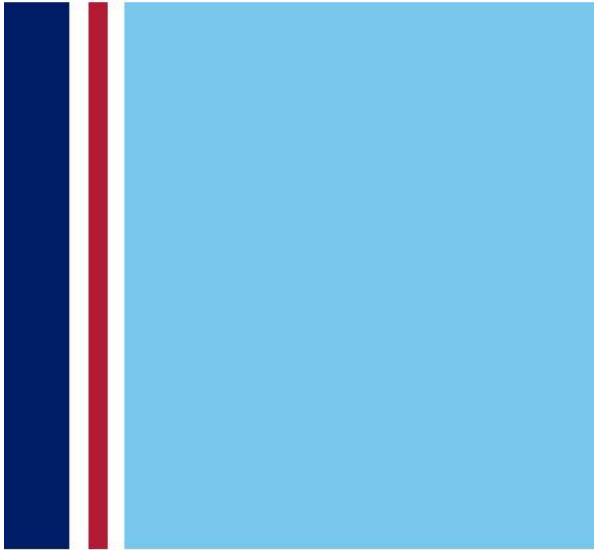


watsonwyatt.com



## Practical Issues in Model Design

### **CAS Special Interest Seminar on Predictive Modeling**

Claudine Modlin, FCAS, MAAA  
October 11, 2007

# 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

# Variates, Polynomials and Splines

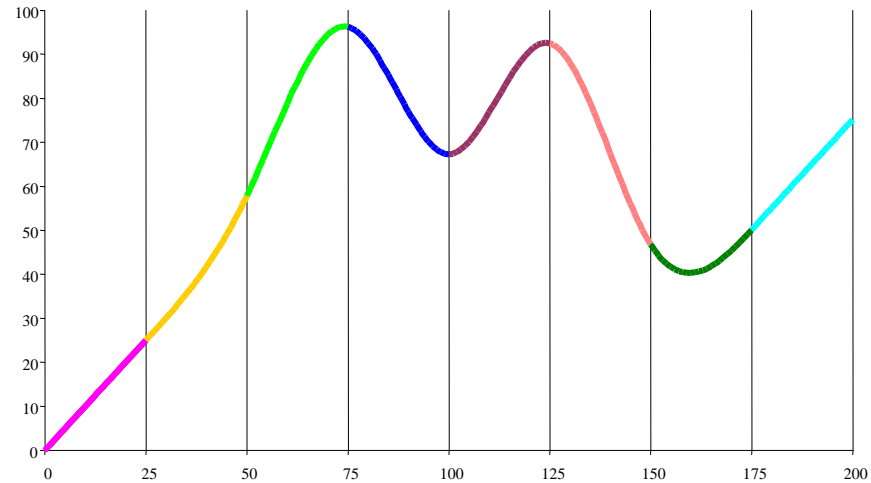
- Common to model treating all variables categorically, with discrete levels
  - allows actual shape to show through
  - produces step changes in genuinely continuous variables (eg AOI)
  - no extrapolation
- Some variables (age) are both continuous and discrete
  - step changes are acceptable
  - some smoothness is desirable

# 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

# Spline definition

- A series of polynomial functions, with each function defined over a short interval
- Intervals are defined by  $k+2$  knots
  - two exterior knots at extremes of data
  - variable number ( $k$ ) of interior knots
- At each interior knot the two functions must join "smoothly"



