

Chapter 7: Categorical responses

7.1 For the degree of crash data, regard the event of interest as the occurrence of injury (fatal or non-fatal) in a crash. Develop a statistical model.

Firstly create a new variable `injured`, with `injured=0` when `degree=1` (non-casualty) and `injured=1` when `degree=2` (injury) or `degree=3` (fatal):

```
data injury;
set act.injury;
if degree=1 then injured=0;
else injured=1;
run;
```

The same variables as are significant for the ordinal regression, are significant here:

```
proc logistic data=injury descending;
class agecat sex roaduserclass/ param=ref;
model injured = roaduserclass agecat sex agecat*sex;
weight number;
run;
```

The LOGISTIC Procedure

Model Information

Data Set	WORK.INJURY	
Response Variable	injured	
Number of Response Levels	2	
Weight Variable	Number	Number
Model	binary logit	
Optimization Technique	Fisher's scoring	

Number of Observations Read	209
Number of Observations Used	209
Sum of Weights Read	76341
Sum of Weights Used	76341

Response Profile

Ordered Value	injured	Total Frequency	Total Weight
1	1	132	32045.000
2	0	77	44296.000

Probability modeled is injured=1.

Class Level Information

Class	Value	Design Variables					
agecat	1	1	0	0	0	0	0
	2	0	1	0	0	0	0
	3	0	0	1	0	0	0
	5	0	0	0	1	0	0
	6	0	0	0	0	1	0
	7	0	0	0	0	0	1
	10	0	0	0	0	0	0

Sex	F	1		
	M	0		
roaduserclass	2	1	0	0
	4	0	1	0
	6	0	0	1
	10	0	0	0

Model Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

Model Fit Statistics

Criterion	Intercept Only	Intercept and Covariates
AIC	103858.56	101068.62
SC	103861.91	101125.44
-2 Log L	103856.56	101034.62

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	2821.9414	16	<.0001
Score	2609.8781	16	<.0001
Wald	1623.5450	16	<.0001

Type 3 Analysis of Effects

Effect	DF	Wald	
		Chi-Square	Pr > ChiSq
roaduserclass	3	1404.9916	<.0001
agecat	6	114.9136	<.0001
Sex	1	27.1318	<.0001
agecat*Sex	6	24.4544	0.0004

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	-0.4683	0.0213	481.1289	<.0001
roaduserclass 2	1	0.1406	0.0270	27.0427	<.0001
roaduserclass 4	1	0.2558	0.0367	48.6248	<.0001
roaduserclass 6	1	2.8668	0.0776	1365.3552	<.0001
agecat 1	1	-0.1872	0.0326	33.0242	<.0001
agecat 2	1	-0.1169	0.0324	12.9869	0.0003
agecat 3	1	-0.0631	0.0368	2.9392	0.0865
agecat 5	1	0.0551	0.0303	3.3163	0.0686
agecat 6	1	0.0667	0.0334	3.9926	0.0457
agecat 7	1	0.1375	0.0347	15.7179	<.0001
Sex F	1	0.1751	0.0336	27.1318	<.0001
agecat*Sex 1 F	1	0.1353	0.0534	6.4174	0.0113
agecat*Sex 2 F	1	0.1236	0.0525	5.5423	0.0186
agecat*Sex 3 F	1	0.0509	0.0605	0.7065	0.4006
agecat*Sex 5 F	1	0.0283	0.0490	0.3338	0.5634
agecat*Sex 6 F	1	0.0162	0.0559	0.0844	0.7714
agecat*Sex 7 F	1	-0.1406	0.0598	5.5360	0.0186

Odds Ratio Estimates

Point	95% Wald
-------	----------

Effect	Estimate	Confidence Limits
roaduserclass 2 vs 10	1.151	1.092 1.214
roaduserclass 4 vs 10	1.292	1.202 1.388
roaduserclass 6 vs 10	17.581	15.101 20.469

Association of Predicted Probabilities and Observed Responses

Percent Concordant	48.1	Somers' D	0.002
Percent Discordant	48.0	Gamma	0.002
Percent Tied	3.9	Tau-a	0.001
Pairs	10164	c	0.501

The striking result is that the odds of an injury for a motor cycle driver (road user class 6) is 17.6 times the odds of an injury for a car driver (road user class 10).

7.2 For the personal injury data, develop a statistical model for the occurrence of legal representation.

It only makes sense to include those variables that are known at the time of claim as predictors for legal representation. One may find, for example, that operational time is significant, but this would be as a result of the fact that claims being legally contested will take longer to settle.

