



Practical issues in model design

CAS Special Interest Seminar
on Predictive Modeling

Boston, October 2006

James Tanser

WWW.WATSONWYATT.COM

W Watson Wyatt
Worldwide



Agenda

- Variates, Polynomials and Splines
- Continuous interactions
- Practical problems with continuous variables
- Over-dispersed Poisson





Agenda

- Variates, Polynomials and Splines
- Continuous interactions
- Practical problems with continuous variables
- Over-dispersed Poisson



Copyright © Watson Wyatt Worldwide. All rights reserved.



Variates, Polynomials and Splines

- Common to model using treating all variables categorically, with discrete levels
 - Allows actual shape to show through
 - Produces step changes in genuinely continuous variables (eg AOI)
 - No extrapolation
- Some variables (age) are both continuous and discrete
 - Step changes are acceptable
 - Some smoothness is desirable



Copyright © Watson Wyatt Worldwide. All rights reserved.



Variates, Polynomials and Splines

- Variates allow each unique data value to have a different effect on the linear predictor, but force some smoothness
- In practice implemented via:
 - Polynomials
 - or
 - Splines



Copyright © Watson Wyatt Worldwide. All rights reserved.



Agenda

- Variates, Polynomials and Splines
- Continuous interactions
- Practical problems with continuous variables
- Over-dispersed Poisson



Copyright © Watson Wyatt Worldwide. All rights reserved.



Continuous interactions

- Interaction = multiplication
- Continuous * Discrete
- Continuous * Continuous



Copyright © Watson Wyatt Worldwide. All rights reserved.



Interaction = Multiplication

- The notation is the clue:
 - A.B or A*B
- For example:
 - Two variates X and Y
 - Model should be linear in X and linear in Y
 - Interaction between X and Y to be included



Copyright © Watson Wyatt Worldwide. All rights reserved.



Interaction = multiplication

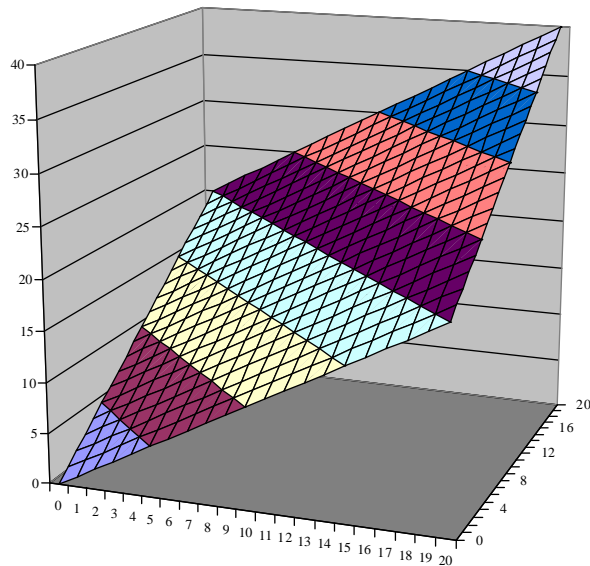
- Model 1: $Z = aX + bY$
 - Linear in X
 - Linear in Y
 - No interaction – "a" does not depend on Y



Copyright © Watson Wyatt Worldwide. All rights reserved.



Interaction = Multiplication $Z = aX + bY$



Copyright © Watson Wyatt Worldwide. All rights reserved.

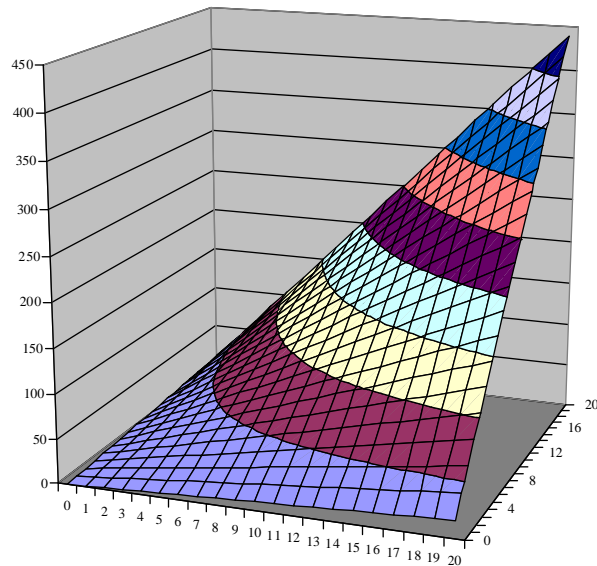
Interaction = multiplication

- Model 1: $Z = aX + bY$
 - Linear in X
 - Linear in Y
 - No interaction – "a" does not depend on Y
- Model 2: $Z = aX + bY + cXY$
 - Linear in X (for any given Y)
 - Linear in Y (for any given X)
 - Interaction present – the gradient for X depends on the value of Y (and vice versa)
 - Quadratic in X=Y direction



Copyright © Watson Wyatt Worldwide. All rights reserved.

