+1) = PHI1 * X(T) + W(T) + CONSTANT
   ## WHERE W(T) IS GAUSSIAN WHITE NOISE WITH MEAN 0 AND ST. DEV. 'STDV'
   ## INITIAL CONDITION IS RANDOMLY DRAWN FROM STATIONARY DISTRIBUTION
   ## INPUT:
        PHI1: AUTOREGRESSIVE PARAMETER
   #
        NTOT: LENGTH OF THE DESIRED TIME SERIES
8
        STDV: STANDARD DEVIATION OF THE NOISE (DEFAULT = 1)
   #
9
        CONSTANT: CONSTANT TERM IN THE AR1 MODEL (DEFAULT = 0)
   #
10
   #
        ISEED: SEED FOR THE RANDOM NUMBER GENERATOR (DEFAULT = 1)
11
12
   #
      OUTPUT:
13
   #
        X: [NTOT] - RANDOM TIME SERIES FROM AR(1) PROCESS
```

The initial condition should be drawn from the *stationary distribution* of the AR process. What is this distribution? State this distribution for general  $k, \phi_1, \sigma_W^2$ . Use this equation in your function.

**Exercise 5.3.** Use the above functions to generate a time series of length 50, and corresponding sample autocorrelation functions, for  $\phi_1 = 0, 0.5, 0.9$ . Make plots of both quantities. In addition, on the plot for the autocorrelation, superimpose the population autocorrelation function, and the approximate 95% confidence interval for zero correlation. Comment on whether the results make sense, or whether they are 'surprising.'

**Exercise 5.4.** Use the above function to generate a time series of length 50 for  $\phi_1 = 0$ . Split this time series into two halves, each of length 25. Use the function mean.equal.test to calculate the t-statistic for testing equality of means. Repeat this 1000 times, thereby generating 1000 t-values. Plot a histogram of these t-values. Superimpose on this histogram the expected distribution of the t-statistic. This can be done using the curve command, as follows:

```
hist(tval,col="grey",freq=FALSE)
curve(dt(x,dof),add=TRUE)
```

(In the above example, x does not need to be pre-defined. See the manual on R graphics or the help pages for curve.) You should find that the histogram is reasonably consistent with the exact t distribution. Repeat this for  $\phi_1 = 0.5, 0.9$ . Are the histograms still consistent with the exact t distribution? Explain the result you find. If someone concludes that there is significant difference in mean based on a standard t-test (i.e., using iid assumption), how might this conclusion be affected if autocorrelation were taken into account? Does the true significance level increase or decrease when autocorrelation is taken into account? Exercise 5.5 (Mean and Variance of an AR(p) Process). Consider the AR(p) process

$$X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + W_t + k,$$
(5.2)

where  $W_t$  is white noise with zero mean and variance  $\sigma^2$ , and k is a constant. Show that the expectation of the stationary AR(p) process is

$$\mathbb{E}[X_t] = \frac{k}{1 - \phi_1 - \phi_2 - \dots - \phi_p}.$$
(5.3)

Show that the variance of the stationary process is

$$\operatorname{var}[X_t] = \frac{\operatorname{var}[W_t]}{1 - \rho_1 \phi_1 - \rho_2 \phi_2 - \dots - \rho_p \phi_p}.$$
(5.4)

(Hint: multiply both sides of (5.2) by  $X_t - \mathbb{E}[X_t]$  and take the expectation.)

Exercise 5.6 (Autocorrelation of a running mean). Consider the running mean

$$Y_t = \frac{1}{2K+1} \sum_{k=-K}^{K} X_{t+k}.$$
(5.5)

Assume  $X_t \sim GWN(0, \sigma^2)$ . Compute the mean and variance of  $Y_t$ . Show that the autocorrelation function of  $Y_t$  is

$$\operatorname{cor}[Y_{t+\tau}, Y_t] = \begin{cases} 1 - \frac{|\tau|}{2K+1} & \text{for } |\tau| \le 2K+1\\ 0 & \text{for} |\tau| > 2K+1 \end{cases}.$$
(5.6)

Note that even though the process  $X_t$  is uncorrelated, running averages of  $X_t$  are correlated for lags less than twice the window size K. Explain why this makes sense.

### EXERCISES ON THE POWER SPECTRUM

In this homework you will determine power spectra from models and data. To do this, you will need to use the Fast Fourier Transform, which is performed by the R function fft(). The output of fft is somewhat technical, so I give you the following program to convert the output into a periodogram:

```
periodogram = function(x) {
1
   ## COMPUTES THE PERIODOGRAM OF A VECTOR X USING THE FFT
2
   ## ASSUMES THE LENGTH OF X IS EVEN AND HIGHLY COMPOSITE
3
   ##
        (E.G., EXPRESSIBLE AS POWERS OF 2,3,5)
4
   # INPUT:
5
   #
        X[NTOT]: DATA VECTOR
7
   # OUTPUT:
   #
        PGRAM[NTOT/2]: VECTOR OF PERIODOGRAM AMPLITUDES AT HARMONIC FREQUENCIES
8
          1/N, 2/N, . . ., 1/2 CYCLES PER TIME STEP (A TOTAL NTOT/2 FREQUENCIES)
   #
9
   # COMMENT: (BIASED) VARIANCE = 2 * SUM(PGRAM[1:(NTOT/2-1)]) + PGRAM[NTOT/2]
10
11
  ntot = length(x)
12
  if ( ntot %% 2 != 0 ) stop('length of x must be even')
13
14
  x.fft = fft(x)
  pgram = Re(x.fft * Conj(x.fft))[1:(ntot/2)+1] / ntot
15
   pgram
16
17
   }
```

You should verify that this function works by setting  $\times$  to be specific sum of sines and cosines and ensuring that the output gives the appropriate amplitudes, when appropriately scaled. The following example shows a successful test of this function:

```
> set.seed(4)
1
  > ntot = 32
2
  > t
          = 1:ntot - 1
          = 1 * cos(2*pi*t*5/ntot) - 3 * sin(2*pi*t*13/ntot)
  > x
4
         + 8 * cos(2*pi*t*(ntot/2)/ntot)
5
  > x.pgram = periodogram(x)
6
  > x.amplt = sqrt(4*x.pgram/ntot)
                                       # convert pgram to fourier amplitude
7
  # special conversion for last harmonic
  > x.amplt[ntot/2] = sqrt(x.pgram[ntot/2]/ntot)
  > x.amplt
10
    [1] 2.519803e-15 2.181560e-15 7.616681e-16 1.930989e-15 1.000000e+00
11
    [6] 1.260945e-15 7.422817e-16 3.202373e-15 2.989367e-15 2.373832e-15
12
  [11] 3.143748e-15 1.223717e-15 3.000000e+00 3.084970e-16 1.271623e-15
13
   [16] 8.000000e+00
14
```

```
The output shows that the absolute amplitude of the 5th, 13th, and 16th harmonics are 1, 3, 8, as they should, while the other amplitudes are effectively zero, as expected. Another check is to verify the formula
```

$$\sum_{t=1}^{N} (x_t - \overline{x})^2 = 2 * \sum_{j=1}^{\frac{N}{2}-1} \hat{P}(j/N) + \hat{P}(N/2),$$
(6.1)

which holds for even N. The following lines perform this check:

```
1 > set.seed(1)
2 > x = rnorm(ntot)
3 > ssx = sum((x-mean(x))^2)
4 >
5 > x.pgram = periodogram(x)
6 > ssx.fft = 2*sum(x.pgram[1:(ntot/2-1)]) + x.pgram[ntot/2]
7 > print(all.equal(ssx,ssx.fft))
8 [1] TRUE
```

Exercise 6.1. Write an R function to compute the spectral average

$$\hat{p}_M(\omega_j) = \frac{1}{2M+1} \sum_{m=-M}^{M} \hat{p}(\omega_j).$$
(6.2)

The preamble of this R function should be as follows:

```
spectrum.smooth = function(pgram,window,alpha=0.1) {
1
   ## COMPUTES SPECTRAL RUNNING MEAN AS:
2
   ##
        SPEC_AVE = SMOOTHED PERIODOGRAM FROM M = -WINDOW TO M = WINDOW
3
   ## INPUT:
   ##
         PGRAM[NTOT/2]: PERIODOGRAM VALUES FROM PERIODOGRAM.R
5
   ##
         WINDOW: HALF-WIDTH OF AVERAGING WINDOW
6
   ##
         ALPHA: THE (1-ALPHA)100% LEVEL FOR THE CONFIDENCE LIMITS
7
   ## OUTPUT: LIST
8
   ##
         PGRAM.AVE[NTOT/2]: AVERAGED PERIODOGRAM
0
   ##
         FREQ[NTOT/2]: FREQUENCIES (CYCLES PER TIME STEP)
10
   ##
         BW[2]: BANDWIDTH C(-WINDOW/NTOT, WINDOW/NTOT)
11
12
   ##
         CI.LIMITS[2]: PROPORTIONALITY CONSTANT FOR LOWER AND UPPER CONFIDENCE LIMITS
13
   ##
            (DOF/CHI-SQUARE[ALPHA/2, DOF], DOF/CHI-SQUARE[1-ALPHA/2,DOF])
   ##
         PGRAM[NTOT/2]: REPEATS PERIODOGRAM (FOR LATER PLOTTING)
14
```

The spectral average near the extreme low and high frequencies, where the smoothing cannot be performed consistently, should be set to NA.  $\hfill \Box$ 

**Exercise 6.2.** Generate a time series of length 64 from a standardized Gaussian white noise process, perform a spectral average using M = 5, and write out the results, as follows:

```
> set.seed(1)
1
  > x.sd
               = 1
2
  > window
              = 5
3
  > ntot
                = 64
4
  > freq
               = 1:(ntot/2) / ntot
5
  > x
                = rnorm(ntot,sd=x.sd)
6
              = periodogram(x)
  > x.pgram
7
  > x.pgram.smooth = spectrum.smooth(x.pgram,window)
8
  > cbind(freq, x.pgram, x.pgram.smooth$pgram.ave)
```

Plot the results using the R function spectrum.plot.R, provided on the class website. Based on this result, would you say that the estimated power spectrum is consistent with a white noise process? Explain why or why not.

**Exercise 6.3.** What is the theoretical (i.e., population) power spectrum for the Gaussian white noise process? Superimpose this theoretical power spectrum on your spectrum plot. Is the estimated power spectrum consistent with the theoretical spectrum?

**Exercise 6.4.** Generate a time series of length 64 from an AR(1) process with parameter  $\phi_1 = 0.8$ , as follows:

1 ntot = 64
2 phil = 0.8
3 x = arl.ts(phil,ntot,iseed=1)

Perform a spectral average of this time series using M = 5, and write out the results, as follows:

24 EXERCISES ON THE POWER SPECTRUM

```
1 > window = 5
2 > x.pgram = periodogram(x)
3 > x.pgram.smooth = spectrum.smooth(x.pgram,window)
4 > cbind(freq, x.pgram, x.pgram.smooth$pgram.ave)
```

Make a plot of the power spectrum, smoothed periodogram, bandwidth, and 90% confidence interval using spectrum.plot.R. Is the estimated spectrum more consistent with white noise or an AR(1) process? Explain why or why not.  $\Box$ 

**Exercise 6.5.** What is the theoretical power spectrum of the above AR(1) process? Superimpose this theoretical spectrum on your spectrum plot. Is your estimated power spectrum consistent with the theoretical spectrum?

**Exercise 6.6.** Estimate the power spectrum of the NINO3.4 index. The NINO3.4 index can be extracted using the following R program:

```
fname
              = paste(dir.enso,'nina34.data', sep='')
2
  line1
               = scan(fname, nlines=1, quiet=TRUE)
  nino34.yr = line1[1]:line1[2]
3
  nino34.nyrs = length(nino34.yr)
4
  cnames = c('YEAR', month.abb[])
5
  data2
              = read.table(fname, skip=1, nrows=nino34.nyrs,
      na.strings=-99.99, col.names=cnames)
  ## CREATE VECTOR WITH NINO34 TIME SERIES AND TIME AXIS
  nino34.ts = as.numeric(t(as.matrix(data2[,1:12+1])))
10
  nino34.time = seq(from=data2[1,1],length.out=length(nino34.ts),by=1/12)
11
```

The resulting time series will have NA's which need to be stripped out before performing spectral analysis. Trim the length of the index to be highly composite (i.e., a product of powers of 2, 3, 5). State the final length that you use. What are the units of "frequency"? Are there any significant peaks? Do these peaks make sense (do you expect certain peaks)? Could an AR-process fit this spectrum? ENSO is often said to be characterized by oscillations around 4-7 years. Does your estimated power spectrum support this characterization? Explore different bandwidths, choose a single one that you believe is most revealing, and show that one.

## EXERCISES FOR VECTORS, MATRICES, AND GEOMETRY

**Exercise 7.1** (Linear combination of random vectors). Let  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$  be independent *M*-dimensional random vectors. However, the vectors are not drawn from the same population. Instead, the populations have individual means  $\mathbb{E}[\mathbf{x}_k] = \boldsymbol{\mu}_k$  and covariance matrices  $\operatorname{cov}[\mathbf{x}_k] = \boldsymbol{\Sigma}_k$ . Because the vectors are independent,

$$\operatorname{cov}[\mathbf{x}_k, \mathbf{x}_{k'}^T] = \begin{cases} \mathbf{0} & k \neq k' \\ \mathbf{\Sigma}_k & k = k' \end{cases}$$
(7.1)

Note that there are K distinct covariance matrices  $\Sigma_1, \Sigma_2, \ldots, \Sigma_K$ . Suppose  $\mathbf{y} = c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + \cdots + c_N \mathbf{x}_N$ , where the c's are known scalars. What is the mean and covariance matrix of  $\mathbf{y}$ ?

**Exercise 7.2** (Matrix combinations of random vectors). Let  $\mathbf{x} \sim N_M(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Suppose  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ , where  $\mathbf{A}$  is a constant matrix and  $\mathbf{b}$  is a constant vector. What is the complete distribution of  $\mathbf{y}$ ?

**Exercise 7.3** (Properties of the Trace and Determinant). Let **A** be a square matrix, but not necessarily symmetric. Square matrices still have eigenvectors and eigenvalues that satisfy

$$\mathbf{As} = \lambda \mathbf{s}.\tag{7.2}$$

As a result, the matrix is close enough for practical purposes to having an eigen decomposition of the form

$$\mathbf{A} = \mathbf{S}\mathbf{\Lambda}\mathbf{S}^{-1},\tag{7.3}$$

25

#### 26 EXERCISES FOR VECTORS, MATRICES, AND GEOMETRY

in which case A is said to diagonalizable and the properties below are easy to prove. If A is symmetric, then we know that there exists a decomposition of the above form where S is an orthogonal matrix and  $\Lambda$  has real diagonal elements, but these special properties are not needed in the following problems.

(a) Show that  $\operatorname{tr} \mathbf{A} = \operatorname{tr} \mathbf{A}^T$  for any matrix  $\mathbf{A}$ .

(b) Show that  $tr{AB} = tr{A^TB^T}$  for any two matrices A and B such that  $tr{AB}$  is defined.

(c) Show that tr  $\mathbf{A} = \sum \lambda_i$ , where the  $\lambda_i$ 's are the eigenvalues of  $\mathbf{A}$ , and the sum is over all eigenvalues.

(d) What is  $tr{A^k}$  in terms of the eigenvalues of A?

(e) Show that  $tr[\mathbf{A}^{-1}] = \sum_{m=1}^{M} \lambda_m^{-1}$ , where the  $\lambda_m$ 's are the eigenvalues of  $\mathbf{A}$ .

(f) Show that  $tr{\mathbf{A}\mathbf{A}^T} = \sum_i \sum_j a_{ij}^2$ , where  $a_{ij}$  is the (i, j) element of **A**.

(g) Show that  $|\mathbf{A}| = \lambda_1 \lambda_2 \dots \lambda_p$ , where  $\mathbf{A}$  is a  $p \times p$  matrix and the  $\lambda_i$ 's are the eigenvalues of  $\mathbf{A}$ .

(h) If A is symmetric and positive definite, show that  $\operatorname{tr} \mathbf{A} > 0$  and  $|\mathbf{A}| > 0$ .

(i) Show that the trace of a positive semi-definite matrix is non-negative. Show that the determinant of a positive semi-definite matrix is non-negative.

**Exercise 7.4** (Idempotent matrices). An *idempotent* matrix is a matrix  $\mathbf{A}$  such that  $\mathbf{A}^2 = \mathbf{A}$ . Show that the eigenvalues of an idempotent matrix can have only two values, namely 0 and 1. Show that the number of eigenvalues equal to 1 is given by tr[ $\mathbf{A}$ ].  $\Box$ 

Exercise 7.5 (Properties of Orthogonal Matrices).

(a) Show that orthogonal transformations preserve the Euclidean length of a vector. That is, if  $\boldsymbol{\xi} = \mathbf{U}\mathbf{x}$ , where U is an orthogonal matrix, then show that  $\boldsymbol{\xi}^T \boldsymbol{\xi} = \mathbf{x}^T \mathbf{x}$ .

(b) Show that orthogonal transformations preserve the cosine-angle between two vectors. That is, show that the cosine-angle between x and y, and the cosine-angle between  $\boldsymbol{\xi} = \mathbf{U}\mathbf{x}$  and  $\boldsymbol{\eta} = \mathbf{U}\mathbf{y}$ , are equal. The cosine-angle between two vectors x and y is

$$\cos \theta = \frac{\mathbf{x}^T \mathbf{y}}{\sqrt{(\mathbf{x}^T \mathbf{x})(\mathbf{y}^T \mathbf{y})}}.$$
(7.4)

(c) Show that the determinant of an orthogonal matrix is +1 or -1.

(d) Show that  $|\mathbf{U}\mathbf{A}\mathbf{U}^T| = |\mathbf{A}|$ , where U is an orthogonal matrix and A is any matrix.

(e) Show that  $tr{\mathbf{U}\mathbf{A}\mathbf{U}^T} = tr{\mathbf{A}}$ , where U is orthogonal and A is a square matrix.

(f) If  $U_1$  and  $U_2$  are orthogonal matrices, show that  $U_1U_2$  is also orthogonal. Hint: remember that orthogonal matrices satisfy *two* properties.

**Exercise 7.6.** Suppose A is a constant matrix and y is a random vector with mean  $\mu_Y$  and covariance matrix  $\Sigma_Y$ . Show that the expected value of the quadratic form  $y^T A y$  is  $tr[A\Sigma_Y] + \mu^T A \mu$ .

Exercise 7.7. Show that the inverse of a transpose equals the transpose of the inverse:

$$\left(\mathbf{A}^{T}\right)^{-1} = \left(\mathbf{A}^{-1}\right)^{T}.$$
(7.5)

Hint: start with  $AA^{-1} = I$ .

**Exercise 7.8** (Frobenius Norm). The "length" of a vector is defined as the square root of the sum square vector elements. A natural generalization of "length" to a matrix **X** is

$$\|\mathbf{X}\|_{F} = \sqrt{\sum_{i} \sum_{j} x_{ij}^{2}}.$$
(7.6)

This measure is called the *Frobenius Norm*. It is clear that the Frobenius norm vanishes if and only if  $\mathbf{X} = \mathbf{0}$ , consistent with a measure of "length."

(a) Show that the Frobenius Norm can be written equivalently as

$$\|\mathbf{X}\|_F^2 = \operatorname{tr}\left[\mathbf{X}\mathbf{X}^T\right]. \tag{7.7}$$

(b) Suppose U and V are orthogonal matrices with dimensions such that the product UAV is defined. Show that

$$\|\mathbf{UAV}\|_F = \|\mathbf{A}\|_F.$$
(7.8)

**Exercise 7.9.** Consider a  $M \times M$  matrix **A** with eigenvalues  $\lambda_1, \lambda_2, \ldots, \lambda_M$ . Show that the eigenvalues of  $a\mathbf{I} + b\mathbf{A}$  are  $a + b\lambda_m, m = 1, 2, \ldots, M$ .

**Exercise 7.10.** If A is positive definite, show that *each* diagonal element of A must be positive.  $\Box$ 

**Exercise 7.11** (Distribution of a random quadratic form). Suppose  $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma}$  is positive definite. In this problem you will determine the distribution of the quadratic form

$$Q = (\mathbf{x} - \boldsymbol{\mu})^T \, \boldsymbol{\Sigma}^{-1} \left( \mathbf{x} - \boldsymbol{\mu} \right) \,. \tag{7.9}$$

(a) Let the eigenvector decomposition of the covariance matrix be

$$\Sigma = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T, \tag{7.10}$$

where U is an orthogonal matrix, and  $\Lambda$  is the diagonal eigenvalue matrix. What is the eigenvector decomposition of  $\Sigma^{-1}$ ?

(b) Using the results from part (a), construct a transformation matrix A such that

$$\boldsymbol{\zeta} = \mathbf{A} \left( \mathbf{x} - \boldsymbol{\mu} \right) \quad \text{and} \quad \boldsymbol{\zeta}^T \boldsymbol{\zeta} = Q \,.$$
 (7.11)

28 EXERCISES FOR VECTORS, MATRICES, AND GEOMETRY

What is **A**?

(c) What is the distribution of  $\boldsymbol{\zeta}$ ?

(d) A theorem in statistics states that if two random variables are uncorrelated *and* distributed as a multivariate Gaussian, then they are independent. Use this theorem to show that  $\zeta_i$  and  $\zeta_j$  are independent if  $i \neq j$ .

(e) What is the distribution of Q?

## EXERCISES FOR LINEAR REGRESSION: LEAST SQUARES ESTIMATION

In this homework assignment you will write an R function that estimates the regression parameter  $\beta$  in the model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},\tag{8.1}$$

where y is an N-dimensional vector of predictands, and X is an  $N \times M$  dimensional matrix of predictors. The core of this function is to compute an estimate of  $\beta$  based on solution of the normal equations

$$\hat{\boldsymbol{\beta}} = \left(\mathbf{X}^T \mathbf{X}\right)^{-1} \mathbf{X}^T \mathbf{y}.$$
(8.2)

In practice, most packages use the SVD or QR algorithm to estimate  $\beta$ . You will not be asked to use these algorithms because we are interested in focusing on *statistics*, not *numerics*. Having said that, you should be aware that you are solving the least squares problem inefficiently and perhaps inaccurately.

To calculate (8.2), we need to calculate  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{X}^T \mathbf{y}$ . In R, the transpose operation is t (), and matrix multiplication is  $\$ \star \$$ . Thus, these terms are obtained by the commands

xtx	= t(x) %*% x
xty	= t(x) %*% y

The inverse of a general matrix can be calculated using the solve command. However, the matrix we want to invert is symmetric, and it is faster and more accurate to invert

a symmetric matrix using the *Cholesky decomposition*. R can invert a matrix using the Cholesky decomposition as follows:

xtx.inv = chol2inv(chol(xtx))

Verify that xtx %\*% xtx.inv equals the identity matrix (to within roundoff error).

In addition to estimating the parameters  $\beta$ , you also should estimate the sum square error SSE and coefficient of determination  $R^2$ . The sum square error is

$$SSE = \left(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\right)^{T} \left(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\right).$$
(8.3)

In the notes,  $R^2$  was estimated from *centered* variables. It is not necessary to center the variables as long as you include the intercept as a predictor. Therefore, your R function should *automatically insert* the intercept as a predictor. That is, you will give the function all the predictors except the intercept, and the function will insert the intercept among the predictors. You can generate a vector of 1s using the command rep, and then join this vector with the predictor matrix using the command cbind.

In general,  $R^2$  can be calculated as

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}},\tag{8.4}$$

where

$$SST = \sum_{n} \left( y_n - \hat{\mu} \right)^2.$$
(8.5)

**Exercise 8.1.** Write an R function that sets up the normal equations for solving (8.1), solves the normal equations to obtain the least squares estimates, and computes  $R^2$ . The function call based on the normal equations should be regress.normal (y, x), where y and x are the appropriate vector and matrix of the regression model (8.1). The preamble of this function should be the following:

```
regress.normal = function(y,x,include.intercept=TRUE) {
  ******
  ## DETERMINES THE LEAST SQUARES ESTIMATE OF B IN THE EQUATION Y = XB + E
  ## BASED ON THE NORMAL EOUATIONS
  ## INPUT:
                  N-DIMENSIONAL VECTOR OF PREDICTANDS
  ##
      Y[NTOT]:
      X[NTOT, MTOT]: N X M DIMENSIONAL MATRIX OF PREDICTORS
  ##
  ##
      INCLUDE.INTERCEPT: INCLUDE THE INTERCEPT? (DEFAULT=TRUE)
  ## OUTPUT:
      BHAT: M-DIMENSIONAL VECTOR OF ESTIMATES OF B
10
  ##
  ##
      R2: R-SQUARED
11
      SSE: SUM SQUARE ERROR OF THE LEAST SQUARES PREDICTION
  ##
12
  ##
      DOF: DEGREES OF FREEDOM OF THE SSE.
13
  ##
      RES.SE: STANDARD ERROR OF THE RESIDUALS
14
  ******
15
```
Exercise 8.2. Apply your function to the following random numbers:

```
1 > set.seed(1)
2 > ntot = 20
3 > y = rnorm(ntot); pred1 = rnorm(ntot); pred2 = rnorm(ntot)
4 > x = cbind(pred1,pred2)
```

After running your function, print out  $\hat{\beta}$ ,  $R^2$ , SSE, and dof produced by your function.

You should check your function by comparing with the R function lm(). To do this, apply the lm function as follows:

```
> xy.lm
              = lm(y^x)
1
  > summary(xy.lm)
2
3
  Call:
4
  lm(formula = y ~ x)
5
6
7
  Residuals:
    Min 1Q Median
                            ЗQ
                                     Max
8
   -2.0455 -0.6782 0.2175 0.5049 1.4532
9
10
  Coefficients:
11
     Estimate Std. Error t value Pr(>|t|)
12
  (Intercept) 0.1494 0.2072 0.721
                                         0.481
13
  xpred1 -0.1516
                          0.2504 -0.605
                                            0.553
14
  xpred2
              0.2894
                         0.2695 1.074
                                            0.298
15
16
  Residual standard error: 0.9119 on 17 degrees of freedom
17
  Multiple R-squared: 0.1078,
                               Adjusted R-squared: 0.002868
18
  F-statistic: 1.027 on 2 and 17 DF, p-value: 0.3791
19
```

In line 1, the function lm is called using the formula notation. Then, in line 2, the summary function is used to extract basic information about the linear model. For the purpose of the present homework, we are interested in only three parts of this summary: the coefficients,  $R^2$ , and residual standard error. The value of the coefficients are listed under Coefficients: Estimate, and are 0.1494, -0.1516, and 0.2894, corresponding to the three predictors: intercept, pred1, and pred2. These should match the values computed from regress.normal above. The  $R^2$  is in the second to last line: 0.1078, and should agree with that calculated from regress.normal. Finally, the residual standard error derived from regress.normal should agree with the result from lm (line 17).  $\Box$ 

**Exercise 8.3.** Apply your regression function to estimate the growth rate of atmospheric  $CO_2$  concentration over the past half-century or so. State the growth rate in units of ppm/yr. The  $CO_2$  concentration data can be downloaded as  $co2\_mm\_mlo.txt$  from the class website (which in turn was downloaded from http://www.esrl.noaa.gov/gmd/ccgg/trends/). This data set can be read into R as follows:

32 EXERCISES FOR LINEAR REGRESSION: LEAST SQUARES ESTIMATION

```
iyst
           = 1960
1
           = 2017
  iynd
2
3
  ****
  ######## GET CO2 DATA
5
  6
  fdata = '/Users/delsole/data/indices/co2_mm_mlo.txt'
7
  nskip = 72
8
  col.names = c('year','month','date','average','interp','trend','#days')
0
  co2.table = read.table(fdata,skip=nskip,col.names=col.names,na.strings=-99.99)
10
11
year.get = co2.table[,'year'] >= iyst & co2.table[,'year'] <= iynd</pre>
vear.say = co2.table[year.get,'date']
14 month = co2.table[year.get,'month']
15 CO2
        = co2.table[year.get,'average']
16 plot(year.say,co2,type="l",col="black",xlab='year',ylab='Parts Per Million')
```

You should change fdata to correspond to the data file on your computer. The resulting plot should reproduce the figure in the notes.

Unfortunately, there exists missing data. Therefore, *inside your R function*, you will have to strip out this missing data before applying the least squares method. I recommend including the following inside your R function:

```
ntot = length(y)
1
  if (length(x) \% ntot != 0) stop('x not dimensioned correctly')
2
  mtot = length(x)/ntot
3
  if ( ntot <= mtot ) stop('regression problem is not over-determined')
4
   ### STRIP MISSING DATA
6
7
  dim(x)
          = c(ntot, mtot)
  is.missing = is.na(y) | is.na(rowSums(x))
8
  x.good = x[!is.missing,]
             = y[!is.missing ]
10 y.good
11 nsamp
           = sum(!is.missing)
```

Note that the correct sample size after missing data has been stripped is nsamp. This is important when you augment the predictor matrix  $\mathbf{X}$  by a column of ones.

After running your function, print out  $\hat{\beta}$ ,  $R^2$ , SSE, and dof produced by your function. You can check your calculations against the built-in R function  $\lim$  (you will need to use the na.action=na.omit option to deal with missing data).

**Exercise 8.4.** Compute the residuals of the regression equation. Make a plot of them, and *state the first 50 values* as a vector (e.g., as.numeric(residuals[1:50]).

**Exercise 8.5.** Use your regression function to estimate the annual cycle of the  $CO_2$  concentration. The annual cycle should be defined as the first two Fourier harmonics of the annual cycle (i.e., sin/cos function with periods of 12 months and 6 months). The predictor matrix for *just these harmonics* can be constructed in R as follows:

```
year = iyst + (1:ntot - 0.5)/12
1
  year.shift = year - 1960
2
  t = seq(year.shift)/12
3
  nharm = 2
4
  x = NULL
5
  for ( n in 1:nharm) x = cbind(x, cos(2*pi*t*n), sin(2*pi*t*n))
6
  colnames(x) = c(paste(rep(c('cos', 'sin'), nharm),
7
     rep(1:nharm,each=2),sep=""))
8
```

State the five coefficients of this fit (i.e, the intercept, the 2 coefficients for the sin function, 2 coefficients for the cosine function). *Print out the resulting* coefficients,  $R^2$ , dof, and SSE.

**Exercise 8.6.** Plot the residuals after the annual cycle has been removed. Superimpose a plot of the actual  $CO_2$  data for comparison. To do this, you will need to add a constant term to the residuals to make them fit on the same figure; state what constant should you use and why.

# EXERCISES FOR LINEAR REGRESSION: INFERENCE

In this homework, you will investigate whether Guam sea level has been rising over the last few decades. To do this, you need to download the following files from the class website:

file name	purpose
540.rlrdata	Guam sea level Data file
detrend.nino34.ascii.txt	SST data file
enso.index.v2.R	R code for reading SST
hw.LinRegInference.student.R	R code for reading sea level

Sea level data was downloaded from https://www.psmsl.org/data/obtaining/stations/540.php. The SST data was downloaded from

https://www.cpc.ncep.noaa.gov/products/analysis\_monitoring/ensostuff/detrend.nino34.ascii.txt. The first few lines of the sea level data file look like this:

1	1948.0417;	6977;	0;000
2	1948.1250;	6992;	0;000
3	1948.2083;	7023;	0;000
4	1948.2917;	7114;	0;000

#### **36** EXERCISES FOR LINEAR REGRESSION: INFERENCE

The first column is the year-month is the decimal form year + (month-0.5)/12.0. The second column is the mean sea level value for the month, in mm, and the remaining numbers pertain to missing data. If there are no values to average, then the value for the month is set to -99999. However, you do not need to know these details. The R code reads the data file and extracts the data needed to do the homework. At the end of the code, you will see the following lines explaining the variables:

```
1 #### AT THIS POINT, THE FOLLOWING VARIABLES ARE DEFINED:
2 #### SEASON: CHARACTER DESCRIBING THE SEASON
3 #### (E.G., 'OND' FOR 'OCT-NOV-DEC')
4 #### IYST: START YEAR
5 #### IYND: END YEAR
6 #### NYRS: TOTAL NUMBER OF YEARS
7 #### Y[NYRS]: GUAM SEA LEVEL (mm) AVERAGED OVER THE SEASON
8 #### X.YEAR[NYRS]: YEAR
9 #### X.SST[NYRS]: NINO3.4 INDEX FOR THE SAME YEAR/SEASON
```

Note: the numbers derived in this homework will not match those in the notes because the data set here has been updated based on new observations (also, the units in the notes were incorrect).

**Exercise 9.1.** Make a plot of the tide gauge data for Guam during JFM, 1973-2006. Be sure that the axis labels are accurate (with units) and the title explains what is shown in the figure, including the station and the calendar period. Comment on whether the figure suggests that the tide gauge measurements have been increasing in the last few decades.  $\Box$ 

**Exercise 9.2.** Use regress.normal to determine the trend of sea level. The trend is obtained by fitting sea level to the model

sea level = 
$$\beta_1 + \beta_2$$
 year +  $\epsilon_Y$ . (9.1)

Print out the output of regress.normal. Report the trend in units of cm per year.

**Exercise 9.3.** Augment regress.normal so that it calculates the t-statistic and associated p-value for <u>each</u> regression coefficient. Call these values tval and pval in the output list. Recall that if the model is defined as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},\tag{9.2}$$

where **X** is  $N \times M$ , then the least squares estimate of  $\beta$  is

$$\hat{\boldsymbol{\beta}} = \left(\mathbf{X}^T \mathbf{X}\right)^{-1} \mathbf{X}^T \mathbf{y},\tag{9.3}$$

and the t-statistic for testing  $\beta_j = 0$  is

$$t = \frac{\hat{\beta}_j}{s_e \sqrt{d_j}},\tag{9.4}$$

where  $s_e = \sqrt{SSE/(N-M)}$  and  $d_j$  is the j'th diagonal element of  $(\mathbf{X}^T \mathbf{X})^{-1}$ . The t-statistic has a t-distribution with N-M degrees of freedom, which allows you to calculate the p-value. You are encouraged to check your calculations using lm.

**Exercise 9.4.** Is the trend significant at the 5% level? To answer this, test the hypothesis  $\beta_2 = 0$  in model (9.1). Print out the output of regress.normal. Check to ensure that the t-value you calculate agrees with the t-value obtained from the R function lm().

**Exercise 9.5.** Augment regress.normal to calculate the  $(1 - \alpha)100\%$  confidence interval of each regression coefficient. Call the interval bhat.ci in the output list (it should be dimensioned [MTOT, 2]). The confidence interval can be calculated from

$$\hat{\beta}_j \pm t_{\alpha/2,N-M} s_e \sqrt{d_j}.$$
(9.5)

Give the 95% confidence interval for the slope and intercept. Is this confidence interval consistent with your answer to the previous question? You are encouraged to check your calculations using confint.

**Exercise 9.6.** Test whether sea level has a significant trend *after the effect of ENSO has been taken into account*. This test is performed by fitting the model

sea level = 
$$\beta_1 + \beta_2 * \text{year} + \beta_3 * \text{NINO3.4} + \epsilon$$
 (9.6)

and testing the hypothesis  $\beta_2 = 0$ . Print out the output of regress.normal. Use this output to form a conclusion about whether sea level has a statistically significant trend.  $\Box$ 

1

### EXERCISES FOR MODEL SELECTION

In this problem set you will estimate the annual cycle of the NINO3.4 index and apply model selection criteria to determine the number of terms in the Fourier series that should be fitted. The NINO3.4 index can be downloaded from

http://www.esrl.noaa.gov/psd/data/correlation/nina34.data

An R-code for reading this data set is ModSel.student.R, which can be downloaded from the class website. You will need to modify the variables dir.rlib and dir.enso to point to the directories containing your R-functions and ENSO time series, respectively. The variable nharm specifies the maximum number of Fourier harmonics to be considered. When you run this code, it should produce a figure showing the *raw* NINO3.4 time series (without annual cycle removed).

The exercises below will ask you to write a function to evaluate a variety of model selection criteria. Each exercise will focus on a single criterion, but in the end you should submit a single function that calculates all these criteria. This function should have the following name and preamble:

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

#### 40 EXERCISES FOR MODEL SELECTION

```
modsel.loo.simple = function(y, x, kmax=dim(x)[2]+1) {
   #### EVALUATES VARIOUS MODEL SELECTION CRITERIA FOR THE MODEL
2
   #### Y = X B + E
   #### NOTE: THE CONSTANT (INTERCEPT) SHOULD NOT BE INCLUDED IN X;
   ####
          IT IS INSERTED BY THIS PROGRAM
   #
6
   # INPUT:
7
   #
       Y[N]: PREDICTAND, WHERE N = NUMBER OF INDEPENDENT SAMPLES
8
   #
       X[N,K]: PREDICTOR MATRIX, K = NUMBER OF PREDICTORS
0
   #
          (NOT INCLUDING CONSTANT TERM)
10
   #
       KMAX: MAXIMUM NUMBER OF PREDICTORS (INCLUDING THE CONSTANT)
11
12
   #
          (DEFAULT: ALL PREDICTORS)
13
   # OUTPUT:
       SE.SQR[KMAX]: UNBIASED ERROR VARIANCE OF EACH MODEL 1, 2, ..., KMAX
   #
14
   #
       NSE.SQR[KMAX]: NORMALIZED UNBIASED ERROR VARIANCE OF EACH MODEL
15
   #
       CVMSE[KMAX]: LOO CROSS-VALIDATED MEAN SQUARE ERROR FOR EACH MODEL
16
   #
       CVSTD[KMAX]: STANDARD ERROR OF LOO CROSS-VALIDATED SQUARED ERRORS
17
   #
       AIC[KMAX], AICC[KMAX], BIC[KMAX]: INFORMATION CRITERIA FOR EACH MODEL
18
   #
       IC.STD[KMAX]: STANDARD ERROR OF THE INFORMATION CRITERIA
19
```

#### Exercise 10.1. We want to fit the regression model

$$y_n = \beta_0 + \sum_{h=1}^{H} \left( c_h \cos\left(2\pi nh/12\right) + s_h \sin\left(2\pi nh/12\right) \right)$$
(10.1)

where  $y_n$  is the NINO3.4 index and n = 1, 2, ..., N. Show that this model can be written in the form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.\tag{10.2}$$

What is X? What is  $\beta$ ? Hint: in R, the predictor matrix containing the Fourier harmonics can be generated as

x.pred = NULL
for ( nh in 1:nharm) x.pred =

3

cbind(x.pred,cos(2\*pi\*nh\*nino34.time),sin(2\*pi\*nh\*nino34.time))

**Exercise 10.2.** Write a function that computes the unbiased error variance of the regression model (10.1) as a function of the number of harmonics H. The unbiased error variance is

$$s_e^2 = \frac{SSE}{N-M}.$$
(10.3)

The input to this function is the time series y and the predictor matrix X without the constant term. Internally, you can insert the constant function as follows:

```
1 ntot = length(y)
2 x.const = cbind(rep(1,ntot),x); # ADD COLUMN OF ONES IN PREDICTOR MATRIX
```

If k is the number of columns of x.const, then the following commands will fit the regression model and produce fitted values and residuals.