In this data set, injury code is therefore the only available explanatory variable for legal representation. We could include accident year as an explanatory variable, in order to assess whether there is a trend in legal representation over time, but legal representation for the most recent claims will not have happened by the time the data was recorded. In practice, other aspects of the accident or claimant may be predictive. We examine legal representation by injury code:

```
proc freq data=act.persinj;
tables legrep*inj1 / norow nopercent;
run;
```

Table of LEGREP by INJ1

LEGREP(LEGREP)	INJ1(INJ1)							
Frequency	1	2	3	4	5	6	9	Total
Col Pct								
0	5571 35.62	1152 34.12	374 33.01	56 29.63	85 45.21	121 47.27	649 51.67	8008
1	10067 64.38	2224 65.88	759 66.99	133 70.37	103 54.79	135 52.73	607 48.33	14028
Total	15638	3376	1133	189	188	256	1256	22036

Overall the percentage of legal representation is 64% (14028/22036). This varies from 48% (injury code 9) to 70% (injury code 4). As injury code 1 is by far the most frequent, it is declared as the base level:

```
proc genmod data=act.persinj descending;
class inj1 (ref="1") / param=ref;
model legrep = inj1 / dist=bin type3;
run;
```

The GENMOD Procedure

Model Information

Data Set	ACT.PERSINJ	
Distribution	Binomial	
Link Function	Logit	
Dependent Variable	LEGREP	LEGREP

Number of Observations Read	22036
Number of Observations Used	22036
Number of Events	14028
Number of Trials	22036

Class Level Information

Class	Value	Design Variables						
INJ1	1	0	0	0	0	0	0	0
	2	1	0	0	0	0	0	0
	3	0	1	0	0	0	0	0
	4	0	0	1	0	0	0	0
	5	0	0	0	1	0	0	0
	6	0	0	0	0	1	0	0
	9	0	0	0	0	0	0	1

Response Profile

Ordered Value	LEGREP	Total Frequency
1	1	14028
2	0	8008

PROC GENMOD is modeling the probability that LEGREP='1'.

Criteria For Assessing Goodness Of Fit

Criterion	DF	Value	Value/DF
Deviance	22E3	28721.3092	1.3038
Scaled Deviance	22E3	28721.3092	1.3038
Pearson Chi-Square	22E3	22035.9999	1.0003
Scaled Pearson X2	22E3	22035.9999	1.0003
Log Likelihood		-14360.6546	

Algorithm converged.

Analysis Of Parameter Estimates

Parameter	DF	Estimate	Standard Error	Wald 95% Confidence Limits		Chi-Square	Pr > ChiSq	
Intercept	1	0.5917	0.0167	0.5590	0.6244	1255.56	<.0001	
INJ1	2	1	0.0661	-0.0122	0.1444	2.74	0.0980	
INJ1	3	1	0.1161	-0.0120	0.2441	3.15	0.0757	
INJ1	4	1	0.2733	-0.0406	0.5872	2.91	0.0879	
INJ1	5	1	-0.3996	0.1475	-0.6887	-0.1105	7.34	0.0067
INJ1	6	1	-0.4822	0.1263	-0.7297	-0.2347	14.58	0.0001
INJ1	9	1	-0.6586	0.0589	-0.7740	-0.5432	125.10	<.0001

```
Scale          0          1.0000          0.0000          1.0000          1.0000
```

NOTE: The scale parameter was held fixed.

```
LR Statistics For Type 3 Analysis

Source          DF          Chi-Square      Pr > ChiSq
INJ1              6          161.38          <.0001
```

Injury code is highly significant (p -value < 0.0001).

7.3 The RTA also publish driver crash statistics by age, sex and blood alcohol concentration. This data set is provided on the companion website. Develop

(a) an ordinal regression model; and

(b) a nominal regression model

for degree of crash.

(a) Ordinal regression model Perform preparatory data manipulations and ordinal regression:

```
data bac;
set act.bac;
if number=0 or number=. then delete; /*take out zero frequencies*/
agecat=age;
if age=8 then agecat=7 ; /* Combine 60-69 and 70+*/
if age=4 then agecat=10; /* Make age 30-39 base level*/
if bac=1 then bac=10; /* Make bac=1 (legal) base level*/
run;

proc logistic data=bac;
class agecat sex bac / param=ref;
model degree = bac agecat sex agecat*sex ;
weight number;
run;
```

The interaction term `bac*agecat` is marginally significant (p -value = 0.0410) and has been omitted.

The LOGISTIC Procedure

Model Information

Data Set	WORK.BAC	
Response Variable	Degree	Degree
Number of Response Levels	3	
Weight Variable	Number	Number
Model	cumulative logit	
Optimization Technique	Fisher's scoring	
Number of Observations Read	174	
Number of Observations Used	174	
Sum of Weights Read	58890	
Sum of Weights Used	58890	

Response Profile

Ordered Value	Degree	Total Frequency	Total Weight
1	1	63	34911.000
2	2	69	23351.000
3	3	42	628.000

Probabilities modeled are cumulated over the lower Ordered Values.

Class Level Information

Class	Value	Design Variables					
agecat	1	1	0	0	0	0	0
	2	0	1	0	0	0	0
	3	0	0	1	0	0	0
	5	0	0	0	1	0	0
	6	0	0	0	0	1	0
	7	0	0	0	0	0	1
	10	0	0	0	0	0	0
Sex	f	1					
	m	0					
BAC	2	1	0	0	0		
	3	0	1	0	0		
	4	0	0	1	0		
	5	0	0	0	1		
	10	0	0	0	0		

Model Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

Score Test for the Proportional Odds Assumption

Chi-Square	DF	Pr > ChiSq
92.2948	17	<.0001

Model Fit Statistics

Criterion	Intercept Only	Intercept and Covariates
AIC	85415.956	84952.633
SC	85422.274	85012.655
-2 Log L	85411.956	84914.633