Interaction = Multiplication $Z = aX + bY + cXY$

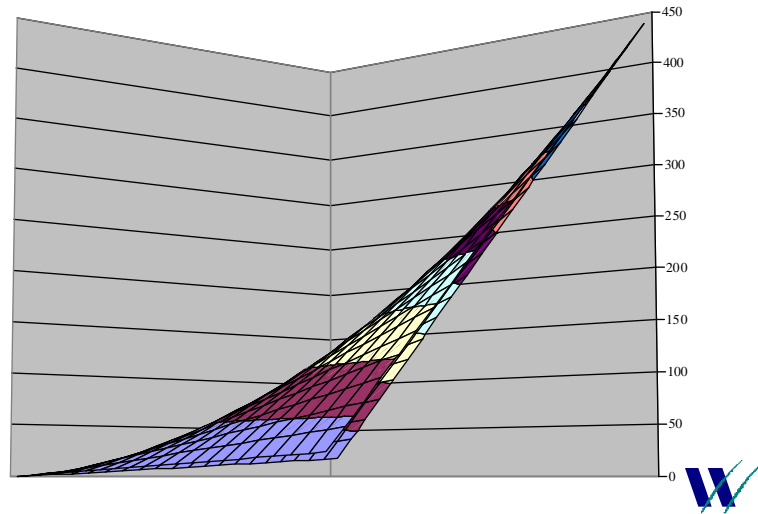


Copyright © Watson Wyatt Worldwide. All rights reserved.



Interaction = Multiplication

$$Z = aX + bY + cXY$$



Copyright © Watson Wyatt Worldwide. All rights reserved.



Continuous interactions

- Interaction = multiplication
- Continuous * Discrete
- Continuous * Continuous



Copyright © Watson Wyatt Worldwide. All rights reserved.

Continuous * Discrete

- Define a new set of variates, one for each factor level, so that variate n is:
 - Variate value if factor at level n
 - Zero otherwise
- Treat each of these variates as usual:
 - Polynomial (same order?)
 - Spline (same knots?)
- Useful to include factor in model for neatness



Copyright © Watson Wyatt Worldwide. All rights reserved.

Design matrix Spline

18	M	1	1	1	0	0	...	1
20	F	1	0.52	0.98	0.02	0	...	0
22	F	1	0.17	0.83	0.17	0	...	0
24	M	1	0.02	0.5	0.48	0.02	...	1
26	M	1	0	0.17	0.67	0.17	...	1
28	M	1	0	0.02	0.48	0.48	...	1
30	F	1	0	0	0.17	0.67	...	0
32	F	1	0	0	0.02	0.48	...	0
34	M	1	0	0	0	0.17	...	1
36	F	1	0	0	0	0.02	...	0



Copyright © Watson Wyatt Worldwide. All rights reserved.



Design matrix Discrete * Spline

1	1	1	0	0	...	0	0	0	0	...
1	0	0	0	0	...	0.52	0.98	0.02	0	...
1	0	0	0	0	...	0.17	0.83	0.17	0	...
1	0.02	0.5	0.48	0.02	...	0	0	0	0	...
1	0	0.17	0.67	0.17	...	0	0	0	0	...
1	0	0.02	0.48	0.48	...	0	0	0	0	...
1	0	0	0	0	...	0	0	0.17	0.67	...
1	0	0	0	0	...	0	0	0.02	0.48	...
1	0	0	0	0.17	...	0	0	0	0	...
1	0	0	0	0	...	0	0	0	0.02	...

Copyright © Watson Wyatt Worldwide. All rights reserved.



Design matrix Discrete * Polynomial

1	18	324	5832	104976	...	0	0	0	0	...
1	0	0	0	0	...	20	400	8000	160000	...
1	0	0	0	0	...	22	484	10648	234256	...
1	24	576	13824	331776	...	0	0	0	0	...
1	26	676	17576	456976	...	0	0	0	0	...
1	28	784	21952	614656	...	0	0	0	0	...
1	0	0	0	0	...	30	900	27000	810000	...
1	0	0	0	0	...	32	1024	32768	1048576	...
1	34	1156	39304	1336336	...	0	0	0	0	...
1	0	0	0	0	...	36	1296	46656	1679616	...

Copyright © Watson Wyatt Worldwide. All rights reserved.





Continuous interactions

- Interaction = multiplication
- Continuous * Discrete
- Continuous * Continuous



Copyright © Watson Wyatt Worldwide. All rights reserved.



Continuous * Continuous

- Simply create $X*Y$ terms!
- Eg Polynomial order 2 for X and Y:
 - X, X^2, Y, Y^2
 - XY
 - XY^2, X^2Y
 - X^2Y^2
- For splines combine together all the basis functions
 - $f_1(x), f_2(x), f_3(x), \dots, g_1(y), g_2(y), g_3(y), \dots$
 - $f_1(x).g_1(y), f_2(x).g_1(y), f_3(x).g_1(y)$
 - ...



Copyright © Watson Wyatt Worldwide. All rights reserved.



Design matrix Spline * Spline

- Not enough room on slide to show!



Copyright © Watson Wyatt Worldwide. All rights reserved.



Agenda

- Variates, Polynomials and Splines
- Continuous interactions
- Practical problems with continuous variables
- Over-dispersed Poisson



Copyright © Watson Wyatt Worldwide. All rights reserved.



Practical problems

- Number of variates
- Missing values
- Edge effects



Copyright © Watson Wyatt Worldwide. All rights reserved.