# Agenda

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

# Continuous interactions

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



# Interaction = Multiplication

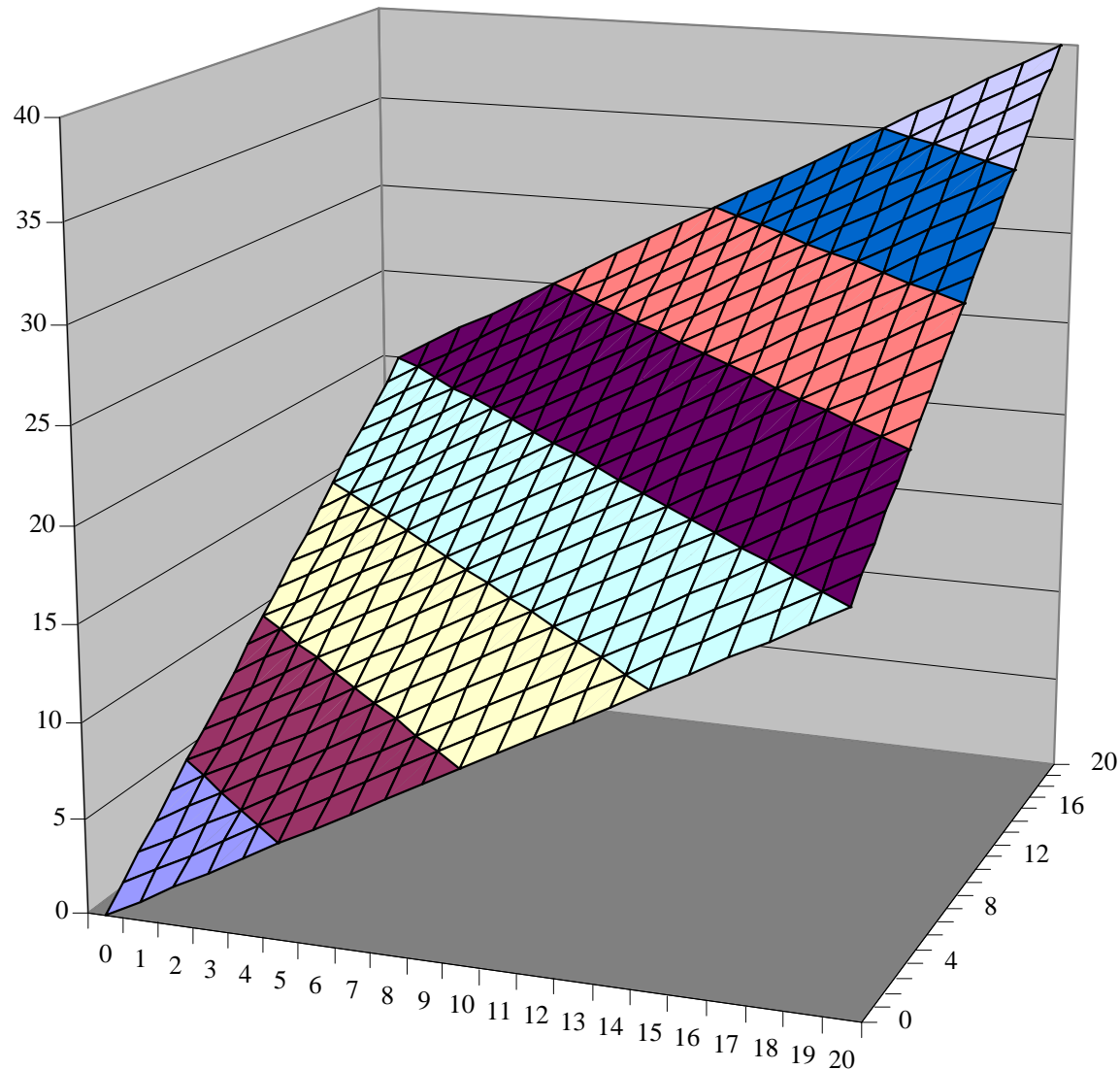
- The notation is the clue:
  - $X.Y$  or  $X*Y$
- For example:
  - two variates  $X$  and  $Y$
  - model should be linear in  $X$  and linear in  $Y$
  - interaction between  $X$  and  $Y$  to be included