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	497.3236	17	<.0001
Score	495.6705	17	<.0001
Wald	484.8243	17	<.0001

Type 3 Analysis of Effects

Effect	DF	Wald Chi-Square	Pr > ChiSq
BAC	4	385.3662	<.0001
agecat	6	88.9992	<.0001
Sex	1	0.3955	0.5294

agecat*Sex 6 23.9044 0.0005

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	0.3970	0.0230	297.5092	<.0001
Intercept	2	4.5721	0.0456	10071.2059	<.0001
BAC	2	-1.0029	0.2782	12.9949	0.0003
BAC	3	-0.9437	0.1381	46.7086	<.0001
BAC	4	-0.6432	0.0730	77.7037	<.0001
BAC	5	-1.2249	0.0764	256.9530	<.0001
agecat	1	0.2212	0.0356	38.5304	<.0001
agecat	2	0.1200	0.0360	11.1252	0.0009
agecat	3	0.0454	0.0410	1.2258	0.2682
agecat	5	-0.0344	0.0339	1.0288	0.3104
agecat	6	-0.0548	0.0375	2.1367	0.1438
agecat	7	-0.0806	0.0385	4.3840	0.0363
Sex	f	0.0240	0.0381	0.3955	0.5294
agecat*Sex 1 f	1	-0.2374	0.0599	15.7105	<.0001
agecat*Sex 2 f	1	-0.1955	0.0598	10.6708	0.0011
agecat*Sex 3 f	1	-0.0591	0.0694	0.7244	0.3947
agecat*Sex 5 f	1	-0.0948	0.0559	2.8748	0.0900
agecat*Sex 6 f	1	-0.1082	0.0637	2.8842	0.0895
agecat*Sex 7 f	1	0.00126	0.0671	0.0004	0.9851

Odds Ratio Estimates

Effect	Point Estimate	95% Wald Confidence Limits
BAC 2 vs 10	0.367	0.213 0.633
BAC 3 vs 10	0.389	0.297 0.510
BAC 4 vs 10	0.526	0.456 0.606
BAC 5 vs 10	0.294	0.253 0.341

Association of Predicted Probabilities and Observed Responses

Percent Concordant	45.2	Somers' D	-.064
Percent Discordant	51.6	Gamma	-.066
Percent Tied	3.2	Tau-a	-.042
Pairs	9891	c	0.468

As the test for proportional odds is rejected (p -value < 0.0001), the partial proportional odds model (Section 7.8) or the nominal regression model should be fitted.

Partial proportional odds model

```
data bac2; set bac;
do i=1 to number;output;end;drop i;run;

data bac2; set bac2;
id=_n_;
do; if degree>2 then y=1; else y=0; degreej=2;output;end;
do; if degree>1 then y=1; else y=0; degreej=1;output;end;
run;

proc sort; by id; run;

proc genmod data=bac2;
class agecat bac sex degreej id;
model y = agecat sex bac agecat*sex degreej agecat*degreej sex*degreej
      bac*degreej agecat*sex*degreej
```

```

      / dist=bin link=logit type3;
repeated subject=id / type=unstr;
run;

```

Examine the Type 3 tests:

Score Statistics For Type 3 GEE Analysis

Source	DF	Chi-Square	Pr > ChiSq
agecat	6	21.97	0.0012
Sex	1	23.25	<.0001
BAC	4	64.51	<.0001
agecat*Sex	6	5.17	0.5222
degreej	1	33.02	<.0001
agecat*degreej	6	9.73	0.1363
Sex*degreej	1	48.52	<.0001
BAC*degreej	4	17.09	0.0019
agecat*Sex*degreej	6	2.34	0.8857

As the interaction term **agecat*Sex*degreej** is not significant, remove it from the model and re-estimate:

```

proc genmod data=bac2;
class agecat bac sex degreej id;
model y = agecat sex bac agecat*sex degreej agecat*degreej sex*degreej bac*degreej
      / dist=bin link=logit type3;
repeated subject=id / type=unstr;
run;

```

Score Statistics For Type 3 GEE Analysis

Source	DF	Chi-Square	Pr > ChiSq
agecat	6	33.10	<.0001
Sex	1	20.92	<.0001
BAC	4	64.55	<.0001
agecat*Sex	6	23.39	0.0007
degreej	1	33.02	<.0001
agecat*degreej	6	14.93	0.0208
Sex*degreej	1	46.41	<.0001
BAC*degreej	4	17.18	0.0018

The proportional odds assumption is violated for age (p -value = 0.0208), sex (p -value < 0.0001) and BAC (p -value = 0.0018). This model is therefore retained. The table of parameter estimates is