Number of variates

- Various tricks used to make modeling with categorical factors quick
 - Only one calculation per factor per row
 - Calculation is addition
- Tricks don't work variates and calculation is multiplication
 - Polynomial with 5 terms is slower than adding 5 new factors



Copyright © Watson Wyatt Worldwide. All rights reserved.



Number of variates

- Splines make it easy to include many variates, slowing down the model
 - Use only at final modelling stages
 - Be parsimonious!
- Interactions with variates creates many (tens or hundreds) of variates very quickly



Copyright © Watson Wyatt Worldwide. All rights reserved.



Practical problems

- Number of variates
- Missing values
- Edge effects



Copyright © Watson Wyatt Worldwide. All rights reserved.



Missing values

- Missing values in a variate often cause entire record to be ignored
 - Replace missing values with zeros
- Care is needed to differentiate "real" zeros and "missing" zeros
 - Create a missing flag and include in all models involving variate
 - Remember spline basis functions transform zero to some other (non-zero) value (extrapolation)



Copyright © Watson Wyatt Worldwide. All rights reserved.



Practical problems

- Number of variates
- Missing values
- Edge effects



Copyright © Watson Wyatt Worldwide. All rights reserved.



Edge effects

- One or two records with extreme variate values can have a disproportionate effect on the model
 - Look at leverage or Cook's distance
 - Understand your data
 - Consider limiting range of variate
 - Be careful when extrapolating



Copyright © Watson Wyatt Worldwide. All rights reserved.



Cautionary example

- Artificial data, loosely based on actual naive analysis
- Retention analysis containing three records with incorrect premium change, all of which renewed
- Problems:
 - Overfitting to edges
 - Knot placement

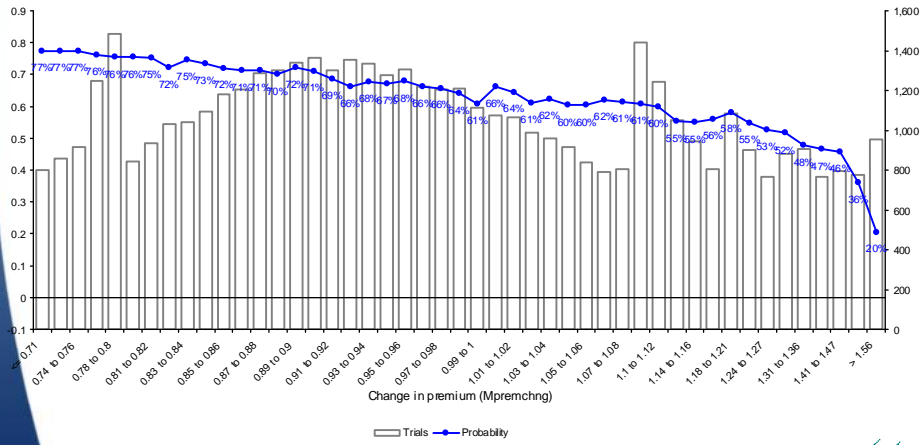


Copyright © Watson Wyatt Worldwide. All rights reserved.

Simple grouped oneway

Retention job

Example of problem factor



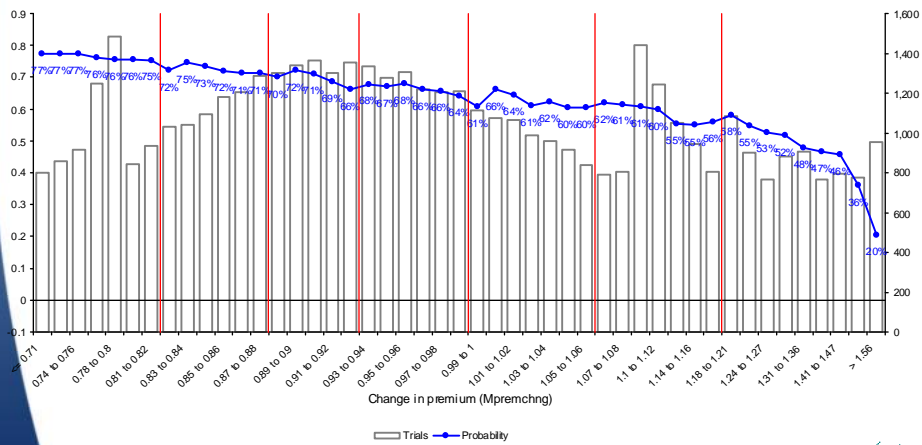
Copyright © Watson Wyatt Worldwide. All rights reserved.