```
for ( k in 1:kmax)
                        {
1
            x.pred
                         = x.const[,1:k]
2
            xy.lm
                         = lm(y^x.pred-1)
3
4
            yhat
                         = fitted(xy.lm)
                         = residuals(xy.lm)
5
            yresid
6
   }
```

Add to this code by computing  $s_e^2$  and the normalized error variance

$$\frac{s_e^2}{\hat{\sigma}_Y^2},\tag{10.4}$$

where  $\hat{\sigma}_Y^2$  is the sample variance of y. Make a plot of the normalized error variance as a function of the number of predictors. Where does the minimum occur? How many Fourier harmonics does this correspond to? How much variance is explained by the annual cycle? How much variance is unexplained by the annual cycle?

**Exercise 10.3.** Write a function that computes the leave-one-out cross validated mean square error. Hint: use "negative" arguments in R arrays to leave data out. For instance:

State the values of the leave-one-out cross validated mean square error (after normalizing by the variance of y).

**Exercise 10.4.** The above approach to leave-one-out cross validation is computationally intensive. A faster numerical method is explained in exercise 11.2. You do not have to prove exercise 11.2, but you will use the result. Write a function that computes the leave-one-out cross validated mean square error based on the numerical approach discussed in exercise 11.2. Verify that the result is identical to what you found in the previous exercise. Discuss how much faster the new function is relative to the old.

**Exercise 10.5.** Compute the standard error of the cross validated mean square error. Indicate the standard error as error bars on the figure produced in the previous problem. Error bars can be plotted using the arrows command as follows:

#### 42 EXERCISES FOR MODEL SELECTION

```
modsel.list = modsel.loo.simple(nino34.ts,x.pred)
x0 = 0:(2*nharm)
y0 = modsel.list$cvmse - modsel.list$cvstd
y1 = modsel.list$cvmse + modsel.list$cvstd
yrange = range(y0,y1)
plot(x0,modsel.list$cvmse,type='b',pch=19,xlab='number of Fourier terms',
ylab='unbiased error variance',ylim=yrange)
arrows(x0,y0,x0,y1,length=0.1,angle=90,code=3,lwd=2)
```

What model does the "one-standard-deviation rule" select?

**Exercise 10.6.** Write a function to evaluate AIC, BIC, and AICC. Plot these values as a function of the number of predictors in the annual cycle model. State the numerical values of these quantities. What model does these criteria select?

**Exercise 10.7.** Write a function that computes the standard error of the information criteria. Show the standard error in the plot generated in the previous exercise. What model does the "one-standard-error" rule select?

**Exercise 10.8.** Based on all of the above results, you should have found that either 1 or 2 harmonics should be used to fit the annual cycle of NINO3.4. Let us decide to use 2 harmonics (so this means you have 5 predictors: 1 intercept term, cos/sine for the first harmonic, and cos/sine for the second harmonic). Now subtract this annual cycle from the NINO3.4 index and plot the residual. This is called the NINO3.4 *anomaly*. You should recognize the 1982, 1998, and 2015 El Niño events.

**Exercise 10.9** (Apparently Different Information Criteria). Different definitions of AIC, AICC, BIC appear in the literature. However, these definitions differ by either a factor of N or an additive constant, neither of which affect the *location* of the minimum value. For example, some papers define  $AIC_c$  as

$$AIC_c = N \log\left(\frac{SSE}{N}\right) + \frac{2KN}{N-K-1},$$
(10.5)

where K is the total number of "parameters" in the regression model *including the error* variance  $\sigma_{\epsilon}^2$ ; that is, K = M + 1 in terms of the model defined in (??). Show that this definition differs from the  $AIC_c$  defined (10.28) by a constant and therefore the two expressions (10.28) and (10.5) are equivalent model selection criteria.

**Exercise 10.10** (Shortcut for Leave-One-Out Cross Validation). Show that the leave-one-out cross validated error of the regression model (??) can be written as

$$\epsilon_k^{LOO} = y_k - \mathbf{x}_k \boldsymbol{\beta}_k = \frac{y_k - \mathbf{x}_k \boldsymbol{\beta}}{1 - \alpha_k},\tag{10.6}$$

where  $\hat{\beta}$  is the least-squares estimate of  $\beta$  using all data and  $\alpha_k$  is the scalar

$$\alpha_k = \mathbf{x}_k \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_k^T.$$
(10.7)

Interestingly, the cross validated error  $y_n - \mathbf{x}_n \boldsymbol{\beta}_n$  is merely an *inflated* version of the residual of the least squares model derived from the entire data set  $y_n - \mathbf{x}_n \hat{\boldsymbol{\beta}}$ . Equation (10.6) also shows that the leave-one-out cross validated mean square error can be determined from the traditional least squares solution, without re-fitting the model N separate times. Although the method requires computing the inverse of  $\mathbf{X}^T \mathbf{X}$  once, this inverse is already available since it is required to compute  $\hat{\boldsymbol{\beta}}$ . Hint: note that

$$\boldsymbol{\beta}_{k} = \left(\mathbf{X}^{T}\mathbf{X} - \mathbf{x}_{k}^{T}\mathbf{x}_{k}\right)^{-1} \left(\mathbf{X}^{T}\mathbf{y} - \mathbf{x}_{k}^{T}y_{k}\right), \qquad (10.8)$$

and apply the Sherman-Morrison-Woodbury formula.

### PITFALLS OF STATISTICAL INFERENCE

**Exercise 11.1** (Selecting Predictors Using Screening). A forecaster wants to predict hurricanes for the season. To do this, the forecaster evaluates the correlation between hurricane counts and 50 variables around the globe that may be relevant. The correlations are computed using 30 years of data. Based on this analysis, the forecaster selects the predictor with the maximum correlation, which has a correlation of 0.6 with hurricane counts, and uses that predictor to make a forecast. Is the predictor statistically significant at the 5% level? For your null hypothesis, you may assume the 50 variables are independent and normally distributed. Does your conclusion change if the variables are not independent?  $\Box$ 

**Exercise 11.2** (Selecting the Best Forecasts). A forecaster has nine forecasts and would like to combine them to make a single, superior forecast. The forecaster decides to select the four forecasts that are most correlated with observations from a 20-year historical record, and then average those forecasts to make a prediction. When the four best forecasts are averaged, the correlation between the mean 4-member forecast and observations is 0.55. For reference, the 5% significance level for a 20-year forecast is about 0.45. Since the actual correlation exceeds the threshold, the forecaster claims that the 4-member mean forecast has a statistically significant correlation. What is "wrong" with this argument? What is the experimentwise error rate for this procedure? Plot a histogram of the correlation skill under the null hypothesis of independent variables.

Hint: generate random numbers to simulate the above scenario and perform the same steps as the forecaster. Repeat this procedure numerous times to estimate an empirical distribution of the correlation under the null hypothesis of no-skill (all variables are in-

dependent). You can use the cor function to determine the correlation between 1 variable and 9 others (e.g., if obs is 20 elements long, and frcs is a  $20 \times 9$  matrix, then as.numeric(cor(frcs, obs)) gives 9 correlations). Also, the order command allows you to identify the forecasts with the four largest correlations.

**Exercise 11.3** (Discovery of a new Predictor). An investigator believes that a particular index, called *DOOM*, is useful for predicting JFM temperature. The 50-year correlation between JFM temperature and *DOOM* is 0.2 in December and 0.1 in November, but 0.3 in October. Since the 5% significance threshold for a correlation coefficient based on 50 independent samples is about 0.28, the investigator promptly publishes a paper proclaiming that the October value of *DOOM* is a good predictor of JFM temperature. Has the significance test been applied properly? If not, why not? Is the predictor significant when screening is taken into account? If not, how large would the correlation need to be to ensure a true 5% significance?

# EXERCISES FOR PRINCIPAL COMPONENT ANALYSIS

In this problem set you will calculate the principal components of December mean sea surface temperature (SST). To do this, you will need to download some R code, data sets, and install appropriate libraries. This homework set will walk you through these steps.

First, you need prepare R to make plots of spatial fields. The libraries for making spatial plots are not automatically included in R, so you must install them. The libraries you need are the following:

- maps
- fields
- rworldmap
- ncdf4
- chron

On the MAC, it is very easy to install libraries:

- 1. Click "Packages & Data."
- 2. Click "Package Installer."

- 3. Click "Get list" (if a list is not already there).
- 4. Click the appropriate library (see above list).

Next, you need to download the data. The data file is sst.mnmean.v4.nc.

Finally, you should download two R programs

- plot\_latlon\_v4.R for plotting spatial fields
- sst.pca.student.R, which reads this data and plots the SST for a single year.

As usual, before running sst.pca.student.R, you should modify the directory names dir.data and dir.lib to correspond to the directories where you save data and R functions, respectively. If you have successfully followed all of the above directions, then sst.pca.student.R will run with no errors or warnings and will generate a plot showing SST in December 1983.

To calculate a principal component, you must first calculate the climatological mean, and then subtract it out from all fields. Remember that R is faster for vector calculations than for for loops. Therefore, we want to calculate the climatological mean using vector calculations. A trick for doing this is to *re-shape* the data array so that the command rowMeans can be used. When converting from vector to matrix or vice versa, R assumes the elements are stored *column-wise*. You can get the idea from the following R session:

```
> dum = letters[1:12]
  > print(dum)
2
    [1] "a" "b" "c" "d" "e" "f" "q" "h" "i" "j" "k" "l"
  > \dim(dum) = c(3, 4)
4
  > print(dum)
        [,1] [,2] [,3] [,4]
   [1,] "a"
             "d"
                   "g"
                        "j"
7
   [2,] "b"
             "e"
                   "h"
                        "k"
8
   [3,] "c"
             "f"
                   "i"
                        ייןיי
```

Taking advantage of this convention, the monthly anomalies can be calculated efficiently as follows:

```
1 dim(sst) = c(nlon*nlat,nyrs)
2 sst.clim = rowMeans(sst)
3 sst.anom = sst.data - sst.clim
```

In addition to subtracting out the mean, you also need to create weights that take into account the fact that the area of a grid point is proportional to the cosine of latitude. A concise way of generating the appropriate vector of weights is the following.

Make sure you understand what both commands do and why!

Note that the array sst.anom.weight is shaped as [space, time], whereas the notes assumes the transpose. There is a good reason why I define the data array the way that I did, but let's not get into that. There are two ways to deal with this inconsistency: take the transpose of the data array before performing PCA, or perform SVD on the data array but switch the labels U and V. The second approach is taken in this homework.

Another pesky issue is that SST is defined only over the ocean- the grid points over land have no SST value. The undefined points need to be "stripped away" before computing the EOFs. The basic idea is to construct a new (smaller) data matrix in which each column contains *only* the defined ocean points, apply PCA to the resulting data matrix, then "fill in" the undefined land points at the end of the calculation. Since this is a programming trick and not statistics, I'll just give you the answer for a generic data set data.array:

```
1
   # IDENTIFY MISSING OR 0-WEIGHTED DATA
   dim(data.array) = c(nlon*nlat, ntot)
2
   lbad = is.na(rowSums(data.array)) | weight == 0
3
Δ
   # COMPRESS DATA BY ELIMINATING MISSING GRID POINTS
5
   data.array = data.array[!lbad,]
6
7
   ndef = sum(!lbad)
   # COMPUTE SVD
9
   mmin = min(ndef, ntot)
10
   data.svd = svd(data.array,nu=mmin,nv=mmin)
11
12
13
   # FILL IN MISSING POINTS IN EOF
14
   eof = array(NA, dim=c(nlon*nlat, mmin))
   eof [!lbad,] = data.svd$u
15
```

**Exercise 12.1.** Note that the data array sst is shaped as [space, time], whereas the lecture notes assumed the data array is [time, space]. This means that if you take the SVD of sst, the identification of U and V differs from that of the notes. Starting from the SVD of the weighted data

$$\mathbf{W} \quad \mathbf{Y}' = \dot{\mathbf{U}} \quad \dot{\mathbf{S}} \quad \dot{\mathbf{V}}^T \\ M \times M \quad M \times N \quad M \times N \quad N \times N \quad N \times N,$$
(12.1)

where  $\mathbf{W}$  is a diagonal matrix whose diagonal elements equal the square root of the cosine of latitude, state the equations for the scaled EOFs  $\mathbf{E}$  and scaled PCs  $\mathbf{F}$  such that

$$\mathbf{Y}' = \mathbf{E}\mathbf{F}^T. \tag{12.2}$$

Show mathematically that the scaled PCs have unit variance.

**Exercise 12.2.** Based on the above hints, write a function called eof.latlon.simple that performs principal component analysis on a data array formatted as [space, time]. The preamble of this function should be the following:

```
eof.latlon.simple <- function(lon, lat, data.array, neof=30) {</pre>
   # COMPUTE PRINCIPAL COMPONENTS OF A DATA ARRAY.
4
   # INPUT:
5
   # LON[NLON]: VECTOR SPECIFYING LONGITUDES
6
   # LAT[NLAT]: VECTOR SPECIFYING LATITUDES
   # DATA.ARRAY[NLON,NLAT,NTOT] OR DATA.ARRAY[NLON*NLAT,NTOT]:
       THE DATA ARRAY IN [SPACE1, SPACE2, TIME] OR [SPACE, TIME] FORMAT
   #
   # NEOF = NUMBER OF SPATIAL EOFS (WITH MASK) TO BE OUTPUTTED. DEFAULT = 30
10
11
   # OUTPUT LIST:
12
   # $EOF[NLON*NLAT, NEOF]: THE FIRST NEOF SCALED EOFS
13
  # $PC[NTOT, NEOF]: THE PCS, NORMALIZED TO UNIT VARIANCE
14
  # $FEXPVAR: FRACTION OF EXPLAINED VARIANCE FOR EACH EOF (ALL OF THEM).
15
  # $FEXPVAR.CI: CONFIDENCE INTERVALS FOR FRACTION OF EXPLAINED VARIANCE.
17
  # $EOFI[NLON*NLAT, NEOF]: PSEUDO INVERSE OF EOF (I.E., T(EOFI) %*% EOF = I)
   # $NEOF: MINIMUM OF (NEOF IN ARGUMENT LIST, RANK OF DATA.ARRAY)
18
```

This function should subtract out the time mean, apply an area weighting, allow for undefined points, and produce a list containing the quantities indicated above.  $\Box$ 

**Exercise 12.3.** Apply your function to compute the principal components of December SST. Show a plot of the spatial structure and time series of the leading component. Also, show a plot of the explained variances of the leading components along with confidence intervals. These plots should be similar to those displayed in the notes.

Exercise 12.4. Consider the linear combination

$$r = a_1 X_1 + a_2 X_2 + \dots + a_M X_M = \mathbf{a}^T \mathbf{x},$$
 (12.3)

where  $a_1, \ldots, a_M$  are constants and  $X_1, \ldots, X_M$  are random variables. A *standardized linear combination* is one in which the coefficients **a** are constrained such that  $\mathbf{a}^T \mathbf{a} = 1$ .

(a) Show that  $\mathbf{b} = \mathbf{a}/\sqrt{\mathbf{a}^T \mathbf{a}}$  satisfies  $\mathbf{b}^T \mathbf{b} = 1$  for any choice of  $\mathbf{a} \neq \mathbf{0}$ . It follows from this that

$$\mathbf{b}^T \mathbf{x} = \frac{\mathbf{a}^T \mathbf{x}}{\sqrt{\mathbf{a}^T \mathbf{a}}} \tag{12.4}$$

is a standardized linear combination for *any* choice of  $\mathbf{a} \neq \mathbf{0}$ .

(b) Write an expression for the variance of r in terms of the covariance matrix of x,  $\Sigma$ .

(c) Consider the new variable

$$\tilde{\mathbf{a}} = \mathbf{U}^T \mathbf{a},\tag{12.5}$$

where U is the orthogonal eigenvector matrix of  $\Sigma$  (i.e,  $\Sigma = U\Lambda U^T$ , where  $U^T U = I$ ). Write the variance of r in terms of  $\tilde{a}$ . Something "special" happens under this transformation. What is it? Write out a *simple* expression for the variance of r without using vector notation (explicitly write out sums and products).

(d) Show that maximum variance of r is  $\lambda_1$  (i.e., the largest eigenvalue of  $\Sigma$ ). What vector  $\tilde{a}$  gives this maximum? What vector a gives this maximum?

(e) Suppose the covariance matrix is estimated from the sample covariance matrix of the *centered* data matrix **X**:

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}.$$
(12.6)

How are the eigenvalues and eigenvectors of  $\hat{\Sigma}$  related to the singular values and singular vectors of  $\mathbf{X}$ ?

(f) Use these results to give an alternative description of the leading principal component. That is, instead of describing the leading component as something that best approximates the data, describe it in terms of standard linear combinations.

### COMPUTATIONAL EXERCISES ON FIELD SIGNIFICANCE

In this homework set, you will perform a series of field significance tests. The first step of this homework is to obtain the following files from the class website. It is good practice to put all of your data in a single directory (say "data"), and all your auxiliary programs into a single directory that can be referred to in future R programs.

fieldsig.student.R	main program: reads data and computes EOFs
cca.data.clim763.DJF.RData	data (will be used in future homeworks)
eof.latlon.R	auxiliary file: calculate EOFs
plot_latlon_v4.R	auxiliary file: plot spatial maps
gev.R	auxiliary file: function to solve eigenvalue problem
index.climate.v2.R	auxiliary file: function to compute NINO3.4 index

Your starting point is the R program fieldsig.student.R, which does the following things:

• It reads the data file cca.data.clim763.DJF.RData and loads them in the arrays tem.data and sst.data. Inside fieldsig.student.R, you need to modify the variable dir.data to correspond to the directory to which you down-loaded the data.

- 54 COMPUTATIONAL EXERCISES ON FIELD SIGNIFICANCE
  - It reads the auxiliary programs (e.g., eof.latlon.R and plot\_latlon\_v4.R). You will need to modify dir.Rlib to correspond to the directory to which you downloaded the auxiliary programs.
  - It computes the nino3.4 index from sst.data.
  - It also computes the EOFs of U.S. temperature using eof.latlon, which creates the list variable tem.eof. The EOFs are tem.eof\$eof[space,neof], and the pc time series are tem.eof\$pc[time,neof].
  - It plots the leading EOF and PC of North American Temperature and SST. The resulting figure should look like fig. 13.1.

Note: temperature over the ocean is set to NA. Also, temperature is read as an array with dimension nlon, nlat, time, but eventually is "reshaped" into an array with dimension nlon\*nlat, nyrs. If you do not understand what this means, please see me!

The purpose of this homework is to test the hypothesis that tropical Pacific SSTs influence North American temperature in DJF. You will begin by looking at the correlation map between temperature and NINO3.4. This is very convenient to plot in R. First, we compute the correlation between NINO3.4 and every grid point of temperature using the command

cor.map = cor(t(tem.data),nino34)

The R-function cor assumes data matrices are in the form TIME x SPACE, so we must take the transpose of tem.data (which is an array with dimensions SPACE x TIME) to use cor. To plot this, we then use the command

plot\_latlon\_v4(lon,lat,cor.map).

To "maskout" insignificant correlations, we use the commands

```
cor.crit = 2/sqrt(nyrs)
cor.map[abs(cor.map)<cor.crit] = NA
plot_latlon_v4(lon,lat,cor.map)
```

The resulting plot should look like fig. 13.1 in the notes.

### Assignment

**Exercise 13.1.** Write an R code that applies the Livezey-Chen field significance test. The preamble of the function should be the following:

```
livezey.chen = function(lon,lat,xdata,ydata,ntrials=1000,alpha=0.05) {
1
    ### PERFORMS THE LIVEZEY-CHEN FIELD SIGNIFICANCE TEST
2
    # INPUT:
3
        LON[NLON]: LONGITUDE OF THE FIELD DATA
4
        LAT[NLAT]: LATITUDE OF THE FIELD DATA
5
        XDATA[NLON, NLAT, NTOT]: DATA ARRAY FOR THE FIELD
6
7
    #
        YDATA[NTOT]: REFERENCE TIME SERIES
       NTRIALS: NUMBER OF MONTE CARLO TRIALS (DEFAULT = 10000)
8
        ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
    # OUTPUT:
10
        HISTOGRAM OF AREAS WITH SIGNIFICANT CORRELATIONS
11
    #
12
       LIST
    #
          SCRIT.NULL: SIGNIFICANCE THRESHOLD FROM LIVEZEY-CHEN TEST
13
14
    #
          $SIG.AREA.OBS: PERCENTAGE AREA WITH SIGNIFICANT CORRELATIONS
```

Apply the Livezey-Chen field significance test to test whether NINO3.4 is related to DJF temperature over North America. Make a histogram of the Monte Carlo results.

Be sure to express the answer in terms of percentage <u>area</u> with significant correlations. Suggestion: create another field that gives the fractional area at each grid point, then sum over those points. For instance:

```
1 area.frac = rep(cos(lat*pi/180),each=nlon)
2 lgood = !is.na(xdata[,1])
3 area.frac[!lgood] = NA
4 area.frac = area.frac/sum(area.frac,na.rm=TRUE)
5
6 cor.map = as.numeric(cor(t(xdata),ydata))
7 sig.area.obs = sum(area.frac[abs(cor.map) > 2/sqrt(ntot)],na.rm=TRUE)*100
```

You should obtain a plot similar to fig. 13.4 of the lecture notes. Make a conclusion regarding whether the field is significant.  $\Box$ 

**Exercise 13.2.** Write an R code that applies the permutation field significance test. The preamble of the function should be the following:

```
fieldsig.permutation = function(lon,lat,xdata,ydata,ntrials=1000,alpha=0.05) {
1
   ### PERFORMS PERMUTATION FIELD SIGNIFICANCE TEST
2
   # INPUT:
3
       LON[NLON]: LONGITUDE OF THE FIELD DATA
4
   #
       LAT[NLAT]: LATITUDE OF THE FIELD DATA
5
   #
       XDATA[NLON, NLAT, NTOT]: DATA ARRAY FOR THE FIELD
6
       YDATA[NTOT]: REFERENCE TIME SERIES
   #
7
      NTRIALS: NUMBER OF TRIALS (DEFAULT = 10000)
   #
8
       ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
   #
   # OUTPUT:
10
   #
       HISTOGRAM OF AREAS WITH SIGNIFICANT CORRELATIONS
11
   #
       LIST:
12
   #
         $CRIT.NULL: SIGNIFICANCE THRESHOLD FROM LIVEZEY-CHEN TEST
13
   #
         $SIG.AREA.OBS: PERCENTAGE AREA WITH SIGNIFICANT CORRELATIONS
14
```

The sample function is very useful for generating a permutation. For instance, sample (N) generates a permutation of the integers 1-N. To ensure there are no accidental matches, the following loop may be helpful

56 COMPUTATIONAL EXERCISES ON FIELD SIGNIFICANCE

Answer the same questions as the previous exercise. How does this result compare to that from the Livezey-Chen test?  $\Box$ 

**Exercise 13.3.** Write a function that evaluates Mutual Information Criterion (MIC) for the field significance problem. MIC is a function of R-square, which can be computed from R functions as follows

r.square = summary(lm(y<sup>x</sup>))\$r.square

Show a plot of MIC versus the number of EOFs of U.S. temperature and identify the minimum. The figure should be similar to fig. 13.5 in the notes.  $\hfill \Box$ 

**Exercise 13.4.** Write a code to evaluate the significance level for MIC. Include this in your function for MIC. Plot the significance level.  $\Box$ 

**Exercise 13.5.** Write an R code that computes the False Discovery Rate. The preamble of this function should be the following:

```
fdr.fieldsig = function(y,x,fdr=0.1) {
1
   ## APPLIES FALSE DISCOVERY RATE TO A CORRELATION MAP
2
   ## INPUT:
3
         Y[NTOT]: REFERENCE TIME SERIES
   #
4
         X[NTOT,MTOT]: FIELD TIME SERIES
   #
5
   #
         FDR: FALSE DISCOVERY RATE
6
   ## OUTPUT:
7
         LREJECT[MTOT]: LOGICAL VECTOR INDICATING REJECTIONS
   #
8
```

Make a plot showing where the FDR implies rejection of the null hypothesis. The plot should be similar to fig. 12.7 of the notes.  $\Box$ 



Figure 13.1 Leading EOF of DJF North American temperature during 1982-2017 (top) and correlation between DJF NINO34 index and DJF temperature during 1982-2017 (bottom).

### EXERCISES FOR MULTIVARIATE LINEAR REGRESSION

**Exercise 14.1.** It is natural to wonder whether a "better" estimate of B can be obtained by weighting the SSE toward certain variables. To explore this possibility, suppose  $w_s$  is the weight for the s'th spatial location. Then, the SSE in (14.20) can be modified as

$$SSE = \sum_{s=1}^{S} \sum_{n=1}^{N} w_s \left( Y_{ns} - \sum_{m=1}^{M} X_{nm} B_{ms} \right)^2.$$
(14.1)

More generally, one can specify a positive-definite weight matrix W and define

$$SSE = \sum_{s'=1}^{S} \sum_{s=1}^{S} \sum_{n=1}^{N} W_{ss'} \left( Y_{ns} - \sum_{m=1}^{M} X_{nm} B_{ms} \right) \left( Y_{ns'} - \sum_{m=1}^{M} X_{nm} B_{ms'} \right), \quad (14.2)$$

which can be written equivalently as

$$SSE_W = \operatorname{tr}\left[\left(\mathbf{Y} - \mathbf{XB}\right)\mathbf{W}\left(\mathbf{Y} - \mathbf{XB}\right)^T\right].$$
 (14.3)

Show that the matrix **B** that minimizes  $SSE_W$  is independent of **W** and equal to the least squares estimate (14.25). This means that there is nothing to gain by weighting different spatial locations differently. Explain why this is so.

### **Computational Exercise for Multivariate Linear Regression**

60 EXERCISES FOR MULTIVARIATE LINEAR REGRESSION

In this homework, you will derive a multivariate linear prediction model for Sea Surface Temperature.

**Getting the Data** The data set will be the ERSSTv5,<sup>2</sup> which is on the COLA servers.<sup>3</sup> However, you might prefer to download the data set directly on your laptop (I would recommend this). To download this data using R, grab the following files on the class website:

download.ersstv5.R

```
download_ftp_file.R
```

```
3 create_directory.R
```

Then type source ('download.ersstv5.R'). The code will take a few minutes to run and will show the file names being downloaded. It also will generate a warning message about creating directories, which can be ignored.

**Reading the Data** To make this assignment easier, I have written an R code called mlinreg.student.R that does the following things:

- 1. read the NetCDF files for ERSSTv5
- 2. load the data into the array sst.full[space,month,year]
- 3. regress out the mean and linear trend from each grid point
- 4. define the Pacific Ocean domain to be between  $30^{\circ}$ S and  $60^{\circ}$ N
- 5. compute the EOFs of monthly SST in the prescribed domain

This code also requires the following files from the class website:

```
eof.latlon.R
```

```
2 mask.define.R
```

It is good practice to put all of your data in a single directory (say "data"), and all your auxiliary programs into a single directory that can be referred to in future R programs. You should change dir.Rlib to point to the directory of your R functions, and dir.ersst to point to the directory of your ERSSTv5 data.

After running mlinreg.student.R, you should see a plot of the leading EOF and PC. If the code crashes or gives an error message, please see me. Otherwise, you are ready! After the code runs to completion, you should have the following data arrays available.

```
# 1) THE GRIDDED DATA SET IS IN SST.FULL[SPACE,MONTH,YEAR]
```

```
2 # 2) EOFS ARE IN EOF.LIST
```

```
# (E.G., EOF.LIST$EOF[SPACE,NEOF], EOF.LIST$PC[TIME,NEOF])
```

<sup>2</sup>https://www.ncdc.noaa.gov/data-access/marineocean-data/extended-reconstructed-sea-surface-temperatureersst-v5

```
3/shared/obs/sst/ncdc/ersst-v5
```

<sup># # 3)</sup> LAT/LON IS LON.SST, LAT.SST

<sup># 4)</sup> MON.INIT IS THE CALENDAR MONTH OF THE INITAL CONDITION

<sup># 5)</sup> MON.TARG IS THE CALENDAR MONTH OF THE TARGET

Temperature over the ocean is set to NA. Also, temperature is read as an array with dimension nlon, nlat, time, but eventually is "reshaped" into an array with dimension nlon\*nlat, 12, nyrs. If you do not understand what this means, please see me!

**Exercise 14.2.** Write a R function to compute MIC given two data arrays. The function also should compute the 5% significance threshold for MIC. The preamble should be as follows

```
mic.gaussian = function(x,y,equal.dim=TRUE,alpha=0.05) {
   ### THIS FUNCTION COMPUTES CORRECTED MUTUAL INFORMATION CRITERION (MIC)
2
   ###
           FOR Y = XB + E
   ### INPUT:
   ##
         X[NSAMP,XDIM]: X DATA ARRAY, OFTEN FORMATTED AS [TIME,EOF]
   ##
         Y[NSAMP, YDIM]: Y DATA ARRAY, OFTEN FORMATTED AS [TIME, EOF]
   ##
         EQUAL.DIM: LOGICAL INDICATING WHETHER
   ##
           (TRUE) EQUAL NUMBER OF X'S, Y'S CHOSEN: MIC[MIN(XDIM, YDIM)];
8
           (FALSE) MIC FOR ALL TRUNCATIONS ARE COMPUTED: MIC[XDIM, YDIM]
   ##
   ### OUTPUT: LIST (DIMENSIONS DEPEND ON EQUAL.DIM)
10
   ##
         $MIC: MIC VALUES
11
   ##
         $PENALTY: THE PENALTY TERM IN MIC
12
   ##
         $CRIT: SIGNIFICANCE THRESHOLD OF MIC
13
```

**Exercise 14.3.** Consider the case of predicting December SST based on June SST. Calculate the MIC for EOF truncations at least for 1-25. The MIC values should be the following:

1	> mic	CC				
2	[1]	-0.93434963	-1.23487398	-1.61938679	-1.67394748	
3	[5]	-1.84103698	-1.50387241	-1.19900423	-0.69182271	
4	[9]	0.05407005	0.33422799	1.48966123	2.50855073	
5	[13]	4.09619446	5.50819524	7.95854904	11.03580560	
6	[17]	14.07453324	18.15159799	23.34799814	30.02343607	
7	[21]	38.16355309	48.57419513	62.29384427	79.60195564	
8	[25]	101.97693930	131.36865590	172.06833168	227.00353490	
9	[29]	312.44543534	452.13622542	710.10017386	1323.74697742	
0	[33]	4435.93164492	NA			

Plot the MIC values and indicate the minimum, which should be 5 EOFs. Next, do the same thing, except for predicting December SST based on September initial condition. Print the values of MIC (as above). Decide how many EOFs should be used and state your answer.

**Exercise 14.4.** Use multivariate regression to make a prediction of December SSTs based on June SSTs using 5 EOFs. Submit your code for doing this. Make sure your results are consistent with the notes. Then, make a prediction for December SSTs using September SSTs and 5 EOFs. Submit plots of your predictions.

# EXERCISES FOR CANONICAL CORRELATION ANALYSIS

cca.student.R	main program
cca.data.clim763.DJF.RData	SST and N. American temperature for DJF
eof.latlon.R	auxiliary file: calculate EOFs
plot_latlon_v4.R	auxiliary file: plot spatial maps

In this homework you will write an R function that performs Canonical Correlation Analysis. You will need to download data and a few R programs. These files are summarized in the above table. The R code cca.student.R reads the data file

cca.data.clim763.DJF.RData and computes the EOFs of SST and U.S. temperature, and plots out the leading EOF and PC time series. In a previous homework you investigated the relation between ENSO and U.S. temperature using the concept of field significance. However, that analysis assumed you knew the NINO3.4 index. Here, you will *derive* an index for predicting U.S. temperature based on SST.

The following exercises break up CCA into discrete steps. However, in the end, you should submit a single function that performs all the calculations. The preamble of this function should be the following:

```
cca.pca = function(x.pca,y.pca,tx=NULL,ty=NULL) {
   *****
   ## PERFORMS CANONICAL CORRELATION ANALYSIS ON X AND Y.
   ## X AND Y ARE ASSUMED TO BE IN THE FOLLOWING FORMS:
   ## X = FX %*% EX^T ### AND ### Y = FY %*% EY^T
   ## WHERE FX AND FY HAVE COVARIANCE MATRICES = I
   ## FOR EXAMPLE: FROM PRINCIPAL COMPONENT ANALYSIS
7
   ## IF TX OR TY = NULL, THEN BOTH (TX, TY) SELECTED USING MIC
   ## INPUT:
   #
       X.PCA: LIST OUTPUT FROM EOF.LATLON[X-DATA]
10
   #
       Y.PCA: LIST OUTPUT FROM EOF.LATLON[Y-DATA]
11
12
   #
       TX: TRUNCATION FOR X (TX <= MX); IF NULL, TX IS SELECTED
13
   #
       TY: TRUNCATION FOR Y (TY <= MY); IF NULL, TY IS SELECTED
   ## OUTPUT LIST:
14
   #
       MIC[MX, MY]: MUTUAL INFORMATION CRITERION
15
       NMIN[1,2]: VALUES OF TX, TY THAT MINIMIZES MIC
   #
16
   #
       CAN.COR[MIN(TX,TY)]: CANONICAL CORRELATIONS
17
   #
       RX[NTOT, MIN(TX, TY)]: CANONICAL VARIATES FOR X
18
   #
       RY[NTOT, MIN(TX, TY)]: CANONICAL VARIATES FOR Y
19
       PX[SX ,MIN(TX,TY)]: CANONICAL LOADING VECTORS FOR X
   #
20
   #
       PY[SY ,MIN(TX,TY)]: CANONICAL LOADING VECTORS FOR Y
21
   #
       QX.TILDE[TX,MIN(TX,TY)]: WEIGHTING VECTORS FOR X-FEATURES
22
   #
       QY.TILDE[TY,MIN(TX,TY)]: WEIGHTING VECTORS FOR Y-FEATURES
23
   #
       TX, TY: SELECTED VALUES OF TX AND TY
24
25
   *****
```

Following the notes, the data sets are assumed to have been decomposed into the form

$$\mathbf{X} = \mathbf{F}_X \quad \mathbf{E}_X^T \\ N \times S_X \quad N \times M_X \quad M_X \times S_X,$$
(15.1)

and

$$\mathbf{Y} = \mathbf{F}_{Y} \quad \mathbf{E}_{Y}^{T}$$
$$N \times S_{Y} \quad N \times M_{Y} \quad M_{Y} \times S_{Y},$$
(15.2)

where

$$\frac{1}{N-1}\mathbf{F}_X^T\mathbf{F}_X = \mathbf{I} \quad \text{and} \quad \frac{1}{N-1}\mathbf{F}_Y^T\mathbf{F}_Y = \mathbf{I}.$$
 (15.3)

This decomposition is accomplished in cca.student.R using principal component analysis. You will be performing CCA on the time series matrices  $\mathbf{F}_X$  and  $\mathbf{F}_Y$ .

**Exercise 15.1.** In a previous homework you wrote a function that computed Mutual Information Criterion (MIC). Augment that function to compute MIC for all truncations  $T_X$  and  $T_Y$ , up to some maximum number (what \*is\* the maximum number?). As a reminder, MIC is defined as

$$MIC = \sum_{i} \log(1 - \hat{\rho}_i^2) + \mathbb{P}, \qquad (15.4)$$

where  $\hat{\rho}_i$  is the *i*'th sample canonical correlation and

$$\mathbb{P} = (N+1) \left( \frac{T_X + T_Y}{N - T_X - T_Y - 2} - \frac{T_X}{N - T_X - 2} - \frac{T_Y}{N - T_Y - 2} \right).$$
(15.5)

Plot MIC for a range of values and identify the value of  $T_X$  and  $T_Y$  that minimizes MIC. Print out the values for the first 5 rows and columns (i.e., print mic[1:5,1:5]). These results should match the example in the notes. Where is the minimum value, and what is it? The location of the minimum can be found using the following R commands:

n mmin = which(mic == min(micc,na.rm=TRUE),arr.ind=TRUE)
tx = nmin[1]
ty = nmin[2]

**Exercise 15.2.** Compute all canonical correlations for the optimum choice of  $T_X$  and  $T_Y$ . State the values. These values should be consistent with those in the notes.

**Exercise 15.3.** Write a function that computes canonical variates. Compute the leading canonical variates between SST and U.S. temperature using the same truncation parameters as above. Make a plot that shows the leading canonical variate for the two data sets. Verify that the sample covariance matrices of the canonical variates equals the identity matrix. The sample covariance matrix of rx can be obtained using the R command cov(rx).

**Exercise 15.4.** Write a function that computes canonical loading vectors. Compute the leading canonical loadings between SST and U.S. temperature using the same truncation parameters as above. Make a plot that shows the leading canonical loadings for the two data sets.

**Exercise 15.5.** Write a function to compute the fraction of variance explained by the canonical component. This is somewhat tricky because (1) the area weighting needs to be included, (2) missing data should be skipped, and (3) the total variance needs to be computed to compute the fraction. All of this information is available from the output of eof.latlon. To help you out, here is the way to compute it for the EOFs of SST:

```
var.x.tot = sum(sst.eof$sval^2)/(nyrs-1)
exp.var.x = rep(NA,dim=tx)
for ( n in 1:tx) exp.var.x[n] = sum(px[!sst.eof$lbad,n]^2 *
sst.eof$weight[!sst.eof$lbad]^2)/var.x.tot
```

Verify that the sum of the fractional variances equals the sum of the fractional variances of the first  $T_X$  EOFs sum(sst.eof\$fexpvar[1:tx]). You might have to compute *all* of the singular vectors using the command svd(cov.mat,nu=tx,nv=ty). State the explained variances of the canonical components for SST and for U.S. Temperature. These variances should match those in the notes.

**Exercise 15.6.** Write a *separate* code that computes the 5% significance levels of the canonical correlations based on 5000 trials of Monte Carlo experiments. State the critical values for each canonical correlation based on the optimum choice of  $T_X$  and  $T_Y$ . Use these results to decide whether the canonical correlations computed from data are significant. Clearly state your conclusion. (Coding advice: first use only 100 trials until you have debugged your code, and then change to 5000 trials when the code is working.)

**Exercise 15.7.** Repeat the above steps, except this time *do not detrend the data*. To prevent detrending, set npoly =0 at the top of the code. How do the results differ? Explain why this makes sense.

**Exercise 15.8.** Show mathematically that the error covariance matrix for the model (15.80) is

$$\tilde{\boldsymbol{\Sigma}}_{\epsilon Y} = \mathbf{I} - \hat{\mathbf{S}}_{\rho}^2. \tag{15.6}$$

This result shows that the error covariance matrix is diagonal, and that the diagonal elements equal  $1 - \hat{\rho}_k^2$ , which can be interpreted as the variance explained by the canonical variate.
# EXERCISES: COVARIANCE DISCRIMINANT ANALYSIS

**Exercise 16.1.** Prove (16.14). Hint, substitute (16.12) into the outer product and expanding gives

$$\left(\mathbf{x}_{t}-\hat{\boldsymbol{\mu}}_{Y}\right)\left(\mathbf{x}_{t}-\hat{\boldsymbol{\mu}}_{Y}\right)^{T}=\mathbf{A}+\mathbf{B}+\mathbf{C},$$
(16.1)

where

$$\mathbf{A} = \left(\mathbf{u}_t^* - \hat{\boldsymbol{\mu}}_Y\right) \left(\mathbf{u}_t^* - \hat{\boldsymbol{\mu}}_Y\right)^T$$
(16.2)

$$\mathbf{B} = \mathbf{f}_t \mathbf{f}_t^T \tag{16.3}$$

$$\mathbf{C} = \left(\mathbf{u}_t^* - \hat{\boldsymbol{\mu}}_Y\right) \mathbf{f}_t^T + \mathbf{f}_t \left(\mathbf{u}_t^* - \hat{\boldsymbol{\mu}}_Y\right)^T.$$
(16.4)

Show that  $\mathbb{E}[\mathbf{C}] = \mathbf{0}$ . What is  $\mathbb{E}[\mathbf{A}]$ ?

**Exercise 16.2.** Suppose  $d(\sigma_X^2, \sigma_Y^2)$  is some measure of the difference between the variances  $\sigma_X^2, \sigma_Y^2$ . Suppose further that the function d(,) is invariant to an invertible linear transformation of X and Y. Prove that d(,) can depend only on the ratio of variances

$$d(\sigma_X^2, \sigma_Y^2) = f\left(\frac{\sigma_X^2}{\sigma_Y^2}\right)$$
(16.5)

**Exercise 16.3** (Discriminant analysis via SVD). Given two data matrices **X** and **Y**, solve discriminant analysis based on the singular value decomposition (SVD) (i.e., without solv-

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

ing an eigenvalue problem). Show how the discriminant ratios, variates, and loading vectors can be derived from the results of the SVD. Hint: compute the SVD of  $\mathbf{Y}$  to construct a whitening transformation, then compute the SVD of the whitened data matrix  $\mathbf{X}$ .

Exercise 16.4 (Loading Patterns). Show that the matrix P that minimizes

$$\gamma_Y = E\left[\|\mathbf{y} - E[\mathbf{y}] - \mathbf{Pr}_Y\|_W^2\right],\tag{16.6}$$

is

$$\mathbf{P} = \operatorname{cov}[\mathbf{y}, \mathbf{r}_Y] \left( \operatorname{cov}[\mathbf{r}_Y, \mathbf{r}_Y] \right)^{-1}, \qquad (16.7)$$

thereby proving (16.64). Also, show that the matrix  $\mathbf{P}$  that minimizes

$$\dot{\gamma}_Y = \|\dot{\mathbf{Y}} - \mathbf{R}_Y \mathbf{P}^T\|_W^2. \tag{16.8}$$

is

$$\dot{\mathbf{P}} = \dot{\mathbf{Y}}^T \mathbf{R}_Y \left( \mathbf{R}_Y^T \mathbf{R}_Y \right)^{-1}, \qquad (16.9)$$

thereby proving (16.94). Note that  $\dot{\mathbf{P}}$  and  $\mathbf{P}$  are merely the sample and population versions of each other.

**Exercise 16.5.** Prove that if  $\Sigma_X \neq \Sigma_Y$ , then there exists a **q** such that  $\lambda \neq 1$ , where

$$\lambda = \frac{\mathbf{q}^T \boldsymbol{\Sigma}_X \mathbf{q}}{\mathbf{q}^T \boldsymbol{\Sigma}_Y \mathbf{q}}.$$
(16.10)

**Exercise 16.6.** Define the sum total variance of  $\dot{\mathbf{Y}}$  as

$$\|\dot{\mathbf{Y}}\|_{W}^{2} = \frac{1}{N_{Y} - 1} \operatorname{tr} \left[ \dot{\mathbf{Y}} \mathbf{W} \dot{\mathbf{Y}}^{T} \right], \qquad (16.11)$$

where  $\mathbf{W}$  is a positive definite matrix defining how different points are weighted. If  $\dot{\mathbf{Y}} = \mathbf{R}_Y \dot{\mathbf{P}}_Y^T$ , use the properties of CDA to show explicitly that the sum total variance can be written as

$$\|\dot{\mathbf{Y}}\|_W^2 = \sum_{k=1}^T \mathbf{p}_k^T \mathbf{W} \mathbf{p}_k.$$
(16.12)

#### Numerical Exercises

discr.student.R	program for reading the data sets
tas_Amon_CCSM4_historical_r1i1p1_185001-200512.nc	20C, 1st member
tas_Amon_CCSM4_historical_r2i1p1_185001-200512.nc	20C, 2nd member
tas_Amon_CCSM4_piControl_r1i1p1_080001-130012.nc	PI simulation

In this homework you will write an R function to perform covariance discriminant analysis. You will need to download data and a few R programs. These files are summarized in the above table. The R code discr.student.R reads data files for the 20C and PI simulations, combines them, and computes EOFs. The following exercises break up the discriminant analysis into discrete steps, but you should submit a single function that performs all the calculations. The preamble of this function should be the following:

```
cda.eof = function(xdata,ydata,eof.list) {
1
   ### PERFORMS COVARIANCE DISCRIMINANT ANALYSIS ON X AND Y
2
   ### INPUT:
   ###
         XDATA[NX, MDIM]
4
   ###
         YDATA[NY, MDIM]
         EOF.LIST: LIST FROM EOF CALCULATION
   ###
   ### OUTPUT:
7
         MIC[NEOF]: MIC AS A FUNCTION OF NUMBER OF PCS
   ###
8
   ###
         NMIN: LOCATION OF MINIMUM MIC
9
   ###
         DISCR.RATIO[NEOF]: DISCRIMINANT RATIOS VS. NUMBER OF PCS
10
         RX[NX,NMIN]: VARIATE TIME SERIES FOR X
   ###
11
   ###
         RY[NY,NMIN]: VARIATE TIME SERIES FOR Y
12
13
   ###
         PMAT[SPACE, NMIN]: LOADING VECTOR
```

In the following calculations, you should include \*both\* ensemble members from 20C. A trick for doing this is to reshape the array so that the time series looks twice as long:

```
1 ### RESHAPE PC.20C[TIME,NENS,NEOF] TO PC.20C[TIME*NENS,NEOF]
2 dim(pc.20c) = c(tdim.20c*nens,neof)
```

Then, when you want individual ensemble members, you can reshape the array back to [time,ensemble,eof]

**Exercise 16.7.** Write a function that evaluates Mutual Information Criterion (MIC) for comparing covariance matrices. MIC is defined as

$$\text{MIC} = \frac{1}{N_T} \log \left( \frac{|\overline{\Sigma}_X|^{N_X} |\overline{\Sigma}_Y|^{N_Y}}{|\overline{\Sigma}_T|^{N_T}} \right) + \mathbb{P}, \tag{16.13}$$

where

$$\mathbb{P} = \frac{T}{N_T} \left( \frac{N_X(N_X+1)}{N_X - P - 2} + \frac{N_Y(N_Y+1)}{N_Y - P - 2} - \frac{N_T(N_T+1)}{N_T - P - 2} \right),$$

70 EXERCISES: COVARIANCE DISCRIMINANT ANALYSIS

and

$$\hat{\Sigma}_T = \frac{N_X \overline{\Sigma}_X + N_Y \overline{\Sigma}_Y}{N_T} \quad \text{and} \quad N_T = N_X + N_Y.$$
(16.14)

The function should evaluate MIC over all possible EOF truncations T. Plot MIC for a range of values and identify the value of T that minimizes MIC. Print out the first 5 values of MIC. These results should match the example in the notes.

**Exercise 16.8.** Compute the discriminant ratios for the optimum choice of T. State the values. These values should be consistent with those in the notes

**Exercise 16.9.** Write a function that computes discriminant variates. Compute the variates for the optimum choice of T and plot them. Verify that the sample covariance matrix of the PI variates equals the identity matrix. Verify that the sample covariance matrix of 20C variates is diagonal, with diagonal elements equal to the discriminant ratios.

**Exercise 16.10.** Write a function that computes the loading vectors. Plot the leading loading vector.

**Exercise 16.11.** Write a *separate* code that computes the 5% significance levels of the discriminant ratios based on 5000 trials of Monte Carlo experiments. State the 95% percentile for all discriminant ratios for the optimum choice of T. Are your discriminant ratios significant or not?

**Exercise 16.12.** What is the 5% significance level of the *univariate* F-test for equality of variances for sample sizes  $N_X = 51$ ,  $N_Y = 51$ ? Using your Monte Carlo code, show a plot of the 5% and 95% percentiles of *the leading discriminant ratio* as a function of the truncation parameter T. In this exercise, let  $N_X = 51$ ,  $N_Y = 51$ , and the number of trials = 1000. Also, let the maximum dimension be 30. What happens to the ratios as T increases? Explain why this happens. The 95% percentile from the Monte Carlo experiments should be close to the univariate F-test for T = 1.

# ANALYSIS OF VARIANCE

Table	17.1	default
Lable	T1.T	ueraun

prca.student.R	main program
[year]11_tref_NCEP-CFSv2.nc	NetCDF data files
land_cover.nc	NetCDF file for the land-sea mask
interesting.points.R	auxiliary file: R function for identifying large values in a field
index.climate.v2.R	auxiliary file: R function for defining N. American land area
eof.latlon.R	auxiliary file: calculate EOFs
plot_latlon_v4.R	auxiliary file: plot spatial maps
gev.R	auxiliary file for solving generalized eigenvalue problem

In this homework set you will use Analysis of Variance (ANOVA) to quantify the predictability of seasonal forecasts from the CFSv2 Retrospective Hindcasts. More precisely, you will determine whether April-mean temperature over North America is predictable by CFSv2 based on November initial conditions.

The above table shows the programs and data sets you need to download to do this homework. Download these files and run prca.student.R. As usual, modify dir.data and dir.Rlib to correspond to your data and library directories, respectively. prca.student.R should run to completion without any error or warning messages (please contact me if this is not true). After completion, the following variables are available to you:

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

```
## VARBL: NAME OF THE VARIABLE (E.G., TEMPERATURE)
1
   ## ICMON: INITIAL CONDITION MONTH (E.G., 'NOV')
2
   ## VFMON: VERIFICATION MONTH (E.G., 'JAN')
   ## EPIC: SELECTED ENSEMBLE MEMBERS (E.G., 1,2,...,24)
   ## NENS: TOTAL NUMBER OF ENSEMBLE MEMBERS (LENGTH (EPIC))
   ## NYRS: NUMBER OF YEARS
6
   ## IYST: FIRST YEAR OF THE DATA SET
  ## LON[NLON]: LONGITUDE VALUES OF THE DOMAIN
   ## LAT[NLAT]: LATITUDE VALUES OF THE DOMAIN
   ## HCST.DATA[NLON,NLAT,NENS,NYRS]: HINDCAST DATA
10
   ## VIEW[NLON, NLAT]: N. AMERICAN MASK
11
```

Note that the data array has the format [space, ensemble, year], which is not the same as in the notes (the notes assume [ensemble, year, space]). There is a good reason why the notes use a different format than the usual format for data, but that is a long story.

In dealing with multivariate data in R, your programs will run much faster if instead of using for loops, you make use of the commands rowMeans, rowSums, colMeans, and colSums. In order to make use of these commands, you will need to "reshape" your arrays using the dim command. Any array can be reshaped into another array with different dimensions, as long as the product of the dimensions are the same. The rule is that R fills up columns first, then rows. For instance, consider the following set of commands:

```
> x = rnorm(20)
   > print(x)
2
     [1] -0.08253025 -0.66934158 -2.55017522 2.19427432 -0.56563243 1.17370178 1.43692367
     [8] -0.63444609 -1.17081896 -2.23559009 -1.83286874 0.76325786 -0.20322325 -0.59558207
    [15] -2.17997749 -0.99505771 0.23259272 1.14159830 -1.63726907 -0.70107359
   > dim(x) = c(4,5)
   > print(x)
                [,1]
                          [,2]
                                     [,3]
                                                [,4]
                                                            [,5]
    [1,] -0.08253025 -0.5656324 -1.1708190 -0.2032233
                                                     0.2325927
10
   [2,] -0.66934158 1.1737018 -2.2355901 -0.5955821 1.1415983
    [3,] -2.55017522 1.4369237 -1.8328687 -2.1799775 -1.6372691
11
   [4,] 2.19427432 -0.6344461 0.7632579 -0.9950577 -0.7010736
12
```

In the above example, a random vector of length 20 is printed. Then, the vector is reshaped into a  $4 \times 5$  array and printed. You can see that R takes the sequence of numbers in the vector and fills the  $4 \times 5$  array, filling *columns first*, then the rows.

To illustrate the above technique, we will calculate a few quantities that are needed in ANOVA. One quantity that is needed is the grand mean at each grid point. This can be calculated as follows:

```
1 ## COMPUTE GRAND MEAN
2 dim(hcst.data) = c(nlon*nlat,nens*nyrs)
3 gmean = rowMeans(hcst.data)
```

The resulting variable gmean is a vector of length nlon\*nlat giving the mean at each grid point. Another quantity that is needed is the ensemble mean, which can be calculated as follows:

```
1 ## COMPUTE ENSEMBLE MEAN
```

```
2 dim(hcst.data) = c(nlon*nlat,nens,nyrs)
```

```
3 emean = array(NA, dim=c(nlon*nlat, nyrs))
```

```
4 for (ny in 1:nyrs) emean[,ny] = rowMeans(hcst.data[,,ny])
```

The first line reshapes the data array. The second line creates an array of the appropriate dimensions, but fills the array with NA. It is good practice to always initialize arrays with NA, so that if you accidentally do not fill up the entire array you will get an NA whenever any algebraic calculation is performed with those elements. The resulting NA is then useful for debugging. The last line calculates the ensemble mean for each year (i.e., the "condition" is "year"). The unbiased estimate of the variance of conditional means is therefore

```
1 ## COMPUTE SIGNAL VARIANCE
2 var.sig = rowSums((emean - gmean)^2) / (nyrs-1)
```

Note that emean is a 2-dimensional array while gmean is a vector. R will automatically repeat gmean until it fills up an array of the same size as emean, so that the subtraction can be performed. After the subtraction, each term in the resulting array is squared, then the sum of the squares is computed. Dividing by the degrees of freedom yields the unbiased estimate of the signal variance.

A key quantity in ANOVA is the critical value of F for testing significance. This is very simple to do in R using the command qf. For the above parameters, the  $\alpha 100\%$ significance threshold is

f.crit = qf(alpha,nyrs-1,nyrs\*(nens-1),lower.tail=FALSE)

## Assignment

**Exercise 17.1.** Write a function called f.anova.array that performs ANOVA for each grid point in an array of data. The output should give the F-value at each grid point and the corresponding significance level. The header of this function should be the following:

```
f.anova.array = function(x,nspace,nens,ncon,alpha=0.05) {
    ### COMPUTES THE F-STATISTIC IN ANOVA FOR EACH ELEMENT IN NSPACE
3
    # INPUT:
4
         X: [NSPACE, NENS, NCON] ARRAY OF DATA
         NSPACE: NUMBER OF SPATIAL ELEMENTS FOR INDIVIDUALLY COMPUTING F
5
6
         NENS: NUMBER OF ENSEMBLE MEMBERS
         NCON: NUMBER OF CONDITIONS
         ALPHA: SIGNIFICANCE LEVEL
    # OUTPUT: LIST WITH THE FOLLOWING VARIABLES
         F: [NSPACE] VECTOR CONTAINING THE F-VALUES
10
11
         F.CRIT: THE CRITICAL F-VALUE FOR F AT THE ALPHA*100% SIGNIFICANCE LEVEL
```

Turn in a copy of your program.

**74** ANALYSIS OF VARIANCE

1

3

4

**Exercise 17.2.** Use your function to calculate the F-statistic for CFSv2 hindcasts of January 2m temperature. Mask out insignificant values, like so:

```
f.list = f.anova.array(hcst.data,nlon*nlat,nens,nyrs)
#### MASK OUT INSIGNIFICANT F'S
fval = f.list$f
fval[fval < f.list$f.crit] = NA</pre>
```

Make a plot of the F values and turn it in.

Also, print out the result of summary(fval), which should be 2 lines giving the minimum, 1st quantile, median, etc.

**Exercise 17.3.** Use interesting.points to identify grid points with large *F* values. For instance, this can be done as follows:

ilist = interesting.points(fval,max,3,nlon,nlat,lon,lat,20,20)

The output of ilist should be self-explanatory (but see me if it doesn't make sense to you). Show a box-whisker plot of the data using the command

```
n npic = 1 # (or 2 or 3)
y = hcst.data[ilist$x.pic[npic],ilist$y.pic[npic],]
boxplot(y,names=year,col='grey')
```

In the title of each figure, state the corresponding value of F and the longitude and latitude of the point. Explain how the resulting figure is consistent with the value of F.

**Exercise 17.4.** Explore what happens to the F values as you reduce the number of ensemble members or change the lead time. These parameters can be changed by changing epic and lead, respectively. For instance, the default settings are

1	lead	=	6
2	epic	=	9:24

lead = 6 tells the code to pick the 6th lead time. Since the data is monthly, starting on November, this corresponds to predicting April (November is the "first lead"). epic = 9:24 tells the code to pick ensemble members 9-24. The ensemble members are stored in reverse order, as indicated in the following table from Saha et al. (2014).

1	8	Oct	at	0000,	0600,	1200,	and	1800	UTC
2	13	Oct	at	0000,	0600,	1200,	and	1800	UTC
3	18	Oct	at	0000,	0600,	1200,	and	1800	UTC
4	23	Oct	at	0000,	0600,	1200,	and	1800	UTC
5	28	Oct	at	0000,	0600,	1200,	and	1800	UTC
6	2	Nov	at	0000,	0600,	1200,	and	1800	UTC
7	7	Nov	at	0000,	0600,	1200,	and	1800	UTC

So, to choose an 8-member ensemble, you would set epic = 17:24. Try a few choices, plot the results, and explain why the results make sense. For instance, what do you think happens to predictability when you decrease the ensemble size? Or decrease the lead time?

# PREDICTABLE COMPONENT ANALYSIS

In this homework set you will apply Predictable Component Analysis (PrCA) to quantify the predictability of seasonal forecasts over North America from the CFSv2 Retrospective Hindcasts. The data set and codes for reading this data set are identical to those you used in the last homework for ANOVA.

To apply predictable component analysis, you must first compute the principal components of the data set. This can be done by pooling all ensemble members and years together as if it were one long time series. This will be done automatically by eof.latlon. Since the data set is dimensioned as

hcst.data[nlon,nlat,nens,nyrs]

On output, the PCs will have dimension [nens\*nyrs, neof]. When desired, the PCs can be re-shaped into an array of dimension [nens, nyrs, neof].

In this homework, you will eventually write an R function called prca that performs predictable component analysis on the PCs. The homework assignment below will direct you to write this code in "pieces." In the end, you should combine all these pieces together to produce a single function that performs all the calculations. The header for this function should be:

#### **78** PREDICTABLE COMPONENT ANALYSIS

```
prca.eof = function(data.eof,nens,ncon,alpha=0.01) {
  ## PERFORMAS PREDICTABLE COMPONENT ANALYSIS (OR EQUIVALENTLY MANOVA)
2
  ## ON PC-TIME SERIES OF AN ENSEMBLE FORECAST DATA SET
   ## INPUT:
   ##
         DATA.EOF: LIST OUTPUT FROM EOF.LATLON
   ##
         NENS: NUMBER OF ENSEMBLE MEMBERS
6
   ##
         NCON: NUMBER OF CONDITIONS
7
   ##
        ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 1%)
8
   ## OUTPUT: LIST$
0
   #
       MIC: VALUES OF MUTUAL INFORMATION CRITERION
10
   #
       MIC.CRIT: SIGNIFICANCE VALUES OF MIC
11
12
  #
      F.MAX[NTRUN]: MAXIMIZED F-RATIOS
13
   #
      R[NENS, NCON, NTRUN]: PREDICTABLE VARIATES
   #
      P[NSPACE,NTRUN]: PREDICTABLE LOADING VECTORS
14
   #
      FEXPVAR[NMIN]: FRACTION OF VARIANCE EXPLAINED
15
   #
       RSQR.ADJ[NMIN]: ADJUSTED R-SQUARE
16
  #
        SNR.ADJ[NMIN]: ADJUSTED SIGNAL-TO-NOISE RATIO
17
   #
       NENS, NCON, NMIN, ALPHA: OTHER PARAMETERS
18
```

**Exercise 18.1.** Write an R program called prca.null.fixdim that estimates the 1% significance level of each eigenvalue from predictable component analysis using ntrun feature vectors. Run this program using ntrun = 6 and state the significance values. The header of this program should be

```
prca.null.fixdim = function(ndim, ncon, nens, alpha=0.01, ntrial=10000)
1
   ## ESTIMATES SIGNIFICANCE LEVEL OF THE ALL EIGENVALUES FROM PRCA/MANOVA
2
   ## FOR FIXED TRUNCATION NDIM
  ## BY MONTE CARLO METHODS
  ## INPUT:
   #
       NDIM: DIMENSION OF THE RANDOM VECTOR
   #
       NCOM: NUMBER OF CONDITIONS (E.G., NUMBER OF YEARS)
7
   #
       NENS: NUMBER OF ENSEMBLE MEMBERS (E.G., NUMBER OF REPITITIONS)
8
       ALPHA: DESIRED SIGNIFICANCE LEVELS (CAN BE MORE THAN ONE)
   #
       NTRIALS: NUMBER OF MONTE CARLO TRIALS
  #
10
11
  ## OUTPUT:
        EVAL.CRIT[2,NDIM]: (alpha,1-alpha) SIGNIFICANCE THRESHOLDS FOR EACH EIGENVALUE
12
   #
```

It is a good idea to start out with ntrial = 100 while you are writing and debugging your code, and then change to ntrial = 10000 at the last step.

**Exercise 18.2.** Write an R function to compute the MIC as a function of the number of PCs. The expression for MIC that is appropriate for PrCA is:

$$\operatorname{MIC} = \log \left| \overline{\Sigma}_N \right| - \log \left| \overline{\Sigma}_T \right| + \frac{2CT + T(T+1)}{CE - C - T - 1} - \frac{2T + T(T+1)}{CE - 1 - T - 1}, \quad (18.1)$$

where T is the number of PCs. State the values of MIC that you obtain, and state the value of T at which the minimum occurs.

**Exercise 18.3.** Write an R function to compute the maximized F-ratios from PrCA, and the associated adjusted signal-to-noise ratios and R-square. This function will need to

compute the signal and noise covariance matrices from the PCs and then solve a generalized eigenvalue problem. State the maximized F-ratios for T = 6. Also state the adjusted signal-to-noise ratios and R-square.

**Exercise 18.4.** Write an R function that computes the predictable variates. Remember to normalize the weight vector  $\mathbf{q}$  so that the total variance equals unity. Run this function for T = 6. Make a plot of the variates for the <u>two</u> most predictable components. If the variate is dimensioned as r[nens,nyrs,ntrun], then the variates for the most predictable component can be plotted using the command

```
boxplot(r[,,1],names=year,col='grey')
```

1

Label the plot as in the lecture notes, so that all information required to interpret the plot is contained in the plot.  $\hfill \Box$ 

**Exercise 18.5.** Verify that the predictable variates are orthogonal. More generally, verify that the sample covariance matrix of the predictable components equals the identity matrix.

**Exercise 18.6.** Write a function that calculates the loading vectors. Plot them.  $\Box$ 

# DATA ASSIMILATION

In this homework set you will apply the Kalman Filter to assimilate observations using a vector autoregressive model. A vector autoregressive model allows probabilities of forecasted variables to be solved exactly, so it illustrates the Kalman Filter in an "ideal" situation. However, this ideal situation is very unrealistic, so there is potential for being mislead.

A vector autoregressive model is of the form

$$\mathbf{y}_t = \mathbf{A} \quad \mathbf{y}_{t-1} + \mathbf{w}_t \\ M \times 1 \qquad M \times M \quad M \times 1 \qquad M \times 1,$$
(19.1)

where  $\mathbf{A}$  is a matrix called the *dynamical operator* and  $\mathbf{w}_t$  is Gaussian white noise with zero mean and covariance matrix  $\mathbf{Q}$ . After specifying the matrices  $\mathbf{A}$  and  $\mathbf{Q}$  and an initial condition  $\mathbf{y}_0$ , the vectors  $\mathbf{y}_1, \mathbf{y}_2, \ldots$  can be solved for a particular realization of the noise  $\mathbf{w}_t$ .

In this homework set, you will consider the following experimental situation. First, you will create a vector time series  $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{100}$  from (19.1) for a particular choice of  $\mathbf{A}$ ,  $\mathbf{Q}$ , and  $\mathbf{y}_0$  and for a particular realization of  $\mathbf{w}_t$ . The resulting vector time series will be the "truth." To mimic reality, we will pretend that we do not know the truth. Instead, you will create "observations" using the model

$$\mathbf{o}_t = \mathbf{H} \quad \mathbf{y}_t + \mathbf{r}_t \\ K \times 1 \quad K \times M \quad M \times 1 \quad K \times 1'$$
(19.2)

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

#### 82 DATA ASSIMILATION

for a particular choice of **H** and covariance matrix **R** of  $\mathbf{r}_t$ . The observations  $\mathbf{o}_t$  are known and our goal is to estimate the truth  $\mathbf{y}_t$  from these observations (pretending that we don't know the truth). You will use the Kalman Filter to assimilate these observations to estimate the true state. Then you will compare this estimate with the true state to see how well the Kalman Filter works.

Download the R program kf\_student.R from the class website. This program generates a time series from (19.1) for the choice

$$\mathbf{A} = \begin{pmatrix} 0.9 & 0.5 & 0.3 \\ 0.0 & 0.5 & 2.0 \\ 0.0 & 0.0 & 0.4 \end{pmatrix}, \quad \mathbf{Q} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \text{and} \quad \mathbf{y}_0 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \tag{19.3}$$

and generates "observations" using the choice

$$\mathbf{H} = \begin{pmatrix} 1 & 0 & 0 \end{pmatrix}$$
 and  $\mathbf{R} = 49.$  (19.4)

Note that this choice implies that only one variable (i.e., the first element of  $y_t$ ) is observed with noise. The program assumes that Q is diagonal, so that the elements of the noise vector  $w_t$  are independent and have variances equal to the diagonal elements of Q. This special form allows the noise vector to be specified very efficiently numerically as

w = rnorm(M, sd=sqrt(diag(Q)))

Similarly,  $\mathbf{R}$  is assumed to be diagonal and hence the observations generated by (19.2) can be specified efficiently in a similar way. This is done in the R code that you will download.

After the R code runs (hopefully without error messages!), you should have the following quantities available:

#	NDIM:	THE DI	MENSIO	N OF THE	E VECTOR AUT	FORE	GRESSIVE MODEL (VAR); DEFAULT = 3
#	NTOT:	THE TO	TAL LE	NGTH OF	THE TIME SH	ERIES	S; DEFAULT = 100
#	NOBS:	NUMBER	OF OB	SERVATIO	ONS PER TIME	E STH	EP; DEFAULT = 1
#	Χ:	[NDIM,	NTOT]	MATRIX	SPECIFYING	THE	"TRUTH" AT EACH TIME STEP
#	OBS:	[NOBS,	NTOT]	MATRIX	SPECIFYING	THE	"OBSERVATIONS" AT EACH TIME STEP
#	DYNOP:	[NDIM,	NDIM]	MATRIX	SPECIFYING	THE	DYNAMICAL OPERATOR OF THE VAR
#	HOP:	[NOBS,	NDIM]	MATRIX	SPECIFYING	THE	INTERPOLATION OPERATOR
#	Q.COV:	[NDIM,	NDIM]	MATRIX	SPECIFYING	THE	NOISE COVARIANCE MATRIX OF THE VAR
#	R.COV:	[NOBS,	NOBS]	MATRIX	SPECIFYING	THE	ERROR COVARIANCE MATRIX OF THE OBS

**Exercise 19.1.** Plot the first element of the true time series. Also plot the corresponding observations. The result should look something like the curve and dots in fig. 19.1.

Exercise 19.2. Write a function to evaluate the Kalman Filter equations

$$\boldsymbol{\mu}_{t}^{A} = \boldsymbol{\mu}_{t}^{B} + \boldsymbol{\Sigma}_{t}^{B} \mathbf{H}_{t}^{T} \left( \mathbf{H}_{t} \boldsymbol{\Sigma}_{t}^{B} \mathbf{H}_{t}^{T} + \mathbf{R}_{t} \right)^{-1} \left( \mathbf{o}_{t} - \mathbf{H}_{t} \boldsymbol{\mu}_{t}^{B} \right)$$
(19.5)

$$\boldsymbol{\Sigma}_{t}^{A} = \boldsymbol{\Sigma}_{t}^{B} - \boldsymbol{\Sigma}_{t}^{B} \mathbf{H}_{t}^{T} \left( \mathbf{H}_{t} \boldsymbol{\Sigma}_{t}^{B} \mathbf{H}_{t}^{T} + \mathbf{R}_{t} \right)^{-1} \mathbf{H}_{t} \boldsymbol{\Sigma}_{t}^{B}.$$
(19.6)

The function call should be

```
kalman.population = function(mub, sigmab, hop, r.cov, obs) {
2
    ##########
    ## EVALUATES THE KALMAN FILTER EQUATIONS FOR THE MEAN AND COVARIANCE MATRIX OF THE ANALYSIS
3
    ## BASED ON *POPULATION* COVARIANCE MATRICES AND MEANS
4
5
    ## INPUT.
                  [NDIM] MEAN OF THE BACKGROUND DISTRIBUTION
6
         MUB:
7
         SIGMAB: [NDIM, NDIM] BACKGROUND COVARIANCE MATRIX
         HOP .
                  [NOBS, NDIM] INTERPOLATION OPERATOR
8
         R.COV:
                  [NOBS, NOBS] COVARIANCE MATRIX OF THE OBSERVATIONAL ERROR
10
         OBS:
                  [NOBS] THE OBSERVATIONS
11
    ##
       OUTPUT
         MUA:
                  [NDIM] VECTOR OF THE ANALYSIS DISTRIBUTION
12
13
         SIGMAA:
                 [NDIM, NDIM] COVARIANCE MATRIX OF THE ANALYSIS DISTRIBUTION
```

As you can see from the equations, you will need to compute the inverse of a matrix. You should invert the matrix using the commands chollinv(chol(MATRIX)), which takes advantage of the fact that the matrix being inverted is symmetric.

**Exercise 19.3.** The background distribution is obtained from the *previous* analysis. In particular, the mean and covariance matrix of the background distribution are

$$\boldsymbol{\mu}_t^B = \mathbf{A} \boldsymbol{\mu}_{t-1}^A \quad \text{and} \quad \boldsymbol{\Sigma}_t^B = \mathbf{A} \boldsymbol{\Sigma}_{t-1}^A \mathbf{A}^T + \mathbf{Q}.$$
 (19.7)

Assume the initial analysis  $\mu_0^A = \mathbf{0}$  and

$$\boldsymbol{\Sigma}_{0}^{A} = \begin{pmatrix} 400 & 0 & 0\\ 0 & 400 & 0\\ 0 & 0 & 400 \end{pmatrix}.$$
 (19.8)

Combine these equations together with your Kalman Filter function to construct an analysis for one hundred time steps. Plot the mean analysis and its uncertainty for the first element of  $\mathbf{y}_t$ . The result should look something like the grey shading in fig. 19.1, which shows  $(\boldsymbol{\mu}^A)_1 \pm \sqrt{(\boldsymbol{\Sigma}^A)_{11}}$ .

**Exercise 19.4.** Plot the standard error of the analysis for each element of  $\mathbf{y}_t$  as a function of time (i.e., plot the square roots of the diagonal elements of  $\Sigma_t^A$ ). What happens to the errors after a long time?

**Exercise 19.5.** In this problem only, set  $\mathbf{R}_t = 1$  and show the analysis, observations (if available), and the truth for the first and second elements of  $\mathbf{y}_t$ . How do these plots differ from the previous case? Explain why this difference makes sense.

**Exercise 19.6.** Suppose that, instead of observing  $(\mathbf{y}_t)_1$ , we can observe only the sum  $z(t) = (\mathbf{y}_t)_1 + (\mathbf{y}_t)_2$ . Furthermore, suppose the error of this observation is normally distributed with zero mean and variance 25. Explain how you would assimilate these observations (e.g., explain how you would change your R code to handle this case). Actually do this and show the standard error of the analysis and the truth for  $(\mathbf{y}_t)_1$ .

**Exercise 19.7** (Least Squares Derivation of the Kalman Filter Equations). The Kalman Filter equations can be derived by the method of least squares. The essential problem is that we are given a set of observations o, and based on this want to estimate the state x. To do this, we assume a linear prediction model

$$\hat{\mathbf{y}} = \mathbf{A}\mathbf{o} + \mathbf{b}.\tag{19.9}$$



**Figure 19.1** A particular realization of the first element  $y_1(t)$  generated by the vector autoregressive model (19.1) (solid curve), corresponding observations (dots), and the analysis mean plus a standard deviation, as estimated state by the Kalman Filter (grey shading).

Show that the method of least squares gives

$$\mathbf{A} = \boldsymbol{\Sigma}_{YO} \boldsymbol{\Sigma}_O^{-1} \quad \text{and} \quad \mathbf{b} = \mathbb{E}[\mathbf{y}] - \mathbf{A} \mathbb{E}[\mathbf{o}], \tag{19.10}$$

where the expectation operator  $E[\cdot]$  is an average over all possible realizations of the observations and background quantities. Under this definition,  $E[\mathbf{y}] = \boldsymbol{\mu}_B$  and  $\operatorname{cov}[\mathbf{y}] = \boldsymbol{\Sigma}_B$ . The observations are assumed to be related to the state  $\mathbf{y}$  through the model

$$\mathbf{o} = \mathbf{H}\mathbf{y} + \mathbf{r},\tag{19.11}$$

where **H** is the interpolation operator and **r** is independent and normally distributed with zero mean and covariance matrix  $\Sigma_R$ . Show that this equation implies

$$\Sigma_{OY} = \mathbf{H}\Sigma_Y = \mathbf{H}\Sigma_B \tag{19.12}$$

and

$$\Sigma_O = \mathbf{H}\Sigma_Y \mathbf{H}^T + \Sigma_R = \mathbf{H}\Sigma_B \mathbf{H}^T + \Sigma_R.$$
(19.13)

and

$$\mathbb{E}[\mathbf{o}] = \mathbf{H}\boldsymbol{\mu}_B \tag{19.14}$$

Substituting these equations into (19.10) gives

$$\mathbf{A} = \boldsymbol{\Sigma}_{B} \mathbf{H}^{T} \left( \mathbf{H} \boldsymbol{\Sigma}_{B} \mathbf{H}^{T} + \boldsymbol{\Sigma}_{R} \right)^{-1}$$
(19.15)

85

and

$$\mathbf{b} = \boldsymbol{\mu}_B - \boldsymbol{\Sigma}_B \mathbf{H}^T \left( \mathbf{H} \boldsymbol{\Sigma}_B \mathbf{H}^T + \boldsymbol{\Sigma}_R \right)^{-1} \mathbf{H} \boldsymbol{\mu}_B.$$
(19.16)

It follows that the linear regression model (19.9) yields

$$\mathbf{y}_{a} = \boldsymbol{\mu}_{B} + \boldsymbol{\Sigma}_{B} \mathbf{H}^{T} \left( \mathbf{H} \boldsymbol{\Sigma}_{B} \mathbf{H}^{T} + \boldsymbol{\Sigma}_{R} \right)^{-1} \left( \mathbf{o} - \mathbf{H} \boldsymbol{\mu}_{B} \right),$$
(19.17)

which is precisely the Kalman Filter equation for the mean analysis.

Show that the error covariance of the regression model (19.9) is

$$\begin{split} \boldsymbol{\Sigma}_{a} &= \mathbb{E}\left[\left(\mathbf{y} - \mathbf{y}_{a}\right)\left(\mathbf{y} - \mathbf{y}_{a}\right)^{T}\right] \\ &= \mathbb{E}\left[\left(\mathbf{y} - \boldsymbol{\mu}_{B} - \mathbf{A}(\mathbf{o} - \mathbf{H}\boldsymbol{\mu}_{B})\right)\left(\mathbf{y} - \boldsymbol{\mu}_{B} - \mathbf{A}(\mathbf{o} - \mathbf{H}\boldsymbol{\mu}_{B})\right)^{T}\right] \\ &= \boldsymbol{\Sigma}_{Y} - \boldsymbol{\Sigma}_{YO}\mathbf{A}^{T} - \mathbf{A}\boldsymbol{\Sigma}_{OY} + \mathbf{A}\boldsymbol{\Sigma}_{O}\mathbf{A}^{T} \\ &= \boldsymbol{\Sigma}_{B} - \boldsymbol{\Sigma}_{B}\mathbf{H}^{T}\left(\mathbf{H}\boldsymbol{\Sigma}_{B}\mathbf{H}^{T} + \boldsymbol{\Sigma}_{R}\right)^{-1}\mathbf{H}\boldsymbol{\Sigma}_{B}, \end{split}$$
(19.18)

which is precisely the Kalman Filter equation for the analysis error covariance.

# ENSEMBLE SQUARE ROOT FILTERS

**Exercise 20.1.** Recall from sec. 20.5 that observations can be assimilated sequentially if their errors are independent. In this exercise, you will derive how this is done for the square root filter. In the case of one observation, (21.18) can be written as

$$\mathbf{D} = \mathbf{I} - \beta \mathbf{w} \mathbf{w}^T, \tag{20.1}$$

where

$$\mathbf{w} = \mathbf{X}^{BT} \mathbf{H}_t^T \quad \text{and} \quad \beta = \left(\mathbf{w}^T \mathbf{w} + \mathbf{R}_t\right)^{-1}.$$
 (20.2)

Since only one observation is under consideration,  $\beta$  is a scalar. The matrix **D** differs from the identity matrix by a rank-1 matrix. It is natural to assume that the square root also differs from the identity by a term proportional to the same rank-1 matrix. Therefore, we seek a symmetric square root of **D** of the form

 $\mathbf{I} - \beta \mathbf{w} \mathbf{w}^{T} = \left(\mathbf{I} - \delta \mathbf{w} \mathbf{w}^{T}\right) \left(\mathbf{I} - \delta \mathbf{w} \mathbf{w}^{T}\right), \qquad (20.3)$ 

where  $\delta$  is a parameter to be determined. Expanding the right side and simplifying yields

$$\mathbf{I} - \beta \mathbf{w} \mathbf{w}^T = \mathbf{I} - 2\delta \mathbf{w} \mathbf{w}^T + \delta^2 \mathbf{w} \mathbf{w}^T \mathbf{w} \mathbf{w}^T$$
(20.4)

$$0 = \left(\mathbf{w}^T \mathbf{w} \delta^2 - 2\delta + \beta\right) \mathbf{w} \mathbf{w}^T.$$
(20.5)

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

### **88** ENSEMBLE SQUARE ROOT FILTERS

This equation is satisfied for arbitrary  $\mathbf{w}$  only if the quadratic equation for  $\delta$  vanishes, which requires

$$\delta = \frac{1 \pm \sqrt{1 - \mathbf{w}^T \mathbf{w} \beta}}{\mathbf{w}^T \mathbf{w}} = \frac{1 \pm \sqrt{\beta \mathbf{R}_t}}{\mathbf{w}^T \mathbf{w}} = \beta \left( 1 \mp \sqrt{\frac{\mathbf{R}}{\mathbf{R} + \mathbf{w}^T \mathbf{w}}} \right)^{-1}.$$
 (20.6)

This gives two solutions, but only the solution

$$\delta = \frac{1 - \sqrt{1 - \mathbf{w}^T \mathbf{w} \beta}}{\mathbf{w}^T \mathbf{w}} = \frac{1 - \sqrt{\beta \mathbf{R}_t}}{\mathbf{w}^T \mathbf{w}} = \beta \left( 1 + \sqrt{\frac{\mathbf{R}}{\mathbf{R} + \mathbf{w}^T \mathbf{w}}} \right)^{-1}.$$
 (20.7)

produces a positive semi-definite square root matrix. Show that this is true.

**Exercise 20.2.** Prove the Schur Product theorem (theorem 21.1). Hint: show that the Schur product between a positive semi-definite matrix and a rank-1 matrix is positive semi-definite. Then, use the spectral factorization theorem to argue that any symmetric, positive semi-definite matrix can be represented by a sum of rank-1 matrices.

# EXERCISES: EXTREME VALUE ANALYSIS

In this homework set, you will use the Generalized Extreme Value Distribution to quantify the return period for extreme cold temperatures in the mid-Atlantic U.S. The data set and R code to read this data are summarized in the table below and can be downloaded from the class website. You also will need to install the extRemes library. Also, inside tele\_index.R, you need to change indices.fname to the full directory path of the data file tele\_index.nh.

Table 21.1 default

GHNC_OBS_T_mid_atlantic_1950_2015.txt	data file for mid-Atlantic temperature
gev.student.R	R code to read the mid-Atlantic temperature
tele_index.R	R code to read climate indices
tele_index.nh	NOAA climate indices

**Exercise 21.1.** Run gev.student.R and ensure it generates a time series/box plot identical to fig. 19.1. Next, fit the cold extremes to a GEV using the commands

1 t.cold = -t.cold # reverse sign to use extreme value functions 2 cold.x = fevd(t.cold)

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

Submit a print out of the output when you type cold.x. This print out should be consistent with table 19.1.

Make a plot using the command plot (cold.x) and submit it. This plot is discussed in Gilleland and Katz (2016). Consult this document and describe what this plot shows.

What are the return levels for 2-year, 20-year, and 200-year events (you will need a special command in Gilleland and Katz (2016) to answer this question)?  $\Box$ 

**Exercise 21.2.** Is the fitted distribution a Gumbel, Fréchet, or Weibull? Explain your answer.

**Exercise 21.3.** Fit the data to a GEV, but this time use NAO as a covariate for the location parameter. The NAO is index in the R code. Is this covariate significant? Is the fitted distribution a Gumbel, Fréchet, or Weibull? Explain your answer. Generate a plot similar to the first exercise above and describe the new information that it conveys.

**Exercise 21.4.** Suppose the probability of exceedance in any randomly selected year is *p*. If the exceedances in different years are independent, then what is the average number of years one will wait before observing the next exceedance? Hint: the average number of years one will wait before observing the next exceedance is computed as

return period = 
$$\sum_{y=1}^{\infty} y * (\text{probability that next exceedance occurs in year } y)$$
. (21.1)

What is the probability that the exceedance occurs in year 1? year 2? etc?

# ANSWERS TO EXERCISES

## Exercise 1.1, page 2

The frequency of  $H_1$  is the sum over all joint frequencies involving  $H_1$ :

$$P(H_1) = P(H_1, H_2) + P(H_1, T_2) = (30 + 21)/100 = 51\%$$

Similarly for  $H_2$ 

$$P(H_2) = P(H_1, H_2) + P(T_1, H_2) = (30 + 22)/100 = 52\%$$

The joint event  $P(H_1, T_2)$  can be read directly from the table:

$$P(H_1, T_2) = 21/100 = 21\%.$$

The conditional frequency  $P(T_2|H_1)$  is the number of counts of  $T_2$  out of all cases in which the first toss is heads:

$$P(T_2|H_1) = 21/(30+21) = 21/51 \approx 41\%.$$

It is readily verified that

$$P(T_2|H_1) = P(T_2, H_1) / P(H_1) = 21/51 \approx 41\%$$
(22.1)

Exercises for Statistical Methods for Climate Scientists. By DelSole and Tippett

## **92** ANSWERS TO EXERCISES

# Exercise 1.2, page 2

Since X and Y are independent, p(x, y) = p(x)p(y). Therefore:

$$E[XY] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xyp(x,y)dxdy$$
  
= 
$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xyp(x)p(y)dxdy$$
  
= 
$$\left(\int_{-\infty}^{\infty} xp(x)dx\right) \left(\int_{-\infty}^{\infty} yp(y)dy\right)$$
  
$$E[XY] = E[X]E[Y].$$
 (22.2)

# Exercise 1.3, page 2

(a) Since X and Y are independent,  $E[XY] = E[X]E[Y] = \mu_X \mu_Y$ .

(b)  $\operatorname{cov}[X, Y] = 0$ , since X and Y are independent.

(c)

$$\operatorname{var}[X - Y] = E[((X - Y) - (\mu_X - \mu_Y))^2]$$
  
=  $E[(X - \mu_X)^2 + (Y - \mu_y)^2 - 2(X - \mu_X)(Y - \mu_Y)]$   
=  $\sigma_X^2 + \sigma_Y^2$ .

(d)

$$\begin{aligned} \operatorname{var}(XY) &= E[(XY - E[XY])^2] \\ &= E[(XY - \mu_X\mu_Y)^2] \\ &= E[X^2Y^2 + \mu_X^2\mu_Y^2 - 2XY\mu_X\mu_Y] \\ &= E[X^2]E[Y^2] + \mu_X^2\mu_Y^2 - 2\mu_X^2\mu_Y^2 \\ &= (\sigma_X^2 + \mu_X^2)(\sigma_Y^2 + \mu_Y^2) - \mu_X^2\mu_Y^2 \\ &= \sigma_X^2\sigma_Y^2 + \sigma_X^2\mu_Y^2 + \sigma_Y^2\mu_X^2. \end{aligned}$$

## Exercise 1.4, page 2

$$var[k] = E[(k - E[k])^{2}] = E[(k - k)^{2}] = E[0] = 0.$$
  

$$var[kX] = E[(kX - E[kX])^{2}] = E[(kX - kE[X])^{2}]$$
  

$$= E[k^{2}(X - E[X])^{2}] = k^{2} var[X]$$
  

$$var[k + X] = E[((k + X) - E[k + X])^{2}] = E[(k + X - k - E[X])^{2}]$$
  

$$= E[(X - E[X])^{2}] = var[X]$$

Exercise 1.5, page 2

$$cov[X, Y] = E [(X - E[X]) (Y - E[Y])]$$
  
=  $E [XY - XE[Y] - YE[X] + E[X]E[Y]]$   
=  $E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y]$   
=  $E[XY] - E[X]E[Y].$  (22.3)

#### Exercise 1.6, page 2

Calculus method: take the first and second derivatives with respect to k:

$$\frac{\partial E[(X-k)^2]}{\partial k} = -2E[X-k] = 0$$
$$\frac{\partial^2 E[(X-k)^2]}{\partial k^2} = 2 \ge 0$$

The first equation gives the solution k = E[X]. The second equation demonstrates that the second derivative is positive everywhere and hence the solution is a minimum.

Second method: add and subtract E[X] and take expectations:

$$E[(X-k)^{2}] = E[(X-E[X] + E[X] - k)^{2}]$$
  
=  $E[(X-E[X])^{2}] + E[(E[X] - k)^{2}] + E[(X-E[X])(E[X] - k)]$ 

The last term vanishes because E[X] - k is a constant and hence can be taken outside the expectation operator, leaving E[X - E[X]] = E[X] - E[X] = 0. It follows that

$$E[(X - k)^{2}] = E[(X - E[X])^{2}] + E[(E[X] - k)^{2}].$$

Both terms on the right hand side are positive, but only the last depends on k. Moreover, the last term vanishes if k = E[X], so this choice minimizes the left hand side.

#### Exercise 1.7, page 2

From the hint,

$$E[(t(X - E[X]) + (Y - E[Y]))^2] = t^2 \operatorname{var}[X] + 2t \operatorname{cov}[X, Y] + \operatorname{var}[Y] \ge 0. \quad (22.4)$$

The minimum value is found by setting the derivative of t to zero and solving, which gives

$$t_{min} = -\operatorname{cov}[X, Y] / \operatorname{var}[X]$$
(22.5)

#### **94** ANSWERS TO EXERCISES

Substituting this value into the t-equation gives

$$0 \leq \left(\frac{-\operatorname{cov}[X,Y]}{\operatorname{var}[X]}\right)^{2} \operatorname{var}[X] - 2\frac{\operatorname{cov}^{2}[X,Y]}{\operatorname{var}[X]} + \operatorname{var}[Y]$$
$$\leq \frac{\operatorname{cov}^{2}[X,Y]}{\operatorname{var}[X]} - 2\frac{\operatorname{cov}^{2}[X,Y]}{\operatorname{var}[X]} + \operatorname{var}[Y]$$
$$\leq \operatorname{var}[Y] - \frac{\operatorname{cov}^{2}[X,Y]}{\operatorname{var}[X]}$$
$$\leq \operatorname{var}[Y] \left(1 - \frac{\operatorname{cov}^{2}[X,Y]}{\operatorname{var}[X]\operatorname{var}[Y]}\right)$$
$$\leq \operatorname{var}[Y] \left(1 - \rho^{2}\right)$$

Hence

$$\rho^2 = \frac{\operatorname{cov}^2[X, Y]}{\operatorname{var}[X] \operatorname{var}[Y]} \le 1,$$

which proves that the correlation coefficient must be between -1 and 1.

Alternatively, one can reason geometrically: the equation is a parabola that is concave upward. For the inequality to be satisfied, there should either no roots or exactly one root (i.e., if there are two roots, then the parabola crosses the x-axis and hence must be negative). The condition for no root or exactly one root is that the discriminant is not positive. That is, for a parabola of the form  $y = Ax^2 + Bx + C$ , the discriminant  $B^2 - 4AC \le 0$ . This condition gives:

$$4\operatorname{cov}^{2}[X,Y] - 4\operatorname{var}[X]\operatorname{var}[Y] \le 0 \quad \to \quad \frac{\operatorname{cov}^{2}[X,Y]}{\operatorname{var}[X]\operatorname{var}[Y]} \le 1.$$
(22.6)

If the squared correlation is *exactly* equal to 1, then the quadratic equation vanishes exactly. But the only way an expression of the form  $E[Z^2]$  can equal zero is if Z itself is exactly zero, which implies

$$Y = -t(X - E[X]) + E[Y],$$

for some constant t; in other words, if X and Y are exactly linearly related.

#### Exercise 1.8, page 3

The answers to the questions are as follows:

$$\hat{\mu}_X \sim N(\mu_X, \sigma^2/N)$$

$$\hat{\mu}_Y \sim N(\mu_Y, \sigma^2/N)$$

$$\hat{\mu}_X - \hat{\mu}_Y \sim N(\mu_X - \mu_Y, 2\sigma^2/N)$$

$$(N-1)\hat{\sigma}_X^2/\sigma^2 \sim \chi^2_{N-1}$$

$$(N-1)\hat{\sigma}_Y^2/\sigma^2 \sim \chi^2_{N-1}$$

$$(N-1)(\hat{\sigma}_X^2 + \hat{\sigma}_Y^2)/\sigma^2 \sim \chi_{2(N-1)}$$

## Exercise 1.9, page 3

Since the die is fair, the probability of each outcome is equal, so  $P(1) = \cdots = P(6) = 1/6$ . Therefore, the expected value is

$$E[X] = 1\left(\frac{1}{6}\right) + 2\left(\frac{1}{6}\right) + \dots + 6\left(\frac{1}{6}\right) = \frac{7}{2}.$$
 (22.7)

Similarly, the expected squared value is

$$E[X^2] = 1^2 \left(\frac{1}{6}\right) + 2^2 \left(\frac{1}{6}\right) + \dots + 6^2 \left(\frac{1}{6}\right) = \frac{91}{6}.$$
 (22.8)

Hence, the variance is

$$\operatorname{var}[X] = E[X^2] - (E[X])^2 = \frac{91}{6} - \left(\frac{7}{2}\right)^2 = \frac{35}{12}.$$
 (22.9)

#### Exercise 1.10, page 3

Let  $X_i$  denote the outcome of the *i*'th die. Then  $X = X_1 + X_2 + X_3$  and

$$E[X] = E[X_1] + E[X_2] + E[X_3] = 3\left(\frac{7}{2}\right) = \frac{21}{2} = 10.5$$
 (22.10)

where the answer from exercise 1.9 has been used. Similarly, since the dice rolls are independent, the variance of the sum is the sum of the variances:

$$\operatorname{var}[X] = \operatorname{var}[X_1] + \operatorname{var}[X_2] + \operatorname{var}[X_3] = 3\frac{35}{12} = \frac{35}{4} = 8.75.$$
 (22.11)

where the answer from exercise 1.9 has been used again.

# Exercise 1.11, page 3

$$\begin{split} E[\hat{\sigma}^2] &= E\left[\frac{1}{N-1}\sum_{n=1}^{N}\left(X_n - \hat{\mu}_X\right)^2\right] & \text{definition of } \hat{\sigma}^2 \\ &= E\left[\frac{1}{N-1}\sum_{n=1}^{N}\left(X_n - \mu + \mu - \hat{\mu}_X\right)^2\right] & \text{insert } 0 = \mu - \mu \\ &= E\left[\frac{1}{N-1}\sum_{n=1}^{N}\left(\left(X_n - \mu\right)^2 + 2\left(X_n - \mu\right)\left(\mu - \hat{\mu}_X\right) + \left(\hat{\mu}_X - \mu\right)^2\right)\right] & \text{algebra} \\ &= \frac{1}{N-1}\sum_{n=1}^{N} E\left[\left(X_n - \mu\right)^2 + 2\left(X_n - \mu\right)\left(\mu - \hat{\mu}_X\right) + \left(\hat{\mu}_X - \mu\right)^2\right] & \text{linearity of expectation} \\ &= \frac{1}{N-1}\sum_{n=1}^{N} \left(\text{var}[X] + 2E\left[\left(X_n - \mu\right)\left(\mu - \hat{\mu}_X\right)\right] + \frac{\text{var}[X]}{N}\right) & (1.20) \text{ and } (1.45) \\ &= \frac{1}{N-1}\left(N \operatorname{var}[X] + 2E\left[\sum_{n=1}^{N}\left(X_n - \mu\right)\left(\mu - \hat{\mu}_X\right)\right] + \operatorname{var}[X]\right) & \text{expectations are constant} \\ &= \frac{1}{N-1}\left(N \operatorname{var}[X] - 2NE\left[\left(\hat{\mu}_X - \mu\right)\left(\mu - \hat{\mu}_X\right)\right] + \operatorname{var}[X]\right) & \text{algebra} \\ &= \frac{1}{N-1}\left(N \operatorname{var}[X] - 2NE\left[\left(\hat{\mu}_X - \mu\right)\left(\mu - \hat{\mu}_X\right)\right] + \operatorname{var}[X]\right) & \text{algebra} \\ &= \frac{1}{N-1}\left(N \operatorname{var}[X] - 2N\frac{\operatorname{var}[X]}{N} + \operatorname{var}[X]\right) & (1.45) \\ &= \frac{1}{N-1}\left(N - 1\right)\operatorname{var}[X] & \text{algebra} \\ &= \operatorname{var}[X] & \text{algebra}. \end{split}$$

# Exercise 1.12, page 3

The basic idea is to use a computer to generate samples (X, Y) from a population with cor[X, Y] = 0.5 and then estimate the probability

$$P(Y > 0 \mid X > 0) = \frac{P(Y > 0 \text{ and } X > 0)}{P(X > 0)}.$$
(22.12)

A code that generates samples (X, Y) from a bivariate normal population with  $cor[X, Y] = \rho$  and then computes relative frequencies is the following:

```
ntrials = 1000000
1
          = 0.5
  rho
2
           = rnorm(ntrials)
  xvar
  nvar
           = rnorm(ntrials)
           = rho * xvar + sqrt(1-rho^2)*nvar
  yvar
          = xvar >= 0
  x.ge.0
  y.ge.0 = yvar >= 0
  prob.yge0.cond.xge0 = sum(x.ge.0 & y.ge.0) / sum(x.ge.0)
  print(prob.yge0.cond.xge0)
```

After running this code, it is found that  $P(Y > 0 | X > 0) \approx 2/3$  when  $\rho = 0.5$ . A figure showing the probabilities for other values of  $\rho$  is shown in fig. 22.1.



**Figure 22.1** Probability that Y > 0 given that X > 0, as a function of the correlation cor[X, Y] for a bivariate normal distribution with zero means

## Exercise 1.13, page 3

The displacement in the x and y directions have the distributions

$$d_x = x_1 + x_2 + \dots + x_N \sim N(0, N\sigma^2)$$
(22.13)

$$d_y = y_1 + y_2 + \dots + y_N \sim N(0, N\sigma^2).$$
 (22.14)

The total squared displacement is  $r^2 = d_x^2 + d_y^2$ . Since  $d_x/\sigma/\sqrt{N}$  and  $d_y/\sigma/\sqrt{N}$  are standardized Gaussians,

$$\frac{r^2}{N\sigma^2} = \left(\frac{d_x}{\sigma\sqrt{N}}\right)^2 + \left(\frac{d_y}{\sigma\sqrt{N}}\right)^2 \sim \chi_2^2.$$
(22.15)

The upper 5% percentile of a chi-squared distribution with 2 degrees of freedom is approximately 6. Therefore, the radius of the circle in which there is a 95% probability that the

particle is in the circle after 10 steps is

$$r_{95} \approx \sigma \sqrt{6N} = 2\sqrt{6*10} \approx 15.5.$$
 (22.16)

# Exercise 1.14, page 4

$$\begin{split} \mathbb{E}[(X' - \alpha \hat{\mu})^2] &= \mathbb{E}[((X' - \mu_X) - (\alpha \hat{\mu}_X - \mu_X))^2] & \text{add } 0 = \mu_X - \mu_X \\ &= \mathbb{E}[(X' - \mu)^2] + \mathbb{E}[(\alpha \hat{\mu}_X - \mu_X)^2] & X' \text{ and } \hat{\mu}_X \text{ are independent} \\ &= \sigma_X^2 + \mathbb{E}[(\alpha \hat{\mu}_X - \alpha \mu_X + \alpha \mu_X - \mu_X)^2] & \text{definition of variance} \\ &= \sigma_X^2 + \mathbb{E}[(\alpha \hat{\mu}_X - \alpha \mu_X)^2] + \mathbb{E}[(\alpha \mu_X - \mu_X)^2] & \text{cross terms vanish} \\ &= \sigma_X^2 + \alpha^2 \mathbb{E}[(\hat{\mu}_X - \mu_X)^2] + (\alpha - 1)^2 \mu_X^2 & \text{basic identities} \\ &= \sigma_X^2 + \alpha^2 \sigma_X^2 / N + (\alpha - 1)^2 \mu_X^2 & \text{var}[\hat{\mu}_X] = \sigma^2 / N \end{split}$$

Differentiating with respect to  $\alpha$  gives

$$\frac{d}{d\alpha}\mathbb{E}[(X'-\alpha\hat{\mu})^2] = 2\alpha\sigma_X^2/N + 2(\alpha-1)\mu_X^2$$

Setting this to zero and solving for  $\alpha$  gives the desired answer.

# Exercise 2.1, page 8

(a) A function that tests for vanishing correlation is the following

```
cor.equal.test = function(data1, data2, alpha=0.05) {
  ## THIS FUNCTION TESTS VANISHING CORRELATION BETWEEN
2
  ## TWO BI-VARIATE, NORMALLY DISTRIBUTED RANDOM VARIABLES
3
   ##
4
   # INPUT:
5
       DATA1: [N]-DIMENSIONAL VECTOR OF DATA
   #
6
   #
       DATA2: [N]-DIMENSIONAL VECTOR OF DATA
7
   #
       ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
8
   # OUTPUT LIST:
0
  #
      RHO: SAMPLE CORRELATION BETWEEN DATA1 AND DATA2
10
  #
      RHO.CRIT: 100*ALPHA% CRITICAL VALUE FOR THE CORRELATION
11
12
  #
     PVAL: P-VALUE OF THE STATISTIC RHO
13
  if (length(data1) != length(data2)) stop('data sets are not the same length')
14
15
  ntot = length(data1)
16
17
  mean1 = sum(data1)/ntot
18
   mean2 = sum(data2)/ntot
19
20
  ssr1 = sum((data1-mean1)^2)
21
ssr2 = sum((data2-mean2)^2)
23 sscx = sum((data1-mean1)*(data2-mean2))
24
25
  rho.hat = sscx / sqrt(ssr1*ssr2)
26
27 t.crit = qt(alpha/2,ntot-2,lower.tail=FALSE)
28 rho.crit = t.crit / sqrt( ntot - 2 + t.crit<sup>2</sup>)
29
  tval = rho.hat * sqrt(ntot-2) / sqrt(1-rho.hat^2)
30
   pval = 2*pt(abs(tval), ntot-2, lower.tail=FALSE)
31
32
   list(rho=rho.hat, rho.crit=rho.crit, pval=pval, tval=tval)
33
34
   }
```

(b) Applying this function to lagged PDO indices can be done as follows:

```
>print('early period')
1
2 >npic
           = 2:length(pdo1)
3 >pdo1.half1 = pdo1[npic-1]
4 >pdo1.half2 = pdo1[npic]
  >print(cor.equal.test(pdo1.half1,pdo1.half2))$rho
5
6
  [1] 0.05495659
7
8
9
  $rho.crit
10 [1] 0.3808629
11
12 $pval
13 [1] 0.7854252
14
15 $tval
16 [1] 0.2751988
```

```
print('later period')
1
  npic
         = 2:length(pdo2)
2
   pdo2.half1 = pdo2[npic-1]
3
  pdo2.half2 = pdo2[npic]
4
   print(cor.equal.test(pdo2.half1,pdo2.half2))
5
6
  $rho
7
   [1] 0.3836199
8
10
  $rho.crit
  [1] 0.3160319
11
12
  $pval
13
   [1] 0.01591744
14
15
16
   $tval
   [1] 2.526791
17
```

The year-to-year correlation for the first and second periods are 0.05 and 0.38. Only the second is large enough to reject the hypothesis of vanishing correlation at the 5% level, but only marginally (p-value is 2%).

(c) Performing the test on the whole period

```
print('full period')
1
          = 2:length(pdo)
  npic
2
   pdo.half1 = pdo[npic-1]
3
   pdo.half2 = pdo[npic]
4
   print(cor.equal.test(pdo.half1,pdo.half2))
5
6
   [1] "full period"
7
   $rho
8
   [1] 0.5081405
9
10
11 $rho.crit
12 [1] 0.1816386
13
14 $pval
  [1] 4.964102e-09
15
16
   $tval
17
   [1] 6.326905
18
19
   $rho.limits
20
   [1] 0.3597967 0.6314279
21
```

For the whole 1950-2017 period, the correlation is 0.50, which is statistically significant at the 5% level.



**Figure 22.2** Scatter plots of the JFM PDO index with itself lagged by one year, for the periods 1950-1977 (a), 1978-2017 (b), and 1950-2017 (c).

Why do the above conclusions differ? Following the hint, we draw scatter diagrams for the three periods. The result is shown in fig. 22.2. There is only a weak positive relation for the two periods individually, shown in (a) and (b), but a much stronger linear relation in the last (c). Closer examination reveals that the values in the later period tend to be above those of the lower period. This "shift in mean" produces an apparent positive relation, in the sense that it exists no matter how the data is sampled (as assumed under stationarity), but instead it is a property of the data as a whole that arises when all the data are pooled together.

(d) The answer to the other questions are

null hypothesis: population correlation vanishes

test statistic:

$$\hat{\rho} = \frac{\sum_{n} (x_n - \hat{\mu}_X) (y_n - \hat{\mu}_Y)}{\sqrt{\left(\sum_{n} (x_n - \hat{\mu}_X)^2\right) \left(\sum_{n} (y_n - \hat{\mu}_X)^2\right)}}$$
(22.17)

value of the test statistic: 0.05 (1950-1977), 0.38 (1978-2017), 0.49 (1950-2017)

- **rejection region:**  $|\hat{\rho}| > \rho_c$  where  $\rho_c = 0.38$  (1950-1977), 0.32 (1978-2017), 0.24 (1950-2017)
- final decision: can reject the null hypothesis  $\rho = 0$  in the period 1950-1977, but not for 1978-2017 (marginally) nor for the whole period 1950-2017.

**assumption:** paired samples (X, Y) are iid from a Gaussian distribution.

(e) A function that tests for equality of means is the following:

#### **102** ANSWERS TO EXERCISES

```
mean.equal.test = function(data1, data2, alpha=0.05) {
  ## THIS FUNCTION TESTS EQUALITY OF MEANS OF TWO IID
2
  ## NORMALLY DISTRIBUTED RANDOM VARIABLES
3
  ##
4
  # INPUT:
5
      DATA1: [N1]-DIMENSIONAL VECTOR OF DATA
   #
6
   #
       DATA2: [N2]-DIMENSIONAL VECTOR OF DATA
7
   #
      ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
8
  # OUTPUT LIST:
0
  #
      DIFF.MEAN: DIFFERENCE IN MEANS (MEAN1 - MEAN2)
10
  #
      DIFF.MEAN.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE DIFFERENCE IN MEANS
11
12
  #
      PVAL: P-VALUE OF THE T-STATISTIC
13
  #
      T: T-STATISTIC FOR DIFFERENCE IN MEANS
      T.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE T STATISTIC
14
  #
  #
     MEAN1: ESTIMATE OF THE MEAN OF DATA1
15
     MEAN2: ESTIMATE OF THE MEAN OF DATA2
  #
16
     SPOOL: POOLED ESTIMATE OF THE STANDARD DEVIATION
  #
17
  #
     DIFF.MEAN.LIMITS: [2] (1-ALPHA)100% CONFIDENCE LIMITS
18
19
  #
       FOR THE DIFFERENCE IN MEANS
20
21 nl
        = length(data1)
22 n2
        = length(data2)
23 \text{ dof} = n1 + n2 - 2
24
25 mean1 = sum(data1)/n1
_{26} mean2 = sum(data2)/n2
27 diff.mean = mean1 - mean2
28
29 ssr1 = sum((data1-mean1)^2)
30 ssr2 = sum((data2-mean2)^2)
31
32
  spool = sqrt((ssr1 + ssr2)/dof)
33
34 tval = diff.mean/spool/sqrt(1/n1 + 1/n2)
35 pval = 2* pt(abs(tval), dof, lower.tail=FALSE)
36 t.crit = qt(alpha/2,dof,lower.tail=FALSE)
37
38 diff.mean.crit = t.crit * spool * sqrt(1/n1 + 1/n2)
39
40 list (diff.mean=diff.mean, diff.mean.crit=diff.mean.crit, pval=pval,
      tval=tval,t.crit=t.crit,mean1=mean1,mean2=mean2,spool=spool)
41
42
  }
```

Applying this test to the PDO index gives:
```
> mean.equal.test(pdo1,pdo2)
1
  > print(mean.equal.test(pdo1,pdo2))
2
  $diff.mean
3
   [1] -1.156262
4
   $diff.mean.crit
6
   [1] 0.4537684
7
8
   $pval
9
   [1] 3.22733e-06
10
11
12
  $tval
13
   [1] -5.08751
14
   $t.crit
15
   [1] 1.996564
16
17
   $mean1
18
19
   [1] -0.7405952
20
   $mean2
21
   [1] 0.4156667
22
23
24
   $spool
25
   [1] 0.9223707
```

(f) The critical value of t is computed in line 38 of the above code.

(g) Answers:

null hypothesis: population means are equal test statistic:  $(\hat{\mu}_X - \hat{\mu}_Y) / s_p / \sqrt{1/N_X + 1/N_Y}$ value of the test statistic: -5.09 rejection region: |t| > 2.0final decision: the null hypothesis is rejected at the 5% level. assumption: samples are iid from a Gaussian distribution.

(h) The  $\rho = 0$  test is relevant to the equality-of-means test because the t-test assumes the samples are independent. The equality-of-variance test is relevant to to the equalityof-means test because the t-test assumes the variance are constant. Specifically, the fact that the year-to-year correlation in each sample is insignificant (or almost insignificant) implies that the independence assumption is reasonable for each sample separately (except possibly for the second period). The fact that the differences in variance are insignificant implies that the assumption of equal variances is reasonable.

Exercise 3.1, page 9

Since  $X_1, \dots, X_N \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$ ,  $\hat{\mu} \sim \mathcal{N}(\mu, \sigma^2/N)$ . (22.18)

Hypothesis tests for the mean are based on

$$\frac{\text{observed} - \text{hypothesized}}{\text{standard error}} = \frac{\hat{\mu} - \mu}{\sigma/\sqrt{N}} \sim \mathcal{N}(0, 1).$$
(22.19)

Replacing  $\sigma^2$  by the sample estimate  $\hat{\sigma}^2$  leads to the t-distribution:

$$t = \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} \sim t_{N-1}.$$
 (22.20)

By definition of  $t_{\alpha,N-1}$ ,

$$1 - \alpha = P\left(-t_{\alpha/2, N-1} \le t < t_{\alpha/2, N-1}\right)$$
(22.21)

$$= P\left(-t_{\alpha/2,N-1} \le \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} < t_{\alpha/2,N-1}\right)$$
(22.22)

$$= P\left(-t_{\alpha/2,N-1} \le \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} \text{ and } \frac{\hat{\mu} - \mu}{\hat{\sigma}/\sqrt{N}} < t_{\alpha/2,N-1}\right)$$
(22.23)

We want to solve for  $\mu$  by itself:

$$1 - \alpha = P\left(\mu \le \hat{\mu} + t_{\alpha/2, N-1}\hat{\sigma}/\sqrt{N} \text{ and } \hat{\mu} - t_{\alpha/2, N-1}\hat{\sigma}/\sqrt{N} < \mu\right)$$
(22.25)

$$= P\left(\hat{\mu} - t_{\alpha/2, N-1}\hat{\sigma}/\sqrt{N} < \mu \le \hat{\mu} + t_{\alpha/2, N-1}\hat{\sigma}/\sqrt{N}\right).$$
(22.26)

An R code for computing the coverage of this confidence interval is

```
ntrials = 100000
1
2 nsamp = 20
           = 10
3 mu
4 stdv
          = 3
         = 0.05
  alpha
5
           = array(rnorm(nsamp*ntrials,mean=mu,sd=stdv),dim=c(ntrials,nsamp))
  Х
7
  x.mean = rowMeans(x)
8
  x.stdv = sqrt(rowSums((x-x.mean)^2)/(nsamp-1)/nsamp)
9
  tcrit = qt(alpha/2,nsamp-1,lower.tail=FALSE)
10
11
12 mean.ci = cbind(x.mean-tcrit*x.stdv,x.mean+tcrit*x.stdv)
13 coverage = sum(mu >= mean.ci[,1] & mu <= mean.ci[,2])/ntrials</pre>
14
15 print(paste('coverage=',coverage))
  [1] "coverage= 0.94961"
16
```

The coverage is very close to 95%.

## Exercise 3.2, page 10

An augmented var.equal.test that calculates a confidence interval is the following:

var.equal.test = function(data1, data2, alpha=0.05) { 1 ### THIS FUNCTION TESTS EQUALITY OF VARIANCE OF TWO 2 ### IID NORMALLY DISTRIBUTED RANDOM VARIABLES 3 4 ### # INPUT: 5 DATA1: [N1]-DIMENSIONAL VECTOR OF DATA # 6 DATA2: [N2]-DIMENSIONAL VECTOR OF DATA 7 # ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%) # 8 # OUTPUT LIST: F.MAX: RATIO OF VARIANCE, CONSTRUCTED TO BE GREATER THAN 1 # 10 F.CRIT: THE UPPER CRITICAL THRESHOLD OF SIGNIFICANCE # 11 PVAL: P-VALUE OF THE F.MAX RATIO # 12 # VAR1: UNBIASED ESTIMATE OF THE VARIANCE OF DATA1 13 # VAR2: UNBIASED ESTIMATE OF THE VARIANCE OF DATA2 14 # RATIO: VAR1/VAR2 15 # RATIO.LOWER: LOWER LIMIT OF ALPHA% CONFIDENCE INTERVAL FOR VARIANCE RATIO 16 17 # RATIO.UPPER: UPPER LIMIT OF ALPHA% CONFIDENCE INTERVAL FOR VARIANCE RATIO 18 n1 = length(data1); n2 = length(data2) 19 20 mean1 = sum(data1)/n1; mean2 = sum(data2)/n2 21 22 var1 =  $sum((data1-mean1)^2)/(n1-1)$ 23  $var2 = sum((data2-mean2)^2)/(n2-1)$ 24 25 if ( var1 > var2) { 26 f.max = var1/var2 27 f.crit = qf(alpha/2,n1-1,n2-1,lower.tail=FALSE) 28 pval = 2\*pf(f.max,n1-1,n2-1,lower.tail=FALSE) 29 } else { 30 f.max = var2/var1 31 f.crit = qf(alpha/2,n2-1,n1-1,lower.tail=FALSE) 32 pval = 2\*pf(f.max,n2-1,n1-1,lower.tail=FALSE) 33 } 34 35 f.upper = qf(alpha/2,n1-1,n2-1,lower.tail=FALSE) 36 f.lower = qf(alpha/2, n1-1, n2-1, lower.tail=TRUE) 37 38 ratio = var1/var2 39 ratio.lower = ratio/f.upper 40 ratio.upper = ratio/f.lower 41 42 43 list(f.max=f.max,f.crit=f.crit,pval=pval,var1=var1,var2=var2, ratio=ratio, ratio.lower=ratio.lower, ratio.upper=ratio.upper) 44 } 45

Note that the critical F values f.upper and f.lower were recomputed after the if statement so that they correspond to the ratio var1/var2 and not its reciprocal.

The output of this function for the PDO index is the following

```
106 ANSWERS TO EXERCISES
```

```
n > nbreak = which( year == 1977)
2 > pdo1 = pdo[1:nbreak]
3 > pdo2 = pdo[(nbreak+1):nyrs]
4 >
  5
  > ## TEST EQUALITY OF VARIANCE
6
  7
8 > source(paste('var.equal.test.R', sep=""))
9 > print(var.equal.test(pdo1,pdo2))
10 $f.max
11 [1] 1.247894
12
13 $f.crit
14 [1] 2.074389
15
16 $pval
17 [1] 0.5522268
18
19
  $var1
20 [1] 0.7420675
21
22 $var2
23 [1] 0.9260218
24
25 $ratio
26 [1] 0.80135
27
28 $ratio.lower
29 [1] 0.4043894
30
31 $ratio.upper
32 [1] 1.662312
```

Or if the data is entered in the opposite direction...

```
> var.equal.test(pdo2,pdo1)
1
   $f.max
2
   [1] 1.247894
3
4
   $f.crit
5
   [1] 2.074389
6
7
   $pval
8
   [1] 0.5522268
9
10
11
  $var1
12
   [1] 0.9260218
13
   $var2
14
   [1] 0.7420675
15
16
   $ratio
17
18
   [1] 1.247894
19
   $ratio.lower
20
   [1] 0.6015718
21
22
  $ratio.upper
23
24
   [1] 2.472864
```

Accordingly, the 95% confidence interval for the ratio of variances is (0.39, 1.63), or (0.61, 2.54), depending on how the two variances are ordered in the ratio. Note that this interval includes one, implying that the ratio is not significantly different from one at the 5% significance level, consistent with the p-value.

Exercise 3.3, page 10

A modified version of cor.equal.test that also computes a confidence interval is the following:

```
cor.equal.test = function(data1, data2, alpha=0.05) {
1
  ## THIS FUNCTION TESTS VANISHING CORRELATION BETWEEN
2
  ## TWO BI-VARIATE, NORMALLY DISTRIBUTED RANDOM VARIABLES
3
  ##
4
  # INPUT:
5
      DATA1: [N]-DIMENSIONAL VECTOR OF DATA
   #
6
   #
      DATA2: [N]-DIMENSIONAL VECTOR OF DATA
7
      ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
  #
8
  # OUTPUT LIST:
0
  #
      RHO: SAMPLE CORRELATION BETWEEN DATA1 AND DATA2
10
  #
     RHO.CRIT: 100*ALPHA% CRITICAL VALUE FOR THE CORRELATION
11
12 # PVAL: P-VALUE OF THE STATISTIC RHO
13 # RHO.LIMITS: [2] (1-ALPHA/2) *100% CONFIDENCE LIMITS OF THE CORRELATION
14
if (length(data1) != length(data2)) stop('data sets are not the same length')
16
  ntot = length(data1)
17
18
  mean1 = sum(data1)/ntot
19
  mean2 = sum(data2)/ntot
20
21
22 ssr1 = sum((data1-mean1)^2)
ssr2 = sum((data2-mean2)^2)
24 sscx = sum((data1-mean1)*(data2-mean2))
25
26 rho.hat = sscx / sqrt(ssr1*ssr2)
27
28 t.crit = qt(alpha/2,ntot-2,lower.tail=FALSE)
29 rho.crit = t.crit / sqrt( ntot - 2 + t.crit<sup>2</sup>)
30
31
  tval = rho.hat * sqrt(ntot-2) / sqrt(1-rho.hat^2)
  pval = 2*pt(abs(tval), ntot-2, lower.tail=FALSE)
32
33
          = 0.5 * \log((1+rho.hat)/(1-rho.hat))
34 Z.
35 zc = qnorm(alpha/2,mean=0,sd=1/sqrt(ntot-3),lower.tail=FALSE)
36 z.limits = c(z - zc, z+zc)
37
_{38} rho.limits = (\exp(2*z.limits)-1)/(\exp(2*z.limits)+1)
39
40 list (rho=rho.hat, rho.crit=rho.crit, pval=pval, tval=tval, rho.limits=rho.limits)
41
  }
```

Applying this function to the PDO data gives the following

#### 108

```
1 > ## test vanishing correlation for first half
2 > 10 = 1:(length(pdo1)-1)
3 > pdo1.10 = pdo1[10]; pdo1.11 = pdo1[10+1]
4 > print(cor.equal.test(pdo1.l0,pdo1.l1))
₅ $rho
  [1] 0.05495659
6
7
8 $rho.crit
9 [1] 0.3808629
10
11 $pval
12 [1] 0.7854252
13
14 $tval
15 [1] 0.2751988
16
17 $rho.limits
18 [1] -0.3319908 0.4260723
```

```
1 > ## test vanishing correlation for 2nd half
_{2} > 10 = 1: (length(pdo2)-1)
3 > pdo2.10 = pdo2[10]; pdo2.11 = pdo2[10+1]
4 > print(cor.equal.test(pdo2.l0,pdo2.l1))
5 $rho
6 [1] 0.3836199
7
8 $rho.crit
9 [1] 0.3160319
10
11 $pval
12 [1] 0.01591744
13
14 $tval
15 [1] 2.526791
16
17 $rho.limits
18 [1] 0.07748101 0.62365106
```

```
> ## test vanishing correlation for whole period
1
   > 10 = 1: (length(pdo)-1)
2
   > pdo.10 = pdo[10]; pdo.11 = pdo[10+1]
3
   > print(cor.equal.test(pdo.10,pdo.11))
4
   $rho
5
   [1] 0.4915623
6
7
   $rho.crit
8
   [1] 0.2404471
9
10
11
  $pval
12
  [1] 2.398144e-05
13
   $tval
14
  [1] 4.550883
15
16
   $rho.limits
17
   [1] 0.2850066 0.6544904
18
```

The first confidence interval includes 0, implying that the correlation is not significantly difference from zero. The second confidence interval does not include 0, but only marginally. The third confidence interval does not include 0, hence the correlation is significantly different from zero. All three conclusions are consistent with the corresponding hypothesis tests.

# Exercise 3.4, page 11

A modified version of mean.equal.test that also computes confidence limits for the difference in means is the following:

```
mean.equal.test = function(data1, data2, alpha=0.05) {
  ## THIS FUNCTION TESTS EQUALITY OF MEANS OF TWO IID
2
  ## NORMALLY DISTRIBUTED RANDOM VARIABLES
3
   ##
4
   # INPUT:
5
       DATA1: [N1]-DIMENSIONAL VECTOR OF DATA
   #
6
   #
       DATA2: [N2]-DIMENSIONAL VECTOR OF DATA
7
   #
       ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
8
   # OUTPUT LIST:
0
   #
      DIFF.MEAN: DIFFERENCE IN MEANS (MEAN1 - MEAN2)
10
  #
      DIFF.MEAN.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE DIFFERENCE IN MEANS
11
12
  #
      PVAL: P-VALUE OF THE T-STATISTIC
      T: T-STATISTIC FOR DIFFERENCE IN MEANS
13
   #
      T.CRIT: 100*ALPHA% LEVEL CRITICAL VALUE OF THE T STATISTIC
   #
14
   #
      MEAN1: ESTIMATE OF THE MEAN OF DATA1
15
      MEAN2: ESTIMATE OF THE MEAN OF DATA2
   #
16
      SPOOL: POOLED ESTIMATE OF THE STANDARD DEVIATION
17
   #
   #
      DIFF.MEAN.LIMITS: [2] (1-ALPHA)100% CONFIDENCE LIMITS
18
   #
       FOR THE DIFFERENCE IN MEANS
19
20
  n1
        = length(data1)
21
22 n 2
         = length(data2)
23 \text{ dof} = n1 + n2 - 2
24
25 mean1 = sum(data1)/n1
_{26} mean2 = sum(data2)/n2
27 diff.mean = mean1 - mean2
28
29 ssr1 = sum((data1-mean1)^2)
  ssr2 = sum((data2-mean2)^2)
30
31
  spool = sqrt((ssr1 + ssr2)/dof)
32
33
  tval
         = diff.mean/spool/sqrt(1/n1 + 1/n2)
34
         = 2* pt(abs(tval),dof,lower.tail=FALSE)
  pval
35
  t.crit = qt(alpha/2, dof, lower.tail=FALSE)
36
37
38
  diff.mean.crit = t.crit * spool * sqrt(1/n1 + 1/n2)
39
40
  diff.mean.limits = diff.mean + c(-1,1) * diff.mean.crit
41
   list(diff.mean=diff.mean,diff.mean.crit=diff.mean.crit,pval=pval,
42
      tval=tval,t.crit=t.crit,mean1=mean1,mean2=mean2,spool=spool,
43
      diff.mean.limits=diff.mean.limits)
44
45
   }
```

Applying this to the PDO data gives

```
> print(mean.equal.test(pdo1,pdo2))
1
   $diff.mean
2
   [1] -1.156262
3
4
   $diff.mean.crit
5
   [1] 0.4537684
6
7
   $pval
8
   [1] 3.22733e-06
9
10
11
  $tval
12
   [1] -5.08751
13
   $t.crit
14
   [1] 1.996564
15
16
   $mean1
17
18
   [1] -0.7405952
19
20
   $mean2
   [1] 0.4156667
21
22
   $spool
23
24
   [1] 0.9223707
25
26
   $diff.mean.limits
   [1] -1.6100303 -0.7024935
27
```

The confidence limits do not include 0, so the difference in means is significantly different from zero. This is consistent with the hypothesis test.

## Exercise 3.5, page 12

An R function for constructing a confidence interval for a correlation coefficient is the following:

```
cor.equal.boot = function(x,y,alpha=0.05,nboot=100000) {
1
  ## ESTIMATES (1-ALPHA)100% CONFIDENCE INTERVAL
2
  ## FOR THE CORRELATION COEFFICIENT USING BOOTSTRAP METHODS
3
  ## INPUT:
4
  ## X[NSAMP]: FIRST DATA SET
5
   ## Y[NSAMP]: SECOND DATA SET
6
   ## ALPHA: EQUIVALENT SIGNIFICANCE LEVEL (DEFAULT = 5%)
7
   ## NBOOT: NUMBER OF BOOTSTRAP SAMPLES (DEFAULT = 100000)
8
  ## OUTPUT:
9
  ## CONFIDENCE LIMITS[2]
10
  nsamp = length(x)
11
12
  if (length(y) != nsamp) stop('x and y must be same length')
13
        = sample(nsamp,nboot*nsamp,replace=TRUE)
14 npic
15 x.boot = x[npic]
  y.boot = y[npic]
16
17
  dim(x.boot) = c(nboot, nsamp)
  dim(y.boot) = c(nboot, nsamp)
18
  x.boot = x.boot - rowMeans(x.boot)
19
  y.boot = y.boot - rowMeans(y.boot)
20
  cor.boot = rowSums(x.boot*y.boot)/sqrt(rowSums(x.boot^2)*rowSums(y.boot^2))
21
22 cor.boot.conf = quantile(cor.boot,probs=c(alpha/2,1-alpha/2))
23
  cor.boot.conf
24
25
  }
```

The result of running this function is:

```
nsamp = 20
1
          = 0.5
2 rho
3 nboot = 100000
  set.seed(310)
4
       = rnorm(nsamp)
5 X
  W
           = rnorm(nsamp)
6
7
  V
           = rho * x + sqrt(1-rho^2)*w
8
  rho.test1 = cor.equal.test(x,y)
9
  rho.test2 = cor.test(x,y)
10
11
  print(paste('Traditional Confidence interval:',
12
   paste(signif(rho.test1$rho.limits,4),collapse=', ')))
13
 print(paste('Bootstrap confidence interval :',
14
    paste(signif(cor.boot.conf,4),collapse=', ')))
15
16
17
18
   [1] "Traditional Confidence interval: 0.3301, 0.8601"
19
20
   [1] "Bootstrap confidence interval : 0.328, 0.8761"
```

The bootstrap interval is nearly identical to the standard interval.

## Exercise 4.1, page 14

An R function to perform the Wilcoxon-Mann-Whitney test is the following:

```
diff.mean.nonparametric = function(data1,data2,alpha=0.05) {
1
   ### PERFORMS MANN-WHITNEY TEST FOR A DIFFERENCE IN MEDIANS
2
   ### INPUT:
3
   ### DATA1[N1]: SAMPLE 1
4
   ###
        DATA2[N2]: SAMPLE 2
5
   ###
        ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 5%)
7
   ### OUTPUT: LIST$
        ZVAL: Z-VALUE FOR THE NORMAL APPROXIMATION TO MANN-WHITNEY TEST
8
   ###
   ### Z.CRIT: CRITICAL VALUE FOR ALPHA*100% SIGNIFICANCE
9
10 ### PVAL: P-VALUE OF THE TEST
11
                 = length(data1)
12 n1
   n2
                 = length(data2)
13

      14
      rank.12
      = rank(c(data1,

      15
      rank.1
      = rank.12[1:n1]

