

# Insurance Risk Analytics

**First Model** : Construct a GLM using only the first **order terms** (linear terms) of the potential predictor variables by assuming **Poisson distribution and log link function**. Please do the following:

1. Assess the significance of each predictor by looking at the global test (i.e. Type 3 test). Use the testing results to perform variable selection.

Full model with all predictors:

Analysis Of Maximum Likelihood Parameter Estimates								
Parameter		DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept		1	-0.6724	0.9737	-2.5809	1.2361	0.48	0.4898
ppc	A	1	-0.4668	0.2307	-0.9189	-0.0147	4.09	0.0430
ppc	B	1	-0.4377	0.2000	-0.8297	-0.0458	4.79	0.0286
ppc	C	1	-0.2229	0.1729	-0.5617	0.1160	1.66	0.1974
cons_typ	fire_re	1	-0.4581	0.2403	-0.9291	0.0130	3.63	0.0567
cons_typ	frame	1	0.0696	0.1201	-0.1659	0.3050	0.34	0.5624
height		1	0.0302	0.0492	-0.0661	0.1266	0.38	0.5385
sqf		1	0.0005	0.0002	0.0001	0.0008	7.72	0.0054
Scale		0	1.0000	0.0000	1.0000	1.0000		

There are three variables insignificant at the 10% significance level: ppc level C, construction type frame, and height. We can first remove height since it has a high p-value and is a single continuous variable, while the other two are dummy variables of other factor variables.

Analysis Of Maximum Likelihood Parameter Estimates								
Parameter		DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept		1	-0.0951	0.2561	-0.5971	0.4068	0.14	0.7102
ppc	A	1	-0.4826	0.2289	-0.9311	-0.0340	4.45	0.0350
ppc	B	1	-0.4516	0.1985	-0.8407	-0.0625	5.17	0.0229
ppc	C	1	-0.2282	0.1725	-0.5663	0.1099	1.75	0.1858
cons_typ	fire_re	1	-0.4764	0.2385	-0.9438	-0.0089	3.99	0.0458
cons_typ	frame	1	0.0762	0.1197	-0.1584	0.3107	0.41	0.5245
sqf		1	0.0006	0.0001	0.0004	0.0007	57.27	<.0001
Scale		0	1.0000	0.0000	1.0000	1.0000		

After removing height, ppc level C and construction type frame are still insignificant. But some levels of the categorical variables are significant, so it's ok keep all the levels of the two categorical variable in the model.

- With all significant variables in the type3 test, further look at the needs of collapsing levels for the categorical variables. Perform contrast tests for the between-level comparisons to determine the needs of reducing levels. Follow the contrast test results to reduce the levels of any categorical variables in the model if needed and refit model.

Contrast Estimate Results										
Label	Mean Estimate	Mean		L'Beta Estimate	Standard Error	Alpha	L'Beta		Chi-Square	Pr > ChiSq
		Confidence Limits					Confidence Limits			
ppc A vs.B	0.9695	0.6804	1.3816	-0.0309	0.1807	0.05	-0.3851	0.3233	0.03	0.8641
Exp(ppc A vs.B)				0.9695	0.1752	0.05	0.6804	1.3816		
ppc B vs.C	0.7998	0.6329	1.0107	-0.2234	0.1194	0.05	-0.4574	0.0106	3.50	0.0613
Exp(ppc B vs.C)				0.7998	0.0955	0.05	0.6329	1.0107		
ppc C vs.D	0.7959	0.5676	1.1161	-0.2282	0.1725	0.05	-0.5663	0.1099	1.75	0.1855
Exp(ppc C vs.D)				0.7959	0.1373	0.05	0.5676	1.1161		
cons_typ fire_re vs.frame	0.5755	0.3718	0.8907	-0.5525	0.2228	0.05	-0.9893	-0.1158	6.15	0.0132
Exp(cons_typ fire_re vs.frame)				0.5755	0.1282	0.05	0.3718	0.8907		