Analysis Of GEE Parameter Estimates
Empirical Standard Error Estimates

Parameter		Estimate	Standard Error	95% Confidence Limits		Z	Pr > Z
Intercept		4.5137	0.0957	4.3261	4.7013	47.16	<.0001
agecat	1	0.2093	0.1463	-0.0775	0.4961	1.43	0.1526
agecat	2	0.2680	0.1493	-0.0246	0.5606	1.80	0.0726
agecat	3	0.0843	0.1634	-0.2360	0.4046	0.52	0.6061
agecat	5	-0.0303	0.1327	-0.2903	0.2297	-0.23	0.8195
agecat	6	-0.2613	0.1401	-0.5359	0.0132	-1.87	0.0621
agecat	7	-0.4658	0.1386	-0.7375	-0.1941	-3.36	0.0008
agecat	10	0.0000	0.0000	0.0000	0.0000	.	.
Sex	f	0.6118	0.1024	0.4111	0.8125	5.97	<.0001
Sex	m	0.0000	0.0000	0.0000	0.0000	.	.
BAC	2	-0.7622	1.0227	-2.7666	1.2422	-0.75	0.4561

BAC	3		-1.5874	0.3279	-2.2300	-0.9447	-4.84	<.0001
BAC	4		-1.1280	0.2218	-1.5627	-0.6932	-5.08	<.0001
BAC	5		-1.8472	0.1608	-2.1623	-1.5321	-11.49	<.0001
BAC	10		0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	1	f	-0.2366	0.0599	-0.3540	-0.1192	-3.95	<.0001
agecat*Sex	1	m	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	2	f	-0.1950	0.0598	-0.3123	-0.0777	-3.26	0.0011
agecat*Sex	2	m	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	3	f	-0.0594	0.0694	-0.1954	0.0765	-0.86	0.3916
agecat*Sex	3	m	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	5	f	-0.0963	0.0560	-0.2061	0.0134	-1.72	0.0853
agecat*Sex	5	m	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	6	f	-0.1115	0.0638	-0.2365	0.0135	-1.75	0.0804
agecat*Sex	6	m	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	7	f	-0.0038	0.0671	-0.1353	0.1276	-0.06	0.9546
agecat*Sex	7	m	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	10	f	0.0000	0.0000	0.0000	0.0000	.	.
agecat*Sex	10	m	0.0000	0.0000	0.0000	0.0000	.	.
degreej	1		-4.1151	0.0949	-4.3010	-3.9291	-43.37	<.0001
degreej	2		0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	1	1	0.0108	0.1447	-0.2729	0.2944	0.07	0.9406
agecat*degreej	1	2	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	2	1	-0.1516	0.1478	-0.4414	0.1381	-1.03	0.3051
agecat*degreej	2	2	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	3	1	-0.0406	0.1613	-0.3568	0.2756	-0.25	0.8014
agecat*degreej	3	2	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	5	1	-0.0041	0.1313	-0.2615	0.2533	-0.03	0.9752
agecat*degreej	5	2	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	6	1	0.2113	0.1384	-0.0600	0.4826	1.53	0.1268
agecat*degreej	6	2	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	7	1	0.3952	0.1368	0.1270	0.6634	2.89	0.0039
agecat*degreej	7	2	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	10	1	0.0000	0.0000	0.0000	0.0000	.	.
agecat*degreej	10	2	0.0000	0.0000	0.0000	0.0000	.	.
Sex*degreej	f	1	-0.5967	0.0961	-0.7851	-0.4083	-6.21	<.0001
Sex*degreej	f	2	0.0000	0.0000	0.0000	0.0000	.	.
Sex*degreej	m	1	0.0000	0.0000	0.0000	0.0000	.	.
Sex*degreej	m	2	0.0000	0.0000	0.0000	0.0000	.	.
BAC*degreej	2	1	-0.2577	1.0296	-2.2756	1.7603	-0.25	0.8024
BAC*degreej	2	2	0.0000	0.0000	0.0000	0.0000	.	.
BAC*degreej	3	1	0.6994	0.3323	0.0481	1.3507	2.10	0.0353
BAC*degreej	3	2	0.0000	0.0000	0.0000	0.0000	.	.
BAC*degreej	4	1	0.5103	0.2227	0.0739	0.9467	2.29	0.0219
BAC*degreej	4	2	0.0000	0.0000	0.0000	0.0000	.	.
BAC*degreej	5	1	0.7083	0.1664	0.3822	1.0344	4.26	<.0001
BAC*degreej	5	2	0.0000	0.0000	0.0000	0.0000	.	.
BAC*degreej	10	1	0.0000	0.0000	0.0000	0.0000	.	.
BAC*degreej	10	2	0.0000	0.0000	0.0000	0.0000	.	.

Model equations are constructed as

$$\ln \frac{\hat{\pi}_j}{1 - \hat{\pi}_j} = \begin{cases} 0.3986 + 0.2201x_1 + \cdots - 1.1389x_{17} & \text{for } j = 1 \text{ (non-casualty)} \\ 4.5137 + 0.2093x_1 + \cdots - 1.8472x_{17} & \text{for } j = 2 \text{ (injury)} \end{cases}$$

The effect of BAC is shown in the table. $\exp(\hat{\beta}_1)$ is the effect of the BAC level (compared to legal), on the odds of a crash being non-casualty; and $\exp(\hat{\beta}_2)$ is the effect of the BAC level (compared to legal), on the odds of a crash being non-casualty or injury.