      16
      rank.2
      = rank.12[1:n2+

                = rank(c(data1,data2))
                = rank.12[1:n2+n1]
17 rank.sum.1 = sum(rank.1)
18 rank.sum.2 = sum(rank.2)
19
20 mu.rank.sum = n1 * (n1+n2+1)/2
var.rank.sum = n1*n2*(n1+n2+1)/12
22 if (rank.sum.1 > mu.rank.sum) continuity.cor = -0.5 else continuity.cor = 0.5
           = (rank.sum.1 - mu.rank.sum + continuity.cor)/sqrt(var.rank.sum)
23 zval
             24 z.crit
                = 2*pnorm(abs(zval),lower.tail=FALSE)
25
   pval
26
   list(zval=zval, z.crit=z.crit, pval=pval,
27
   rank.sum.1=rank.sum.1,rank.sum.2=rank.sum.2)
28
29
   }
```

## Exercise 4.2, page 14

The result of applying this function to the PDO time series is:

```
> print(diff.mean.nonparametric(data1,data2))
1
   $zval
2
   [1] -4.243011
3
4
   $z.crit
5
   [1] 1.959964
6
7
   $pval
8
   [1] 2.205404e-05
9
10
11
  $wval
12
  [1] 219
13
  $rank.sum.1
14
   [1] 625
15
16
  $rank.sum.2
17
18
   [1] 1721
```

The z-value is significant, with a p-value of 0.022%, which is much less than 5%, hence the null hypothesis of equal medians should be rejected. This conclusion is consistent with the t-test from previous homeworks.

Exercise 4.3, page 14

An R function that computes Spearman's rank correlation and performs a hypothesis test is the following

```
cor.test.nonparametric = function(data1, data2, alpha=0.05) {
1
  ### PERFORMS RANK CORRELATION TEST
2
  ### INPUT:
3
         DATA1[N1]: SAMPLE 1
  ###
  ###
         DATA2[N2]: SAMPLE 2
5
   ###
        ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 5%)
6
   ### OUTPUT: LIST$
7
  ###
        RHO.SPEARMAN = SPEARMAN'S RANK CORRELATION
8
  ###
       ZVAL: Z-VALUE FOR THE NORMAL APPROXIMATION TO MANN-WHITNEY TEST
9
10 ###
       Z.CRIT: CRITICAL VALUE FOR ALPHA*100% SIGNIFICANCE
11 ### PVAL: P-VALUE OF THE TEST
12
if (length(data1) != length(data2)) stop('data sets not of equal length')
14 ntot = length(data1)
15 r.mean = (ntot+1)/2
16 r.var = ntot*(ntot^2-1)/12
17 rank.1 = rank(data1)
  rank.2 = rank(data2)
18
19
20
  rho.spearman = sum((rank.1-r.mean)*(rank.2-r.mean))/r.var
21
22 zval = rho.spearman * sqrt(ntot - 1)
23 z.crit = qnorm(alpha/2,lower.tail=FALSE)
24 pval = 2*pnorm(abs(zval),lower.tail=FALSE)
25
26 list(rho.spearman=rho.spearman,zval=zval,z.crit=z.crit,pval=pval)
27
28
   }
```

## Exercise 4.4, page 15

The result of testing correlations for the 3 data sets are as follows:

```
1 > ## test vanishing correlation for first half
2 > 10 = 1:(length(pdol)-1)
3 > pdol.10 = pdol[10]; pdol.11 = pdol[10+1]
4 > print(cor.test.nonparametric(pdol.10,pdol.11))
5 $rho.spearman
6 [1] 0.05494505
7
8 $zval
9 [1] 0.2801659
10
11 $z.crit
12 [1] 1.959964
13
14 $pval
15 [1] 0.7793502
```

```
1 > ## test vanishing correlation for 2nd half
  > 10 = 1: (length(pdo2)-1)
2
  > print(cor.test.nonparametric(pdo2[10],pdo2[10+1]))
3
  $rho.spearman
4
   [1] 0.3487854
5
6
   $zval
7
   [1] 2.150058
8
9
  $z.crit
10
11
  [1] 1.959964
12
13 $pval
14 [1] 0.03155065
```

```
> ## test vanishing correlation for whole period
1
2 > 10 = 1: (length (pdo12)-1)
  > print(cor.test.nonparametric(pdo12[10],pdo12[10+1]))
3
  $rho.spearman
4
   [1] 0.4570496
5
6
   $zval
7
   [1] 3.713089
8
9
10
  $z.crit
11
   [1] 1.959964
12
13
  $pval
   [1] 0.000204745
14
```

The null hypothesis is rejected in the first half, but not in the second half, or in the whole period. This result is consistent with the parametric correlation test.

Exercise 4.5, page 15

A function that tests difference in dispersions is

```
118 ANSWERS TO EXERCISES
```

```
diff.dispersion = function(data1, data2, alpha=0.05) {
1
   ## PERFORMS WILCOXON SQUARED-RANK TEST ON ABSOLUTE DEVIATIONS |X - MEDIAN(X)|
2
   ### INPUT:
3
   ###
         DATA1[N1]: SAMPLE 1
4
   ###
         DATA2[N2]: SAMPLE 2
5
   ###
         ALPHA: SIGNIFICANCE LEVEL (DEFAULT = 5%)
6
   ### OUTPUT: LIST$
7
   ###
         ZVAL: Z-VALUE FOR THE NORMAL APPROXIMATION TO MANN-WHITNEY TEST
8
   ###
        Z.CRIT: CRITICAL VALUE FOR ALPHA*100% SIGNIFICANCE
0
10 ### PVAL: P-VALUE OF THE TEST
11
12 n1
                   = length(data1)
13 n2
                   = length(data2)
14 ntot
                   = n1 + n2
15
16 median1 = median(data1)
17 median2 = median(data2)
18 absdev = c(abs(data1-median1),abs(data2-median2))
19 rank.12 = rank(absdev)
19 rank.12 = rank(absdev)
                  = rank.12[1:n1]
20 rank.1 = rank.12[1:n1]
21 rank.2 = rank.12[1:n2+n1]
22 rank.sqr.sum.1 = sum(rank.1^2)
23 rank.sqr.sum.2 = sum(rank.2<sup>2</sup>)
24
25 mu.rank.sum = n1*(ntot+1)*(2*ntot+1)/6
26 var.rank.sum = n1*(ntot-n1)/(ntot-1)/ntot*sum((((1:ntot)^2-mu.rank.sum/n1)^2))
= 2*pnorm(abs(zval),lower.tail=FALSE)
29 pval
30
31 list(zval=zval,z.crit=z.crit,pval=pval,
rank.sqr.sum.1=rank.sqr.sum.1,rank.sqr.sum.2=rank.sqr.sum.2)
33
34 }
```

The result of running this function on the PDO data set is

```
> print(diff.dispersion(data1,data2))
2 $zval
  [1] -0.9649911
3
4
  $z.crit
5
  [1] 1.959964
6
7
  $pval
8
9 [1] 0.3345493
10
11 $mu.rank.sum
12 [1] 44114
13
14 $var.rank.sum
15 [1] 32644360
16
17 $rank.sqr.sum.1
18
  [1] 38600.5
19
20
  $rank.sqr.sum.2
  [1] 68532
21
```

The null hypothesis is not rejected, and this is consistent with previous homeworks.

# Exercise 5.1, page 17

An R function that calculates the sample autocorrelation function is the following:

```
acf.brute = function(x, lag.max=NULL) {
1
   ## COMPUTES THE AUTOCORRELATION FUNCTION OF A TIME SERIES
2
   ## INPUT:
3
        X: [NTOT]-LENGTH NUMERICAL VECTOR OF THE TIME SERIES
   ##
   ##
        LAG.MAX: MAXUMIM LAG TO COMPUTE
5
   ##
          (DEFAULT = MAX(1, FLOOR(10 * log10(LENGTH(X)))))
6
   ## OUTPUT:
7
   ##
       X.ACF: [1:(LAG.MAX+1)] AUTOCORRELATION FUNCTION OF X
8
   ##
          FOR LAGS 0 TO LAG.MAX
9
10
  ntot = length(x)
11
12
if ( is.null(lag.max)) lag.max = max(1,floor(10*log10(ntot)))
14 x.anom = x - mean(x)
  x.ssa = sum(x.anom<sup>2</sup>)
15
16
   x.acf = numeric(lag.max+1)
17
   for ( lag in 0:lag.max) {
18
   npic = 1:(ntot-lag)
19
   x.acf[lag+1] = sum(x.anom[npic]*x.anom[npic+lag])/x.ssa
20
21
   }
22
23 x.acf
24
   }
```

Result of applying this function to x = rnorm(20) is:

```
1 > set.seed(1)
_{2} > x = rnorm(20)
  > acf.mine = acf.brute(x)
           = acf(x)
  > acf.R
4
  > print(acf.mine)
5
    [1] 1.00000000 -0.12228595 -0.18521142 -0.04940401 0.14719994
6
   [6] -0.28278013 -0.25514538 0.21242725 0.09722174 -0.11993122
  [11] -0.18075942 0.28550651 -0.06259741 0.09444605
  > print(as.numeric(acf.R$acf))
10
    [1] 1.00000000 -0.12228595 -0.18521142 -0.04940401 0.14719994
   [6] -0.28278013 -0.25514538 0.21242725 0.09722174 -0.11993122
11
  [11] -0.18075942 0.28550651 -0.06259741 0.09444605
12
```

### Exercise 5.2, page 18

The stationary distribution of the AR1 process is normal with mean  $k/(1-\phi_1)$  and variance  $\sigma_W^2/(1-\phi^2)$ . An R function that generates a time series from an AR1 model is the following

```
ar1.ts = function(phi1, ntot, stdv=1, constant=0, iseed=1) {
1
  ## GENERATE A REALIZATION FROM THE AR1 MODEL
2
  ## X(T+1) = PHI1 * X(T) + W(T) + CONSTANT
3
   ## WHERE W(T) IS GAUSSIAN WHITE NOISE WITH MEAN 0 AND ST. DEV. 'STDV'
   ## INITIAL CONDITION IS RANDOMLY DRAWN FROM STATIONARY DISTRIBUTION
5
   ## TNPUT:
6
   #
        PHI1: AUTOREGRESSIVE PARAMETER
7
   #
        NTOT: LENGTH OF THE DESIRED TIME SERIES
8
   #
        STDV: STANDARD DEVIATION OF THE NOISE (DEFAULT = 1)
0
   #
        CONSTANT: CONSTANT TERM IN THE AR1 MODEL (DEFAULT = 0)
10
  #
        ISEED: SEED FOR THE RANDOM NUMBER GENERATOR (DEFAULT = 1)
11
12
  # OUTPUT:
13
  #
        X: [NTOT] - RANDOM TIME SERIES FROM AR(1) PROCESS
14
  set.seed(iseed)
15
16
          = numeric(ntot)
   Х
17
         = rnorm(1,mean=constant/(1-phi),sd=stdv/sqrt(1-phi^2))
18
   x[1]
   W
          = rnorm(ntot, sd=stdv)
19
  for (n in 2:ntot) x[n] = phi1 * x[n-1] + w[n] + constant
20
  x }
21
```

### Exercise 5.3, page 18

Time series and corresponding autocorrelation functions can be generated as follows:

```
pdf('ts.acf.pdf')
2 par(mfrow=c(3,2),mar=c(5,5,3,2),cex.axis=1.5,cex.lab=1.5,cex.main=1.5)
  ntot = 50
3
  for ( phi1 in c(0,0.5,0.9)) {
4
        = ar1.ts(phi1,ntot,iseed=2)
     Х
     plot(x,type="l",col='blue',xlab='time',lwd=2,
6
       main=paste('Time Series for phi1=', signif(phi1, 2)))
     acf.mine = acf.brute(x)
8
     lag.max = length(acf.mine)-1
     yrange
             = c(-0.4, 1)
10
     plot(0:lag.max,acf.mine,xlab="lag",ylab="ACF",type="l",ylim=yrange,lwd=2,
11
        main=paste('Autocorrelation for phi1=', signif(phi1,2)))
12
     acf.exact = phi1^(1:ntot-1)
13
14
     par(new=TRUE)
     plot(1:ntot-1,acf.exact,type="1",lty="solid",col="red",xlim=c(0,lag.max),
15
       ylim=yrange,xlab="",ylab="",axes=FALSE,lwd=2)
16
     abline(h=2/sqrt(ntot), col="blue", lty="dashed", lwd=2)
17
     abline (h=-2/sqrt(ntot), col="blue", lty="dashed", lwd=2)
18
     legend('topright',legend=c('sample','population'),col=c('black','red'),lwd=2)
19
   1
20
   dev.off()
21
```

A sample plot is shown in fig. 22.3.



Figure 22.3 Sample time series and corresponding autocorrelation function from an AR1 model.

As  $\phi_1 \rightarrow 1$ , the time series get "smoother" and the ACF decays more slowly. The sample ACF for  $\phi_1 = 0.9$  is "surprising" because it looks very different from the true ACF: the sample ACF oscillates, changes sign, and exceeds the 95% significance level. However, these results are not inconsistent with theory, because we expect an average of 0.75 autocorrelations out of 15 to exceed the 95% level. Since the true ACF is relatively high for all lags less than 15, any high correlation due to sampling variability will exhibit a smooth oscillation. The ACF is especially misleading because the trough and peak occur at multiples of 7, suggesting the existence of a oscillation with a period of 14, when in fact the process is strictly AR(1).

## Exercise 5.4, page 18

An R code that generates time series, splits it into 2, and performs a t-test, is the following:

```
source('/Users/delsole/R/delsole_tools/mean.equal.test.R')
1
   ntrials = 10000
2
           = 25
   nx
3
            = 25
4
   ny
           = numeric(ntrials)
5
   tval
   yrange = c(0, 0.4)
6
   for ( phi1 in c(0, 0.5, 0.9)) {
7
     for ( nt in 1:ntrials) {
8
       z = ar1.ts(phi1,nx+ny,iseed=nt)
9
       x = z[1:nx]
10
       y = z [(nx+1):(nx+ny)]
11
       tval[nt] = mean.equal.test(x,y)$tval
12
     }
13
     t1 = 1 + 2 * sum(phi1^(1:length(x)) * (1-(1:length(x)/ntot)))
14
     ftitle=paste('Histogram of T for phil=', signif(phil, 2))
15
     hist(tval,col="grey",freq=FALSE,main=ftitle,ylim=yrange)
16
     dof = nx + ny - 2
17
     curve(dt(x,dof),add=TRUE)
18
19
   }
```

The plot for a particular realization is shown in fig. 22.4. The histograms show that as the autocorrelation parameter  $\phi_1$  increases, the variance of t-values increases. This implies that performing the t-test on serially correlated data will tend to incorrectly reject the null hypothesis at a greater rate than the chosen significance level. Equivalently, a sample that is deemed "significant" by the standard decision rule may actually be insignificant when autocorrelation is taken into account.

### Exercise 5.5, page 19

Taking the expectations of both sides of (5.2) gives

$$\mathbb{E}[X_t] = \phi_1 \mathbb{E}[X_{t-1}] + \phi_2 \mathbb{E}[X_{t-2}] + \dots + \phi_p \mathbb{E}[X_{t-p}] + \mathbb{E}[W_t] + k.$$
(22.27)



**Figure 22.4** Histograms of the t-statistic for testing equality of means calculated from time series generated by an AR1 model with  $\phi_1 = 0, 0.5, 0.9$ . The time series are each of length 25.

Since the process is stationary, all expectations of  $X_{t+\tau}$  are identical. Solving for the expectation therefore gives

$$\mathbb{E}[X_t] = \frac{k}{1 - \phi_1 - \phi_2 - \dots - \phi_p},$$
(22.28)

where  $\mathbb{E}[W_t] = 0$  has been used.

To calculate the variance, first subtract the mean from both sides of (5.2), which yields

$$X_t - \mu = \phi_1 \left( X_{t-1} - \mu \right) + \phi_2 \left( X_{t-2} - \mu \right) + \dots + \phi_p \left( X_{t-p} - \mu \right) + W_t, \quad (22.29)$$

where the k term vanishes because of (22.28). Multiply both sides by  $X_t - \mu$  and take expectations:

$$\operatorname{var}[X_t] = \phi_1 \operatorname{cov}[X_t, X_{t-1}] + \phi_2 \operatorname{cov}[X_t, X_{t-2}] + \dots + \phi_p \operatorname{cov}[X_t, X_{t-p}] + \operatorname{cov}[X_t, W_t].$$
(22.30)

The last term can be evaluated by substituting (5.2) for  $X_t$ , which yields

$$\operatorname{cov}[X_t, W_t] = \operatorname{cov}[\phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + W_t + k, W_t] = \operatorname{var}[W_t], \ (22.31)$$

since  $W_t$  is independent of  $X_{t-\tau}$  for all  $\tau \ge 1$ . Substituting the relation  $\operatorname{cov}[X_t, X_{t-\tau}] = \operatorname{var}[X_t]\rho_{\tau}$  and solving for the variance gives

$$\operatorname{var}[X_t] = \frac{\operatorname{var}[W_t]}{1 - \rho_1 \phi_1 - \rho_2 \phi_2 - \dots - \rho_p \phi_p}.$$
(22.32)

[ . . . . ]

### Exercise 5.6, page 19

The expectation of  $Y_t$  is

$$\mathbb{E}[Y_t] = \frac{1}{2K+1} \sum_{k=-K}^{K} \mathbb{E}[X_{t+k}] = 0.$$
(22.33)

The variance of  $Y_t$  is

$$\operatorname{var}[Y_t] = \frac{1}{(2K+1)^2} \sum_{k=-K}^{K} \sum_{j=-K}^{K} \mathbb{E}\left[X_{t+k} X_{t+j}\right].$$

Since  $X_t$  is white noise,  $\mathbb{E}[X_{t+k}X_{t+j}]$  vanishes except when k = j. Thus, the double sum collapses to a single sum:

$$\operatorname{var}[Y_t] = \frac{1}{(2K+1)^2} \sum_{k=-K}^{K} \sigma^2 \\ = \frac{\sigma^2}{(2K+1)}.$$

The autocovariance function is

$$\operatorname{cov}[Y_{t+\tau}, Y_t] = \frac{1}{(2K+1)^2} \sum_{k=-K}^{K} \sum_{j=-K}^{K} \mathbb{E}\left[X_{t+\tau+k} X_{t+j}\right]$$
(22.34)

Since  $X_t$  is white noise,  $\mathbb{E}[X_{t+\tau+k}X_{t+j}]$  vanishes except when  $k + \tau = j$ . When  $\tau = 0$ , this identity holds 2K + 1 times. When  $\tau = 1$ , this identity holds 2K times. For general  $\tau$ , the identity holds  $2K + 1 - |\tau|$  times. Thus

$$\operatorname{cov}[Y_{t+\tau}, Y_t] = \frac{\sigma^2}{(2K+1)^2} \left(2K+1 - |\tau|\right) \quad \text{for } |\tau| <= 2K+1$$
(22.35)

The autocorrelation is therefore

$$\operatorname{cor}[Y_{t+\tau}, Y_t] = \frac{\operatorname{cov}[Y_{t+\tau}, Y_t]}{\operatorname{var}[Y_t]} = 1 - \frac{|\tau|}{2K+1} \quad \text{for } |\tau| <= 2K+1$$

For  $|\tau| > 2K + 1$ , the autocorrelation of  $Y_t$  vanishes. Thus, the autocorrelation function of  $Y_t$  decreases linearly from 1 until it vanishes, and then remains 0 thereafter.

The fact that running means are correlated even when the underlying process is not follows from the fact that a running mean is computed from the time mean within a window, so when the two windows intersect they contain the same data, and hence are correlated.

Exercise 6.1, page 23

A spectral averaging function is the following:

```
spectrum.smooth = function(pgram,window,alpha=0.1) {
1
  ## COMPUTES SPECTRAL RUNNING MEAN AS:
2
  ##
        SPEC_AVE = SMOOTHED PERIODOGRAM FROM M = -WINDOW TO M = WINDOW
3
   ## INPUT:
4
         PGRAM[NTOT/2]: PERIODOGRAM VALUES FROM PERIODOGRAM.R
   ##
5
   ##
         WINDOW: HALF-WIDTH OF AVERAGING WINDOW
6
   ##
         ALPHA: THE (1-ALPHA)100% LEVEL FOR THE CONFIDENCE LIMITS
7
   ## OUTPUT: LIST
8
   ##
         PGRAM.AVE[NTOT/2]: AVERAGED PERIODOGRAM
0
   ##
         FREQ[NTOT/2]: FREQUENCIES (CYCLES PER TIME STEP)
10
   ##
        BW[2]: BANDWIDTH C(-WINDOW/NTOT,WINDOW/NTOT)
11
12
  ##
        CI.LIMITS[2]: PROPORTIONALITY CONSTANT FOR LOWER AND UPPER CONFIDENCE LIMITS
            (DOF/CHI-SQUARE[ALPHA/2, DOF], DOF/CHI-SQUARE[1-ALPHA/2, DOF])
13
   ##
   ##
       PGRAM[NTOT/2]: REPEATS PERIODOGRAM (FOR LATER PLOTTING)
14
15
    nhalf
               = length(pgram)
16
              = nhalf \star 2
17
    ntot
     freq
              = 1:(ntot/2) / ntot
18
     dof
               = 2 * (2 * window + 1)
19
20
     pgram.ave = as.numeric(rep(NA, nhalf))
21
     for ( n in (window+1): (nhalf-window-1)) pgram.ave[n] =
22
         mean(pgram[(n-window):(n+window)])
23
24
25
     bw = c(-window, window) / ntot
26
     ci = dof/qchisq(c(alpha/2,1-alpha/2),dof)
27
     list (pgram.ave=pgram.ave, freq=freq, bw=bw, ci=ci, pgram=pgram)
28
29
   }
```

## Exercise 6.2, page 23

The output from the spectral estimation is:

1	> set.seed(1)				
2	> x.sd = 1				
3	> window = 5				
4	> ntot = 64				
5	> freq = 1:(ntot/	2) / ntot			
6	> x = rnorm(nt	ot,sd=x.sd)			
7	> x.pgram = periodogram(x)				
8	<pre>&gt; x.pgram.smooth = spectrum.smooth(x.pgram,window)</pre>				
9	<pre>&gt; cbind(freq,x.pgram,x.p</pre>	gram.smooth\$pgram.ave)			
10	freq x.pgra	n			
11	[1,] 0.015625 0.4695402	8 NA			
12	[2,] 0.031250 0.0867886	9 NA			
13	[3,] 0.046875 0.2007800	7 NA			
14	[4,] 0.062500 0.0881236	2 NA			
15	[5,] 0.078125 0.7340169	0 NA			
16	[6,] 0.093750 2.1562158	1 0.8100061			
17	[7,] 0.109375 0.4095593	7 0.7723600			
18	[8,] 0.125000 0.2681699	6 0.8389093			
19	[9,] 0.140625 1.1434929	7 0.8596130			
20	[10,] 0.156250 1.0305373	9 0.9965911			
21	[11,] 0.171875 2.3228423	4 1.0368931			
22	[12,] 0.187500 0.0554327	6 0.9298071			
23	[13,] 0.203125 0.8188309	0 1.0700647			
24	[14,] 0.218750 0.4285212	3 1.0775075			
25	[15,] 0.234375 1.5948823	8 1.0317232			
26	[16,] 0.250000 1.1773389	7 0.9613558			
27	[17,] 0.265625 0.9782700	6 0.7704560			
28	[18,] 0.281250 1.9523923	4 0.8123769			
29	[19,] 0.296875 0.3500410	3 0.8457826			
30	[20,] 0.312500 0.6398652	3 0.9243640			
31	[21,] 0.328125 0.2564968	5 0.8015068			
32	[22,] 0.343750 0.2229443	0 0.7269105			
33	[23,] 0.359375 0.5165629	7 0.7570746			
34	[24,] 0.375000 1.1862929	2 0.6136065			
35	[25,] 0.390625 1.2929173	6 0.6251425			
36	[26,] 0.406250 0.2434529	2 0.5862000			
37	[27,] 0.421875 0.3567790	5 NA			
38	[28,] 0.43/500 1.3100758	8 NA			
39	[29,] 0.453125 0.3742433	/ NA			
40	[30,] 0.468750 0.4769365	3 NA			
41	[31,] 0.484375 0.2114980	8 NA			
42	[32,] 0.500000 1.9593428	9 NA			

# The plot produced by

```
spectrum.plot(x.pgram.smooth)
```

```
abline(h=x.sd,lty="dashed",lwd=2)
```

```
3 ftitle = paste("White Noise (N=",ntot,")",sep="")
```

```
4 title(main=ftitle)
```

is shown in fig. 22.5. To judge whether the spectrum is consistent with a white noise process, we must to decide whether the power spectrum is "flat," or constant. Technically, to do this, we would compare the ratio of spectral estimates with an appropriate threshold from an F-distribution. Instead of doing this, we crudely use the confidence interval by checking whether the *differences* between spectral estimates exceed the confidence interval. It can be seen that no two spectral estimates differ by more than a confidence interval, so we conclude the spectrum is consistent with a white noise process.



**Figure 22.5** Periodogram and spectral average of a white noise process. The error bars in the upper right show the 90% confidence interval and the bandwidth. The horizontal dashed line shows the theoretical power spectrum of the process.

## Exercise 6.3, page 23

The theoretical power spectrum for  $GWN(0,\sigma^2)$  is derived from the Discrete Fourier Transform

$$p(\omega) = c_0 + 2\sum_{\tau=1}^{\infty} c_\tau \cos(\omega\tau)$$
(22.36)

For GWN,  $c_{\tau} = 0$  for all  $\tau \ge 1$ , so the power spectrum is easily evaluated to be

$$p(\omega) = c_0 = \sigma^2.$$
 (22.37)

This power spectrum can be superimposed on the plot using the commands:

```
n exact.pgram = x.sd^2
```

abline(h=exact.pgram,lwd=2,lty="dashed",col="black")

The difference between the theoretical and estimated power spectra lies within a confidence interval, so we conclude the estimated spectra is consistent with the theoretical power spectrum.

# Exercise 6.4, page 24

The output from the spectral estimation is:

1	> window	= 5		
2	> x.pgram	= periodogra	am(x)	
3	> x.pgram.smoo	th = spectru	um.smooth(x	.pgram,window)
4	> cbind(freq,x	.pgram,x.pg	ram.smooth\$	Spgram.ave)
5	freq	x.pgram		
6	[1,] 0.015625	4.12171567	NA	
7	[2,] 0.031250	4.28688721	NA	
8	[3,] 0.046875	2.18587695	NA	
9	[4,] 0.062500	0.80868492	NA	
10	[5,] 0.078125	5.00694613	NA	
11	[6,] 0.093750	8.68718662	3.0192022	
12	[7,] 0.109375	0.45912976	2.6531974	
13	[8,] 0.125000	0.69869530	2.3055704	
14	[9,] 0.140625	2.36259034	2.1226395	
15	[10,] 0.156250	2.03132280	2.1816207	
16	[11,] 0.171875	2.56218866	1.8167134	
17	[12,] 0.187500	0.09566224	1.1014335	
18	[13,] 0.203125	0.46299053	1.1621120	
19	[14,] 0.218750	0.17363724	1.1047530	
20	[15,] 0.234375	1.45747803	0.9041859	
21	[16,] 0.250000	0.99296599	0.7397831	
22	[17,] 0.265625	0.81910720	0.5225284	
23	[18,] 0.281250	1.12659371	0.5375536	
24	[19,] 0.296875	0.06774579	0.5196011	
25	[20,] 0.312500	0.15635229	0.5302911	
26	[21,] 0.328125	0.22289282	0.4021078	
27	[22,] 0.343750	0.17238675	0.3240738	
28	[23,] 0.359375	0.26093924	0.2755012	
29	[24,] 0.375000	0.26551307	0.1774613	
30	[25,] 0.390625	0.29122684	0.1816354	
31	[26,] 0.406250	0.04746208	0.1769702	
32	[27,] 0.421875	0.13459156	NA	
33	[28,] 0.437500	0.28480903	NA	
34	[29,] 0.453125	0.04815466	NA	
35	[30,] 0.468750	0.11366153	NA	
36	[31,] 0.484375	0.10503442	NA	
37	[32,] 0.500000	0.81701638	NA	

A plot of the estimated power spectra is shown in fig. 22.6.



**Figure 22.6** Periodogram and spectral average of an AR(1) process with  $\phi = 0.8$ . The error bars in the upper right show the 90% confidence interval and the bandwidth. The horizontal dashed line shows the theoretical power spectrum of the process.

### Exercise 6.5, page 24

The theoretical power spectrum for an AR(1) process was derived in the text as

$$p(\omega) = \frac{\sigma_{\epsilon}^2}{1 + \phi_1^2 - 2\phi_1 \cos(2\pi f)}.$$
(22.38)

R commands for evaluating this power spectrum are as follows:

```
1 freq = 1:(ntot/2) / ntot
2 x.spec.exact = 1/(1+phi1^2-2*phi1*cos(2*pi*freq))
3 lines(freq,x.spec.exact,lty="dashed",lwd=3)
```

The final plot is shown in fig. 22.6.

## Exercise 6.6, page 24

The power spectrum for NINO3.4 can be estimated using the following:

```
fname
              = paste(dir.enso,'nina34.data', sep='')
1
              = scan(fname, nlines=1, quiet=TRUE)
2 line1
3 nino34.yr = line1[1]:line1[2]
   nino34.nyrs = length(nino34.yr)
              = read.table(fname,skip=1,nrows=nino34.nyrs,na.strings=-99.99,col.names=cnames)
   data2
6
   ## CREATE VECTOR WITH NINO34 TIME SERIES AND TIME AXIS
7
  nino34.ts = as.numeric(t(as.matrix(data2[,1:12+1])))
8
9 nino34.time = seq(from=data2[1,1],length.out=length(nino34.ts),by=1/12)
10
nomit = is.na(nino34.ts)
nino34.ts = nino34.ts[!nomit]
nino34.time = nino34.time[!nomit]
14
15 ### FIND HIGHLY COMPOSITE NUMBER
16 chunk = 10
17
  length.composite = length(nino34.time)
   while ( nextn(length.composite) > length(nino34.time)) length.composite =
18
    length.composite - chunk
19
   length.composite = nextn(length.composite)
20
21
22 ### INCLUDE ONLY HIGHLY COMPOSITE NUMBER OF SAMPLES
23 nino34.ts = nino34.ts [1:length.composite + length(nino34.ts) - length.composite]
24 nino34.time = nino34.time[1:length.composite + length(nino34.ts) - length.composite]
25
26 nino34.pgram
                    = periodogram(nino34.ts)
27 nino34.pgram.smooth = spectrum.smooth(nino34.pgram,window)
  spectrum.plot(nino34.pgram.smooth)
28
```

The resulting power spectra is shown in fig. 22.7.

There are significant peaks near  $1/12 \approx 0.083$  and  $1/6 \approx 0.17$ , which are obviously associated with the annual and sub-annual cycles.

The largest peak occurs near  $1/48 \approx 0.02$ , which is near 4-years. However, the peak does not rise outside the confidence interval relative to its neighboring values, so it is not objectively clear that the peak is at 4-years or at even lower frequencies. The characterization that ENSO has dominant oscillatory periods around 4-7 years is not unambiguous.

Exercise 7.1, page 25

$$\mathbb{E}[\mathbf{y}] = \mathbb{E}\left[\sum_{j=1}^{N} c_j \mathbf{x}_j\right] = \sum_{j=1}^{N} c_j \mathbb{E}[\mathbf{x}_j] = c_1 \boldsymbol{\mu}_1 + \dots c_N \boldsymbol{\mu}_N$$



**Figure 22.7** Periodogram and spectral average of the raw NINO3.4 index. The error bars in the upper right show the 90% confidence interval and the bandwidth.

$$\operatorname{cov}[\mathbf{y}] = \mathbb{E}\left[\left(\sum_{j=1}^{N} c_j \left(\mathbf{x}_j - \boldsymbol{\mu}_j\right)\right) \left(\sum_{k=1}^{N} c_k \left(\mathbf{x}_k - \boldsymbol{\mu}_k\right)\right)^T\right]$$
$$= \sum_{j=1}^{N} \sum_{k=1}^{N} c_j c_k \mathbb{E}\left[\left(\mathbf{x}_j - \boldsymbol{\mu}_j\right) \left(\mathbf{x}_k - \boldsymbol{\mu}_k\right)^T\right]$$
$$= \sum_{j=1}^{N} c_j^2 \boldsymbol{\Sigma}_j \qquad \text{due to independence of } \mathbf{x}_j \text{ and } \mathbf{x}_k$$

## Exercise 7.2, page 25

Since y is a linear combination of Gaussian variables, it is itself Gaussian and hence only the mean and covariance need be computed. These are  $\mathbb{E}[y] = A\mu + b$ , and

$$cov(\mathbf{y}) = \mathbb{E}\left[ (\mathbf{y} - \mathbb{E}[\mathbf{y}]) (\mathbf{y} - \mathbb{E}[\mathbf{y}])^T \right]$$
$$= \mathbb{E}\left[ (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{A}\boldsymbol{\mu} - \mathbf{b}) (\mathbf{A}\mathbf{x} + \mathbf{b} - \mathbf{A}\boldsymbol{\mu} - \mathbf{b})^T \right]$$
$$= \mathbb{E}\left[ \mathbf{A} (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{A}^T \right]$$
$$= \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$$

Thus,  $\mathbf{y} \sim N(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$ 

### Exercise 7.3, page 26

(a) A and  $A^T$  have the same diagonal, so they have the same trace.

(b)  $\operatorname{tr}[\mathbf{AB}] = \operatorname{tr}\left[(\mathbf{AB})^T\right] = \operatorname{tr}\left[\mathbf{B}^T\mathbf{A}^T\right] = \operatorname{tr}[\mathbf{A}^T\mathbf{B}^T]$ . The first equality follows from part (a), the others from properties of transpose and trace discussed in the chapter.

(c) 
$$\operatorname{tr}[\mathbf{A}] = \operatorname{tr}[\mathbf{S}\mathbf{\Lambda}\mathbf{S}^{-1}] = \operatorname{tr}[\mathbf{\Lambda}\mathbf{S}^{-1}\mathbf{S}] = \operatorname{tr}[\mathbf{\Lambda}\mathbf{I}] = \operatorname{tr}[\mathbf{\Lambda}].$$

(d)

$$tr[\mathbf{A}^{k}] = tr[(\mathbf{S}\mathbf{\Lambda}\mathbf{S}^{-1})^{k}]$$
definition  
= tr[(\boldsymbol{S}\mathbf{\Lambda}\mathbf{S}^{-1}) (\boldsymbol{S}\mathbf{\Lambda}\mathbf{S}^{-1}) ... (\boldsymbol{S}\mathbf{\Lambda}\mathbf{S}^{-1})] expand the power  
= tr[\boldsymbol{S}\mathbf{\Lambda}^{k}\mathbf{S}^{-1}] the intermediate terms  $\mathbf{S}^{-1}\mathbf{S} = \mathbf{I}$  cancel out  
= tr[\boldsymbol{\Lambda}^{k}]

Alternatively,  $\mathbf{A}^k \mathbf{x} = \lambda^k \mathbf{x}$  implies that  $\lambda^k$  is an eigenvalue of  $\mathbf{A}^k$ , in which case the result follows from part c.

(e)

$$\operatorname{tr}[\mathbf{A}^{-1}] = \operatorname{tr}\left[\left(\mathbf{S}\mathbf{\Lambda}\mathbf{S}^{-1}\right)^{-1}\right] = \operatorname{tr}\left[\mathbf{S}^{-1}\mathbf{\Lambda}^{-1}\mathbf{S}\right] = \operatorname{tr}\left[\mathbf{\Lambda}^{-1}\left(\mathbf{S}\mathbf{S}^{-1}\right)\right] = \operatorname{tr}\left[\mathbf{\Lambda}^{-1}\right] = \sum_{m=1}^{M} \lambda_m^{-1}.$$

The last equality follows from the fact that the inverse of a *diagonal* matrix is diagonal with diagonal elements equal to reciprocals of the original matrix.

(f) tr[AB] = 
$$\sum_{i} \sum_{j} \mathbf{A}_{ij} \mathbf{B}_{ji}$$
 for any two matrices A and B. Substituting  $\mathbf{B} = \mathbf{A}^{T}$  gives tr[AA<sup>T</sup>] =  $\sum_{i} \sum_{j} \mathbf{A}_{ij} \mathbf{A}_{ij} = \sum_{i} \sum_{j} \mathbf{A}_{ij}^{2}$ 

(g)  $|\mathbf{A}| = |\mathbf{S}\mathbf{A}\mathbf{S}^{-1}| = |\mathbf{S}||\mathbf{A}||\mathbf{S}^{-1}| = |\mathbf{S}||\mathbf{S}^{-1}||\mathbf{A}| = |\mathbf{S}\mathbf{S}^{-1}||\mathbf{A}| = |\mathbf{I}||\mathbf{A}| = \lambda_1\lambda_2\dots\lambda_P$ . We use the fact that the determinant is a scalar so the order does not matter, and the determinant of a diagonal matrix equals the product of diagonals.

(h) By part (c), the trace equals the sum of eigenvalues values. Since the eigenvalues of a symmetric, positive definite matrix are positive, it follows the trace is positive. Similarly, by part (f), the determinant equals the product of eigenvalues. Since the eigenvalues are positive, the product is positive.

(i) A positive semi-definite matrix has only non-negative eigenvalues. Since the trace of a matrix equals the sum of eigenvalues, the sum of eigenvalues of a positive semi-definite matrix must also be non-negative, and hence the trace is non-negative. Since the determinant of a matrix equals the product of eigenvalues, the product of eigenvalues of a positive semi-definite matrix must be non-negative, and hence the determinant is non-negative.

### Exercise 7.4, page 26

Consider the eigenvalue problem

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}.$$

Multiply both sides of this equation by A, and invoke the fact that A is idempotent:

$$\mathbf{A}^2 \mathbf{x} = \lambda \mathbf{A} \mathbf{x} \quad \Rightarrow \quad \mathbf{A} \mathbf{x} = \lambda^2 \mathbf{x}.$$

Eliminating Ax from the above equations gives

$$\lambda \mathbf{x} = \lambda^2 \mathbf{x}.$$

Since the eigenvector is non-zero (i.e.,  $\mathbf{x} \neq \mathbf{0}$ ), this equation can be satisfied in general only if  $\lambda = \lambda^2$ , which has the two solutions  $\lambda = 0$  and  $\lambda = 1$ .

The trace of any matrix equals the sum of eigenvalues. Since the eigenvalues of A are either 0 or 1, the sum of eigenvalues equals the number that are equal to one.

## Exercise 7.5, page 26

(a)  $\boldsymbol{\xi}^T \boldsymbol{\xi} = (\mathbf{U}\mathbf{x})^T (\mathbf{U}\mathbf{x}) = \mathbf{x}^T \mathbf{U}^T \mathbf{U}\mathbf{x} = \mathbf{x}^T \mathbf{x}$  since U is orthogonal. (b)  $\boldsymbol{\xi}^T \boldsymbol{\eta} = (\mathbf{U}\mathbf{x})^T (\mathbf{U}\mathbf{y}) = \mathbf{x}^T \mathbf{U}^T \mathbf{U}\mathbf{y} = \mathbf{x}^T \mathbf{y}$ . Thus, we have the identities  $\boldsymbol{\xi}^T \boldsymbol{\xi} = \mathbf{x}^T \mathbf{x}$ ,  $\boldsymbol{\eta}^T \boldsymbol{\eta} = \mathbf{y}^T \mathbf{y}$ , and  $\boldsymbol{\xi}^T \boldsymbol{\eta} = \mathbf{x}^T \mathbf{y}$ . Since the angle depends only on these three quantities, it is preserved by the transformation.

(c) Since U is orthogonal,  $UU^T = I$ . Taking the determinant of both sides gives  $|I| = 1 = |UU^T| = |U||U^T| = |U|^2$ . Thus,  $|U| = \pm 1$ .

(d) 
$$|\mathbf{U}\mathbf{A}\mathbf{U}^{T}| = |\mathbf{U}||\mathbf{A}^{T}||\mathbf{U}^{T}| = |\mathbf{U}||\mathbf{U}^{T}||\mathbf{A}| = |\mathbf{U}\mathbf{U}^{T}||\mathbf{A}| = 1|\mathbf{A}|.$$

(e) 
$$\operatorname{tr}[\mathbf{U}\mathbf{A}\mathbf{U}^T] = \operatorname{tr}[\mathbf{U}\mathbf{U}^T\mathbf{A}] = \operatorname{tr}[\mathbf{A}] = \operatorname{tr}[\mathbf{A}].$$

(f) Show that  $(\mathbf{U}_1\mathbf{U}_2)(\mathbf{U}_1\mathbf{U}_2)^T = \mathbf{I}$  and  $(\mathbf{U}_1\mathbf{U}_2)^T(\mathbf{U}_1\mathbf{U}_2) = \mathbf{I}$  (need to show both).

### Exercise 7.6, page 26

By definition,  $\Sigma_Y = \mathbb{E}[(\mathbf{y} - \boldsymbol{\mu}_Y) (\mathbf{y} - \boldsymbol{\mu}_Y)^T] = \mathbb{E}[\mathbf{y}\mathbf{y}^T] - \boldsymbol{\mu}_Y \boldsymbol{\mu}_Y^T$ . Therefore  $\mathbb{E}[\mathbf{y}^T \mathbf{A}\mathbf{y}] = \mathbb{E}[\sum_i \sum_j y_i y_j A_{ij}] = \sum_i \sum_j \mathbb{E}[y_i y_j] A_{ij} = \sum_i \sum_j (\Sigma_{ij} + \mu_i \mu_j) A_{ij}$  $= \operatorname{tr}[\mathbf{A}\Sigma_Y] + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}.$ 

### Exercise 7.7, page 27

Take the transpose of the hint:

$$(\mathbf{A}^{-1})^T \mathbf{A}^T = \mathbf{I}$$
 transpose of the hint (22.39)  
 $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$  multiply both sides by the inverse of transpose. (22.40)

Exercise 7.8, page 27

(a)

$$\|\mathbf{X}\|_{F}^{2} = \operatorname{tr}\left[\mathbf{X}\mathbf{X}^{T}\right] = \sum_{i} \sum_{j} \left(\mathbf{X}\right)_{ij} \left(\mathbf{X}^{T}\right)_{ji} = \sum_{i} \sum_{j} \left(\mathbf{X}\right)_{ij} \left(\mathbf{X}\right)_{ij} = \sum_{i} \sum_{j} \left(\mathbf{X}\right)_{ij}^{2}.$$

Taking the square root of both sides gives (7.6).

(b)

$$\begin{aligned} \|\mathbf{U}\mathbf{A}\mathbf{V}\|_{F}^{2} &= \operatorname{tr}[(\mathbf{U}\mathbf{A}\mathbf{V})(\mathbf{U}\mathbf{A}\mathbf{V})^{T}] = \operatorname{tr}[\mathbf{U}\mathbf{A}\mathbf{V}\mathbf{V}^{T}\mathbf{A}^{T}\mathbf{U}^{T}] \\ &= \operatorname{tr}[\mathbf{A}\mathbf{A}^{T}\mathbf{U}^{T}\mathbf{U}] = \operatorname{tr}[\mathbf{A}\mathbf{A}^{T}] = \|\mathbf{A}\|_{F}^{2}. \end{aligned}$$

# Exercise 7.9, page 27

Consider the eigenvalue equation

$$\mathbf{A}\mathbf{x} = \lambda \mathbf{x}.$$

Multiplying this equation by *b* gives

 $b\mathbf{A}\mathbf{x} = b\lambda\mathbf{x}.$ 

Of course,

 $a\mathbf{Ix} = a\mathbf{x}.$ 

Adding these two equations together gives:

$$a\mathbf{x} + b\mathbf{A}\mathbf{x} = a\mathbf{x} + b\lambda\mathbf{x}$$
  
 $(a\mathbf{I} + \mathbf{A})\mathbf{x} = (a + b\lambda)\mathbf{x}$ 

The last equation shows that the eigenvalue of  $a\mathbf{I} + \mathbf{A}$  is  $a + b\lambda$ .

## Exercise 7.10, page 27

Positive definite means  $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$  for all  $\mathbf{x}$ . Choose

$$\mathbf{x} = \begin{pmatrix} 1\\0\\\vdots\\0 \end{pmatrix}.$$
 (22.41)

For this choice,  $\mathbf{x}^T \mathbf{A} \mathbf{x} = A_{11}$ . Hence, the first diagonal element must be positive. Next, let x be a vector of all zeros except for a one in the second row. Then in this case  $\mathbf{x}^T \mathbf{A} \mathbf{x} =$  $A_{22}$ , which must be positive. Repeating this procedure for each row of x reveals that each and every diagonal element must be positive.

### Exercise 7.11, page 28

(a) 
$$\boldsymbol{\Sigma}^{-1} = (\mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T)^{-1} = (\mathbf{U}^T)^{-1}\boldsymbol{\Lambda}^{-1}\mathbf{U}^{-1} = \mathbf{U}\boldsymbol{\Lambda}^{-1}\mathbf{U}^T$$
, since  $\mathbf{U}^T = \mathbf{U}^{-1}$ 

(b)

$$Q = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{U} \boldsymbol{\Lambda}^{-1} \mathbf{U}^T (\mathbf{x} - \boldsymbol{\mu})$$
$$= \left( (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{U} \boldsymbol{\Lambda}^{-1/2} \right) \left( \boldsymbol{\Lambda}^{-1/2} \mathbf{U}^T (\mathbf{x} - \boldsymbol{\mu}) \right) = \boldsymbol{\xi}^T \boldsymbol{\xi}$$

where  $\boldsymbol{\xi} = \mathbf{A} (\mathbf{x} - \boldsymbol{\mu})$  and  $\mathbf{A} = \boldsymbol{\Lambda}^{-1/2} \mathbf{U}^T$ . Since the covariance matrix is positive definite, all diagonal elements of  $\boldsymbol{\Lambda}$  are positive and hence have an inverse and a positive square root.

(c)  $\boldsymbol{\xi}$  is a linear combination of constants and Gaussians, so it is itself Gaussian. The mean is

$$\mathbb{E}[\boldsymbol{\xi}] = \mathbb{E}[\mathbf{A}(\mathbf{x} - \boldsymbol{\mu})] = \mathbf{A}(\mathbb{E}[\mathbf{x}] - \boldsymbol{\mu}) = 0$$

and the covariance is

$$cov[\boldsymbol{\xi}] = \mathbb{E}[(\boldsymbol{\xi}) (\boldsymbol{\xi})^{T}]$$
  
=  $\mathbb{E}[\mathbf{A} (\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{T} \mathbf{A}^{T}]$   
=  $\mathbf{A}\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}) (\mathbf{x} - \boldsymbol{\mu})^{T}]\mathbf{A}^{T}$   
=  $\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^{T} = \boldsymbol{\Lambda}^{-1}\mathbf{U}^{T} (\mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^{T}) \mathbf{U}\boldsymbol{\Lambda}^{-1}$   
=  $\mathbf{I}$ 

Thus,  $\boldsymbol{\xi} \sim N(0, \mathbf{I})$ .

(d) The covariance matrix of  $\boldsymbol{\xi}$  is diagonal, so  $\operatorname{cov}[\boldsymbol{\xi}_i, \boldsymbol{\xi}_j] = 0$  if  $i \neq j$ . That is, each element of  $\boldsymbol{\xi}$  is uncorrelated with every other element of  $\boldsymbol{\xi}$ . Since  $\boldsymbol{\xi}$  also is Gaussian, uncorrelatedness implies independence.

(e) Since  $\mathbf{Q} = \xi_1^2 + \xi_2^2 + \dots \xi_P^2$  is a sum of squared independent Gaussians, each with zero mean and unit variance, it must have the chi-squared distribution with P degrees of freedom.

### Exercise 8.1, page 30

An R function that solves the least squares problem using the normal equations is the following:

```
regress.normal = function(y,x,include.intercept=TRUE,alpha=0.05) {
1
  *****
2
  ## DETERMINES THE LEAST SQUARES ESTIMATE OF B IN THE EQUATION Y = XB + E
3
  ## BASED ON THE NORMAL EQUATIONS
  ## INPUT:
5
  ##
       Y[NTOT]:
                  N-DIMENSIONAL VECTOR OF PREDICTANDS
6
  ##
      X[NTOT,MTOT]: N X M DIMENSIONAL MATRIX OF PREDICTORS
7
  ##
      INCLUDE.INTERCEPT: INCLUDE THE INTERCEPT? (DEFAULT=TRUE)
8
  ## OUTPUT:
  ## BHAT: M-DIMENSIONAL VECTOR OF ESTIMATES OF B
10
11 ## R2: R-SQUARED
12 ## SSE: SUM SQUARE ERROR OF THE LEAST SQUARES PREDICTION
13 ## DOF: DEGREES OF FREEDOM OF THE SSE.
14 ## RES.SE: STANDARD ERROR OF THE RESIDUALS
  *****
15
16
17 ntot = length(y)
  if ( length(x) %% ntot != 0) stop('x not dimensioned correctly')
18
  mtot = length(x)/ntot
19
  if ( ntot <= mtot ) stop('regression problem is not over-determined')
20
21
22 ### STRIP MISSING DATA
23 \dim(x) = c(ntot, mtot)
24 is.missing = is.na(y) | is.na(rowSums(x))
25 x.good = x[!is.missing,]
26 y.good
           = y[!is.missing ]
27 nsamp
           = sum(!is.missing)
28
if (include.intercept) x.good = cbind(rep(1,nsamp),x.good)
_{30} mtot = dim(x.good)[2]
31
32 Xt.X
           = t(x.good) %*% x.good
           = t(x.good) %*% y.good
33 xty
  xtx.inv = chol2inv(chol(xtx))
34
           = as.numeric(xtx.inv %*% xty)
35 bhat
36
          = sum((y.good - x.good %*% bhat)^2)
37 SSE
38 r2
          = 1 - sse/sum((y.good-mean(y.good))^2)
39 dof
          = nsamp - mtot
40 res.se = sqrt(sse/dof)
41
42 list(bhat=bhat,r2=r2,sse=sse,dof=dof,res.se=res.se)
43
  }
```

### Exercise 8.2, page 31

Using this 'data,' we obtain least squares estimates of the model (8.1) as follows:
```
_1 > set.seed(1)
2 > ntot = 20
3 > y
          = rnorm(ntot); pred1 = rnorm(ntot); pred2 = rnorm(ntot)
        = cbind(pred1,pred2)
4 > X
  > xy.normal = regress.normal(y,x)
5
  > print(xy.normal)
6
  $bhat
7
              [,1]
8
9 [1,] 0.1493806
10 [2,] -0.1516176
11 [3,] 0.2893589
12
13 $r2
14 [1] 0.1078292
15
  $sse
16
  [1] 14.13789
17
18
19
  $dof
20 [1] 17
21
22 $res.se
  [1] 0.9119432
23
```

# Exercise 8.3, page 32

The result of fitting a trend line to CO<sub>2</sub> concentration is:

```
> ntot = length(co2)
1
_{2} > year = iyst + (1:ntot - 0.5)/12
  > year.shift = year - 1960
3
4 > year.sqr = year.shift^2
s > co2.trend.normal = regress.normal(co2,year.shift)
6 > print(co2.trend.normal)
7
  $bhat
8
  [1] 308.698631 1.565872
10 $r2
11 [1] 0.9792493
12
13
  $sse
  [1] 10031.19
14
15
16 $dof
17 [1] 689
18
19 $res.se
20
  [1] 3.815633
```

Thus, the growth rate is 1.6 ppm/year.

Using the lm function gives:

```
> co2.year.lm = lm(co2~year.shift,na.action=na.omit)
1
2
  > print(summary(co2.year.lm))
  Call:
4
  lm(formula = co2 ~ year.shift, na.action = na.omit)
5
7
  Residuals:
     Min
                1Q Median
                                3Q
                                       Max
8
  -7.4921 -2.6494 -0.3407 2.3416 11.1495
10
  Coefficients:
11
               Estimate Std. Error t value Pr(>|t|)
12
  (Intercept) 3.087e+02 2.917e-01 1058.4
                                              <2e-16 ***
13
  year.shift 1.566e+00 8.684e-03
                                     180.3
                                              <2e-16 ***
14
15
16
  Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1
                                                      1
17
  Residual standard error: 3.816 on 689 degrees of freedom
18
   (5 observations deleted due to missingness)
19
20 Multiple R-squared: 0.9792, Adjusted R-squared: 0.9792
21 F-statistic: 3.251e+04 on 1 and 689 DF, p-value: < 2.2e-16
```

The coefficients,  $R^2$ , and dof are consistent.

### Exercise 8.4, page 32

The residuals are:

```
> options (width=70)
1
 > x = cbind(rep(1, ntot), year.shift)
2
  > y.pred = x %*% co2.trend.normal$bhat
3
  > resid = co2 - y.pred
4
  > print(as.numeric(resid[1:50]))
5
        7.666124 8.075635 8.555146 9.864656 10.744167 10.173678
    [1]
                                     3.891720 4.931231 5.990742
                 6.232699 4.352210
7
    [7]
        8.633188
                  7.239763 7.949274
   [13]
        6.600252
                                     8.758784 9.728295 8.787806
8
  [19]
        7.457316
                 5.546827 3.426338
                                     3.875848 4.465359
                                                         5.244870
9
        6.044380 6.533891 7.523402 8.342912
                                               8.592423
                                                         8.001933
10
  [25]
        6.901444 4.590955 3.320465 2.349976 3.489487
11 [31]
                                                         4.358997
        5.278508 5.488019 6.137529 7.537040 8.266551 7.356061
12 [37]
        5.495572 3.395083 1.704593 1.354104 2.353615 3.413125
13
  [43]
14
  [49]
        4.542636
                        NA
```

A plot of the residuals is shown in fig. 22.8.



Figure 22.8 Residuals of  $CO_2$  data after linear trend has been removed.

# Exercise 8.5, page 33

The annual cycle can be estimated as follows:

```
> nharm = 2
1
   >
2
   > t = seq(year.shift)/12
3
   > x = NULL
4
   > for ( n in 1:nharm) x = cbind(x, cos(2*pi*t*n), sin(2*pi*t*n))
5
   > colnames(x) = c(paste(rep(c('cos', 'sin'), nharm), rep(1:nharm, each=2), sep=""))
6
   > co2.cycle.normal = regress.normal(co2,x)
7
   > print(co2.cycle.normal)
8
   $bhat
9
   [1] 354.3263100 -1.4843479
                                    2.1128289
                                                  0.7733454 -0.3005125
10
11
12
   $r2
   [1] 0.005265222
13
14
15
   $sse
   [1] 480868.2
16
17
   $dof
18
   [1] 686
19
20
   $res.se
21
   [1] 26.47591
22
```

Exercise 8.6, page 33

The residuals after the annual cycle has been removed can be computed as

x = cbind(rep(1,ntot),x)
y.pred = x %\*% co2.cycle.normal\$bhat
resid = co2 - y.pred

A plot of the residuals and corresponding raw data is produced by

```
resid
          = co2 - y.pred + co2.cycle.normal$bhat[1]
1
2
   fout = 'co2.cycle.residuals.eps'
3
   if ( lplotfile) postscript(fout, horizontal=FALSE,
4
      onefile=FALSE, height=6, width=8, pointsize=12)
5
   par(mfcol=c(1,1), mar=c(5,5,3,1))
6
   par(cex.lab=1.2,cex.axis=1.2,cex.main=1.2)
7
   yrange = range(co2,resid,na.rm=TRUE)
8
   plot(year.say, resid, type='l', xlab='year',
9
      ylab='Parts Per Million',ylim=yrange,col='red',lwd=2)
10
  par(new=TRUE)
11
12 plot(year.say,co2,type='1',xlab='',ylab='',ylim=yrange,axes=FALSE)
  if ( lplotfile) dev.off()
13
```

The resulting plot is shown in fig. 22.9.



Figure 22.9 CO<sub>2</sub> (black curve) and CO<sub>2</sub> after annual cycle has been removed (red).

### Exercise 9.1, page 36

A code for generating a plot is the following. The resulting figure is shown in fig. 22.10.

```
plot(x.year,y,type='b',pch=19,lwd=2,xlab='year',ylab='sea level (mm)')
```

```
2 ftitle = paste(station,'Sea Level During', season)
```

```
3 title(main=ftitle,line=0.5)
```

```
abline(lm(y~x.year),lty='dashed',lwd=3,col='grey50')
```



Figure 22.10 Sea level at Guam during JFM.

The figure suggests a slight increase in sea level over the past few decades.

# Exercise 9.2, page 36

The result of fitting JFM Guam sea level to year is the following:

```
> print(regress.normal(y,x.year))
1
   $bhat
2
   [1] -560.046422
                         3.805933
3
4
   $r2
5
   [1] 0.135494
   $sse
8
   [1] 302447.4
9
10
   $dof
11
   [1] 32
12
13
   $res.se
14
   [1] 97.21873
15
```

# Exercise 9.3, page 36

An R function that performs regression and computes the t-value, p-value, and confidence interval is the following:

regress.normal = function(y,x,include.intercept=TRUE,alpha=0.05) { \*\*\*\*\* 2 ## DETERMINES THE LEAST SQUARES ESTIMATE OF B IN THE EQUATION Y = XB + E 3 ## BASED ON THE NORMAL EQUATIONS 4 ## INPUT: 5 ## Y[NTOT]: N-DIMENSIONAL VECTOR OF PREDICTANDS 6 ## X[NTOT,MTOT]: N X M DIMENSIONAL MATRIX OF PREDICTORS 7 ## INCLUDE.INTERCEPT: INCLUDE THE INTERCEPT? (DEFAULT=TRUE) 8 ## OUTPUT: 0 ## BHAT: M-DIMENSIONAL VECTOR OF ESTIMATES OF B 10 11 ## R2: R-SQUARED 12 ## SSE: SUM SQUARE ERROR OF THE LEAST SQUARES PREDICTION 13 ## DOF: DEGREES OF FREEDOM OF THE SSE. ## RES.SE: STANDARD ERROR OF THE RESIDUALS 14 ## TVAL[MTOT]: T-STATISTIC FOR EACH REGRESSION COEFFICIENT 15 ## PVAL[MTOT]: P-VALUE FOR EACH REGRESSION COEFFICIENT 16 ## BHAT.CI[MTOT,2]: CONFIDENCE INTERVAL FOR EACH REGRESSION COEF. 17 \*\*\*\*\* 18 19 ntot = length(y)20 if (length(x) % ntot != 0) stop('x not dimensioned correctly') 22 mtot = length(x)/ntot 23 if (ntot <= mtot ) stop('regression problem is not over-determined') 24 25 ### STRIP MISSING DATA  $26 \dim(x) = c(ntot, mtot)$ 27 is.missing = is.na(y) | is.na(rowSums(x)) 28 x.good = x[!is.missing,] y.good = y[!is.missing ] 29 = sum(!is.missing) 30 nsamp 31 if (include.intercept) x.good = cbind(rep(1,nsamp),x.good) 32 33 mtot = dim(x.good)[2]34 35 xtx = t(x.good) %\*% x.good 36 xty = t(x.good) %\*% y.good 37 xtx.inv = chol2inv(chol(xtx)) 38 bhat = as.numeric(xtx.inv %\*% xty) 39 =  $sum((y.good - x.good %*% bhat)^2)$ 40 SSE  $= 1 - \text{sse/sum}((y.\text{good}-\text{mean}(y.\text{good}))^2)$ 41 r2 42 dof = nsamp - mtot 43 res.se = sqrt(sse/dof) 44 ### INFERENCE 45 st.err = sqrt(sse/dof\*diag(xtx.inv)) 46 = bhat/st.err 47 tval = 2\*pt(abs(tval),dof,lower.tail=FALSE) 48 pval 49 t.crit = qt(alpha/2,dof,lower.tail=FALSE) 50 bhat.ci = cbind(bhat - t.crit\*st.err, bhat + t.crit\*st.err) 51 s2 list(bhat=bhat,r2=r2,sse=sse,dof=dof,res.se=res.se,tval=tval,pval=pval,bhat.ci=bhat.ci) 53 }

### Exercise 9.4, page 36

The output is identical to earlier problems, except for the following added items

```
> print(regress.normal(y,x.year))
1
2
  . . .
  $tval
3
  [1] -0.1656398 2.2394999
4
5
  $pval
6
  [1] 0.86948250 0.03220094
7
  $bhat.ci
9
                        [,2]
10
                 [,1]
  [1,] -7447.1429574 6327.050113
11
  [2,] 0.3442531
                       7.267614
12
```

The trend is significant at the 5% level, but just barely. The results agree with lm:

```
> trend.lm = lm(y~x.year)
1
2
  > trend.sum = summary(trend.lm)
  > print(trend.sum)
3
4
5 Call:
6 lm(formula = y ~ x.year)
7
8 Residuals:
     Min 1Q Median 3Q Max
9
  -204.45 -100.43 24.59 78.49 145.99
10
11
12 Coefficients:
             Estimate Std. Error t value Pr(>|t|)
13
  (Intercept) -560.046 3381.110 -0.166 0.8695
14
                                 2.239 0.0322 *
15
  x.year
               3.806
                        1.699
16
  ___
  Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
17
18
19 Residual standard error: 97.22 on 32 degrees of freedom
20 Multiple R-squared: 0.1355, Adjusted R-squared: 0.1085
21 F-statistic: 5.015 on 1 and 32 DF, p-value: 0.0322
```

And with confint:

Exercise 9.5, page 37

The output is given above.

### Exercise 9.6, page 37

The result of fitting sea level to both year and NINO3.4 is the following:

```
> print(regress.normal(y,cbind(x.year,x.sst)))
1
  $bhat
2
  [1] -676.326254 3.865792 -77.019967
3
4
5
  $r2
  [1] 0.6148428
6
7
  $sse
8
  [1] 134747.3
9
10
  $dof
11
  [1] 31
12
13
14 $res.se
  [1] 65.9294
15
16
17
  $tval
  [1] -0.2949533 3.3541617 -6.2113731
18
19
  $pval
20
  [1] 7.699959e-01 2.112795e-03 6.733968e-07
21
22
  $bhat.ci
23
                [,1] [,2]
24
  [1,] -5352.919193 4000.266686
25
  [2,] 1.515179 6.216405
26
  [3,] -102.309586 -51.730348
27
```

Trend is significant: p-value is less than 5% and confidence interval does not include 0. This result agrees with lm:

```
v > trend.enso.lm = lm(y<sup>*</sup>x.year+x.sst)
2 > trend.enso.sum = summary(trend.enso.lm)
3 > print(trend.enso.sum)
4
5 Call:
6 lm(formula = y ~ x.year + x.sst)
7
8 Residuals:
    Min 1Q Median
                        3Q Max
9
10 -99.02 -49.59 9.31 28.70 201.20
11
12 Coefficients:
13
    Estimate Std. Error t value Pr(>|t|)
14 (Intercept) -676.326 2292.994 -0.295 0.77000
                        1.153 3.354 0.00211 **
15 x.year 3.866
16 x.sst
              -77.020
                         12.400 -6.211 6.73e-07 ***
  ____
17
18 Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
19
20 Residual standard error: 65.93 on 31 degrees of freedom
21 Multiple R-squared: 0.6148, Adjusted R-squared: 0.59
22 F-statistic: 24.74 on 2 and 31 DF, p-value: 3.779e-07
```

# And with confint:

# Exercise 10.1, page 40

$$X_{n,j} = \begin{cases} 1 & \text{if } j = 1\\ \cos(2\pi n(j/2)/12) & \text{if j is even}\\ \sin(2\pi n(j-1)/2/12) & \text{if j is odd and } j > 1 \end{cases}$$

and

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ c_1 \\ s_1 \\ c_2 \\ s_2 \\ \vdots \\ c_H \\ s_H \end{pmatrix}$$
(22.42)

# Exercise 10.2, page 41

A complete code for all parts of the exercises is the following:

```
modsel.loo.simple = function(y,x,kmax=dim(x)[2]+1) {
2 ntot = length(y)
3 x.const = cbind(rep(1,ntot),x); # ADD COLUMN OF ONES IN PREDICTOR MATRIX
  if ( \dim(x)[1] = ntot ) stop("inconsistent dimension of y and x in modsel.loo")
  if (kmax > dim(x.const)[2]) stop("kmax exceeds 1 + number of predictors")
6 yvar
         = var(y)
7 ypss
           = sum((y-mean(y))^2)
8
9 cvmse = numeric(kmax); cvstd
                                    = numeric(kmax); se.sqr = numeric(kmax)
nse.sqr = numeric(kmax); ic.std = numeric(kmax)
n aic
          = numeric(kmax); aicc = numeric(kmax); bic = numeric(kmax)
12
13 for ( k in 1:kmax) {
14 x.pred = x.const[,1:k]
            = lm(y~x.pred-1)
15 xy.lm
             = fitted(xy.lm)
16 yhat
17 yresid = residuals(xy.lm)
  se.sqr[k] = sum(yresid<sup>2</sup>)/(ntot-k)
18
nse.sqr[k] = se.sqr[k]/yvar
20
21 ## BRUTE FORCE LOO
22 # yhat.loo = numeric(ntot)
23 # for ( n in 1:ntot) {
_{24} # y.pred = y[-n]
25 # x.pred = x.const[-n,1:k]
26 # yhat.loo[n] = sum(coefficients(lm(y.pred~x.pred-1)) * x.const[n,1:k])
27 # }
28 # cverr
               = (y-yhat.loo)^2
29
30 ## FAST LOO
\dim(x.pred) = c(ntot, k)
32 xtx.inv = chol2inv(chol(t(x.pred) %*% x.pred))
33 alpha
             = numeric(ntot)
34 for (n in 1:ntot) alpha[n] = as.numeric(t(x.pred[n,]) %*% xtx.inv %*% x.pred[n,])
_{35} cverr = (y - yhat)^2/(1-alpha)^2
36
37 cvmse[k] = mean(cverr)
38 cvstd[k] = sd(cverr)/sqrt(ntot)
39
40 nparm
            = k + 1
41 aic[k]
             = ntot *log(2*pi) + ntot*log(mean(yresid^2)) + ntot + 2*nparm
42 aicc[k] = ntot *log(2*pi) + ntot*log(mean(yresid^2)) + ntot +
       2*nparm*(nparm+1)/(ntot-nparm-1)
43
44 bic[k] = ntot *log(2*pi) + ntot*log(mean(yresid^2)) + ntot + nparm*log(ntot)
   ic.std[k] = ntot * sqrt(psigamma(ntot/2, deriv=1))
45
46
   }
47
48 cvmse = cvmse/yvar
49 cvstd = cvstd/yvar
50
s1 list(cvmse=cvmse,cvstd=cvstd,k.min=k.min,k.sel=k.sel,se.sqr=se.sqr,nse.sqr=nse.sqr,
   npcs.min=npcs.min,npcs.sel=npcs.sel,aic=aic,aicc=aicc,bic=bic,ic.std=ic.std) }
52
```

A plot of the normalized unbiased error variance as a function of the number of predictors is shown in fig. 22.11. "0" means only the intercept term is included, and in this case the normalized error variance is 1 because the model error variance equals the total variability. The minimum error variance occurs when 4 Fourier terms are included, which corresponds to 2 harmonics (since the sine and cosine count as one harmonic). The minimum normalized error variance is about 78%. This means the annual cycle explains about 22% of the variance of the NINO3.4 index, leaving 78% unexplained by the annual cycle. The values of the normalized unbiased error variance are

```
1 [1] "unbiased error variance"
```

```
2 [1] 1.0000000 0.9431272 0.8116455 0.7987252 0.7796379 0.7806310
```

```
3 [7] 0.7815282 0.7824943 0.7832677 0.7842617 0.7852547
```



Figure 22.11 Normalized unbiased error variance of annual cycle model for NINO3.4.

#### Exercise 10.3, page 41

The normalized cross validated mean square error of annual cycle model for NINO3.4 is shown in fig. 22.12. The numerical values are

```
      1
      [1] "unbiased error variance"

      2
      [1] 1.0012658 0.9454002 0.8147378 0.8027723 0.7846022 0.7865786

      3
      [7] 0.7885135 0.7905072 0.7922921 0.7943716 0.7963414
```

Exercise 10.4, page 41



# Normalized Cross Validated Mean Square Error of Annual Cycle Model

Figure 22.12 Normalized cross validated mean square error of annual cycle model for NINO3.4. Error bars show the standard error.

### Exercise 10.5, page 42

See figure 22.12 for the error bars. The "one-standard-deviation rule" would select the model with 1 Fourier harmonic (but both sine and cosine).

### Exercise 10.6, page 42

The numerical values of the information criteria are

[1] "AIC" [1] 2113.715 2068.396 1950.634 1938.937 1920.799 1922.799 1924.699 2 [8] 1926.667 1928.438 1930.429 1932.416 3 [1] "AICC" 4 [1] 2109.730 2062.427 1942.685 1929.013 1908.906 1908.942 1908.883 5 [8] 1908.897 1908.720 1908.767 1908.817 6 [1] "BIC" 7 [1] 2123.061 2082.416 1969.328 1962.304 1948.839 1955.512 1962.086 [8] 1968.727 1975.171 1981.835 1988.496

### Exercise 10.7, page 42

See figure. 22.13



Figure 22.13 Information criteria for various annual cycle models for NINO3.4.

### Exercise 10.8, page 42

The NINO3.4 anomaly can be computed and plotted as follows

```
y = nino34.ts
   nh.sel = 4
   x = x.pred[,1:nh.sel]
   xy.lm = lm(y^x)
   fout = 'nino34.anomaly.pdf'
6
   if (lplotfile) pdf(fout,width=8,height=5)
7
   par(mfcol=c(1,1),mar=c(5,5,5,1))
8
   plot(nino34.time, residuals(xy.lm), type='l', col='blue', xlab='year', ylab='NINO3.4 Anomaly')
9
   ftitle.top = paste('NINO3.4 Anomaly')
10
   ftitle.bot = paste(nh.sel/2,'Annual Harmonics Removed')
11
   title(main=ftitle.top,line=2.0)
12
   title(main=ftitle.bot,line=0.5)
13
   if (lplotfile) dev.off()
14
```

The resulting plot is shown in fig. 22.14

# Exercise 10.9, page 42

Substituting K = M + 1 in (10.5) gives

2

$$AIC_c = N \log\left(\frac{SSE}{N}\right) + \frac{2(M+1)N}{N-M-2}.$$
 (22.43)



Figure 22.14 NINO3.4 anomaly after 2 annual harmonics have been removed.

Taking the difference of the two expressions gives

$$\Delta AIC_c = \frac{2(M+1)N}{N-M-2} - \left(2(M+1) + \frac{2(M+1)(M+2)}{N-M-2}\right) = 0$$

This shows that the two expressions differ by a constant. However, shifting a function upward or downward by a constant does not change the location of the minimum, so both expressions would be minimized by the same model and hence give identical model selections.

# Exercise 10.10, page 43

Using the notation (??), it follows that

$$\mathbf{X}^{T}\mathbf{X} = \begin{pmatrix} \mathbf{x}_{1}^{T} & \mathbf{x}_{2}^{T} & \dots & \mathbf{x}_{N}^{T} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{1} \\ \mathbf{x}_{2} \\ \vdots \\ \mathbf{x}_{N} \end{pmatrix} = \sum_{n=1}^{N} \mathbf{x}_{n}^{T} \mathbf{x}_{n}, \quad (22.44)$$

and

$$\mathbf{X}^{T}\mathbf{y} = \begin{pmatrix} \mathbf{x}_{1}^{T} & \mathbf{x}_{2}^{T} & \dots & \mathbf{x}_{N}^{T} \end{pmatrix} \begin{pmatrix} y_{1} \\ y_{2} \\ \vdots \\ y_{N} \end{pmatrix} = \sum_{n=1}^{N} \mathbf{x}_{n}^{T} y_{n}.$$
(22.45)

Let  $\beta_k$  denote the least squares estimate of  $\beta$  based on all samples except the k'th. This estimate can be expressed as

$$\boldsymbol{\beta}_{k} = \left(\mathbf{X}^{T}\mathbf{X} - \mathbf{x}_{k}^{T}\mathbf{x}_{k}\right)^{-1} \left(\mathbf{X}^{T}\mathbf{y} - \mathbf{x}_{k}^{T}y_{k}\right).$$
(22.46)

To simplify this expression, we invoke the Sherman-Morrison-Woodbury formula:

$$(\mathbf{A} - \mathbf{BDC})^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B} \left(\mathbf{D}^{-1} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}\right)^{-1}\mathbf{C}\mathbf{A}^{-1},$$
 (22.47)

where it is assumed that the matrices are conformable and all matrix inverses are defined (i.e., A and D are not singular). This formula implies

$$\left(\mathbf{X}^{T}\mathbf{X} - \mathbf{x}_{k}^{T}\mathbf{x}_{k}\right)^{-1} = \left(\mathbf{X}^{T}\mathbf{X}\right)^{-1} + \left(\mathbf{X}^{T}\mathbf{X}\right)^{-1}\mathbf{x}_{k}^{T}\left(\mathbf{I} - \mathbf{x}_{k}\left(\mathbf{X}^{T}\mathbf{X}\right)^{-1}\mathbf{x}_{k}^{T}\right)^{-1}\mathbf{x}_{k}\left(\mathbf{X}^{T}\mathbf{X}\right)^{-1}$$
(22.48)

Substituting this expression in (22.46), and noting that  $\mathbf{I} - \mathbf{x}_k \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_k^T$  is a scalar, gives

$$\boldsymbol{\beta}_{k} = \left( \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} + \frac{\left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T} \mathbf{x}_{k} \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1}}{1 - \mathbf{x}_{k} \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T}} \right) \left( \mathbf{X}^{T} \mathbf{y} - \mathbf{x}_{k}^{T} y_{k} \right)$$
$$= \hat{\boldsymbol{\beta}} + \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T} \left( -\mathbf{y}_{k} + \frac{\mathbf{x}_{k} \hat{\boldsymbol{\beta}} - \mathbf{x}_{k} \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T} y_{k}}{1 - \mathbf{x}_{k} \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T}} \right)$$
$$= \hat{\boldsymbol{\beta}} + \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T} \left( \frac{\mathbf{x}_{k} \hat{\boldsymbol{\beta}} - y_{k}}{1 - \mathbf{x}_{k} \left( \mathbf{X}^{T} \mathbf{X} \right)^{-1} \mathbf{x}_{k}^{T}} \right).$$
(22.49)

The difference between the withheld sample  $y_k$  and its least squares prediction  $\mathbf{x}_k \boldsymbol{\beta}_k$  is therefore

$$y_{k} - \mathbf{x}_{k}\boldsymbol{\beta}_{k} = y_{k} - \mathbf{x}_{k}\hat{\boldsymbol{\beta}} - \frac{\mathbf{x}_{k} \left(\mathbf{X}^{T}\mathbf{X}\right)^{-1} \mathbf{x}_{k}^{T} \left(\mathbf{x}_{k}\hat{\boldsymbol{\beta}} - y_{k}\right)}{1 - \mathbf{x}_{k} \left(\mathbf{X}^{T}\mathbf{X}\right)^{-1} \mathbf{x}_{k}^{T}}$$
$$= y_{k} - \mathbf{x}_{k}\hat{\boldsymbol{\beta}} + \frac{\alpha_{k}}{1 - \alpha_{k}} \left(y_{k} - \mathbf{x}_{k}\hat{\boldsymbol{\beta}}\right)$$
$$= \frac{y_{k} - \mathbf{x}_{k}\hat{\boldsymbol{\beta}}}{1 - \alpha_{k}}, \qquad (22.50)$$

where we define the scalar

$$\alpha_k = \mathbf{x}_k \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{x}_k^T.$$
(22.51)

# Exercise 11.1, page 45

If the comparisons are independent, then the experimentwise error rate is related to the per-comparison significance level by

$$\alpha_{pc} = 1 - (1 - \alpha_E)^{1/M}.$$
(22.52)

An experimentwise error rate of  $\alpha_E = 5\%$  therefore corresponds to a per-comparison significance level of 0.1%:

<sup>1 &</sup>gt; nsamp = 30; nvarb = 50; alpha.e = 0.05; rho.obs = 0.6

 $_2$  > alpha.pc = 1 - (1-alpha.e)^(1/nvarb)

<sup>3 &</sup>gt; alpha.pc

<sup>4 [1] 0.00102534</sup> 

Under the null hypothesis of independent and normally distributed variables, the statistic

$$T = \frac{\hat{\rho}}{\sqrt{1 - \hat{\rho}^2}} \sqrt{N - 2},$$
 (22.53)

has a t distribution with N - 2 degrees of freedom. The question is ambiguous because either the maximum value or maximum *absolute* value are appropriate, so we compute critical t-values and critical correlations for both interpretations:

```
> t.1tail = qt(alpha.pc ,nsamp-2,lower.tail=FALSE)
> t.2tail = qt(alpha.pc/2,nsamp-2,lower.tail=FALSE)
> rho.1tail = t.1tail / sqrt(nsamp-2+t.1tail^2)
> rho.2tail = t.2tail / sqrt(nsamp-2+t.2tail^2)
> rho.1tail
[1] 0.5403947
> rho.2tail
[1] 0.5693206
```

Both critical correlations (i.e, 0.54, 0.57) are less than the observed correlation of 0.6, so we conclude that the predictor is significant at the 5% significance level, regardless of whether the forecaster is looking at the maximum, or maximum absolute, correlation.

Alternatively, we can compute the actual experimentwise error rate and compare with 5%. The observed correlation of 0.6 corresponds to a t-value of

```
1 > tval = rho.obs*sqrt(nsamp-2)/sqrt(1-rho.obs^2)
2 > tval
3 [1] 3.968627
```

The area to the right of this value is

```
1 > alpha.pc.1tail = pt(tval,nsamp-2,lower.tail=FALSE)
2 > alpha.pc.1tail
3 [1] 0.0002285276
```

Therefore, the experimentwise error rate is

```
1 > alpha.pc.2tail = 2*alpha.pc.1tail
2 > alpha.e.1tail = 1 - (1-alpha.pc.1tail)^nvarb
3 > alpha.e.2tail = 1 - (1-alpha.pc.2tail)^nvarb
4 > alpha.e.1tail
5 [1] 0.01136264
6 > alpha.e.2tail
7 [1] 0.02259872
```

Both of these error rates are less than 5%, so the predictor is significant at the 5% level even after taking screening into account, in agreement with the previous analysis.

If the variables are not independent, then we can use the Bonferroni bound  $\alpha_E \leq M \alpha_{pc}$ :

```
1 > alpha.e.bonferroni.ltail = nvarb * alpha.pc.ltail
2 > alpha.e.bonferroni.2tail = nvarb * alpha.pc.2tail
3 > alpha.e.bonferroni.ltail
4 [1] 0.01142638
```

```
s > alpha.e.bonferroni.2tail
```

```
5 [1] 0.02285276
```

Both values of  $\alpha_E$  are less than 5% and hence imply significance.

### Exercise 11.2, page 46

The significance test is invalid because the selection of the four forecasts is based on the same data as used to test significance. The experimentwise error rate cannot be estimated from  $\alpha_E = 1 - (1 - \alpha_{pc})^M$  because this equation pertains to selecting a *single maximum* and determining the probability of achieving *that value*. In contrast, we are choosing the best *four out of nine variables*, then we are applying some *other* calculation on those variables (i.e., computing their mean, and then computing the correlation using that mean).

To estimate the experimentwise error rate, we perform a numerical simulation. First, we generate random numbers to simulate 20 years of 9 forecasts, plus 20 years of observations. By construction, the random numbers are independent, so the true correlation between any pair of variables is zero. The correlation between each of the 9 "forecasts" and "observations" are computed and ranked, the four forecasts with largest correlation are averaged to construct a single forecast, and the correlation between this 4-member mean and "observations" is computed. This procedure is repeated 10,000 times to produce an empirical distribution for the correlation skill. An R code for this analysis is the following:

```
> set.seed(1)
  > ntot = 20
2
  > ntrials = 10000
              = 9
  > nmod
  > nsel
              = 4
  > obs
               = rnorm(ntot)
6
  > cor.skill = numeric(ntrials)
7
  > for ( nt in 1:ntrials) {
     frcs = array(rnorm(ntot*nmod),dim=c(ntot,nmod))
  +
     cor.xy = as.numeric(cor(frcs, obs))
  +
10
  + npic = order(cor.xy,decreasing=TRUE)
11
12
  + frcs.mean = rowMeans(frcs[,npic[1:nsel]])
  + cor.skill[nt] = cor(frcs.mean,obs)
13
  + }
14
  > cor.95 = quantile(cor.skill,probs=0.95)
15
  > cor.95
16
         95%
17
  0.6145697
18
  > pval.e = sum( cor.skill >= 0.55)/ntrials
19
  > pval.e
20
  [1] 0.1166
21
```

A histogram of the resulting correlations of the 4-member mean is shown in fig. 22.15. First, we see that the correlations are skewed above zero, even though none of the forecasts have skill. Second, the ninety-fifth percentile, indicated by the vertical dashed line, shows that the correlation would have to exceed 0.61 in order to reject the null hypothesis of zero correlation at the 5% level. Consistent with this, the experimentwise p-value for this procedure is 11.7% (i.e., greater than 5%). Therefore, the correlation skill of 0.55 is not significant when the selection process has been taken into account.

### Exercise 11.3, page 46

The significance test is not valid because it did not account for the fact that the investigator looked at 3 correlations and chose the largest. The probability of finding *at least one* out of three independent correlations to exceed the 5% level is

$$p_e = 1 - (1 - 0.05)^3 \approx 14\%.$$
 (22.54)

This calculation is not realistic because the correlations are not independent (e.g., the 3 correlations are computed using exactly the same observations). The Bonferroni bound, which can be applied even if the variables are dependent, implies  $\alpha_E \leq M \alpha_{pc}$ , or  $\alpha_E \leq 15\%$ . Both calculations suggest that the experimentwise error rate for this analysis could be as high as 14-15%, so there is little basis to assume the variable is significant.

To *ensure* an experimentwise error rate of 5%, the Bonferroni bound implies that the comparisonwise error rate should be

$$p_c = 0.05/3 \approx 1.7\%.$$
 (22.55)

The critical value of the correlation for a 1.7% significance level is



Figure 22.15 Histogram of correlation coefficients of a multi-model forecast derived from averaging the four most correlated "no-skill" forecasts out of nine no-skill forecasts. The vertical dashed line indicates the ninety-fifth percentile.