Simple grouped oneway

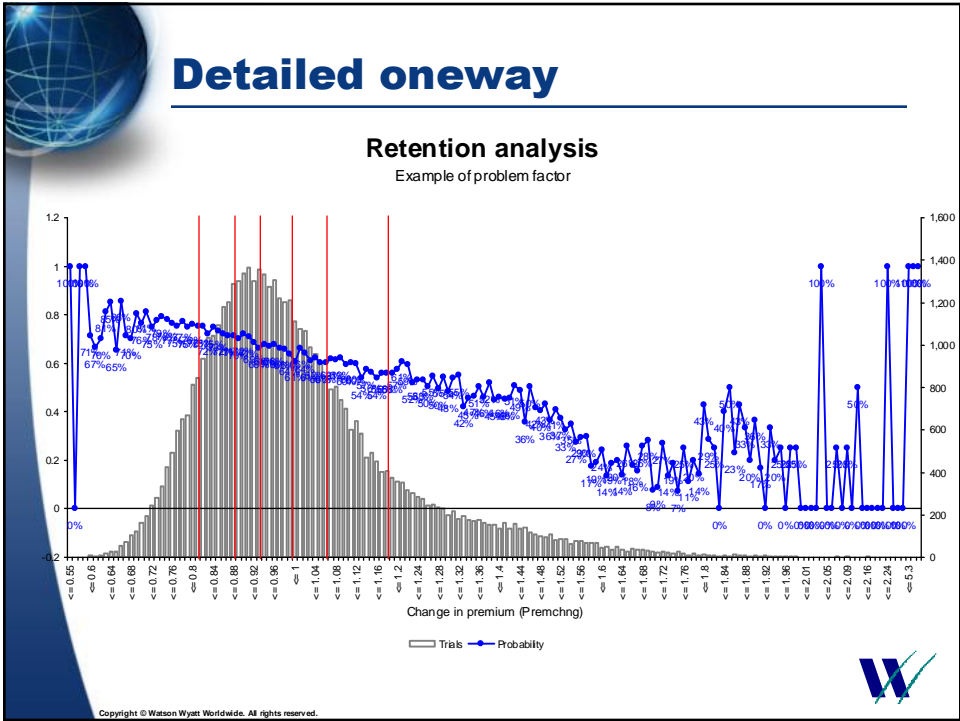
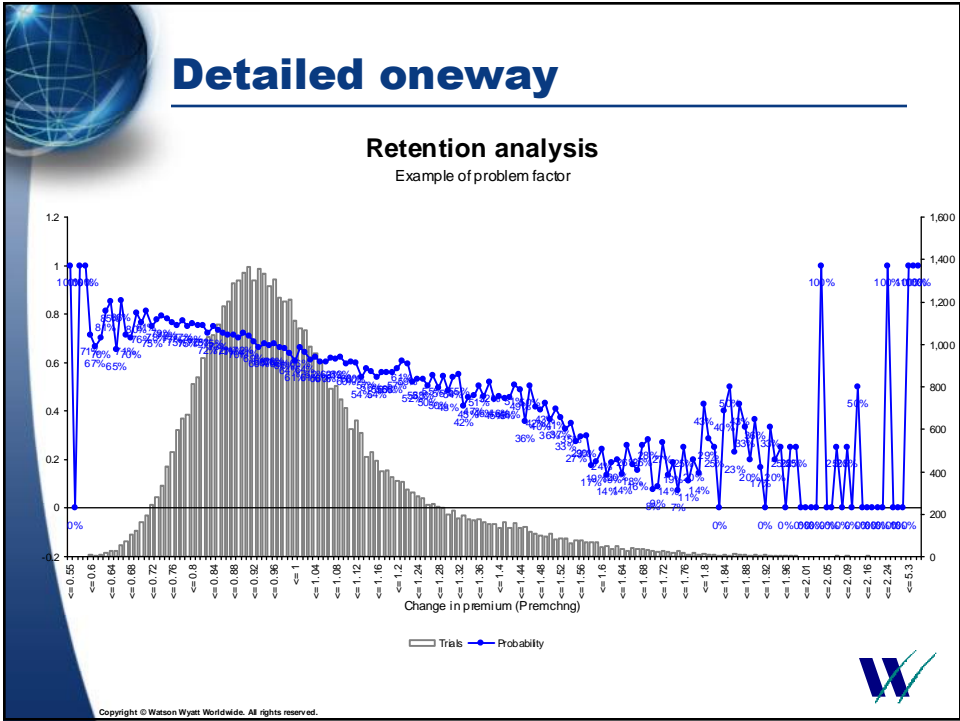
Retention job