BAC	$\hat{\beta}_1$	$\hat{\beta}_2$	$\exp(\hat{\beta}_1)$	$\exp(\hat{\beta}_2)$
.020-.049	-0.7622	-1.0199	0.47	0.77
.050-.079	-1.5874	-0.888	0.20	2.01
.080-.149	-1.128	-0.6177	0.32	1.67
$\geq .150$	-1.8472	-1.1389	0.16	2.03

(b) Nominal regression model

```
proc logistic data=bac;
class agecat sex bac / param=ref;
model degree (ref=first) = bac agecat sex agecat*sex / link=glogit;
weight number;
run;
```

The same explanatory variables as for ordinal regression are significant.

The LOGISTIC Procedure

Model Information

Data Set	WORK.BAC	
Response Variable	Degree	Degree
Number of Response Levels	3	
Weight Variable	Number	Number
Model	generalized logit	
Optimization Technique	Fisher's scoring	

Number of Observations Read	174
Number of Observations Used	174
Sum of Weights Read	58890
Sum of Weights Used	58890

Response Profile

Ordered Value	Degree	Total Frequency	Total Weight
1	1	63	34911.000
2	2	69	23351.000
3	3	42	628.000

Logits modeled use Degree=1 as the reference category.

Class Level Information

Class	Value	Design Variables					
agecat	1	1	0	0	0	0	0
	2	0	1	0	0	0	0
	3	0	0	1	0	0	0
	5	0	0	0	1	0	0
	6	0	0	0	0	1	0
	7	0	0	0	0	0	1
	10	0	0	0	0	0	0
Sex	f	1					
	m	0					
BAC	2	1	0	0	0		
	3	0	1	0	0		
	4	0	0	1	0		
	5	0	0	0	1		
	10	0	0	0	0		

Model Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

Model Fit Statistics

Criterion	Intercept Only	Intercept and Covariates
-----------	----------------	--------------------------

AIC	85415.956	84897.512
SC	85422.274	85011.238
-2 Log L	85411.956	84825.512

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	586.4445	34	<.0001
Score	686.1817	34	<.0001
Wald	601.4194	34	<.0001

Type 3 Analysis of Effects

Effect	DF	Wald Chi-Square	Pr > ChiSq
BAC	8	434.6839	<.0001
agecat	12	103.7009	<.0001
Sex	2	6.2950	0.0430
agecat*Sex	12	24.5469	0.0171

Analysis of Maximum Likelihood Estimates

Parameter	Degree	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	2	1	-0.4245	0.0233	332.3616	<.0001
Intercept	3	1	-4.0295	0.1054	1462.7506	<.0001
BAC 2	2	1	1.0145	0.2858	12.6004	0.0004
BAC 2	3	1	1.2716	1.0277	1.5311	0.2160
BAC 3	2	1	0.8369	0.1421	34.6840	<.0001
BAC 3	3	1	1.9881	0.3389	34.4142	<.0001
BAC 4	2	1	0.5892	0.0746	62.3733	<.0001
BAC 4	3	1	1.3942	0.2257	38.1527	<.0001
BAC 5	2	1	1.0735	0.0790	184.6167	<.0001
BAC 5	3	1	2.4065	0.1691	202.4469	<.0001
agecat 1	2	1	-0.2205	0.0360	37.4249	<.0001
agecat 1	3	1	-0.2789	0.1686	2.7370	0.0980
agecat 2	2	1	-0.1125	0.0364	9.5614	0.0020
agecat 2	3	1	-0.2995	0.1715	3.0524	0.0806
agecat 3	2	1	-0.0495	0.0415	1.4250	0.2326
agecat 3	3	1	-0.00579	0.1809	0.0010	0.9744
agecat 5	2	1	0.0341	0.0343	0.9856	0.3208
agecat 5	3	1	0.0440	0.1530	0.0829	0.7734
agecat 6	2	1	0.0414	0.0380	1.1862	0.2761
agecat 6	3	1	0.2912	0.1597	3.3248	0.0682
agecat 7	2	1	0.0532	0.0391	1.8565	0.1730
agecat 7	3	1	0.5194	0.1554	11.1748	0.0008
Sex f	2	1	-0.00531	0.0385	0.0190	0.8903
Sex f	3	1	-0.5297	0.2112	6.2922	0.0121
agecat*Sex 1 f	2	1	0.2370	0.0604	15.3813	<.0001
agecat*Sex 1 f	3	1	0.2596	0.3397	0.5840	0.4448
agecat*Sex 2 f	2	1	0.1921	0.0604	10.1274	0.0015
agecat*Sex 2 f	3	1	0.2269	0.3455	0.4312	0.5114
agecat*Sex 3 f	2	1	0.0748	0.0700	1.1445	0.2847
agecat*Sex 3 f	3	1	-0.4521	0.4388	1.0612	0.3029
agecat*Sex 5 f	2	1	0.0950	0.0565	2.8293	0.0926
agecat*Sex 5 f	3	1	0.1332	0.3036	0.1925	0.6608
agecat*Sex 6 f	2	1	0.1167	0.0644	3.2815	0.0701
agecat*Sex 6 f	3	1	0.0927	0.3282	0.0797	0.7777
agecat*Sex 7 f	2	1	0.0165	0.0679	0.0590	0.8080
agecat*Sex 7 f	3	1	-0.1309	0.3381	0.1500	0.6985

Odds Ratio Estimates

Effect	Degree	Point Estimate	95% Wald Confidence Limits
--------	--------	-------------------	-------------------------------

BAC	2	vs 10	2	2.758	1.575	4.829
BAC	2	vs 10	3	3.567	0.476	26.731
BAC	3	vs 10	2	2.309	1.748	3.051
BAC	3	vs 10	3	7.302	3.758	14.187
BAC	4	vs 10	2	1.802	1.557	2.086
BAC	4	vs 10	3	4.032	2.590	6.275
BAC	5	vs 10	2	2.926	2.506	3.416
BAC	5	vs 10	3	11.095	7.965	15.456

The effect of blood alcohol concentration is clear from the odds ratio estimates.