Contrast Estimate Results										
Label	Mean Estimate	Mean		L'Beta Estimate	Standard Error	Alpha	L'Beta		Chi-Square	Pr > ChiSq
		Confidence Limits					Confidence Limits			
ppc AB vs.CD	0.7825	0.6400	0.9566	-0.2453	0.1025	0.05	-0.4463	-0.0443	5.72	0.0167
Exp(ppc AB vs.CD)				0.7825	0.0802	0.05	0.6400	0.9566		
cons_typ ao vs.fire_re	1.6796	1.0911	2.5855	0.5185	0.2201	0.05	0.0872	0.9499	5.55	0.0185
Exp(cons_typ ao vs.fire_re)				1.6796	0.3697	0.05	1.0911	2.5855		

The p value for ppc A vs. B, and ppc C vs. D are larger than 0.1, indicating that they're not significantly different. Since ppc is ordinal variable, it makes sense to combine level A and B, level C and D. On the other hand, the p value for contrast of cons\_type fire and cons\_type frame is smaller than 0.1, indicating that these two levels are significantly different. Furthermore, since frame is insignificant in the previous test, we can collapse frame and masonry together as base level.

Analysis Of Maximum Likelihood Parameter Estimates								
Parameter		DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept		1	-0.2090	0.2009	-0.6028	0.1849	1.08	0.2983
ppc_new	AB	1	-0.2453	0.1025	-0.4463	-0.0443	5.72	0.0167
cons_new	fire_re	1	-0.5185	0.2201	-0.9499	-0.0872	5.55	0.0185
sqf		1	0.0005	0.0001	0.0004	0.0007	61.37	<.0001
Scale		0	1.0000	0.0000	1.0000	1.0000		

Now, we have three variables remaining in the model and all are significant.

The final model is:  $E(\text{claims}) = -0.209 - 0.2453 \cdot \text{ppc\_new} - 0.5185 \cdot \text{cons\_new} + 0.0005 \cdot \text{sqf}$

if ppc = A or B, then ppc\_new = 1, otherwise ppc\_new = 0;

if construction type is fire\_re, then cons\_new = 1, otherwise cons\_new = 0.

3. Please write a small paragraph to summarize and explain the model results to the Chief Actuary including interpreting all parameter estimates of the model.

The expected value of claim of a house is linearly related to its public protection class, construction type and the square footage of the house.

If the house has a Public Protection Class level A or B, its expected claim will be 0.2453 less than level C or D when other conditions are the same; If the construction type is Fire Resistant, the house's expected claim will be - 0.5185 less than masonry and frame when other conditions are the same; one unit increase of the square footage of the house will increase claims by 0.0005.

4. Now, the Chief Actuary wants to know for a particular house with Level B in Public Protection Class, Masonry Construction Type, 25 ft tall, and 2850 sqf, what is the expected loss frequency based on this first model?

$$E(\text{claims}) = -0.209 - 0.2453 * \text{ppc\_new} - 0.5185 * \text{cons\_new} + 0.0005 * \text{sqf}$$

if ppc = A or B, then ppc\_new = 1, otherwise ppc\_new = 0;

if construction type is fire\_re, then cons\_new = 1, otherwise cons\_new = 0.

In the above case, ppc level = B, construction type = Masonry, then ppc\_new = 1, cons\_new = 0.

$$y = -0.209 - 0.2453 + 0.0005 * 2850$$

$$= 0.9707$$

If the exposure for this particular house is 1, then:

The expected loss frequency = Expected claims / Exposure = 0.9707

5. Assess the goodness of fit for the model by looking at:
- Overdispersion from specific statistics or residual analysis, and propose 2 different ways to correct overdispersion if it exists