```
alpha = 0.05
   >
            = 50
   >
     nyrs
     nsel
            = 3
   >
   >
            = 1 - (1-alpha)^nsel
   >
     pe
              1 - (1-alpha)^(1/nsel)
   >
     рс
            =
6
   >
     рс
7
   [1] 0.01695243
8
            = qt(pc/2,nyrs-2,lower.tail=FALSE)
   > t.c
9
10
   > rho.c
              = t.c/sqrt(nyrs-2+t.c^2)
11
   > rho.c
   [1] 0.3362835
12
```

Therefore, the observed correlation of 0.3 is not significant, in an experimentwise sense, when screening is taken into account.

Also, one normally expects a predictor to become less useful as the lead time increases. However, in this example, the correlation is strong in October and weak in the months in closer proximity (November and December). Since this behavior contradicts expectations based on past experience and physical models, the investigator is obliged to explain how the value in October is a better predictor than the values in November or December.

# Exercise 12.1, page 49

$$\mathbf{Y}' = \mathbf{W}^{-1} \dot{\mathbf{U}} \dot{\mathbf{S}} \dot{\mathbf{V}}^T \tag{22.56}$$

$$= \left(\frac{1}{\sqrt{N-1}}\mathbf{W}^{-1}\dot{\mathbf{U}}\dot{\mathbf{S}}\right)\left(\sqrt{N-1}\dot{\mathbf{V}}^{T}\right)$$
(22.57)

$$= \mathbf{E}\mathbf{F}^T \tag{22.58}$$

where

$$\mathbf{E} = \frac{1}{\sqrt{N-1}} \mathbf{W}^{-1} \dot{\mathbf{U}} \dot{\mathbf{S}}$$
(22.59)  
$$\mathbf{F} = \sqrt{N-1} \dot{\mathbf{V}}$$
(22.60)

The PCs satisfy

$$\frac{1}{N-1}\mathbf{F}^T\mathbf{F} = \mathbf{I}.$$
(22.61)

# Exercise 12.2, page 50

A function that performs PCA is the following:

```
eof.latlon.simple <- function(lon, lat, data.array, neof=30) {</pre>
1
    # COMPUTE PRINCIPAL COMPONENTS OF A DATA ARRAY.
2
    # INPUT:
3
    # LON[NLON]: VECTOR SPECIFYING LONGITUDES
4
    # LAT[NLAT]: VECTOR SPECIFYING LATITUDES
5
    # DATA.ARRAY[NLON,NLAT,NTOT] OR DATA.ARRAY[NLON*NLAT,NTOT]:
6
    # THE DATA ARRAY IN [SPACE1, SPACE2, TIME] OR [SPACE, TIME] FORMAT
7
   # NEOF = NUMBER OF SPATIAL EOFS (WITH MASK) TO BE OUTPUTTED. DEFAULT = 30
8
0
    # OUTPUT LIST:
   # $EOF[NLON*NLAT, NEOF]: THE FIRST NEOF SCALED EOFS
10
    # $PC[NTOT, NEOF]: THE PCS, NORMALIZED TO UNIT VARIANCE
11
   # $FEXPVAR: FRACTION OF EXPLAINED VARIANCE FOR EACH EOF (ALL OF THEM).
12
13
    # $FEXPVAR.CI: CONFIDENCE INTERVALS FOR FRACTION OF EXPLAINED VARIANCE.
    # $EOFI[NLON*NLAT, NEOF]: PSEUDO INVERSE OF EOF (I.E., T(EOFI) %*% EOF = I)
14
15
    # $NEOF: MINIMUM OF (NEOF IN ARGUMENT LIST, RANK OF DATA.ARRAY)
16
17
    # DEFINE PARAMETERS
   nlon = length(lon)
nlat = length(lat)
18
19
    ntot = length(data.array)/(nlon*nlat)
20
21
    if ( length(data.array) %% nlon*nlat != 0 )
       stop("data.array dimension not integral multiple of nlon*nlat")
22
23
24
    # DEFINE WEIGHTING
25
    weight =rep(sqrt(cos(lat*pi/180)),each=nlon)
26
27
    # IDENTIFY MISSING OR 0-WEIGHTED DATA
   dim(data.array) = c(nlon*nlat,ntot)
28
    lbad = is.na(rowSums(data.array)) | weight == 0
29
30
31
    data.array = data.array - rowMeans(data.array) # REMOVE CLIMATOLOGY
    data.array = data.array * weight # NORMALIZE/WEIGHT THE VARIABLES
32
33
34
    # COMPRESS DATA BY ELIMINATING MISSING GRID POINTS
35
    data.array = data.array[!lbad,]
    ndef = sum(!lbad)
36
37
    # COMPUTE SVD
38
    mmin = min(ndef, ntot, neof)
39
40
   data.svd = svd(data.array,nu=mmin,nv=mmin)
41
   fexpvar = data.svd$d^2/sum(data.svd$d^2) # COMPUTE FEXPVAR
42
43
    # COMPUTE CONFIDENCE INTERVAL OF FEXPVAR
44
    fexpvar.ci = cbind(fexpvar * ( 1 + sqrt(2/ntot)) , fexpvar * (1 - sqrt(2/ntot)))
45
46
47
    # COMPUTE PCS
    pc = sqrt(ntot-1)*data.svd$v[,(1:mmin)]
48
49
   # FILL IN MISSING POINTS IN EOF
50
    eof = array(NA, dim=c(nlon*nlat, mmin))
51
    eofi = array(NA, dim=c(nlon*nlat, mmin))
52
    for ( n in 1:mmin ) eof [!lbad,n] = data.svd$u[,n]/weight[!lbad]*data.svd$d[n]/sqrt(ntot-1)
for ( n in 1:mmin ) eofi[!lbad,n] = data.svd$u[,n]*weight[!lbad]/data.svd$d[n]*sqrt(ntot-1)
53
54
55
56
    list(eof=eof,pc=pc,fexpvar=fexpvar,fexpvar.ci=fexpvar.ci,eofi=eofi,neof=mmin)
57
    }
```

# Exercise 12.3, page 50

The plots should be similar to those shown in the notes.

### Exercise 12.4, page 51

(a)  $\mathbf{b}^T \mathbf{b} = \mathbf{a}^T \mathbf{a} / \mathbf{a}^T \mathbf{a} = 1.$ 

(b)

$$\operatorname{var}[r] = \operatorname{var}\left[\frac{\mathbf{a}^T \mathbf{x}}{\sqrt{\mathbf{a}^T \mathbf{a}}}\right] = \left(\frac{1}{\sqrt{\mathbf{a}^T \mathbf{a}}}\right)^2 \operatorname{var}[\mathbf{a}^T \mathbf{x}] = \frac{\mathbf{a}^T \Sigma \mathbf{a}}{\mathbf{a}^T \mathbf{a}}.$$
 (22.62)

(c) Substitute  $\mathbf{a} = \mathbf{U}\tilde{\mathbf{a}}$  in the above expression:

$$\operatorname{var}[r] = \frac{\mathbf{a}^T \Sigma \mathbf{a}}{\mathbf{a}^T \mathbf{a}} = \frac{\mathbf{a}^T \mathbf{U} \Lambda \mathbf{U}^T \mathbf{a}}{\mathbf{a}^T \mathbf{a}} = \frac{\tilde{\mathbf{a}}^T \Lambda \tilde{\mathbf{a}}}{\tilde{\mathbf{a}}^T \mathbf{U}^T \mathbf{U} \tilde{\mathbf{a}}} = \frac{\tilde{\mathbf{a}}^T \Lambda \tilde{\mathbf{a}}}{\tilde{\mathbf{a}}^T \tilde{\mathbf{a}}}.$$
 (22.63)

 $\Lambda$  is diagonal, so variance of r can be written as

$$\operatorname{var}[r] = \frac{\tilde{a}_{1}^{2}\lambda_{1} + \dots + \tilde{a}_{M}^{2}\lambda_{M}}{\tilde{a}_{1}^{2} + \dots + \tilde{a}_{M}^{2}}$$
(22.64)

(d) Because  $\lambda_1 > \lambda_2 > \cdots > \lambda_M$ , we have the inequality

$$\frac{\tilde{a}_1^2\lambda_1 + \dots + \tilde{a}_M^2\lambda_M}{\tilde{a}_1^2 + \dots + \tilde{a}_M^2} < \frac{\tilde{a}_1^2\lambda_1 + \dots + \tilde{a}_M^2\lambda_1}{\tilde{a}_1^2 + \dots + \tilde{a}_M^2} = \lambda_1 \left(\frac{\tilde{a}_1^2 + \dots + \tilde{a}_M^2}{\tilde{a}_1^2 + \dots + \tilde{a}_M^2}\right) = \lambda_1.$$
(22.65)

This maximum is obtained when

$$\tilde{\mathbf{a}} = \begin{pmatrix} 1\\0\\\vdots\\0 \end{pmatrix}.$$
 (22.66)

This solution corresponds to the vector  $\mathbf{a}$  equal to the leading eigenvector of  $\boldsymbol{\Sigma}$ .

(e) Let the SVD be  $\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ . Then, the sample covariance matrix becomes

$$\frac{1}{N-1}\mathbf{X}^T\mathbf{X} = \frac{1}{N-1}\mathbf{V}\mathbf{S}^T\mathbf{S}\mathbf{V}^T.$$
(22.67)

This implies that the orthogonal eigenvector matrix of  $\Sigma$  is V, and the eigenvalues are related to the singular values as

$$\hat{\lambda}_i = \frac{1}{N-1} s_i^2.$$
(22.68)

(f) The leading principal component maximizes the variance of a standardized linear combination.

### Exercise 13.1, page 55

An R code for doing the Livezey-Chen test is the following:

```
livezey.chen = function(lon,lat,xdata,ydata,ntrials=1000,alpha=0.05) {
1
    ### PERFORMS THE LIVEZEY-CHEN FIELD SIGNIFICANCE TEST
2
    # INPUT:
3
         LON[NLON]: LONGITUDE OF THE FIELD DATA
4
    #
         LAT[NLAT]: LATITUDE OF THE FIELD DATA
5
    #
        XDATA[NLON,NLAT,NTOT]: DATA ARRAY FOR THE FIELD YDATA[NTOT]: REFERENCE TIME SERIES
6
    #
7
    #
        NTRIALS: NUMBER OF MONTE CARLO TRIALS (DEFAULT = 10000)
8
    #
        ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
0
    # OUTPUT:
10
        HISTOGRAM OF AREAS WITH SIGNIFICANT CORRELATIONS
11
    #
    #
        LIST:
12
           $CRIT.NULL: SIGNIFICANCE THRESHOLD FROM LIVEZEY-CHEN TEST
13
    #
14
    #
           $SIG.AREA.OBS: PERCENTAGE AREA WITH SIGNIFICANT CORRELATIONS
15
16
   nlon
                  = length(lon)
17
    nlat
                  = length(lat)
18
    ntot
                  = length(ydata)
    dim(xdata) = c(nlon*nlat, ntot)
19
20
                = rep(cos(lat*pi/180),each=nlon)
= !is.na(xdata[,1])
21
    area.frac
22
    lgood
23
    area.frac[!lgood] = NA
24
    area.frac
                = area.frac/sum(area.frac,na.rm=TRUE)
25
                 = as.numeric(cor(t(xdata),ydata))
26
    cor.map
27
    sig.area.obs = sum(area.frac[abs(cor.map) > 2/sqrt(ntot)],na.rm=TRUE)*100
28
29
    perc.mc = numeric(ntrials)
    for ( n in 1:ntrials) {
30
31
    signif.goodonly = abs(cor(t(xdata[lgood,]),rnorm(ntot))) > 2/sqrt(ntot)
    perc.mc[n] = sum(area.frac[lgood][signif.goodonly])*100
32
33
34
    crit.null = quantile(perc.mc,probs=1-alpha)
    perc.hist = hist(perc.mc,breaks=hbreaks,plot=FALSE)
xrange = range(perc.hist$breaks,sig.area.obs)
35
36
                  = paste('Percent Area with Significant Correlations (p<',alpha,')',sep='')
37
    xlab.say
    hist(perc.mc, col='grey70', xlab=xlab.say, ylab='count', xlim=xrange, main="")
38
    abline(v=crit.null,lwd=2,col='grey30',lty='dashed')
text(crit.null,mean(par()$yaxp[1:2]),'95%\npercentile',pos=2,cex=1.2)
39
40
41
    abline(v=sig.area.obs,lwd=2)
    text(sig.area.obs,mean(par()$yaxp[1:2]),paste('obs fraction= ',round(sig.area.obs),'%',sep=''),pos=4,cex=1.2)
42
    ftitle.top = 'Percent Area with Significant Correlations'
43
    title(ftitle.top,line=2.0)
44
45
    list (crit.null=crit.null,sig.area.obs=sig.area.obs)
46
    }
```

The figure should look like fig. 13.4 in the lecture notes. State the values of crit.null and sig.area.obs that you obtain.

### Exercise 13.2, page 55

An R code that performs the permutation test is the following:

```
fieldsig.permutation = function(lon,lat,xdata,ydata,ntrials=1000,alpha=0.05) {
1
    ### PERFORMS PERMUTATION FIELD SIGNIFICANCE TEST
2
    # INPUT:
3
4
        LON[NLON]: LONGITUDE OF THE FIELD DATA
    #
        LAT[NLAT]: LATITUDE OF THE FIELD DATA
5
    #
        XDATA[NLON,NLAT,NTOT]: DATA ARRAY FOR THE FIELD
YDATA[NTOT]: REFERENCE TIME SERIES
6
    #
7
    #
        NTRIALS: NUMBER OF TRIALS (DEFAULT = 10000)
8
    #
0
        ALPHA: SIGNIFICANCE LEVEL OF THE TEST (DEFAULT = 5%)
    # OUTPUT:
10
        HISTOGRAM OF AREAS WITH SIGNIFICANT CORRELATIONS
11
    #
12
       LIST:
    #
          $CRIT.NULL: SIGNIFICANCE THRESHOLD FROM LIVEZEY-CHEN TEST
13
    #
14
    #
         $SIG.AREA.OBS: PERCENTAGE AREA WITH SIGNIFICANT CORRELATIONS
15
16
   nlon
                  = length(lon)
17
    nlat
                 = length(lat)
    ntot
18
                 = length(ydata)
    dim(xdata) = c(nlon*nlat, ntot)
19
20
                = rep(cos(lat*pi/180),each=nlon)
= !is.na(xdata[,1])
21
    area.frac
    lgood
22
23
    area.frac[!lgood] = NA
24
   area.frac = area.frac/sum(area.frac,na.rm=TRUE)
25
                = as.numeric(cor(t(xdata),ydata))
26
    cor.map
27
    sig.area.obs = sum(area.frac[abs(cor.map) > 2/sqrt(ntot)],na.rm=TRUE)*100
28
29
    perc.mc = numeric(ntrials)
    for ( n in 1:ntrials) {
30
    temp = TRUE
31
    while(temp) {
32
    nperm = sample(nyrs)
temp = any(nperm == 1:nyrs)
33
34
35
36
    signif.goodonly = abs(cor(t(xdata[lgood,]),ydata[nperm])) > 2/sqrt(ntot)
    perc.mc[n] = sum(area.frac[lgood][signif.goodonly])*100
37
38
    crit.null
                = quantile(perc.mc,probs=1-alpha)
39
                = hist(perc.mc,plot=FALSE)
40
    perc.hist
41
    xrange
                 = range(perc.hist$breaks,sig.area.obs)
                 = paste ('Percent Area with Significant Correlations (p<', alpha,')', sep='')
    xlab.say
42
    hist(perc.mc, col='grey70', xlab=xlab.say, ylab='count', xlim=xrange, main="")
43
    abline(v=crit.null,lwd=2,col='grey30',lty='dashed')
text(crit.null,mean(par()$yaxp[1:2]),'95%\npercentile\nUnder H0',pos=2,cex=1.2)
44
45
    abline(v=sig.area.obs,lwd=2)
46
    text(sig.area.obs,mean(par()$yaxp[1:2]),paste('obs fraction= ',round(sig.area.obs),'%',sep=''),pos=4,cex=1.2)
47
    ftitle.top = 'Percent Area with Significant Correlations Under Null Hypothesis'
48
    title(ftitle.top,line=2.0)
49
    list(crit.null=crit.null,sig.area.obs=sig.area.obs,ntrials=ntrials)
50
51
    }
```

The figure should look like fig. 13.4 in the lecture notes.

Exercise 13.3, page 56

MIC can be computed as follows:

```
r.sqr = as.numeric(rep(NA,tem.eof$neof))
  for ( n in 1:tem.eofpreof) if (nyrs -n - 1 - 1 > 0) r.sqr[n] =
2
     summary(lm(nino34<sup>tem.eof$pc[,1:n]))$r.square</sup>
3
4
  n = 1:tem.eof$neof
5
  mic.rsqr = log(1-r.sqr) +
6
           (nyrs+1) * (n+1) / (nyrs-(n+1)-2) -
7
            (nyrs+1)/(nyrs-1-2) -
8
           (nyrs+1) * (n) / (nyrs-n-2)
9
```

Exercise 13.4, page 56

# Exercise 13.5, page 56

An R code for calculating FDR is the following:

```
1
    fdr.fieldsig = function(y,x,fdr=0.1) {
    ## APPLIES FALSE DISCOVERY RATE TO A CORRELATION MAP
2
    ## INPUT:
3
           Y[NTOT]: REFERENCE TIME SERIES
4
           X[NTOT,MTOT]: FIELD TIME SERIES
5
          FDR: FALSE DISCOVERY RATE
6
    #
    ## OUTPUT:
7
          LREJECT[MTOT]: LOGICAL VECTOR INDICATING REJECTIONS
8
    #
    ## DEFINE SHAPE OF ARRAYS
10
    ntot = length(y)
11
    if ( length(x) %% ntot != 0 ) stop('x not dimensioned correctly')
12
           = length(x)/ntot
13
   mtot
   \dim(x) = c(ntot, mtot)
14
15
    ## COMPUTE CORRELATION MAP AND ASSOCIATED P-VALUES
16
17
   dof = ntot - 2
    cor.map = as.numeric(cor(x,y))
18
   tval = abs(cor.map) * sqrt( dof / (1-cor.map<sup>2</sup>) )
pval = 2 * pt(tval,dof,lower.tail=FALSE)
19
   pval
20
21
   ## SCREEN OUT MISSING DATA
22
23
    lgood = !is.na(x[1,])
   ngood = sum(lgood)
24
25
   n.order = order(pval)
26
27
   p.upper = (1:ngood)/ngood*alpha
28
29
                      = rep(NA,mtot)
    lreject
30
    lreject[n.order[1:ngood]] = pval[n.order[1:ngood]] <= p.upper</pre>
31
32
33
    lreject
34
35
    }
```

The plot should be similar to fig. 13.8 of the notes.

# Exercise 14.1, page 59

To minimize  $SSE_W$ , differentiate (14.2) with respect to **B** and set the result to zero. The derivative is

$$\frac{\partial SSE_W}{\partial B_{ab}} = (-2) \sum_{s'=1}^{S} \sum_{n=1}^{N} W_{bs'} (X_{na}) \left( Y_{ns'} - \sum_{m=1}^{M} X_{nm} B_{ms'} \right)$$
$$= (-2) \mathbf{W} (\mathbf{Y} - \mathbf{XB})^T \mathbf{X}$$

Because  ${\bf W}$  is positive definite, the derivative vanishes if and only if

$$\left(\mathbf{Y} - \mathbf{XB}\right)^T \mathbf{X} = \mathbf{0}.$$
 (22.69)

The solution is independent of  $\mathbf{W}$  and identical to the least squares estimate.

Exercise 14.2, page 61

An R function that computes MIC is the following:

```
micc.gaussian = function(x,y,equal.dim=TRUE,alpha=0.05) {
  ### THIS FUNCTION COMPUTES CORRECTED MUTUAL INFORMATION CRITERION (MICC)
2
  ###
          FOR Y = XB + E
3
  ### INPUT:
4
         X[NSAMP,XDIM]: X DATA ARRAY, OFTEN FORMATTED AS [TIME,EOF]
  ##
5
   ##
         Y[NSAMP, YDIM]: Y DATA ARRAY, OFTEN FORMATTED AS [TIME, EOF]
6
   ##
        EQUAL.DIM: LOGICAL INDICATING WHETHER
7
   ##
         (TRUE) EQUAL NUMBER OF X'S, Y'S CHOSEN: MIC[MIN(XDIM, YDIM)];
8
  ##
         (FALSE) MIC FOR ALL TRUNCATIONS ARE COMPUTED: MIC[XDIM, YDIM]
0
  ### OUTPUT: LIST (DIMENSIONS DEPEND ON EQUAL.DIM)
10
11 ##
       $MICC: MIC VALUES
12
 ##
        $PENALTY: THE PENALTY TERM IN MIC
13
  ##
      $CRIT: SIGNIFICANCE THRESHOLD OF MIC
14
  if (is.null(dim(x))) dim(x) = c(length(x), 1)
15
  if (is.null(dim(y))) dim(y) = c(length(y), 1)
16
17
  nsamp = dim(x)[1]
18
  nx.dim = dim(x)[2]
19
  ny.dim = dim(y)[2]
20
21
  if (nsamp != dim(y)[1]) stop('x and y have inconsistent leading dimensions')
22
23
  cov.xy = cov(cbind(x,y))
24
25
26 if (equal.dim) {
27 max.dim = min(nsamp-1,nx.dim,ny.dim)
28 micc
              = as.numeric(rep(NA,max.dim))
29 micc.penalty = as.numeric(rep(NA,max.dim))
30 micc.crit = as.numeric(rep(NA,max.dim))
  for (n in 1:max.dim) if (nyrs -2*n-2>0) {
31
  nall = c(1:n, 1:n+nx.dim)
32
  micc.penalty[n] = (nyrs+1) * (2*n/(nyrs-2*n-2)-2*n/(nyrs-n-2))
33
34 micc.crit[n] = - qchisq(alpha,n^2,lower.tail=FALSE)/(nyrs-(2*n+3)/2) + micc.penalty[n]
35 micc[n] = log(det(cov.xy[nall,nall])) -
             log(det(cov.xy[1:n,1:n,drop=FALSE])) -
36
             log(det(cov.xy[1:n+nx.dim,1:n+nx.dim,drop=FALSE])) +
37
38
            micc.penalty[n]
39 }
40 } else {
41 [LEAVE BLANK]
  }
42
43
  list(micc=micc,penalty=micc.penalty,crit=micc.crit)
44
45
46
   }
```

### Exercise 14.3, page 61

The MIC values for prefix= "Sep2Dec" are

```
168 ANSWERS TO EXERCISES
```

1	> mic	cc				
2	[1]	-2.4045042	-3.0197617	-3.1685965	-3.5064433	
3	[5]	-3.8591420	-3.8506078	-3.6629287	-3.3292924	
4	[9]	-2.6104376	-2.2448622	-1.3491189	-0.4359491	
5	[13]	1.2171204	3.1035460	5.3236336	8.2142499	
6	[17]	11.2348781	15.7603783	20.8051143	27.8231220	
7	[21]	36.3598677	46.9349010	59.7489209	76.9114337	
8	[25]	98.6326493	128.3214484	169.2004940	226.5889051	
9	[29]	310.8455055	450.1993069	707.8745492	1318.4198146	
10	[33]	4431.5898013	NA			

### Exercise 14.4, page 61

A code for making predictions and plotting results is the following:

```
****
2
      ######### PLOT PREDICTED AND OBSERVED FIELDS
      ****
3
                     = chol2inv(chol(cov(pc[mon.init,year.keep,l:nmin]))) %*% cov(pc[mon.init,year.keep,l:nmin],pc[mon.targ,yea
4
     beta.hat
     ypred.tilde = pc[mon.init,,1:nmin] %*% beta.hat
     cbar.breaks = c(-3.5, -3.0, -2.5, -2.0, -1.5, -1.0, -0.5, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 4.5)
cbar.cols = c("white","white","white","white","white","white","#BEBEBE","#A7A7A7","#909090","#7A7A7A","#63636
cbar.list = list(plot.breaks=cbar.breaks,plot.cols=cbar.cols)
10
11
     if (leave.data.out) yhat.say = 'Prediction' else yhat.say = 'Regression Fit'
12
     for ( year in c(1981,1982,1988)) {
13
     nypic = year - iyst + 1
14
     fout = paste('ersstv5.pred', year, prefix, sep='')
15
     if (lplotfile) pdf.eps(fout, 'pdf', height=3)
16
     ypred = as.numeric(eof.list$eof[,1:nmin] %*% ypred.tilde[nypic,1:nmin])
17
     if (year == 1981) cbar.list$plot.breaks = cbar.breaks/2 else cbar.list$plot.breaks = cbar.breaks
18
19
     par(mfrow=c(1,2),mar=c(0.5,0.5,2,0.5))
     pai(minow-c(1,2),marc(0,0,0,0,2,0,0))
plot_latlon_bw(lon.sst,lat.sst,sst.full[,mon.targ,nypic],lmaskzero=FALSE,shrinkdomain=TRUE,suppress.axislab=TRUE,plot.l
title(main=paste('Observed',month.abb[mon.targ],year),line=0.5)
plot_latlon_bw(lon.sst,lat.sst,ypred,lmaskzero=FALSE,shrinkdomain=TRUE,suppress.axislab=TRUE,plot.legend=FALSE,cbar.lis
title(main=paste(yhat.say,' for ',year,' (',month.abb[mon.init],' IC; ',nmin,'EOFs)',sep=''),line=0.5)
20
21
22
23
     if (lplotfile) dev.off()
24
25
     }
```



**Figure 22.16** Comparison between the predicted (left) and observed (right) December anomaly SSTs for a neutral (top; 1981), warm (middle; 1982), and cold (bottom; 1988) ENSO year. The contour label is in units of degrees Celcius. The prediction is based on 22 EOFs with September initial conditions.

# Exercise 15.1, page 65

The plot of MIC should be consistent with fig. 15.5 from the notes. The first 5 rows and columns are

1	> mic[1:5,1:5]				
2	[,1]	[,2]	[,3]	[,4]	[ <b>,</b> 5]
3	[1,] -0.001874504	-0.05766078	-0.2613093	-0.2377176	-0.2327597
4	[2,] -0.198177256	-0.42024605	-0.6405989	-0.5908572	-0.5438949
5	[3,] -0.185517754	-0.51092279	-0.6917943	-0.6108790	-0.6101217
6	[4,] -0.149274572	-0.65337786	-0.9198444	-0.8214891	-0.7746863
7	[5,] -0.111265453	-0.61204103	-0.8646096	-0.7910013	-0.6999580

The minimum occurs at  $T_X = 4, T_Y = 3$  and is -0.9198444.

# Exercise 15.2, page 65

An R function that performs all of CCA is given below. For  $T_X = 4, T_Y = 3$ , the canonical correlations are

1 > cca.list\$can.cor

2 [1] 0.73858363 0.66234594 0.05366594

```
cca.pca = function(fx,fy,ex,ey,tx=NULL,ty=NULL) {
1
    ****
2
    ## PERFORMS CANONICAL CORRELATION ANALYSIS ON X AND Y.
3
    ## X AND Y ARE ASSUMED TO BE IN THE FOLLOWING FORMS:
## X = FX %*% EX^T ### AND ### Y = FY %*% EY^T
4
5
    ## WHERE FX AND FY HAVE COVARIANCE MATRICES = I
6
    ## FOR EXAMPLE: FROM PRINCIPAL COMPONENT ANALYSIS
7
    ## IF TX OR TY = NULL, THEN BOTH (TX,TY) SELECTED USING MIC
8
0
    ## INPUT:
         FX[NTOT,MX]: TIME SERIES FOR X-COMPONENTS: ORTHOGONAL
FY[NTOT,MY]: TIME SERIES FOR Y-COMPONENTS: ORTHOGONAL
10
11
         EX[SX,MX]: SPATIAL STRUCTURES OF THE X-COMPONENTS
12
13
         EY[SY,MY]: SPATIAL STRUCTURES OF THE Y-COMPONENTS
         TX: TRUNCATION FOR X (TX <= MX); IF NULL, TX IS SELECTED TY: TRUNCATION FOR Y (TY <= MY); IF NULL, TY IS SELECTED
14
15
    #
16
    ## OUTPUT LIST:
17
         MIC[MX,MY]: MUTUAL INFORMATION CRITERION
    #
18
    #
         NMIN[1,2]: VALUES OF TX, TY THAT MINIMIZES MIC
19
         CAN.COR[MIN(TX,TY)]: CANONICAL CORRELATIONS
    #
         RX[NTOT, MIN(TX, TY)]: CANONICAL VARIATES FOR X
20
21
         RY[NTOT, MIN(TX, TY)]: CANONICAL VARIATES FOR Y
    #
         PX[SX ,MIN(TX,TY)]: CANONICAL LOADING VECTORS FOR X
PY[SY ,MIN(TX,TY)]: CANONICAL LOADING VECTORS FOR Y
22
23
    #
24
         QX.TILDE[TX,MIN(TX,TY)]: WEIGHTING VECTORS FOR X-FEATURES
25
    #
         QY.TILDE[TY, MIN(TX, TY)]: WEIGHTING VECTORS FOR Y-FEATURES
         TX, TY: SELECTED VALUES OF TX AND TY
26
27
    *****
28
    ntot = dim(fx)[1]
29
    mx = dim(fx)[2]
30
31
          = dim(fy)[2]
    my
   if (ntot != dim(fy)[1]) stop('fx and fy have inconsistent time dimension')
32
33
    micc = array(NA, dim=c(mx, my))
34
    for (ty in 1:mx) for (tx in 1:my) if (ntot-tx-ty-2 > 0) {
35
    penalty = (ntot+1) * ((tx+ty)/(ntot-tx-ty-2) - tx/(ntot-tx-2) - ty/(ntot-ty-2))
36
                  = svd(cov(fx[,1:tx],fy[,1:ty]))
37
    xy.svd
   micc [tx,ty] = sum(log(1-xy.svd$d^2)) + penalty
38
39
    }
40
41
    nmin
            = which (micc == min (micc, na.rm=TRUE), arr.ind=TRUE)
            = nmin[1]
42
    tx
            = nmin[2]
43
   tv
   cov.xy = cov(fx[,1:tx],fy[,1:ty])
xy.svd = svd(cov.xy)
44
45
   can.cor = xy.svd$d
46
47
   qx
            = xy.svd$u
            = xy.svd$v
48
   qy
            = fx[,1:tx] %*% qx
49
   rx
           = fy[,1:ty] %*% qy
50
   ry
51
            = ex[,1:tx] %*% qx
    рх
           = ey[,1:ty] %*% qy
52
    ру
53
   list(can.cor=can.cor,rx=rx,ry=ry,px=px,py=py,tx=tx,ty=ty,mic=micc)
54
55
56
    }
```

### Exercise 15.3, page 65

see fig. ?? of the lecture notes.

#### Exercise 15.4, page 65

see figs. 15.1 and 15.2 of the lecture notes.

Exercise 15.5, page 65

```
1 > cca.pca.list$fexpvar.x
2 [1] 0.1678503 0.2410055 0.1668347 0.1111007
3 > cca.pca.list$fexpvar.y
4 [1] 0.1997934 0.3069071 0.2975915
5 >
6 [1] "sum exp. var. x= 0.68679126577992 from EOFs= 0.686791265779918"
7 [1] "sum exp. var. y= 0.804291929099929 from EOFs= 0.804291929099928"
```

### Exercise 15.6, page 65

An R code that computes the 5% significance levels of the canonical correlations based on 5000 trials of Monte Carlo experiments is the following:

```
## SIGNIFICANCE THRESHOLDS FOR CORRELATIONS
1
  ntrials = 5000
2
  alpha = 0.05
3
  rho.mc = array(NA, dim=c(min(nx.min, ny.min), ntrials))
4
  for (nt in 1:ntrials) {
  x = array(rnorm(nyrs*nx.min),dim=c(nx.min,nyrs))
  y = array(rnorm(nyrs*ny.min),dim=c(ny.min,nyrs))
7
  rho.mc[,nt] = svd( t(svd(x - rowMeans(x))$v) %*% svd(y - rowMeans(y))$v )$d
8
  }
10 rho.crit = numeric(min(nx.min, ny.min))
  for (n in 1:min(nx.min,ny.min)) rho.crit[n] = quantile(rho.mc[n,],probs=1-alpha)
11
```

The resulting values are

```
1 > rho.crit
2 [1] 0.4619982 0.3083790 0.1772026
```

The first two canonical correlations exceed these 5% significance thresholds.

# Exercise 15.7, page 65

The figures for no detrending are shown below. The big change is that the leading component contains a strong trend component. This makes sense because trends are highly correlated, and CCA maximizes correlation.





# Exercise 15.8, page 66

$$\tilde{\boldsymbol{\Sigma}}_{\epsilon Y} = \left(\mathbf{R}_{Y} - \mathbf{R}_{X}\mathbf{S}_{\rho}\right)^{T} \left(\mathbf{R}_{Y} - \mathbf{R}_{X}\mathbf{S}_{\rho}\right) / (N-1)$$

$$= \left(\mathbf{R}_{Y}^{T}\mathbf{R}_{Y} - \mathbf{S}_{\rho}\mathbf{R}_{X}^{T}\mathbf{R}_{Y} - \mathbf{R}_{Y}^{T}\mathbf{R}_{X}\mathbf{S}_{\rho} + \mathbf{S}_{\rho}\mathbf{R}_{X}^{T}\mathbf{R}_{X}\mathbf{S}_{\rho}\right) / (N-1)$$

$$= \mathbf{I} - \mathbf{S}_{\rho}^{2}, \qquad (22.70)$$

where we have used

$$\frac{\mathbf{R}_Y^T \mathbf{R}_Y}{N-1} = \mathbf{I}, \quad \frac{\mathbf{R}_X^T \mathbf{R}_X}{N-1} = \mathbf{I}, \quad \frac{\mathbf{R}_Y^T \mathbf{R}_X}{N-1} = \mathbf{S}_{\rho}.$$
 (22.71)

# Exercise 16.1, page 67

$$\mathbb{E}[\mathbf{C}] = \mathbf{0} \text{ because } \mathbb{E}[\mathbf{u}_t^* - \hat{\boldsymbol{\mu}}_Y] = \mathbf{0}. \text{ Also}$$
$$\mathbb{E}[\mathbf{A}] = \operatorname{cov}[\mathbf{u}_t^* - \hat{\boldsymbol{\mu}}_Y] = \operatorname{cov}[\mathbf{u}_t^*] + \operatorname{cov}[\hat{\boldsymbol{\mu}}_Y] = \boldsymbol{\Sigma}_U + \frac{1}{N_Y}\boldsymbol{\Sigma}_U = \left(1 + \frac{1}{N_Y}\right)\boldsymbol{\Sigma}_U.$$

# Exercise 16.2, page 67

By the invariance property, X and Y may be re-scaled without affecting the measure. Accordingly, normalize each variable by the standard deviation of Y:

$$X^* = X/\sigma_Y$$
 and  $Y^* = Y/\sigma_Y$ . (22.72)
Then,  $\operatorname{var}[Y^*] = 1$ , and  $\operatorname{var}[X^*] = \sigma_X^2 / \sigma_Y^2$ , hence

$$d(\sigma_X^2, \sigma_Y^2) = d(\sigma_X^2 / \sigma_Y^2, 1).$$
(22.73)

The last function clearly depends only on the variance ratio.

## Exercise 16.3, page 67

First compute the economy SVD of Y:

$$\mathbf{Y} = \mathbf{U}_Y \quad \dot{\mathbf{S}}_Y \quad \dot{\mathbf{V}}_Y^T$$

$$N_Y \times S \qquad N_Y \times T \quad T \times T \quad T \times S'$$
(22.74)

where we assume  $T < \min(N_Y, S)$ . Now define the transformation matrix

$$\mathbf{L} = \dot{\mathbf{V}}_Y \dot{\mathbf{S}}_Y^{-1}.$$
 (22.75)

The transformed variables become

$$\mathbf{Y}' = \mathbf{Y}\mathbf{L} = \mathbf{U}_Y \tag{22.76}$$

$$\mathbf{X}' = \mathbf{X}\mathbf{L}.\tag{22.77}$$

Thus the covariance matrices for the transformed variables become

$$\tilde{\boldsymbol{\Sigma}}_{Y} = \frac{1}{N_{Y} - 1} \mathbf{Y}^{'T} \mathbf{Y}^{'} = \frac{1}{N_{Y} - 1} \mathbf{I} \quad \text{and} \quad \tilde{\boldsymbol{\Sigma}}_{X} = \frac{1}{N_{X} - 1} \mathbf{X}^{'T} \mathbf{X}^{'}.$$
(22.78)

Substituting this into the generalized eigenvalue problem gives

$$(N_Y - 1)\tilde{\boldsymbol{\Sigma}}_X \tilde{\mathbf{q}} = \lambda \tilde{\mathbf{q}}.$$
(22.79)

Note, this is a standard eigenvalue problem. Thus, we compute the SVD

$$\sqrt{N_Y - 1}\mathbf{X}\mathbf{L} = \mathbf{U}_X \mathbf{S}_X \mathbf{V}_X^T.$$
(22.80)

It is straightforward to show that the eigenvectors  $\tilde{\mathbf{q}}$  are the columns of  $\mathbf{V}_X$ , and the eigenvalues are the squared singular values. Thus, ignoring the normalization convention, the variates are

$$\mathbf{R}_X = \mathbf{X}\mathbf{L}\mathbf{V}_X$$
 and  $\mathbf{R}_Y = \mathbf{Y}\mathbf{L}\mathbf{V}_X$ . (22.81)

The loading vectors can be found by regression.

### Exercise 16.4, page 68

As usual, taking derivatives is straightforward when the quantity is written in index notation. For notational simplicity, we temporarily suppress the X and Y subscripts. Thus, the sum total variance is

$$\gamma_Y = E\left[\left(\mathbf{y}_s - E[\mathbf{y}_s] - \mathbf{P}_{sk}\mathbf{r}_k\right)\mathbf{W}_{ss'}\left(\mathbf{y}_{s'} - E[\mathbf{y}_{s'}] - \mathbf{P}_{s'k'}\mathbf{r}_{k'}\right)\right].$$
(22.82)

Differentiating with respect to P yields

$$\frac{\partial \gamma_Y}{\partial \mathbf{P}_{ab}} = (-2)E\left[\mathbf{r}_b \mathbf{W}_{as'} \left(\mathbf{y}_{s'} - E[\mathbf{y}_{s'}] - \mathbf{P}_{s'k'} \mathbf{r}_{k'}\right)\right]$$
(22.83)

$$= (-2)\mathbf{W} \left( \operatorname{cov}[\mathbf{y}, \mathbf{r}] - \mathbf{P} \operatorname{cov}[\mathbf{r}, \mathbf{r}] \right).$$
(22.84)

Since W is positive definite, the derivative vanishes only if the term in parentheses vanishes. Setting this term to zero and solving for P yields the desired solution.

In the finite sample case, the sum total variance is

$$\dot{\gamma}_{Y} = \left(\dot{\mathbf{Y}}_{ns} - \mathbf{R}_{nk}\mathbf{P}_{sk}\right)\mathbf{W}_{ss'}\left(\dot{\mathbf{Y}}_{ns'} - \mathbf{R}_{nk'}\mathbf{P}_{s'k'}\right),$$
(22.85)

where products of repeated indices are implicitly summed. Differentiating with respect to  $\mathbf{P}$  yields

$$\frac{\partial \dot{\gamma}_Y}{\partial P_{ab}} = (-2)\mathbf{R}_{nb}\mathbf{W}_{as'} \left( \dot{\mathbf{Y}}_{ns'} - \mathbf{R}_{nk'}\mathbf{P}_{s'k'} \right)$$
(22.86)

$$= (-2)\mathbf{W} \left( \dot{\mathbf{Y}} - \mathbf{R}\mathbf{P}^T \right)^T \mathbf{R}$$
 (22.87)

$$= (-2)\mathbf{W}\left(\dot{\mathbf{Y}}^T\mathbf{R} - \mathbf{P}\mathbf{R}^T\mathbf{R}\right).$$
(22.88)

Since W is positive definite, the derivative vanishes only if the term in parentheses vanishes. Setting this term to zero and solving for P yields the desired solution.

### Exercise 16.5, page 68

A useful trick is to re-write  $\lambda$  as

$$\lambda = \frac{\mathbf{q}^T \mathbf{\Delta} \mathbf{q}}{\mathbf{q}^T \mathbf{\Sigma}_Y \mathbf{q}} + 1$$
(22.89)

where  $\mathbf{\Delta} = \mathbf{\Sigma}_X - \mathbf{\Sigma}_Y$ . If  $H_0$  is true, then  $\mathbf{\Delta} = \mathbf{0}$  and  $\mathbf{q}^T \mathbf{\Delta} \mathbf{q}$  vanishes for all  $\mathbf{q}$ . On the other hand, if  $H_0$  is false,  $\mathbf{\Delta}$  must differ from zero in at least one element. Suppose this non-zero element occurs at the *i*'th diagonal element. Then, choosing  $\mathbf{q}$  to be vector with a one at the *i*'th element and zero everywhere else yields  $\mathbf{q}^T \mathbf{\Delta} \mathbf{q} = \Delta_{ii}$ , which is non-zero by assumption. Moreover, this expression depends only on  $\Delta_{ii}$ , so the value of the other elements of  $\mathbf{\Delta}$  are irrelevant in this argument. Now, suppose all the diagonal element. Suppose  $\Delta_{ij} \neq 0$ . Then, choosing  $\mathbf{q}$  to be a vector that has ones at the *i*'th and *j*'th elements and zero everywhere else yields  $\mathbf{q}^T \mathbf{\Delta} \mathbf{q} = \Delta_{ij}$ , which is non-zero by assumption. Moreover, this expression depends only on  $\Delta_{ij}$ , so the value of the other offdiagonal elements are irrelevant. This proof exhausts all possible cases, therefore if the hypothesis  $H_0$  is false, then  $\lambda$  must differ from one for at least one choice of  $\mathbf{q}$ .

# Exercise 16.6, page 68

Substituting  $\dot{\mathbf{Y}} = \mathbf{R}_Y \dot{\mathbf{P}}_Y^T$  and using standard properties of the trace gives

$$\|\dot{\mathbf{Y}}\|_{W}^{2} = \frac{1}{N_{Y} - 1} \operatorname{tr} \left[ \mathbf{R}_{Y} \mathbf{P}^{T} \mathbf{W} \mathbf{P} \mathbf{R}_{Y}^{T} \right] = \operatorname{tr} \left[ \mathbf{P}^{T} \mathbf{W} \mathbf{P} \left( \frac{\mathbf{R}_{Y}^{T} \mathbf{R}_{Y}}{N_{Y} - 1} \right) \right]$$
$$= \operatorname{tr} \left[ \mathbf{P}^{T} \mathbf{W} \mathbf{P} \right] = \sum_{k=1}^{T} \mathbf{p}_{k}^{T} \mathbf{W} \mathbf{p}_{k}.$$
(22.90)

The last equation shows that the sum total variance equals the sum of inner products of the loading vectors with no cross terms. Thus,  $\mathbf{p}_k^T \mathbf{W} \mathbf{p}_k$  may be interpreted as the variance explained by the k'th component.

Exercise 16.7, page 69

A function that computes MIC and performs CDA is the following