7.4 *In the Enterprise Miner data set, develop a statistical model for the occurrence of a claim.*

Variables in this data set are examined thoroughly in the document “SAS Miner preliminary analysis”. The following transformations are carried out prior to analysis:

```
data claims;
set exercise.claims_sas_miner;

    bluebk = bluebook/1000-14.2; /* more stable for computation*/

    if npolicy=1 then npolicy1=1; else npolicy1=0; /*dichotomise loyalty variable*/

* Categorize income because of zero spike;
  if income=. then incomecat=.;
  else if income=0 then incomecat=99; /* base level*/
  else if income<=25000 then incomecat=2;
  else if income<=50000 then incomecat=3;
  else if income<=75000 then incomecat=4;
  else if income<=100000 then incomecat=5;
  else incomecat=6;

* Categorize oldclaim because of zero spike;
  if oldclaim=. then oldclaimcat=.;
  else if oldclaim=0 then oldclaimcat=99; /* base level*/
  else if oldclaim<=5000 then oldclaimcat=2;
  else if oldclaim<=10000 then oldclaimcat=3;
  else oldclaimcat=4;

* Categorize mvr_pts;
  if mvr_pts=. then mvrcat=.;
  else if mvr_pts=0 then mvrcat=99; /* base level*/
  else if mvr_pts=1 then mvrcat=1;
  else if mvr_pts=2 then mvrcat=2;
  else if mvr_pts=3 then mvrcat=3;
  else if mvr_pts<=6 then mvrcat=4;
  else mvrcat=5;

run;
```

After observing the diagrams and tables, we select the variables which look promising as explanatory variables for `claim_flag`. Fitting these explanatory variables in a binomial model with logit link, we then observe the Type 3 (LR) test results. If the variable is not significant in the type 3 test, it is eliminated from the model. The continuous variables `kidsdriv`, `bluebk`, `retained`, `travtime` and `age` all enter the model in quadratic form. An alternative strategy is to band these variables.

Note that we have observed that home value is associated with income. Even though those two variables are significant in type 3 test, we eliminate home value from the model. The final model is:

```
proc logistic data=claims descending;
class  car_use car_type oldclaimcat revoked married parent1 max_educ density
      mvrcat incomecat / param=glm;
model clm_flag =
  kidsdriv kidsdriv*kidsdriv
  car_use
  bluebk bluebk*bluebk
  retained retained*retained
  car_type
  oldclaimcat
  revoked
  mvrcat
  incomecat
  married
  parent1
  max_educ
  density
  travtime travtime*travtime
  age age*age
  / outroc=claimsroc;
run;
```

The LOGISTIC Procedure

Model Information

Data Set	WORK.CLAIMS	
Response Variable	CLM_FLAG	Claim Happens
Number of Response Levels	2	
Model	binary logit	
Optimization Technique	Fisher's scoring	

Number of Observations Read	10303
Number of Observations Used	9727

Response Profile

Ordered Value	CLM_FLAG	Total Frequency
1	Yes	2601
2	No	7126

Probability modeled is CLM_FLAG='Yes'.

NOTE: 576 observations were deleted due to missing values for the response or explanatory variables.

Class Level Information

Class	Value	Design Variables					
CAR_USE	Commercial	1	0				
	Private	0	1				
CAR_TYPE	Panel Truck	1	0	0	0	0	0
	Pickup	0	1	0	0	0	0
	SUV	0	0	1	0	0	0

	Sedan	0	0	0	1	0	0
	Sports Car	0	0	0	0	1	0
	Van	0	0	0	0	0	1
oldclaimcat	2	1	0	0	0		
	3	0	1	0	0		
	4	0	0	1	0		
	99	0	0	0	1		
REVOKED	No	1	0				
	Yes	0	1				
MARRIED	No	1	0				
	Yes	0	1				
PARENT1	No	1	0				
	Yes	0	1				
MAX_EDUC	<High School	1	0	0	0	0	
	Bachelors	0	1	0	0	0	
	High School	0	0	1	0	0	
	Masters	0	0	0	1	0	
	PhD	0	0	0	0	1	
DENSITY	Highly Rural	1	0	0	0		
	Highly Urban	0	1	0	0		
	Rural	0	0	1	0		
	Urban	0	0	0	1		
mvrcat	1	1	0	0	0	0	0
	2	0	1	0	0	0	0
	3	0	0	1	0	0	0
	4	0	0	0	1	0	0
	5	0	0	0	0	1	0
	99	0	0	0	0	0	1
incomecat	2	1	0	0	0	0	0
	3	0	1	0	0	0	0
	4	0	0	1	0	0	0
	5	0	0	0	1	0	0
	6	0	0	0	0	1	0
	99	0	0	0	0	0	1

Model Convergence Status

Convergence criterion (GCONV=1E-8) satisfied.