# No Interaction

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

# No Interaction

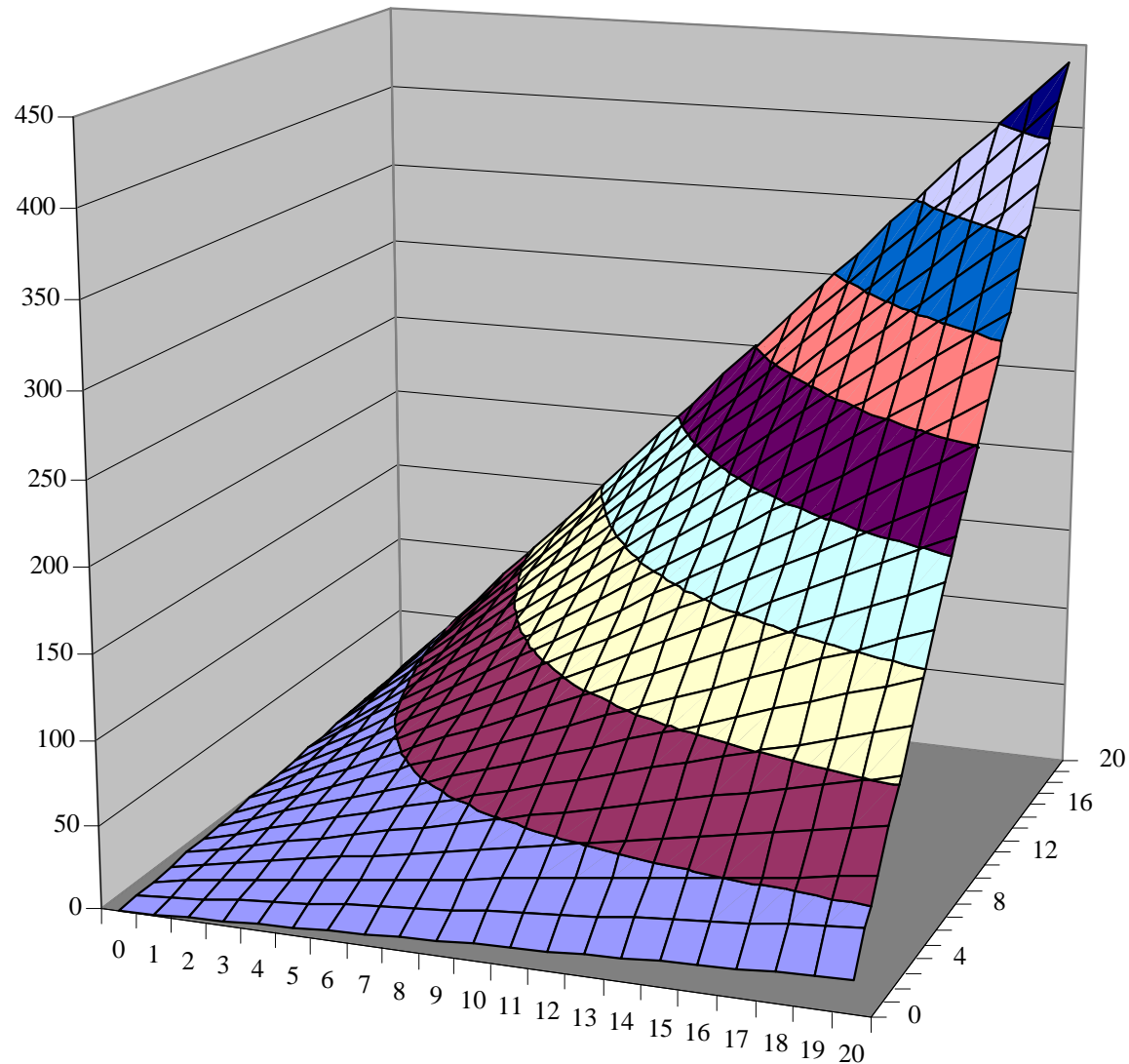
$Z = aX + bY$



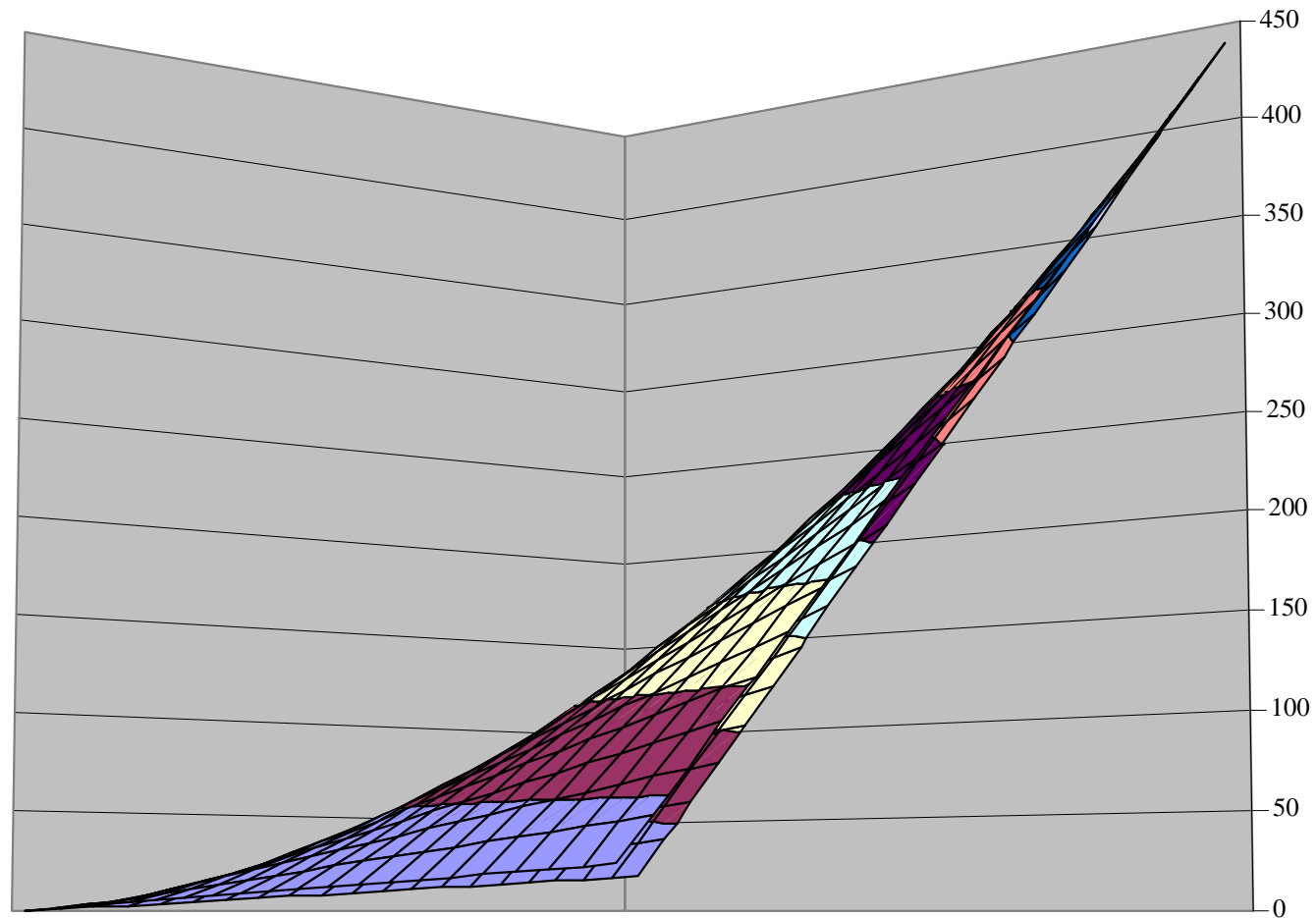
# Interaction = Multiplication

- 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

# Interaction = Multiplication

$$Z = aX + bY + cXY$$


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



# Continuous interactions

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

## Continuous \* Discrete

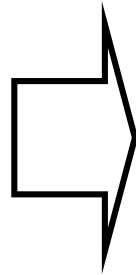
- Define a new set of variates, one for each factor level, so that variate  $n$  is:
  - variate value if factor at level  $n$
  - zero otherwise
- Treat each of these variates as usual:
  - polynomial (same order?)
  - spline (same knots?)
- Useful to include factor in model for neatness



# Design matrix

## Spline and Discrete (no interaction)

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



# 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	...

# 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	...

# Continuous interactions

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

# 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
  - $f_1(x), f_2(x), f_3(x), \dots, g_1(y), g_2(y), g_3(y), \dots$
  - $f_1(x).g_1(y), f_2(x).g_1(y), f_3(x).g_1(y)$
  - $\dots$

# Design matrix Spline \* Spline

- Not enough room on slide to show!

# Agenda

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

# Practical problems

- Number of variates
- Missing values
- Edge effects



## 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 for variates and calculation is multiplication
  - polynomial with 5 terms is slower than adding 5 new factors