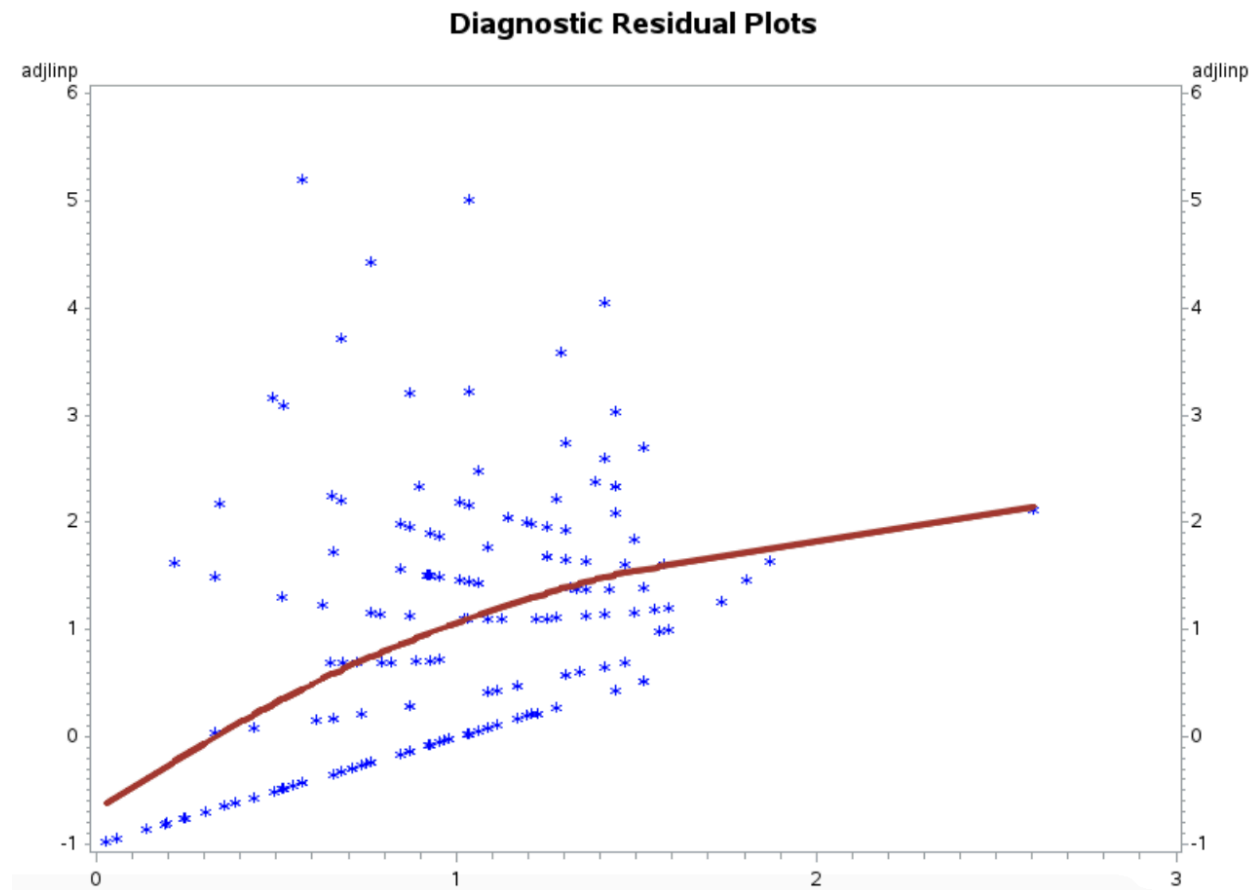
Criteria For Assessing Goodness Of Fit			
Criterion	DF	Value	Value/DF
Deviance	169	543.0998	3.2136
Scaled Deviance	169	543.0998	3.2136
Pearson Chi-Square	169	526.0316	3.1126
Scaled Pearson X2	169	526.0316	3.1126
Log Likelihood		71.6769	
Full Log Likelihood		-447.9070	
AIC (smaller is better)		903.8141	
AICC (smaller is better)		904.0522	
BIC (smaller is better)		916.4272	

- Since the scale of the first model is 1, the scaled deviance is the same with deviance (543.0998), and the scaled Pearson X2 is the same with Pearson Chi-Square (526.0316). The DF is 169 (number of observation = 173, number of parameters = 4). Pearson Chi-

Square Value/DF = 3.1126, which is far larger than 1, indicating that there is a problem of over-dispersion.