Model Fit Statistics

Criterion	Intercept Only	Intercept and Covariates
AIC	11298.075	8254.849
SC	11305.258	8542.155
-2 Log L	11296.075	8174.849

Testing Global Null Hypothesis: BETA=0

Test	Chi-Square	DF	Pr > ChiSq
Likelihood Ratio	3121.2266	39	<.0001
Score	2739.9660	39	<.0001
Wald	1913.4263	39	<.0001

Type 3 Analysis of Effects

Effect	DF	Wald Chi-Square	Pr > ChiSq
--------	----	--------------------	------------

KIDSDRIV	1	48.5730	<.0001
KIDSDRIV*KIDSDRIV	1	9.2754	0.0023
CAR_USE	1	153.5241	<.0001
bluebk	1	40.2729	<.0001
bluebk*bluebk	1	9.9935	0.0016
RETAINED	1	38.5766	<.0001
RETAINED*RETAINED	1	10.3953	0.0013
CAR_TYPE	5	105.0976	<.0001
oldclaimcat	3	79.9446	<.0001
REVOKED	1	105.6912	<.0001
mvrcat	5	39.8971	<.0001
incomecat	5	60.4470	<.0001
MARRIED	1	136.3249	<.0001
PARENT1	1	7.5347	0.0061
MAX_EDUC	4	54.8914	<.0001
DENSITY	3	884.2249	<.0001
TRAVTIME	1	21.0231	<.0001
TRAVTIME*TRAVTIME	1	5.4753	0.0193
AGE	1	156.4409	<.0001
AGE*AGE	1	152.6391	<.0001

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Wald Chi-Square	Pr > ChiSq
Intercept	1	5.2730	0.5584	89.1585	<.0001
KIDSDRIV	1	0.9878	0.1417	48.5730	<.0001
KIDSDRIV*KIDSDRIV	1	-0.1898	0.0623	9.2754	0.0023
CAR_USE Commercial	1	0.8716	0.0703	153.5241	<.0001
CAR_USE Private	0	0	.	.	.
bluebk	1	-0.0313	0.00493	40.2729	<.0001
bluebk*bluebk	1	0.000878	0.000278	9.9935	0.0016
RETAINED	1	-0.1215	0.0196	38.5766	<.0001
RETAINED*RETAINED	1	0.00417	0.00129	10.3953	0.0013
CAR_TYPE Panel Truck	1	-0.2070	0.1384	2.2387	0.1346
CAR_TYPE Pickup	1	-0.1674	0.1159	2.0854	0.1487
CAR_TYPE SUV	1	-0.0449	0.1146	0.1535	0.6952
CAR_TYPE Sedan	1	-0.7259	0.1156	39.4019	<.0001
CAR_TYPE Sports Car	1	0.1615	0.1315	1.5086	0.2194
CAR_TYPE Van	0	0	.	.	.
oldclaimcat 2	1	0.6003	0.0786	58.2903	<.0001
oldclaimcat 3	1	0.5823	0.0809	51.8207	<.0001
oldclaimcat 4	1	0.1340	0.1021	1.7213	0.1895
oldclaimcat 99	0	0	.	.	.
REVOKED No	1	-0.9005	0.0876	105.6912	<.0001
REVOKED Yes	0	0	.	.	.
mvrcat 1	1	0.1243	0.0870	2.0387	0.1533
mvrcat 2	1	0.2238	0.0917	5.9607	0.0146
mvrcat 3	1	0.2131	0.0967	4.8570	0.0275
mvrcat 4	1	0.2068	0.0824	6.3018	0.0121
mvrcat 5	1	0.8251	0.1344	37.7044	<.0001
mvrcat 99	0	0	.	.	.
incomecat 2	1	-0.4491	0.1156	15.0914	0.0001
incomecat 3	1	-0.5661	0.1079	27.5012	<.0001
incomecat 4	1	-0.5248	0.1124	21.8021	<.0001
incomecat 5	1	-0.8474	0.1286	43.4473	<.0001
incomecat 6	1	-0.9402	0.1325	50.3431	<.0001
incomecat 99	0	0	.	.	.
MARRIED No	1	0.7769	0.0665	136.3249	<.0001
MARRIED Yes	0	0	.	.	.
PARENT1 No	1	-0.2631	0.0959	7.5347	0.0061
PARENT1 Yes	0	0	.	.	.
MAX_EDUC <High School	1	0.6230	0.1431	18.9405	<.0001
MAX_EDUC Bachelors	1	0.1166	0.1252	0.8670	0.3518
MAX_EDUC High School	1	0.5942	0.1314	20.4396	<.0001
MAX_EDUC Masters	1	0.1162	0.1240	0.8780	0.3488
MAX_EDUC PhD	0	0	.	.	.
DENSITY Highly Rural	1	-1.8899	0.2046	85.3342	<.0001
DENSITY Highly Urban	1	1.2062	0.0597	408.5895	<.0001
DENSITY Rural	1	-1.7266	0.1260	187.7742	<.0001