```
cda.eof = function(xdata,ydata,eof.list) {
1
    ### PERFORMS COVARIANCE DISCRIMINANT ANALYSIS ON X AND Y
2
    ### INPUT:
3
4
    ###
          XDATA[NX, MDIM]
           YDATA[NY, MDIM]
5
    ###
    ### EOF.LIST: LIST FROM EOF CALCULATION
6
7
    ### OUTPUT:
    ###
          MIC[NEOF]: MIC AS A FUNCTION OF NUMBER OF PCS NMIN: LOCATION OF MINIMUM MIC
8
9
    ###
    ### DISCR.RATIO[NEOF]: DISCRIMINANT RATIOS VS. NUMBER OF PCS
#### RX[NX,NMIN]: VARIATE TIME SERIES FOR X
10
11
    ### RY[NY,NMIN]: VARIATE TIME SERIES FOR Y
### PMAT[SPACE,NMIN]: LOADING VECTOR
12
13
14
15
    neof
                = eof.list$neof
16
    if (length(xdata) %% neof != 0 ) stop('xdata not a multiple of neof')
    if (length(ydata) %% neof != 0 ) stop('ydata not a multiple of neof')
nx = length(xdata) / neof
17
   nx
18
                = length(ydata) / neof
19
    ny
    dim(xdata) = c(nx, neof)
20
21
    dim(ydata) = c(ny, neof)
22
23
    nmax
                = min(nx-1, ny-1, neof)
24
25
    nt
                = nx + ny
    cov.x
                = cov(xdata) * (nx-1)/nx
26
    cov.y
27
                = cov(ydata) * (ny-1)/ny
                 = (nx * cov.x + ny * cov.y)/nt
28
    cov.t
29
30
    mic
                 = as.numeric(rep(NA,nmax))
31
    penalty
                 = as.numeric(rep(NA, nmax))
    for (ne in 1:nmax) penalty[ne] = ne * ( nx*(nx+1)/(nx-ne-2) + ny*(ny+1)/(ny-ne-2) - nt*(nt+1)/(nt-ne-2))
32
33
    for (ne in 1:nmax) {
34
    mic[ne] = nx * log(det(cov.x[1:ne,1:ne,drop=FALSE])) +
35
               ny * log(det(cov.y[1:ne,1:ne,drop=FALSE]))
               nt * log(det(cov.t[1:ne,1:ne,drop=FALSE])) + penalty[ne]
36
37
    }
38
    penalty = penalty / nt
mic = mic / nt
39
40
41
    nmin
             = which.min(mic)
42
    ### rescale covariances to be unbiased
43
    cov.x = cov.x * nx / (nx-1)
44
45
            = cov.y * ny / (ny-1)
    cov.y
46
47
    gev.list
                 = gev(cov.x[1:nmin,1:nmin],cov.y[1:nmin,1:nmin])
    discr.ratio = gev.list$lambda
48
              = xdata[,1:nmin] %*% gev.list$q
= ydata[,1:nmin] %*% gev.list$q
49
    rx
50
    rv
    pmat
                 = eof.list$eof[,1:nmin] %*% cov.y[1:nmin,1:nmin] %*% gev.list$q
51
52
    list (mic=mic, nmin=nmin, discr.ratio=discr.ratio, rx=rx, ry=ry, pmat=pmat)
53
54
    }
```

The MIC values for the data is

1 > mic 2 [1] -0.19404378 -0.60912246 -0.60089462 -0.77768582 -0.97252532

Exercise 16.8, page 70

```
sdiscr.ratio
1 $discr.ratio
2 [1] 21.8064448 1.1885098 1.0114935 0.8586228 0.8169091
```

## Exercise 16.9, page 70

The variates should look like those in fig. 16.3. The diagonal elements of the 20C variates match the discriminant ratios:

```
1 > diag(cov(cda.eof.list$rx))
2 [1] 21.8064448 1.1885098 1.0114935 0.8586228 0.8169091
3 > diag(cov(cda.eof.list$ry))
4 [1] 1 1 1 1 1
```

# Exercise 16.10, page 70

The loading vector should look like that in fig. 16.3.

# Exercise 16.11, page 70

An R code for Monte Carlo experiments is the following:

```
= 1000
1 ntrials
  cov.identity = diag(1, nrow=nmin, ncol=nmin)
2
   lambda.mc = array(NA, dim=c(nmin, ntrials))
3
  for ( nt in 1:ntrials) {
4
5 cov.x = rWishart(1, dof.20c-1, cov.identity)[,,1]/(dof.20c-1)
  cov.y = rWishart(1, dof.ctr-1, cov.identity)[,, 1]/(dof.ctr-1)
6
  lambda.mc[,nt] = gev(cov.x,cov.y)$lambda
7
  }
  lambda.crit = array(NA,dim=c(nmin,2))
9
  for (n in 1:nmin) lambda.crit[n,] = quantile(lambda.mc[n,],probs=c(alpha/2,1-alpha/2))
10
```

The percentiles are

1 > lambda.crit
2 [,1] [,2]
3 [1,] 1.1030110 1.4559969
4 [2,] 0.9907935 1.2569631
5 [3,] 0.8839713 1.1269708
6 [4,] 0.7872901 1.0173159
7 [5,] 0.6806821 0.8910988

Exercise 16.12, page 70

The 95% percentiles as a function of T are

```
[1] "5% significance for univariate F-test= 1.59949546683544"
1
2
  > lambda.crit
       1.603912 2.029011 2.235836 2.562755 2.863984 3.197075
   [1]
       3.511517 3.839611 4.128509 4.473972 4.912822 5.326827
   [7]
4
        5.847943 6.402706 6.933168 7.538059 8.366698 9.042997
  [13]
5
        9.557897 10.486632 11.228660 12.446608 13.606677 14.982855
  [19]
6
  [25] 16.511167 18.575636 19.861070 22.171482 24.300277 27.020040
```

The critical values increase with dimension because each model is nested within the next, and the algorithm maximizes the discriminant ratio for each model. Looking in the reverse direction, as the dimension decreases, some of the parameters that were available for tuning become constrained, so then the maximum must decrease.

### Exercise 17.1, page 73

An R function that performs ANOVA on a spatial field is the following:

```
1
    f.anova.array = function(x,nspace,nens,ncon,alpha=0.05) {
2
    ### COMPUTES THE F-STATISTIC IN ANOVA FOR EACH ELEMENT IN NSPACE
    # INPUT:
3
         X: [NSPACE, NENS, NCON] ARRAY OF DATA
4
    #
         NSPACE: NUMBER OF SPATIAL ELEMENTS FOR INDIVIDUALLY COMPUTING F
5
         NENS: NUMBER OF ENSEMBLE MEMBERS
6
        NCON: NUMBER OF CONDITIONS
8
         ALPHA: SIGNIFICANCE LEVEL
    # OUTPUT: LIST WITH THE FOLLOWING VARIABLES
10
        F: [NSPACE] VECTOR CONTAINING THE F-VALUES
11
       F.CRIT: THE CRITICAL F-VALUE FOR F AT THE ALPHA*100% SIGNIFICANCE LEVEL
12
   ## COMPUTE GRAND MEAN
13
14
   dim(x) = c(nspace, nens*ncon)
   gmean = rowMeans(x)
15
16
   ## COMPUTE ENSEMBLE MEAN
17
18
   dim(x) = c(nspace, nens, ncon)
19
   emean = array(NA, dim=c(nspace, ncon))
   for ( nc in 1:ncon) emean[,nc] = rowMeans(x[,,nc,drop=FALSE])
20
21
   ## COMPUTE SIGNAL VARIANCE
22
23
   var.sig = rowSums((emean - gmean)^2) / (ncon-1)
24
   ### COMPUTE NOISE
25
26
   noise = arrav(NA, dim=c(nspace, nens, ncon))
   for ( nc in 1:ncon) noise[,,nc] = x[,,nc] - emean[,nc]
27
28
   ### COMPUTE NOISE VARIANCE
29
30
   dim(noise) = c(nspace, nens*ncon)
             = rowSums (noise^2) /ncon/ (nens-1)
31
   var.nos
32
33
   ### COMPUTE F
    f = nens * var.sig / var.nos
34
   f.crit = qf(alpha,ncon-1,ncon*(nens-1),lower.tail=FALSE)
35
36
   list(f=f,f.crit=f.crit)
37
38
39
   }
```

# Exercise 17.2, page 74

The F-map should look like fig. 18.1 in the notes. The summary is

1	>	• summary(fval)								
2		Min.	1st	Qu.	Median	Mean	3rd Qu.	Max.	NA's	
3		1.512	2.	698	3.700	4.064	4.918	32.843	3720	

## Exercise 17.3, page 74

The results should like figure 18.1 from the notes.

# Exercise 17.4, page 75

Generally, F values decay as the number ensemble members decreases and as the lead time decreases.

# Exercise 18.1, page 78

An R code for estimating significance levels of the eigenvalues from PrCA is the following:

```
prca.null.fixdim = function(ndim,ncon,nens,alpha=0.01,ntrial=10000) {
## ESTIMATES SIGNIFICANCE LEVEL OF THE ALL EIGENVALUES FROM PRCA/MANOVA
## FOR FIXED TRUNCATION NDIM
 1
2
3
    ## BY MONTE CARLO METHODS
4
     ## INPUT:
5
           NDIM: DIMENSION OF THE RANDOM VECTOR
6
     #
          NDIM: DIMENSION OF THE KANDOM VECTOR
NCOM: NUMBER OF CONDITIONS (E.G., NUMBER OF YEARS)
NENS: NUMBER OF ENSEMBLE MEMBERS (E.G., NUMBER OF REPITITIONS)
ALPHA: DESIRED SIGNIFICANCE LEVELS (CAN BE MORE THAN ONE)
7
     #
8
     #
9
     #
          NTRIALS: NUMBER OF MONTE CARLO TRIALS
10
    #
    ## OUTPUT:
11
         EVAL.CRIT[2,NDIM]: (alpha,1-alpha) SIGNIFICANCE THRESHOLDS FOR EACH EIGENVALUE
12
    #
13
    eval.max = array(NA, dim=c(ndim, ntrial))
14
15
16
    dof.sig = ncon - 1
17
    dof.nos = ncon*(nens-1)
    for ( nt in 1:ntrial ) {
18
19
              = array(rnorm(ndim*ncon*nens),dim=c(ndim*ncon,nens))
    х
    x.emean = rowMeans(x)
20
    x.noise = x - x.emean
21
22
23
    dim(x.emean) = c(ndim, ncon)
24
    cov.sig
                   = nens * cov(t(x.emean))
25
26
    dim(x.noise) = c(ndim, ncon*nens)
27
    cov.nos = x.noise %*% t(x.noise) / dof.nos
28
29
     eval.max[,nt] = gev(cov.sig,cov.nos)$lambda
30
    }
31
    eval.crit = array(NA,dim=c(2,ndim))
32
33
    for ( n in 1:ndim) eval.crit[,n] = quantile(eval.max[n,],probs=c(alpha,1-alpha))
34
    rownames(eval.crit) = paste(c(alpha,1-alpha)*100,'%',sep='')
35
36
     eval.crit
37
     }
```

#### The resulting 1% significance thresholds for ndim = 6 are

1	> f1	> fval.crit							
2		[,1]	[,2]	[,3]	[,4]	[,5]	[,6]		
3	1%	1.274289	0.9646709	0.7218789	0.5190632	0.3530367	0.1903565		
4	99%	2.728496	1.9477712	1.4799000	1.1425884	0.8754198	0.6339855		

### Exercise 18.2, page 78

1	> mic	2				
2	[1]	-0.2422784468	-0.1859868967	-0.2475338825	-0.2474193296	
3	[5]	-0.3759129623	-0.5164772016	-0.4870610152	-0.4322602904	
4	[9]	-0.4153250531	-0.4605784420	-0.3959601760	-0.3631133884	
5	[13]	-0.3060199465	-0.2622818018	-0.1989315625	-0.1834919804	
6	[17]	-0.1291759564	-0.1230023572	-0.0768296814	-0.0253331366	
7	[21]	0.0003693148	-0.1127739471	-0.0523793555	0.0607204885	
8	[25]	0.1454873891	0.2427903034	0.2756050003	0.3125827453	
9	[29]	0.2887010746	0.2935311458			

# Exercise 18.3, page 78

An R code that performs all steps in PrCA is the following:

```
prca.eof = function(data.eof,nens,ncon,alpha=0.01) {
1
   [PREAMBLE OMITTED]
2
4
   neof = data.eof$neof
   ntot = dim(data.eof$pc)[1]
5
   if ( nens * ncon != ntot) stop('time dimension inconsistent with nens, ncon')
6
   ########## COMPUTE MLE OF COVARIANCE MATRICES
8
0
   fx
              = data.eof$pc
   dim(fx)
              = c(nens,ncon*neof)
10
11
   emean
              = colMeans(fx)
   noise
             = t(t(fx)-emean)
12
   dim(emean) = c(
13
                       ncon, neof)
   dim(noise) = c(nens*ncon, neof)
14
15
   dim(fx)
              = c(nens*ncon, neof)
16
   COVIS
              = cov(emean) * (ncon-1)/ncon
17
   cov.n
              = cov(noise) * (ncon*nens-1)/(ncon*nens)
             = cov(fx ) * (ncon*nens-1)/(ncon*nens)
18
   cov.t
19
   ######### COMPUTE MIC
20
21
              = as.numeric(rep(NA, neof))
   mic
   penalty
              = as.numeric(rep(NA, neof))
22
23
   for ( n in 1:neof) if (nens*ncon-n-ncon-2>0) {
24
   penalty.random = (nens*ncon+1)*((n+ncon)/(nens*ncon-n-ncon-2) - n/(nens*ncon-n-2) - ncon/(nens*ncon-ncon-2))
25
   penalty.fixed = 2*(ncon-1)*n/(nens*ncon)
   penalty.fixedc = (2*ncon*n + n*(n+1))/(ncon*nens-ncon-n-1) - (2*n+n*(n+1))/(ncon*nens-1-n-1)
26
27
   penalty[n]
                  = penalty.fixedc
                  = log(det(cov.n[1:n,1:n,drop=FALSE])) -
28
   mic[n]
                    log(det(cov.t[1:n,1:n,drop=FALSE])) +
29
30
                    penalty[n]
31
   nmin = which.min(mic)
32
33
34
    ######### COMPUTE SIGNIFICANCE THRESHOLD FOR MIC
35
   nu.e
            = ncon*(nens-1)
            = ncon-1
36
   nu.h
   mic.crit = as.numeric(rep(NA, neof))
37
   for (n in 1:neof) if (nens*ncon-n-ncon-2>0) mic.crit[n] = -qchisq(1-alpha,n*nu.h)/(nu.e-(n-nu.h+1)/2) + penalty[n]
38
39
40
   ######### SOLVE PRCA/MANOVA
41
   dof.s
           = ncon - 1
   dof.n
            = ncon * (nens - 1)
42
   dof.t
            = ncon * nens - 1
43
            = cov.s * ncon
                                  / dof.s
   cov.s
44
45
   cov.n
            = cov.n * ncon * nens / dof.n
            = cov.t * ncon * nens / dof.t
46
   cov.t
   gev.list = gev(cov.s[1:nmin,1:nmin],cov.n[1:nmin,1:nmin])
47
           = gev.list$q
48
    α
   f.max
            = gev.list$lambda * nens
49
   rsqr.adj = (f.max-1)/(f.max + ncon*(nens-1)/(ncon-1))
50
   snr.adj = (f.max-1)/nens
51
52
   ### NORMALIZE O
53
   for ( n in 1:nmin) q[,n] = q[,n]/sqrt(as.numeric(t(q[,n])%*% cov.t[1:nmin,1:nmin] %*% q[,n]))
54
55
56
   ### COMPUTE VARIATES AND LOADING PATTERN
57
   dim(fx)
              = c(nens*ncon, neof)
              = fx[,1:nmin] %*% q
58
   r
              = data.eof$eof[,1:nmin] %*% cov.t[1:nmin,1:nmin] %*% q
59
   р
60
   ### FLTP SIGN
61
   for ( n in 1:nmin) if (sum(p[,n],na.rm=TRUE) < 0) \{p[,n] = -p[,n]; r[,n] = -r[,n]; q[,n] = -q[,n] \}
62
63
64
    ### COMPUTE EXPLAINED VARIANCES
   total.var = sum(data.eof$sval^2)/(ntot-1)
65
66
   laood
             = !data.eof$lbad
67
   fexpvar = numeric(nmin)
    for ( n in 1:nmin) fexpvar[n] = sum(data.eof$weight[lgood]^2 * p[lgood,n]^2) / total.var
68
69
    70
71
    ######### WRITE OUT RESULTS
72
    ****
   list(mic=mic,mic.crit=mic.crit,f.max=f.max,r=r,p=p,fexpvar=fexpvar,rsqr.adj=rsqr.adj,snr.adj=snr.adj
73
74
       nmin=nmin, ncon=ncon, nens=nens, alpha=alpha)
75
   }
```

The maximized F-ratios are

```
> prca.list$f.max
1
  [1] 19.4746501 3.1956229 2.7206396 1.0809605 0.6636562 0.6385614
2
  > prca.list$rsqr.adj
3
        0.527391997 0.117092530 0.094146488 0.004866437 -0.020737369
  [1]
4
  [6] -0.022319127
5
  > prca.list$snr.adj
6
                    0.137226433 0.107539975 0.005060034 -0.021021486
  [1]
       1.154665632
7
  [6] -0.022589909
```

## Exercise 18.4, page 79

See fig. 22.18.

### Exercise 18.5, page 79

The following shows that the sample covariance matrix of the predictable components equals the identity matrix.

```
> signif(cov(r),4)
1
2
                 [,1]
                              [,2]
                                           [,3]
                                                        [,4]
                                                                     [,5]
     [1,] 1.000e+00 6.578e-18 -2.864e-16 8.649e-17
                                                              1.190e-17
3
4
     [2,] 6.578e-18 1.000e+00 -1.558e-16 2.514e-16
                                                              1.173e-16
     [3,] -2.864e-16 -1.558e-16 1.000e+00 8.631e-17 -1.065e-16
5
    [4,] 8.649e-17 2.514e-16 8.631e-17 1.000e+00 -2.067e-16
[5,] 1.190e-17 1.173e-16 -1.065e-16 -2.067e-16 1.000e+00
6
7
    [6,] -1.353e-16 -4.863e-17 -8.512e-17 1.338e-16 -3.679e-16
8
                 [,6]
     [1,] -1.353e-16
10
     [2,] -4.863e-17
11
12
     [3,] -8.512e-17
     [4,] 1.338e-16
13
    [5,] -3.679e-16
[6,] 1.000e+00
14
15
```

## Exercise 18.6, page 79

See fig. ??

# Exercise 19.2, page 83

A function that solves the Kalman Filter equations is the following:



CFSv2 Forecasts 2m-Temp; Nov start; Apr target; E= 16; 1982-2009 Predictable Component 1 6EOFs 22%

Figure 22.17 First predictable component.

1995

```
kalman.population = function(mub, sigmab, hop, r.cov, obs) {
1
   ##########
   ## EVALUATES THE KALMAN FILTER EQUATIONS FOR THE MEAN AND COVARIANCE MATRIX OF THE ANALYSIS
      BASED ON *POPULATION* COVARIANCE MATRICES AND MEANS
   ##
   ## INPUT:
                 [NDIM] MEAN OF THE BACKGROUND DISTRIBUTION
        MUB:
   #
        SIGMAB: [NDIM, NDIM] BACKGROUND COVARIANCE MATRIX
        HOP:
                 [NOBS, NDIM] INTERPOLATION OPERATOR
        R.COV: [NOBS, NOBS] COVARIANCE MATRIX OF THE OBSERVATIONAL ERROR
   #
                 [NOBS] THE OBSERVATIONS
   #
        OBS:
10
   ## OUTPUT
11
                 [NDIM] VECTOR OF THE ANALYSIS DISTRIBUTION
        MUA:
12
   #
        SIGMAA: [NDIM,NDIM] COVARIANCE MATRIX OF THE ANALYSIS DISTRIBUTION
   #
13
14
              = sigmab %*% t(hop)
   xo.cov
15
              = hop %*% xo.cov + r.cov
   o.cov
16
   o.cov.inv = chol2inv(chol(o.cov))
17
              = xo.cov %*% o.cov.inv
18
   kgain
19
              = mub
                       + kgain %*% ( obs - hop %*% mub )
20
   mua
   sigmaa
              = sigmab - kgain %*% t(xo.cov)
21
22
   list (mua=mua, sigmaa=sigmaa)
23
24
   }
```

2000

2005

2010



Figure 22.18 Second predictable component

# Exercise 19.3, page 83

A loop that solves for the analysis based on observations and background distribution is the following:



Figure 22.19 Time series of the standard error of the analysis.

```
### ASSIMILATE OBSERVATIONS
1
     mua = matrix(rep(0,ndim*ntot),nrow=ndim)
pa = array(rep(0,ndim^2*ntot),dim=c(ndim,ndim,ntot))
2
3
     for ( n in 1:ndim) pa[n, n, 1] = 400
4
5
     for ( n in 2:ntot) {
6
       muf = dynop %*% mua[,n-1]
pf = dynop %*% pa[,,n-1] %*% t(dynop) + q.cov
8
9
        kf = kalman.population(muf,pf,hop,r.cov,obs[,n])
10
11
       mua[,n] = kf$mua
pa[,,n] = kf$sigmaa
12
13
     }
```

# Exercise 19.4, page 83

A plot of the analysis errors for the three elements of  $\mathbf{x}_t$ , as given by the square root of the diagonal elements of  $\Sigma_t^B$ , is shown in fig. 22.19. The errors asymptote to a constant after a long time.

Plots for the case  $\mathbf{R}_t = 1$  are shown in fig. 22.20. The standard error of the analysis is much smaller in this case, because the observational uncertainty is much smaller and hence provides more information about the true state than in the previous case. The analysis errors are always smaller than either the background or observational uncertainties.

## Exercise 19.6, page 83

To assimilate the linear combination  $z(t) = (\mathbf{x}_t)_1 + (\mathbf{x}_t)_2$ , we would use the following operators in the observation equation

$$\mathbf{H} = \begin{pmatrix} 1 & 1 & 0 \end{pmatrix}$$
 and  $\mathbf{R}_t = 25.$  (22.91)

The result of assimilating these observations to estimate  $(\mathbf{x}_t)_1$  is shown in fig. 22.21.

### Exercise 19.7, page 85

The mean square error of the prediction model  $\hat{\mathbf{y}} = \mathbf{A}\mathbf{o} + \mathbf{b}$  is

$$MSE = \mathbb{E}\left[ (\mathbf{y} - \mathbf{Ao} - \mathbf{b})^T (\mathbf{y} - \mathbf{Ao} - \mathbf{b}) \right]$$
(22.92)

Taking the derivative with respect to A and b gives

$$\frac{\partial MSE}{\partial \mathbf{A}} = -2\mathbb{E}\left[\left(\mathbf{y} - \mathbf{A}\mathbf{o} - \mathbf{b}\right)\mathbf{o}^{T}\right] = -2\mathbb{E}\left[\mathbf{y}\mathbf{o}^{T} - \mathbf{A}\mathbf{o}\mathbf{o}^{T} - \mathbf{b}\mathbf{o}^{T}\right]$$
(22.93)

$$\frac{\partial MSE}{\partial \mathbf{b}} = -2\mathbb{E}\left[\mathbf{y} - \mathbf{A}\mathbf{o} - \mathbf{b}\right].$$
(22.94)

Setting the second equation to zero and solving yields

$$\mathbf{b} = \mathbb{E}[\mathbf{y}] - \mathbf{A}\mathbb{E}[\mathbf{o}]. \tag{22.95}$$

Substituting this into the first equation yields

$$\frac{\partial \mathsf{MSE}}{\partial \mathbf{A}} = -2\mathbb{E}\left[\mathbf{y}\mathbf{o}^{T} - \mathbf{A}\mathbf{o}\mathbf{o}^{T} - (\mathbb{E}[\mathbf{y}] - \mathbf{A}\mathbb{E}[\mathbf{o}])\mathbf{o}^{T}\right]$$
(22.96)

$$= (-2) \left( \mathbb{E}[\mathbf{y}\mathbf{o}^{T}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{o}]^{T} \right) + 2 \left( \mathbf{A}\mathbb{E}[\mathbf{o}\mathbf{o}^{T}] - \mathbf{A}\mathbb{E}[\mathbf{o}]\mathbb{E}[\mathbf{o}^{T}] \right)$$
(22.97)  
$$= (-2) \left( \mathbf{\Sigma}_{\mathbf{v}} - \mathbf{A}\mathbf{\Sigma}_{\mathbf{v}} \right)$$
(22.98)

$$= (-2) \left( \boldsymbol{\Sigma}_{YO} - \mathbf{A} \boldsymbol{\Sigma}_{OO} \right), \qquad (22.98)$$

where we have used the fact that for any random vectors x and y,

$$\boldsymbol{\Sigma}_{XY} = \mathbb{E}\left[\left(\mathbf{x} - \mathbb{E}[\mathbf{x}]\right)\left(\mathbf{y} - \mathbb{E}[\mathbf{y}]\right)^{T}\right] = \mathbb{E}[\mathbf{x}\mathbf{y}^{T}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]^{T}.$$
(22.99)

Setting the derivative to zero and solving yields

$$\mathbf{A} = \boldsymbol{\Sigma}_{YO} \boldsymbol{\Sigma}_{OO}^{-1}, \tag{22.100}$$

as desired.

Given  $\mathbf{o} = \mathbf{H}\mathbf{y} + \mathbf{r}$ , it follows that

$$\mathbb{E}[\mathbf{o}] = \mathbf{H}\mathbb{E}[\mathbf{y}] + \mathbb{E}[\mathbf{r}] = \mathbf{H}\boldsymbol{\mu}_B$$
(22.101)

$$\mathbb{E}[\mathbf{0}] = \mathbf{H}\mathbb{E}[\mathbf{y}] + \mathbb{E}[\mathbf{1}] = \mathbf{H}\mu_B \qquad (22.101)$$
  

$$\operatorname{cov}[\mathbf{0}, \mathbf{y}] = \operatorname{cov}[\mathbf{H}\mathbf{y} + \mathbf{r}, \mathbf{y}] = \mathbf{H}\operatorname{cov}[\mathbf{y}] + \operatorname{cov}[\mathbf{r}, \mathbf{y}] = \mathbf{H}\Sigma_B \qquad (22.102)$$
  

$$\operatorname{cov}[\mathbf{0}] = \operatorname{cov}[\mathbf{H}\mathbf{y} + \mathbf{r}, \mathbf{H}\mathbf{y} + \mathbf{r}] = \mathbf{H}\operatorname{cov}[\mathbf{y}]\mathbf{H}^T + \operatorname{cov}[\mathbf{r}] = \mathbf{H}\Sigma_B\mathbf{H}^T + \mathbf{R}.$$

$$\frac{1}{(22.103)} = \frac{1}{(22.103)} = \frac{1}$$

where we used  $\mathbb{E}[\mathbf{r}] = \mathbf{0}$  and  $\operatorname{cov}[\mathbf{r}, \mathbf{y}] = \mathbf{0}$ . Substituting these expressions into the least squares estimates for A gives

$$\mathbf{A} = \left(\mathbf{H}\boldsymbol{\Sigma}_B\right)^T \left(\mathbf{H}\boldsymbol{\Sigma}_B\mathbf{H}^T + \mathbf{R}\right)^{-1} = \boldsymbol{\Sigma}_B\mathbf{H}^T \left(\mathbf{H}\boldsymbol{\Sigma}_B\mathbf{H}^T + \mathbf{R}\right)^{-1}.$$
 (22.104)

Substituting the expressions into the least squares estimate for b gives

$$\mathbf{b} = \boldsymbol{\mu}_B - \boldsymbol{\Sigma}_B \mathbf{H}^T \left( \mathbf{H} \boldsymbol{\Sigma}_B \mathbf{H}^T + \mathbf{R} \right)^{-1} \mathbf{H} \boldsymbol{\mu}_B$$
(22.105)

Therefore, the least squares prediction  $\hat{\mathbf{y}}$  is

$$\hat{\mathbf{y}} = \mathbf{A}\mathbf{o} + \mathbf{b} = \boldsymbol{\mu}_B + \boldsymbol{\Sigma}_B \mathbf{H}^T \left(\mathbf{H}\boldsymbol{\Sigma}_B \mathbf{H}^T + \mathbf{R}\right)^{-1} \left(\mathbf{o} - \mathbf{H}\boldsymbol{\mu}_B\right)$$
 (22.106)

## Exercise 20.1, page 88

The symmetric square root requires

$$\mathbf{I} - \beta \mathbf{w} \mathbf{w}^{T} = \left(\mathbf{I} - \delta \beta \mathbf{w} \mathbf{w}^{T}\right) \left(\mathbf{I} - \delta \beta \mathbf{w} \mathbf{w}^{T}\right).$$
(22.107)

Expanding the right hand side yields

$$\mathbf{I} - \beta \mathbf{w} \mathbf{w}^{T} = \mathbf{I} - 2\delta \beta \mathbf{w} \mathbf{w}^{T} + \delta^{2} \beta^{2} \left( \mathbf{w}^{T} \mathbf{w} \right) \mathbf{w} \mathbf{w}^{T}.$$
 (22.108)

Simplifying this equation yields

$$\left(\delta^2 \beta \left(\mathbf{w}^T \mathbf{w}\right) - 2\delta\beta + \beta\right) \mathbf{w} \mathbf{w}^T = 0.$$
(22.109)

This equation is satisfied only for the choices

$$\delta = \frac{1 \pm \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}{\beta \mathbf{w}^T \mathbf{w}}.$$
(22.110)

Multiplying the top and bottom by the conjugate gives

$$\delta = \frac{1 \pm \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}{\beta \mathbf{w}^T \mathbf{w}} \frac{1 \mp \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}{1 \mp \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}$$
$$= \frac{1 - (1 - \beta \mathbf{w}^T \mathbf{w})}{\beta \mathbf{w}^T \mathbf{w}} \frac{1}{1 \mp \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}$$
$$= \frac{1}{1 \mp \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}$$
(22.111)

To obtain a positive semi-definite square root matrix, we chose  $\delta$  to render

$$\mathbf{x}^{T} \left( \mathbf{I} - \delta \beta \mathbf{w} \mathbf{w}^{T} \right) \mathbf{x} \ge 0$$
 for all  $x$ . (22.112)

Expanding the left hand side yields

$$\mathbf{x}^{T} \left( \mathbf{I} - \delta \beta \mathbf{w} \mathbf{w}^{T} \right) \mathbf{x} = \mathbf{x}^{T} \mathbf{x} - \delta \beta \left( \mathbf{x}^{T} \mathbf{w} \right)^{2}.$$
 (22.113)

The Cauch-Schwartz inequality implies

$$\epsilon = \frac{\left(\mathbf{x}^T \mathbf{w}\right)^2}{\left(\mathbf{x}^T \mathbf{x}\right) \left(\mathbf{w}^T \mathbf{w}\right)} \le 0.$$
(22.114)

Therefore, we have

$$\mathbf{x}^{T}\mathbf{x} - \delta\beta \left(\mathbf{x}^{T}\mathbf{w}\right)^{2} = \left(\mathbf{x}^{T}\mathbf{x}\right) \left(\mathbf{w}^{T}\mathbf{w}\right) \left(\frac{1}{\mathbf{w}^{T}\mathbf{w}} - \delta\beta\epsilon\right)$$
$$= \left(\mathbf{x}^{T}\mathbf{x}\right) \left(1 - \left(1 \pm \sqrt{1 - \beta\mathbf{w}^{T}\mathbf{w}}\right)\epsilon\right)$$
$$= \left(\mathbf{x}^{T}\mathbf{x}\right) \left(1 - \epsilon \mp \epsilon\sqrt{1 - \beta\mathbf{w}^{T}\mathbf{w}}\right)$$
(22.115)

Since  $\beta \leq \mathbf{w}^T \mathbf{w}$ , the term in the radical is between 0 and 1. Since  $\epsilon$  also is between 0 and 1, the quadratic form is positive for all  $\epsilon$ , including  $\epsilon = 1$ , only if  $\delta$  is chosen as

$$\delta = \frac{1 - \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}{\beta \mathbf{w}^T \mathbf{w}} = \frac{1}{1 + \sqrt{1 - \beta \mathbf{w}^T \mathbf{w}}}$$
(22.116)

### Exercise 20.2, page 88

*Proof.* First consider the case in which **B** is rank-1, which implies that  $\mathbf{B} = \mathbf{u}\mathbf{u}^T$  for some vector **u**. Then

$$\mathbf{C} \circ \mathbf{B} = \mathbf{C} \circ \mathbf{u}\mathbf{u}^T = C_{ij}u_iu_j = \mathbf{D}_u\mathbf{C}\mathbf{D}_u, \qquad (22.117)$$

where  $D_u$  is a diagonal matrix whose diagonal elements equal the vector elements of **u**. The last matrix on the right hand side is positive semi-definite, as shown below:

$$\mathbf{x}^{T} \left( \mathbf{C} \circ \mathbf{B} \right) \mathbf{x} = \mathbf{x}^{T} \mathbf{D}_{u} \mathbf{C} \mathbf{D}_{u} \mathbf{x} = \mathbf{y}^{T} \mathbf{C} \mathbf{y} \ge 0,$$
(22.118)

where we have defined  $\mathbf{y} = \mathbf{D}_u \mathbf{x}$ . Since the above inequality holds for arbitrary  $\mathbf{x}$ ,  $\mathbf{C} \circ \mathbf{B}$  is positive semi-definite. Now consider rank greater than 1. By the spectral factorization theorem,  $\mathbf{B} = \sum \mathbf{u}_i \mathbf{u}_i^T$  for some set of vectors  $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N$ . Then

$$\mathbf{C} \circ \mathbf{B} = \mathbf{C} \circ \sum_{i} \mathbf{u}_{i} \mathbf{u}_{i}^{T} = \mathbf{C} \circ \mathbf{u}_{1} \mathbf{u}_{1}^{T} + \mathbf{C} \circ \mathbf{u}_{2} \mathbf{u}_{2}^{T} + \dots + \mathbf{C} \circ \mathbf{u}_{N} \mathbf{u}_{N}^{T}.$$
 (22.119)

Each term is positive semi-definite, so therefore the sum is too. QED

Exercise 21.1, page 90

```
> cold.x
1
2
  fevd(x = t.cold)
3
4
   . . .
  Estimated parameters:
5
    location scale
6
                            shape
   6.5115237 3.5662767 -0.2347719
7
8
9
   Standard Error Estimates:
10
    location scale
                          shape
  0.58044433 0.40333521 0.09324019
11
12
  Estimated parameter covariance matrix.
13
14
              location scale
                                         shape
  location 0.33691562 0.01740978 -0.020276953
15
16
  scale 0.01740978 0.16267929 -0.019456804
  shape
         -0.02027695 -0.01945680 0.008693734
17
```

The return levels are

```
> return.level(cold.x)
1
  fevd(x = t.cold)
2
  get(paste("return.level.fevd.", newcl, sep = ""))(x = x, return.period = return.period)
3
4
   GEV model fitted to t.cold
5
  Data are assumed to be stationary
6
  [1] "Return Levels for period units in years"
7
    2-year level 20-year level 100-year level
8
        7.763954
                      14.138334
                                     16.543253
```

Exercise 21.2, page 90

The shape parameter is

$$\xi = -0.23 \pm 0.09, \tag{22.120}$$

so it is negative even for 2 standard errors, which means it is a Weibull distribution.

Exercise 21.3, page 90

```
> cold.x.tele
1
2
  fevd(x = t.cold, location.fun = ~index)
3
4
  [1] "Estimation Method used: MLE"
5
6
7
   . . .
   Estimated parameters:
8
    mu0 mu1 scale shape
9
   6.4879640 -1.4483469 3.0933949 -0.1545121
10
11
12
   Standard Error Estimates:
13
    mu0 mu1 scale shape
  0.5231300 0.4817077 0.3851407 0.1244832
14
15
   Estimated parameter covariance matrix.
16
    mu0 mu1 scale
17
                                           shape
        0.27366503 0.03729961 0.05315883 -0.02825529
  mu0
18
  mul 0.03729961 0.23204226 0.04143058 -0.02199465
19
  scale 0.05315883 0.04143058 0.14833340 -0.02640878
20
  shape -0.02825529 -0.02199465 -0.02640878 0.01549606
21
```

The coefficient for the location parameter is

$$\mu = -1.45 \pm 0.48. \tag{22.121}$$

This is non-zero even for 2 standard errors, so it is significant. A more accurate test is the likelihood ratio test:

```
1 > lr.test(cold.x,cold.x.tele)
2
3 Likelihood-ratio Test
4
4
4
5 data: t.coldt.cold
6 Likelihood-ratio = 8.2806, chi-square critical value = 3.8415,
7 alpha = 0.0500, Degrees of Freedom = 1.0000, p-value = 0.004007
8 alternative hypothesis: greater
```

The p-value is less than 5%, hence significant.

The shape parameter is not significantly different from zero or positive at the 5% level, so it could be any of the 3 distributions.

## Exercise 21.4, page 90

The probability that an exceedance occurs in year 1 is p. The probability that an exceedance occurs in year 2 is the probability that it did not occur in year 1, times the probability that it did occur in year 2, which gives (1-p)p. Similarly, exceedance in year 3 has probability  $(1-p)^2p$ , and exceedance in year y has probability  $p(1-p)^{y-1}$ . Thus, the average waiting

time is

return period = 
$$p \sum_{y=1}^{\infty} y (1-p)^{y-1}$$
. (22.122)

To evaluate this sum, recall the sum of an infinite geometric series is

$$\sum_{k=0}^{\infty} q^k = \frac{1}{1-q} \quad \text{for } |q| < 1.$$
(22.123)

Differentiating this equation with respect to q yields the identity

$$\sum_{k=0}^{\infty} kq^{k-1} = \frac{1}{(1-q)^2} \quad \text{for } |q| < 1.$$
(22.124)

This identity gives the sum in (22.122) in the case q = 1 - p, thus

return period = 
$$p \frac{1}{(1 - (1 - p))^2} = \frac{1}{p}$$
. (22.125)

Thus, if the probability of exceeding  $x_p$  in any given year is p, then the return period is 1/p years.



**Figure 22.20** A particular realization of the first (top) and second (bottom) elements of  $\mathbf{y}(t)$  generated by the vector autoregressive model (solid curve), corresponding observations (dots), and the analysis mean plus a standard deviation, as estimated state by the Kalman Filter (grey shading), but using  $\mathbf{R}_t = 1$ 



**Figure 22.21** A particular realization of the first (top) and second (bottom) elements of  $\mathbf{y}(t)$  generated by the vector autoregressive model (solid curve), corresponding observations (dots), and the analysis mean plus a standard deviation, as estimated state by the Kalman Filter (grey shading), but using  $\mathbf{R}_t = 25$ , and observing the sum of the first two elements.