## Number of variates

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

# Practical problems

- Number of variates
- Missing values
- Edge effects

# 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)

# Practical problems

- Number of variates
- Missing values
- Edge effects

## 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

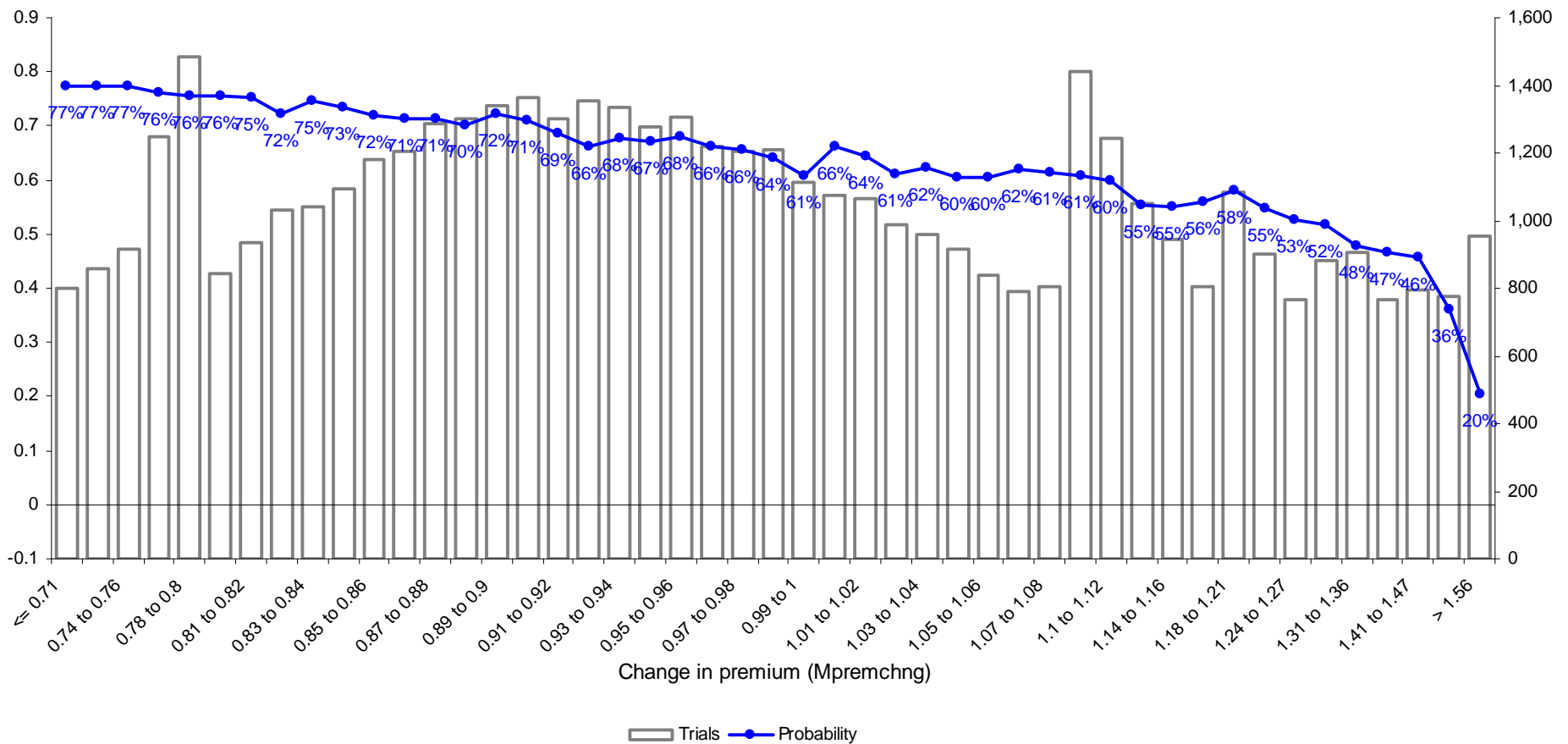
## Cautionary example

- 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