DENSITY	Urban	0	0	.	.	.
TRAVTIME		1	0.0281	0.00613	21.0231	<.0001
TRAVTIME*TRAVTIME		1	-0.00019	0.000081	5.4753	0.0193
AGE		1	-0.2994	0.0239	156.4409	<.0001
AGE*AGE		1	0.00327	0.000265	152.6391	<.0001

Odds Ratio Estimates

Effect		Point Estimate	95% Wald Confidence Limits	
CAR_USE	Commercial vs Private	2.391	2.083	2.744
CAR_TYPE	Panel Truck vs Van	0.813	0.620	1.066
CAR_TYPE	Pickup vs Van	0.846	0.674	1.062
CAR_TYPE	SUV vs Van	0.956	0.764	1.197
CAR_TYPE	Sedan vs Van	0.484	0.386	0.607
CAR_TYPE	Sports Car vs Van	1.175	0.908	1.521
oldclaimcat	2 vs 99	1.823	1.562	2.126
oldclaimcat	3 vs 99	1.790	1.528	2.098
oldclaimcat	4 vs 99	1.143	0.936	1.397
REVOKED	No vs Yes	0.406	0.342	0.482
mvrcat	1 vs 99	1.132	0.955	1.343
mvrcat	2 vs 99	1.251	1.045	1.497
mvrcat	3 vs 99	1.238	1.024	1.496
mvrcat	4 vs 99	1.230	1.046	1.445
mvrcat	5 vs 99	2.282	1.754	2.970
incomecat	2 vs 99	0.638	0.509	0.800
incomecat	3 vs 99	0.568	0.459	0.702
incomecat	4 vs 99	0.592	0.475	0.738
incomecat	5 vs 99	0.429	0.333	0.551
incomecat	6 vs 99	0.391	0.301	0.506
MARRIED	No vs Yes	2.175	1.909	2.478
PARENT1	No vs Yes	0.769	0.637	0.928
MAX_EDUC	<High School vs PhD	1.864	1.408	2.468
MAX_EDUC	Bachelors vs PhD	1.124	0.879	1.436
MAX_EDUC	High School vs PhD	1.812	1.400	2.344
MAX_EDUC	Masters vs PhD	1.123	0.881	1.432
DENSITY	Highly Rural vs Urban	0.151	0.101	0.226
DENSITY	Highly Urban vs Urban	3.341	2.972	3.755
DENSITY	Rural vs Urban	0.178	0.139	0.228

Association of Predicted Probabilities and Observed Responses

Percent Concordant	84.1	Somers' D	0.683
Percent Discordant	15.7	Gamma	0.684
Percent Tied	0.2	Tau-a	0.268
Pairs	18534726	c	0.842

Diagnostics

- **Area under curve**

The area under the curve is 0.842, which indicates good predictive ability. The ROC curve is shown in Figure 1.

- **Hat matrix diagonals**

These are computed in SAS Insight. The threshold value for the hat matrix diagonals is: $h_{ii} > 2p/n = 2 \cdot 40/9727 = 0.0082$. 6.6% of the observations have h_{ii} higher than the threshold. It is worthwhile checking the validity of these cases.

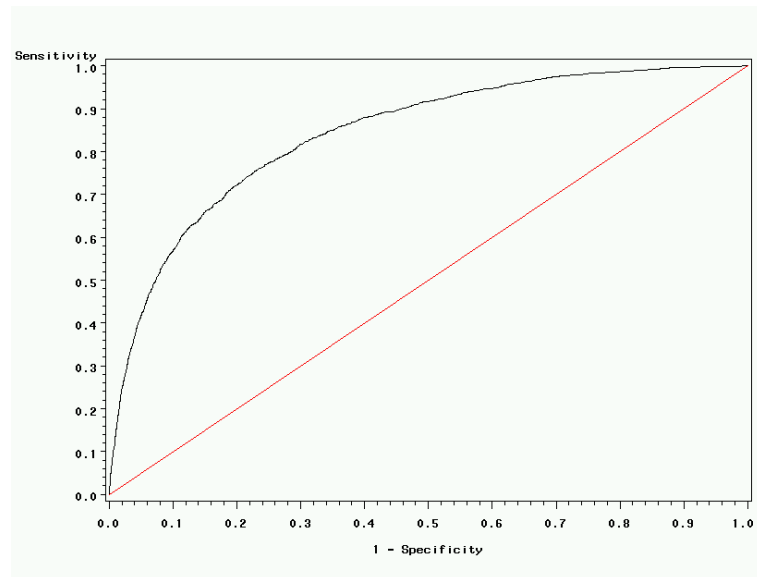


Figure 1: ROC curve, model for claim occurrence

Parameter interpretation

Odds ratio estimates are given in the `proc logistic` output. (This is an advantage of using `proc logistic` rather than `proc genmod`, for logistic regression.) For example:

- The effect of car use being commercial compared with private, is a 139.1% increase in the odds of a claim;
- The effect of a vehicle being a panel truck compared with van, is a $100 \times (1 - 0.813) = 18.7\%$ decrease in the odds of a claim; etc.