There are two ways to address over-dispersion: 1) to run the same model as a negative binomial regression, or by estimate the scale parameter = (Pearson Chi / D.F.) 2) to correct the standard errors of the estimates by removing outliers etc.

#### c. Linearity Assumption



From the residual plot above, we can see that the assumption of linearity is not well satisfied.

#### d. Identifying outliers

We can use Cook's D statistics to identify outliers. If the value  $D_i \geq 4/(n-p-1)$ , then observation  $i$  is considered as a model outlier. In our case,  $4/(n-p-1) = 4/(173-4-1) = 0.0238$ .

For each observation, Cook's D statistics,  $D_i = r_i^2 h_i / (p(1-h_i))$ , where  $p$  is the number of parameters in the model,  $r$  is the standardized Pearson Residual,  $h$  = Hessian weight \* variance of the linear predictor.

Thus, the outliers identified by Cook's D statistics are listed as above.

### Potential Influential Observations

Obs	claims	Pred	ppc_new	cons_new	sqf	CookD	Stresdev
3	9	2.8152508	CD	ao	2300	0.0403045	2.9410694
13	11	4.2237024	CD	ao	3050	0.0377311	2.7583004
15	14	2.8152508	CD	ao	2300	0.1318139	4.7756725
22	3	1.676165	CD	fire_re	2300	0.023969	0.9572842
33	0	1.4251031	CD	fire_re	2000	0.0275188	-1.748016
34	8	2.3935726	CD	ao	2000	0.0484219	2.8655336
45	9	3.6824428	AB	ao	3250	0.0793404	2.380591
56	15	4.1110071	CD	ao	3000	0.0940618	4.1560936
67	10	4.5806653	CD	ao	3200	0.0277841	2.2039848
71	7	2.8868903	AB	ao	2800	0.0344953	2.0664153
77	3	1.2424788	AB	fire_re	2200	0.0440457	1.3754511
117	12	3.6329839	AB	ao	3225	0.1922641	3.5215628
121	6	1.6808488	AB	ao	1800	0.0430397	2.5947135
126	0	1.676165	CD	fire_re	2300	0.0384251	-1.906769
131	3	1.387079	CD	fire_re	1950	0.035278	1.2252474
134	10	2.1440474	AB	ao	2250	0.1154397	3.9146049
139	6	1.6314422	CD	fire_re	2250	0.2601534	2.7307187
141	7	13.512378	CD	ao	5200	0.9201355	-2.541894
146	8	1.9769657	AB	ao	2100	0.0710824	3.2368949
149	10	1.7742664	AB	ao	1900	0.1465648	4.2903481
160	4	2.509927	AB	fire_re	3500	0.0469409	0.9395594
173	0	1.4251031	CD	fire_re	2000	0.0275188	-1.748016

**Second Model:** With the model you have in Step 1 (i.e. your first model), try **Negative Binomial distribution and log link function** and see if you have get a better fit for modeling the dispersion.

1. Assess the significance of each predictor by looking at the global test (i.e. Type 3 test). Use the testing results to perform variable selection. (drop one insignificant variable at a time if you have multiple of them.)

Analysis Of Maximum Likelihood Parameter Estimates								
Parameter		DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept		1	-0.6450	0.4743	-1.5746	0.2846	1.85	0.1738
ppc_new	AB	1	-0.2633	0.1933	-0.6422	0.1155	1.86	0.1731
cons_new	fire_re	1	-0.5401	0.3560	-1.2380	0.1577	2.30	0.1293
sqf		1	0.0007	0.0002	0.0004	0.0011	16.12	<.0001
Dispersion		1	1.0212	0.1873	0.7128	1.4630		