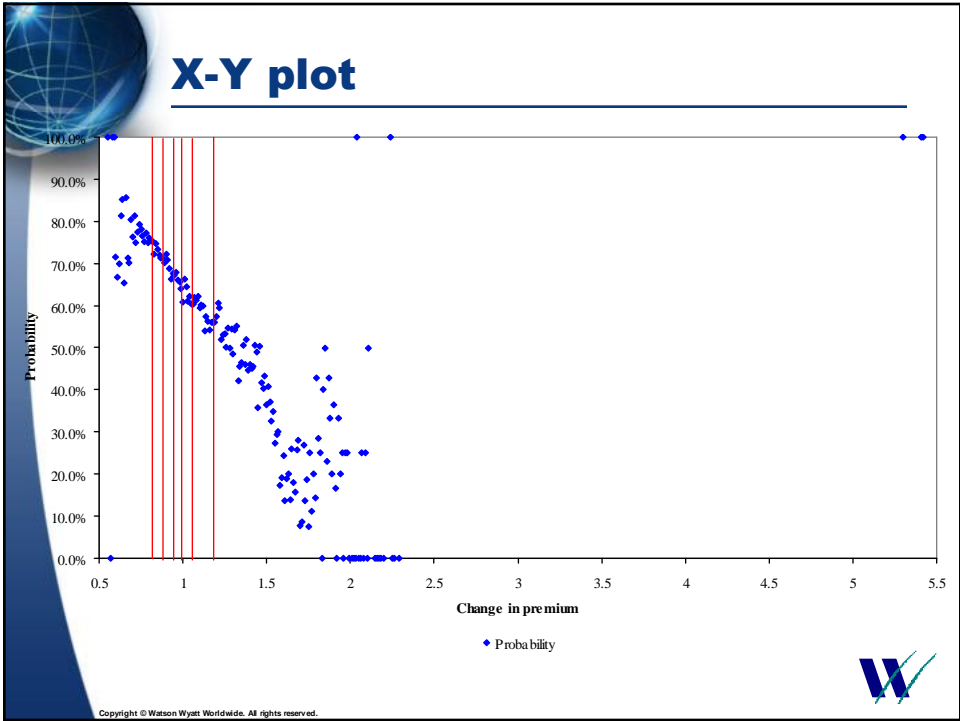
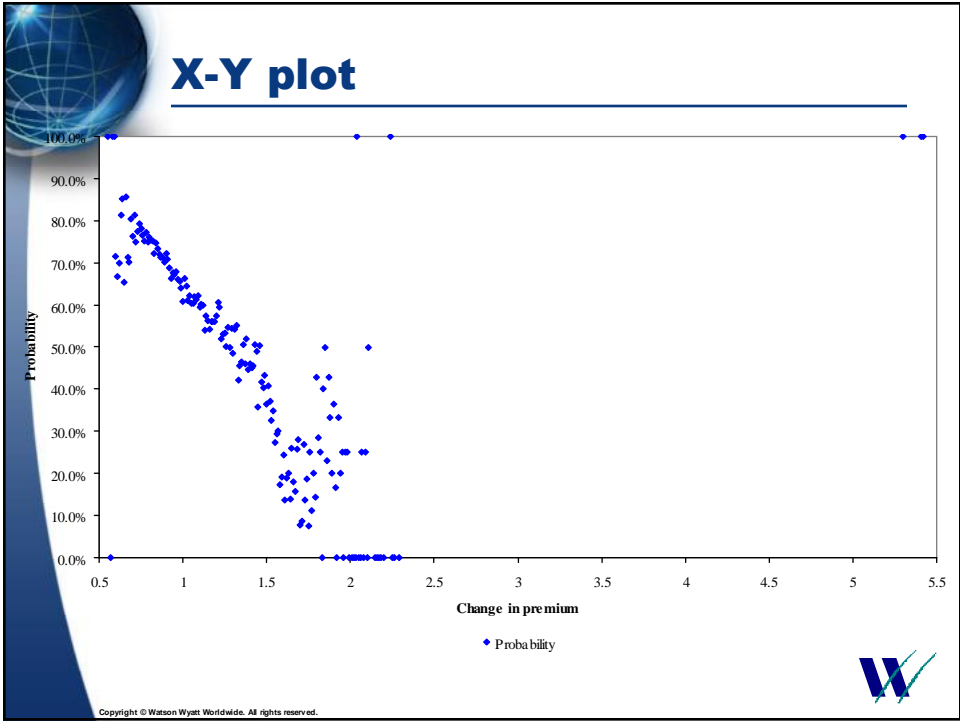
Example of problem factor