# Simple grouped oneway

## Retention job

Example of problem factor

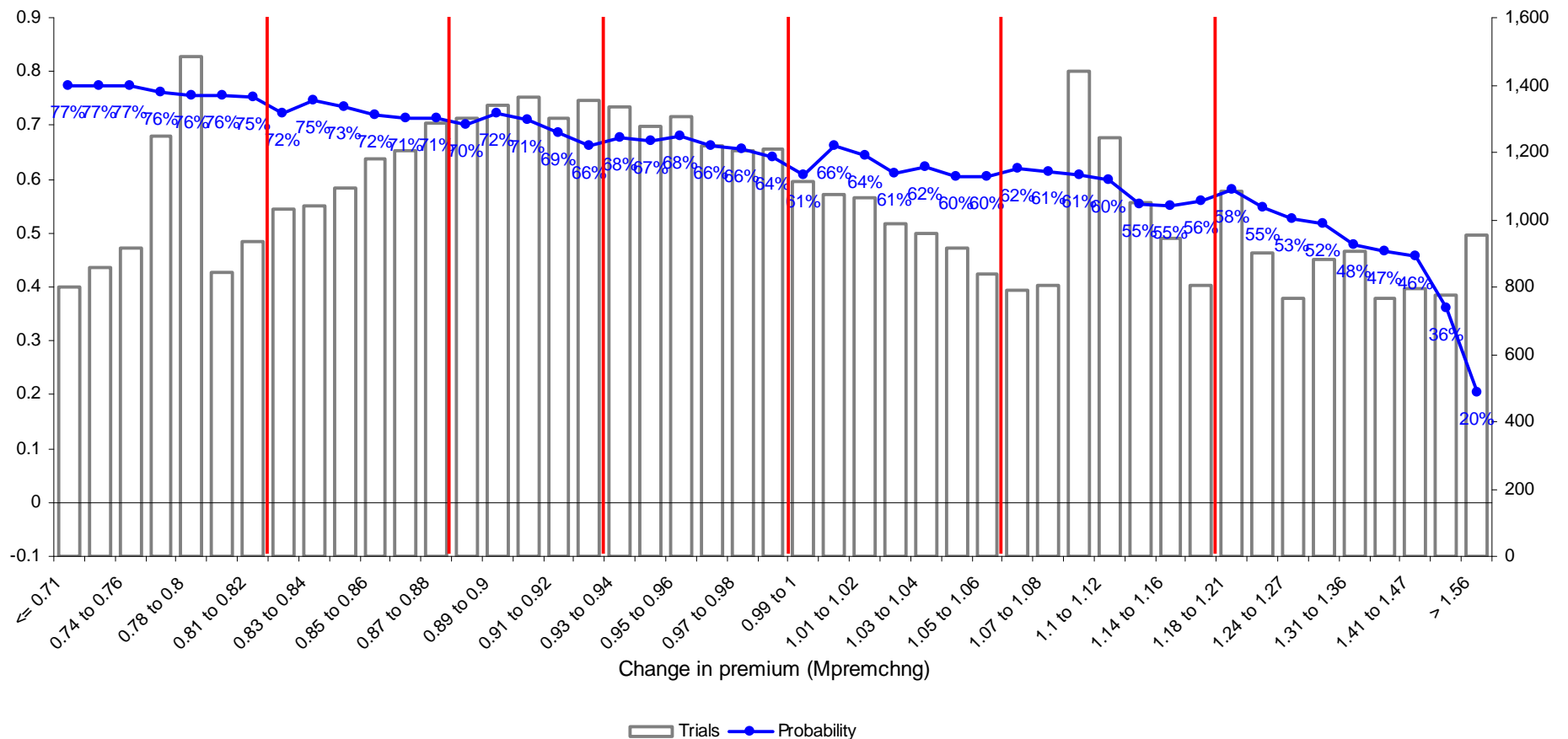




# Simple grouped oneway

## Retention job

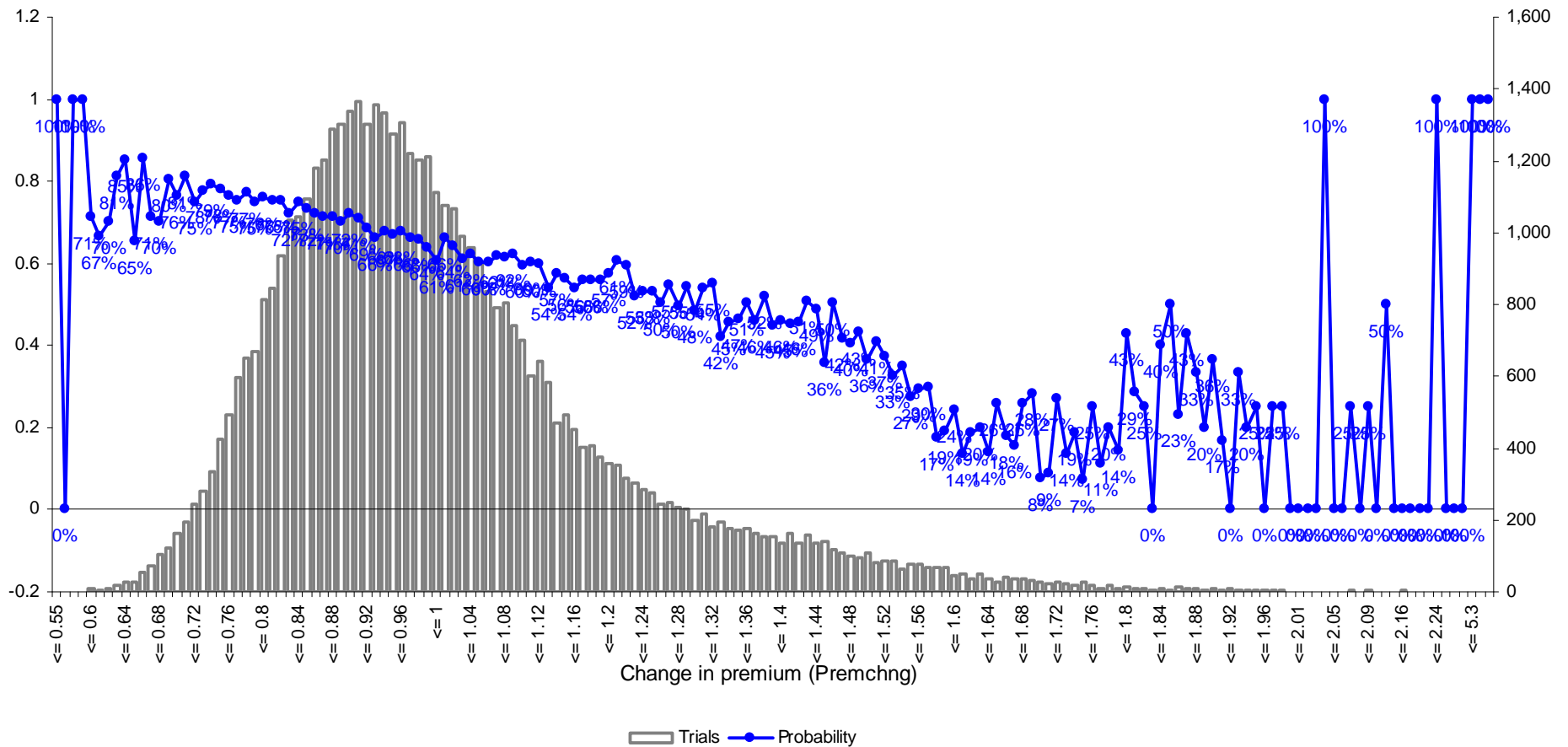
Example of problem factor