Contrast Estimate Results										
Label	Mean Estimate	Mean		L'Beta Estimate	Standard Error	Alpha	L'Beta		Chi-Square	Pr > ChiSq
		Confidence Limits					Confidence Limits			
ppc AB vs.CD	0.7685	0.5261	1.1224	-0.2633	0.1933	0.05	-0.6422	0.1155	1.86	0.1731
Exp(ppc AB vs.CD)				0.7685	0.1485	0.05	0.5261	1.1224		
cons_typ ao vs.fire_re	0.5827	0.2900	1.1708	-0.5401	0.3560	0.05	-1.2380	0.1577	2.30	0.1293
Exp(cons_typ ao vs.fire_re)				0.5827	0.2075	0.05	0.2900	1.1708		

If we run the first model using negative binomial distribution and log link function, we see that only 1 independent variable sqf is significant. Both ppc\_new and cons\_new are insignificant under 10% S.L. The contrast test shows the same p-value result for two categorical variables we have because there are only 2 levels for each.

We first drop ppc\_new with the highest p-value. As we can see from the result below, cons\_new is still insignificant under 10% S.L.

Analysis Of Maximum Likelihood Parameter Estimates								
Parameter		DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept		1	-0.8583	0.4534	-1.7471	0.0304	3.58	0.0584
sqf		1	0.0008	0.0002	0.0004	0.0011	18.64	<.0001
cons_new	fire_re	1	-0.5180	0.3593	-1.2223	0.1862	2.08	0.1494
Dispersion		1	1.0421	0.1897	0.7294	1.4890		

The final model only contains sqf which has a p value less than 0.0001.

Analysis Of Maximum Likelihood Parameter Estimates							
Parameter	DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept	1	-0.9684	0.4513	-1.8528	-0.0839	4.60	0.0319
sqf	1	0.0008	0.0002	0.0004	0.0011	20.01	<.0001
Dispersion	1	1.0637	0.1923	0.7463	1.5161		

**Note:** The negative binomial dispersion parameter was estimated by maximum likelihood.

LR Statistics For Type 3 Analysis			
Source	DF	Chi-Square	Pr > ChiSq
sqf	1	20.29	<.0001

2. compare the differences between Model 1 and 2

Model 1:  $E(\text{claims}) = -0.209 - 0.2453 \cdot \text{ppc\_new} - 0.5185 \cdot \text{cons\_new} + 0.0005 \cdot \text{sqf}$   
 if  $\text{ppc} = A$  or  $B$ , then  $\text{ppc\_new} = 1$ , otherwise  $\text{ppc\_new} = 0$ ;  
 if construction type is  $\text{fire\_re}$ , then  $\text{cons\_new} = 1$ , otherwise  $\text{cons\_new} = 0$ .