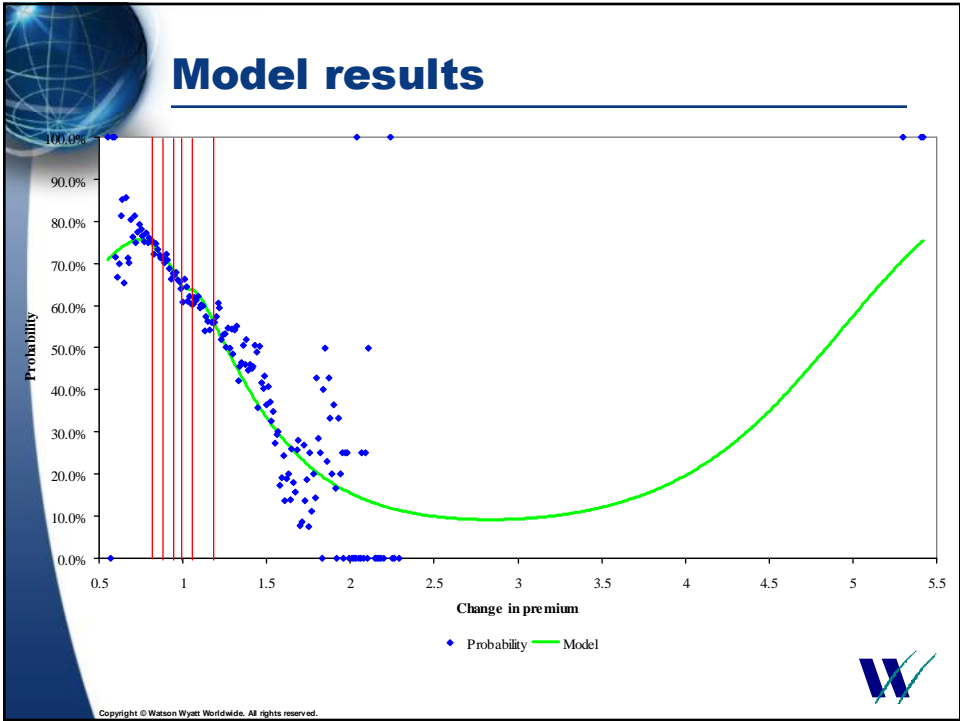
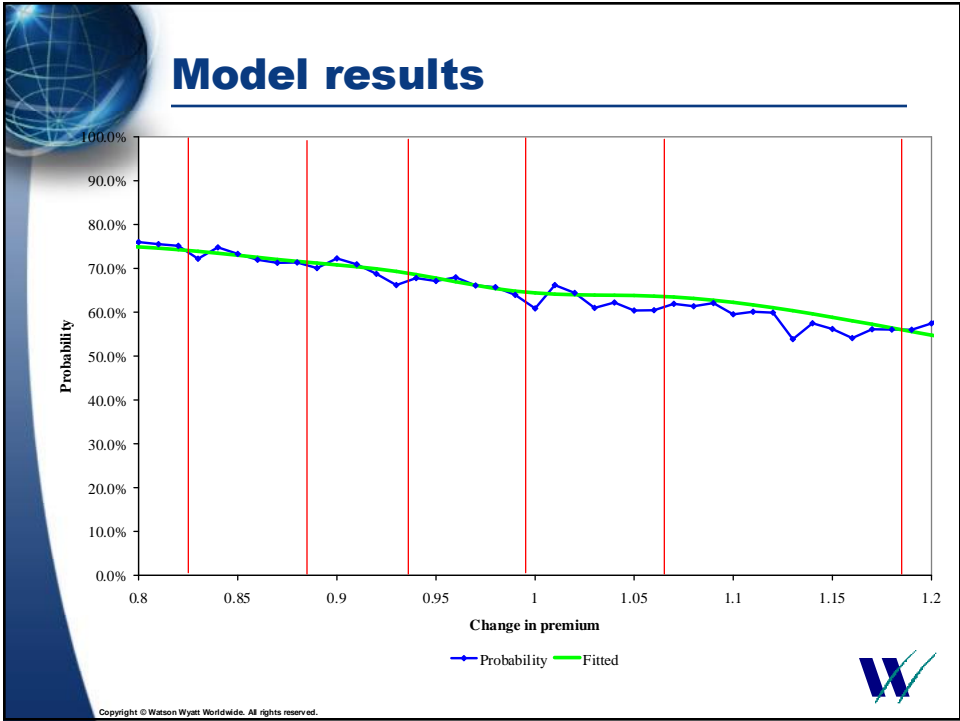


Copyright © Watson Wyatt Worldwide. All rights reserved.











Agenda

- Variates, Polynomials and Splines
- Continuous interactions
- Practical problems with continuous variables
- Over-dispersed Poisson



Copyright © Watson Wyatt Worldwide. All rights reserved.



Interdependence of claim events

- Sometimes claim events are not individually independent
- A real life example:
 - Health insurer recorded data such that individual PMI claim payments could not be matched to a particular single medical event
 - this meant that each transactional claim payment was recorded as a new claim regardless of the event
 - claim numbers therefore appear in multiples per policy in the data
- This invalidates assumptions underlying the GLM framework



Copyright © Watson Wyatt Worldwide. All rights reserved.

Mathematical implications

- Poisson model used for numbers assumes

$$E[Y] = \mu \quad \text{Var}[Y] = \mu$$

- Replacing every one claim with K claims gives

$$E[Y] = K.\mu \quad \text{Var}[Y] = K^2.\mu \quad \checkmark$$

- But the Poisson GLM modelling process applies the Poisson assumptions which are, in this case wrong!

$$E[Y] = K.\mu \quad \text{Var}[Y] = K.\mu \quad \times$$

Copyright © Watson Wyatt Worldwide. All rights reserved.

Mathematical implications

- Fitting a Poisson GLM to claims data that is not independent does not effect the parameter estimates
- But it does effect the standard errors!



Copyright © Watson Wyatt Worldwide. All rights reserved.

Generalised linear models

$$E[\underline{Y}] = \underline{\mu} = g^{-1}(\underline{X} \cdot \underline{\beta} + \underline{\xi})$$

$$\text{Var}[\underline{Y}] = \phi \cdot V(\underline{\mu}) / \underline{\omega}$$

scale parameter

- Inclusion of a scale parameter adjusts the variance assumed in the model



Copyright © Watson Wyatt Worldwide. All rights reserved.

Estimating the scale parameter

- Deviance scale for Poisson

$$\phi = D / (n - p)$$

- Pearson scale

$$\chi^2 = \sum \omega_i (Y_i - \mu_i)^2 / V(\mu_i)$$

$$\phi = \chi^2 / (n - p)$$



Copyright © Watson Wyatt Worldwide. All rights reserved.