Model 2:  $E(\text{claims}) = -0.9684 + 0.0008 \cdot \text{sqf}$

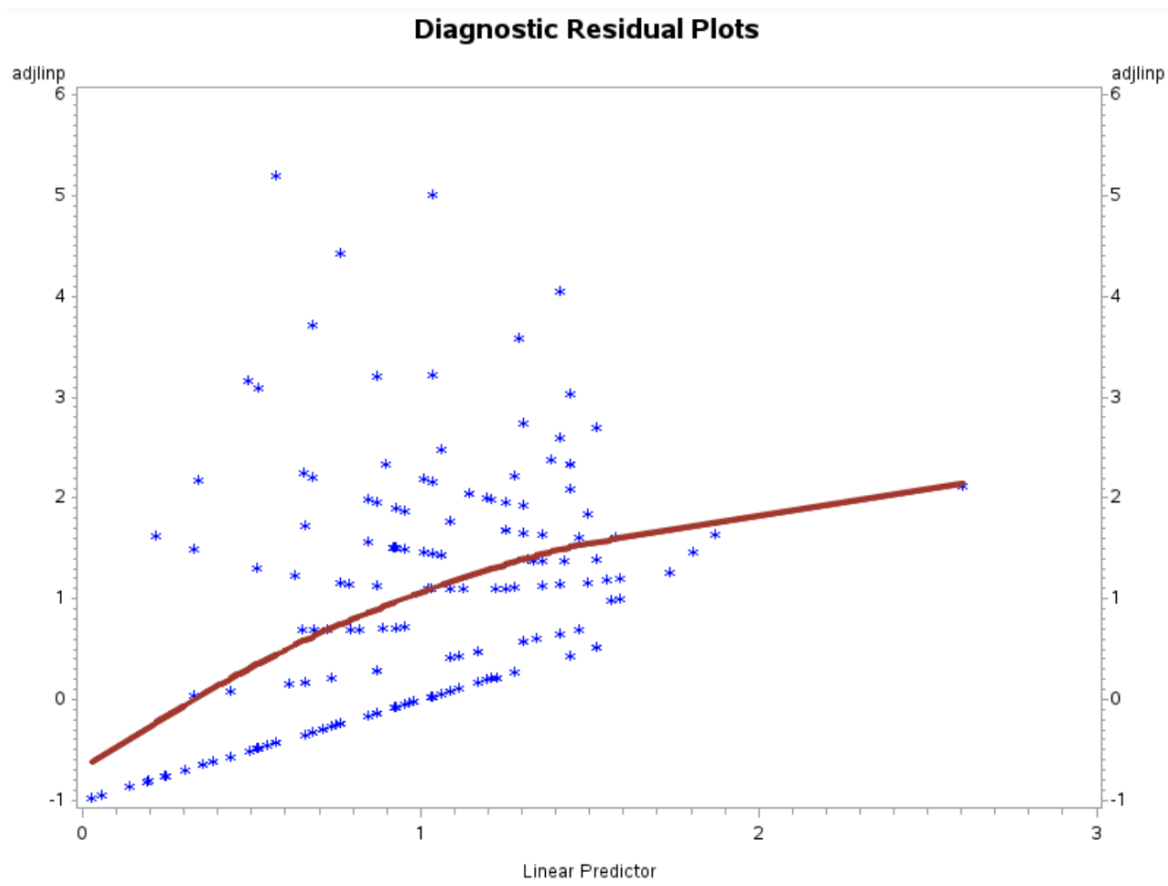
The first model consists 3 significant independent variables while the second model only contains 1 significant independent variable. It is because the first model is a regular Poisson model, assuming variable's mean equal to its variance. It fails to address the problem of over-dispersion and leads to over-estimation of the significance of the predictors, making those insignificant variables significant aka. Type I Error. The second model uses a negative binomial model which successfully addresses the over-dispersion problem.

3. Assess the goodness of fit for the model by looking at:
  - a. Overdispersion

Criteria For Assessing Goodness Of Fit			
Criterion	DF	Value	Value/DF
Deviance	171	195.4597	1.1430
Scaled Deviance	171	195.4597	1.1430
Pearson Chi-Square	171	152.3088	0.8907
Scaled Pearson X2	171	152.3088	0.8907
Log Likelihood		148.4597	
Full Log Likelihood		-371.1243	
AIC (smaller is better)		748.2485	
AICC (smaller is better)		748.3905	
BIC (smaller is better)		757.7084	

The Deviance Value/DF = 1.1430 and the Pearson Chi Value/DF = 0.8907, which are both close to 1, indicating that the over-dispersion issue is addressed.

b. Linearity Assumption



From the residual plot above, we see that the assumption of linearity is still not satisfied.



c. Identifying outliers

The outlier criteria for Cook's D statistic is  $4/(n-p-1) = 4/170 = 0.0235$ . Observations with a higher CookD is listed as below:

Potential Influential Observations					
Obs	claims	Pred	sqf	CookD	Stresdev
3	9	2.8152508	2300	0.0806089	2.9410694
8	0	1.0563759	1900	0.0310926	-1.493504
13	11	4.2237024	3050	0.0754622	2.7583004
14	0	1.0281901	1850	0.0294054	-1.472396
15	14	2.8152508	2300	0.2636279	4.7756725
16	8	4.0013188	2950	0.0246651	1.7678733
22	3	1.676165	2300	0.0479379	0.9572842
29	0	2.5908959	2600	0.0253191	-2.298078
33	0	1.4251031	2000	0.0550375	-1.748016
34	8	2.3935726	2000	0.0968438	2.8655336
45	9	3.6824428	3250	0.1586808	2.380591
51	5	1.9769657	2100	0.0358136	1.8116334
56	15	4.1110071	3000	0.1881235	4.1560936
61	0	2.5217665	2550	0.0237114	-2.26642
62	5	1.9242171	2050	0.0378736	1.8571982
67	10	4.5806653	3200	0.0555682	2.2039848
71	7	2.8868903	2800	0.0689905	2.0664153
74	6	2.7401354	2250	0.0237829	1.7088771
76	0	1.1478513	1600	0.036541	-1.559995
77	3	1.2424788	2200	0.0880914	1.3754511
81	0	1.2765389	2250	0.0467586	-1.6517
82	0	3.0473367	2900	0.0400583	-2.50018
87	0	4.2237024	3050	0.0293178	-2.926276
88	0	2.8098633	2750	0.0314323	-2.396406
90	0	1.2116461	1700	0.040341	-1.604725
94	0	4.5806653	3200	0.0397001	-3.052453
117	12	3.6329839	3225	0.3845282	3.5215628
120	6	2.4544817	2500	0.0464927	1.9235205
121	6	1.6808488	1800	0.0860793	2.5947135
126	0	1.676165	2300	0.0768502	-1.906769
128	4	1.4098913	1475	0.0394365	1.7926542
131	3	1.387079	1950	0.0705559	1.2252474
134	10	2.1440474	2250	0.2308794	3.9146049
138	5	2.3252499	2400	0.0263071	1.5314683
139	6	1.6314422	2250	0.5203067	2.7307187
141	7	13.512378	5200	1.840271	-2.541894
146	8	1.9769657	2100	0.1421648	3.2368949
149	10	1.7742664	1900	0.2931296	4.2903481
150	9	4.1110071	3000	0.0379233	2.0932624
157	0	4.5806653	3200	0.0397001	-3.052453
160	4	2.509927	3500	0.0938817	0.9395594
172	0	2.6261679	2625	0.0261962	-2.314111
173	0	1.4251031	2000	0.0550375	-1.748016

**Third Model** : Test high order terms (i.e. quadratic terms etc.) for any continuous variables in your Second Model and assess the variable significances and goodness of the model.

Criteria For Assessing Goodness Of Fit			
Criterion	DF	Value	Value/DF
Deviance	170	195.9132	1.1524
Scaled Deviance	170	195.9132	1.1524
Pearson Chi-Square	170	156.6727	0.9216
Scaled Pearson X2	170	156.6727	0.9216
Log Likelihood		150.1493	
Full Log Likelihood		-369.4346	
AIC (smaller is better)		746.8693	
AICC (smaller is better)		747.1074	
BIC (smaller is better)		759.4824	

Analysis Of Maximum Likelihood Parameter Estimates							
Parameter	DF	Estimate	Standard Error	Wald 95% Confidence Limits		Wald Chi-Square	Pr > ChiSq
Intercept	1	-2.8940	1.0809	-5.0126	-0.7755	7.17	0.0074
sqf	1	0.0023	0.0008	0.0008	0.0038	8.84	0.0029
sqf*sqf	1	-0.0000	0.0000	-0.0000	-0.0000	4.28	0.0385
Dispersion	1	1.0244	0.1881	0.7147	1.4682		

Both sqf and sqf^2 are significant in the model. But the estimate of sqf \* sqf is very close to 0, we should consider divide sqf\*sqf by 1000 to see the relationship more clearly.

The goodness of fit indicates that this model is better than the second one. 1) the Value/DF is closer to 1; 2) The AIC, AICC is smaller than those of the second model (748.25, 748.39), although BIC is smaller in the second model. The linearity plot also looks better.