Theoretical case study

- Fitted numbers model to the TPPD claims numbers
- Data was adjusted by
 - multiplying exposure by 1000
 - multiplying claims numbers by 10
- Tried fitting models
 - Poisson
 - overdispersed Poisson (with Deviance scale)
 - overdispersed Poisson (with Pearson scale)



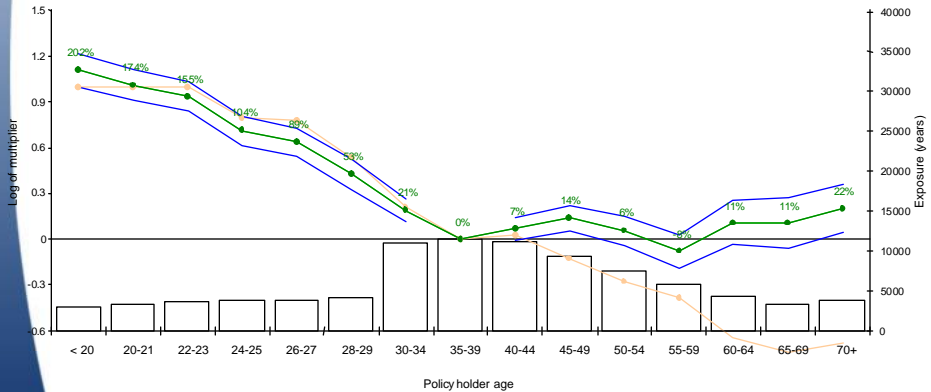
Copyright © Watson Wyatt Worldwide. All rights reserved.



The correct answer

Fully worked example of the tutorial job

Run 8 Model 1 - Final models with analysis - TPPD numbers



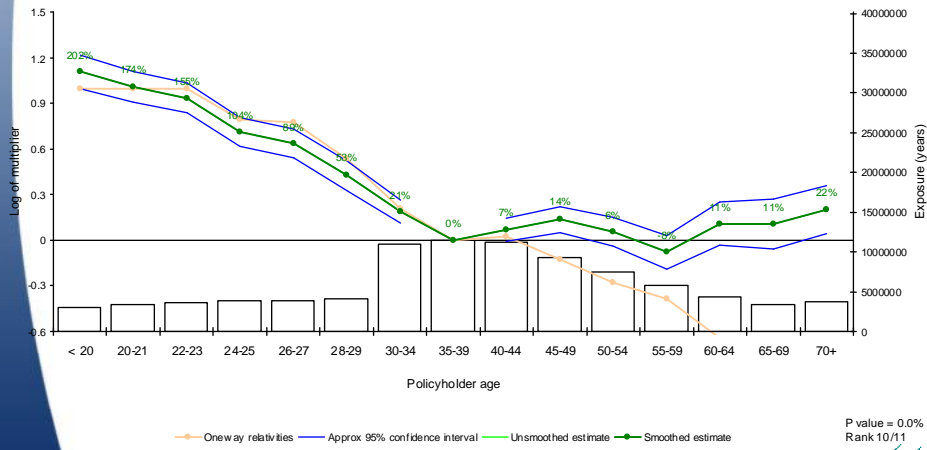
P value = 0.0%
Rank 10/11

Copyright © Watson Wyatt Worldwide. All rights reserved.

Exposure * 1,000

Fully worked example of the tutorial job

Run 8 Model 2 - Final models with analysis - TPPD2 numbers



P value = 0.0%
Rank 10/11

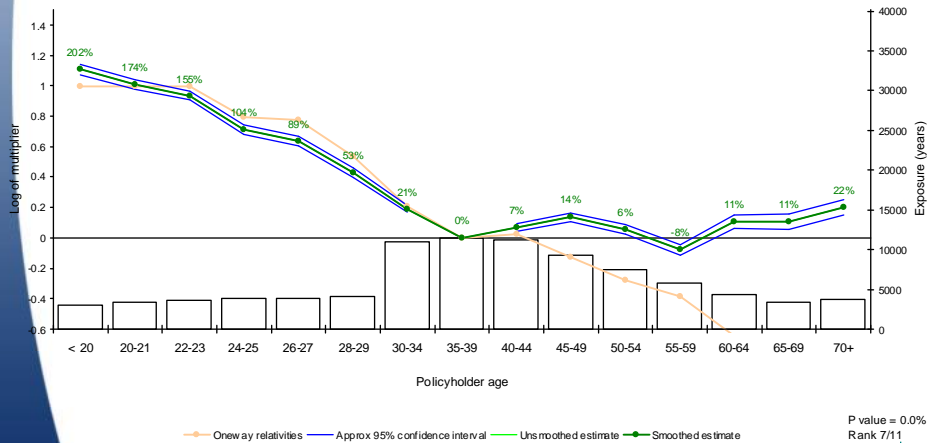


Copyright © Watson Wyatt Worldwide. All rights reserved.

Num * 10

Fully worked example of the tutorial job

Run 8 Model 3 - Final models with analysis - TPPD3 numbers (log poisson)



P value = 0.0%
Rank 7/11

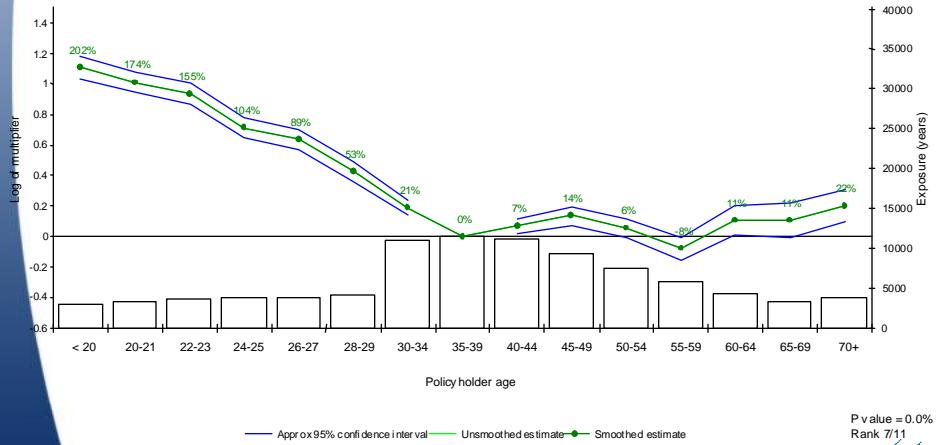


Copyright © Watson Wyatt Worldwide. All rights reserved.

Overdispersed Poisson (deviance)

Fully worked example of the tutorial job

Run 8 Model 4 - Final models with analysis - TPPD3 numbers (log over-dispersed deviance poisson)



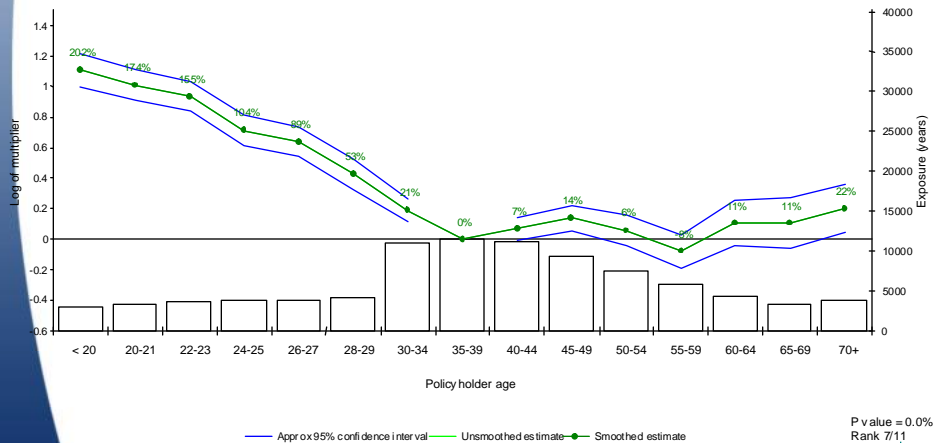
P value = 0.0%
Rank 7/11

Copyright © Watson Wyatt Worldwide. All rights reserved.

Overdispersed Poisson (Pearson)

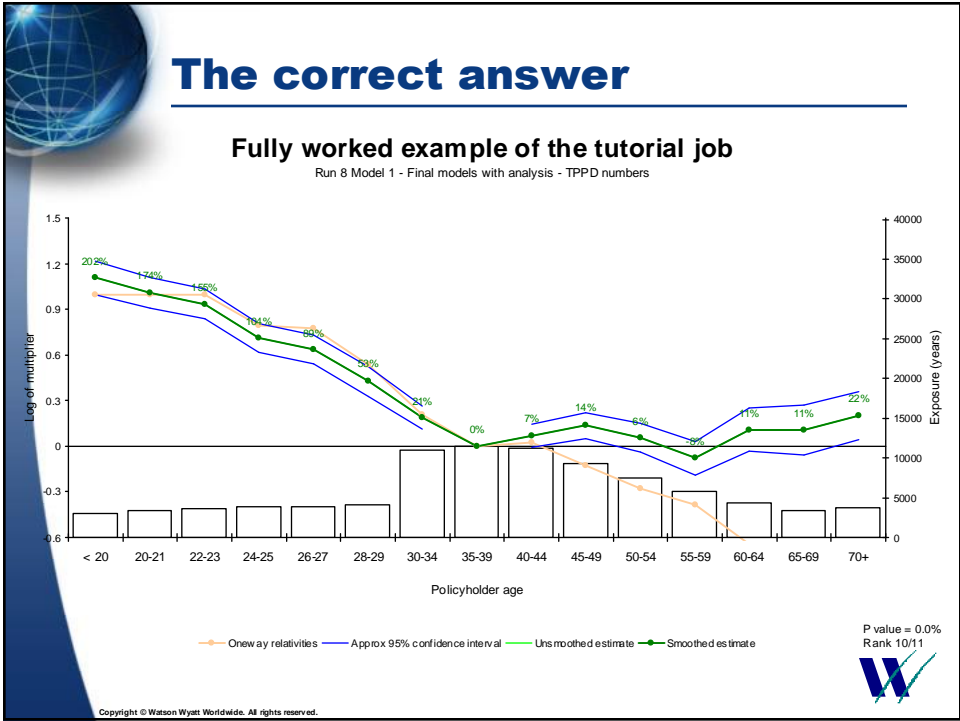
Fully worked example of the tutorial job

Run 8 Model 5 - Final models with analysis - TPPD3 numbers (log over-dispersed pearson poisson)



P value = 0.0%
Rank 7/11

Copyright © Watson Wyatt Worldwide. All rights reserved.



Practical issues in model design

CAS Special Interest Seminar on Predictive Modeling

Boston, October 2006

James Tanser

WWW.WATSONWYATT.COM

W Watson Wyatt
Worldwide