# Detailed oneway

## Retention analysis

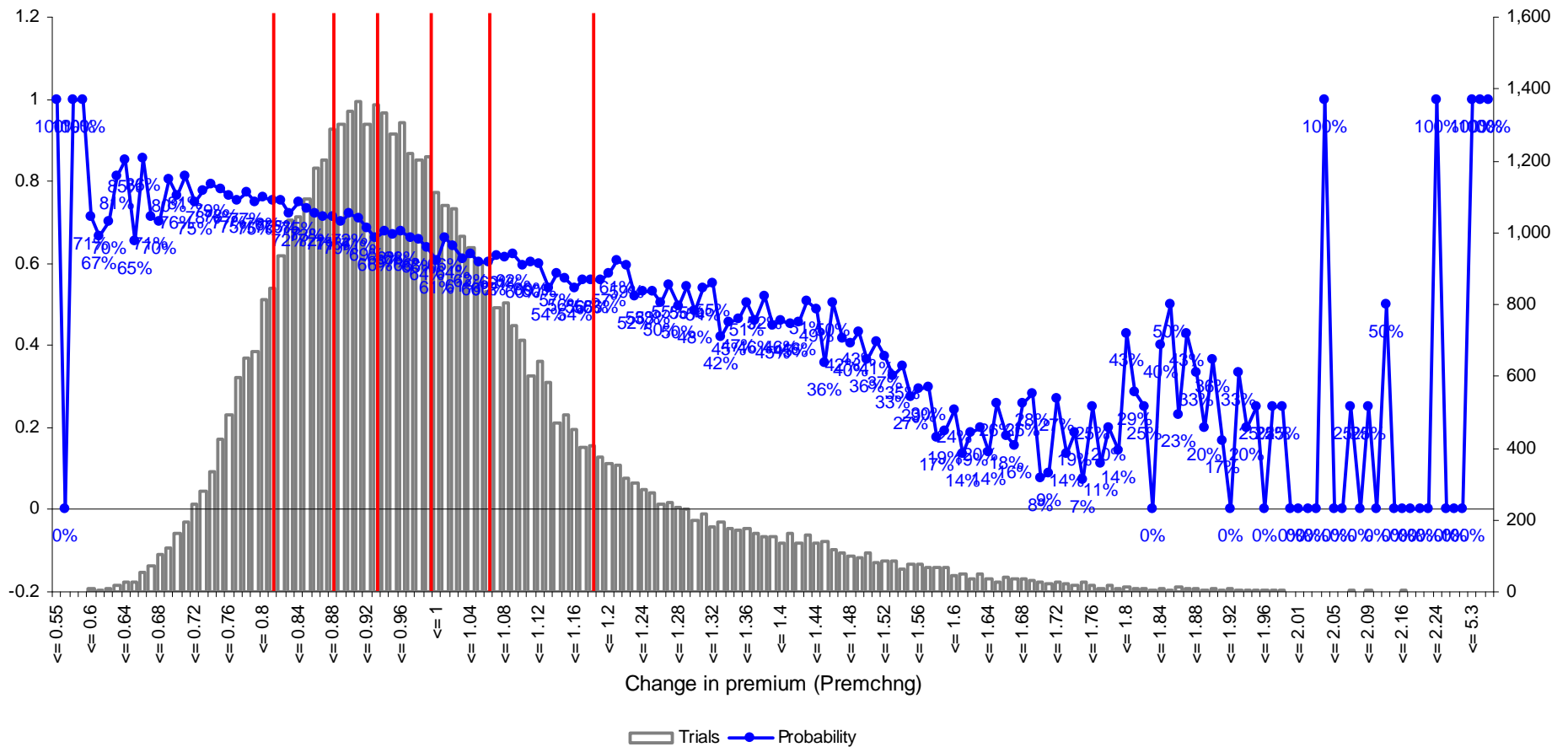
Example of problem factor



# Detailed oneway

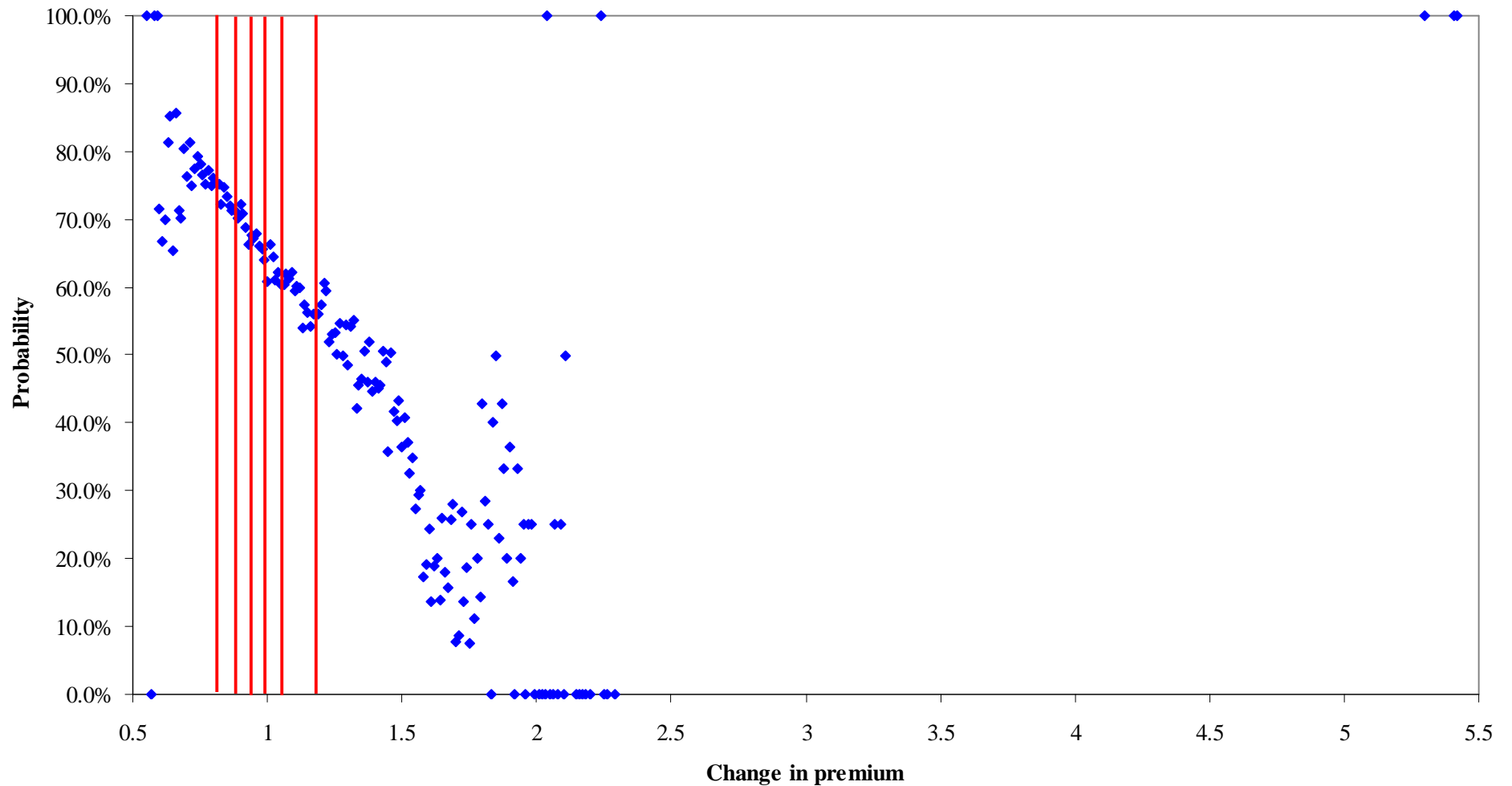
## Retention analysis

Example of problem factor



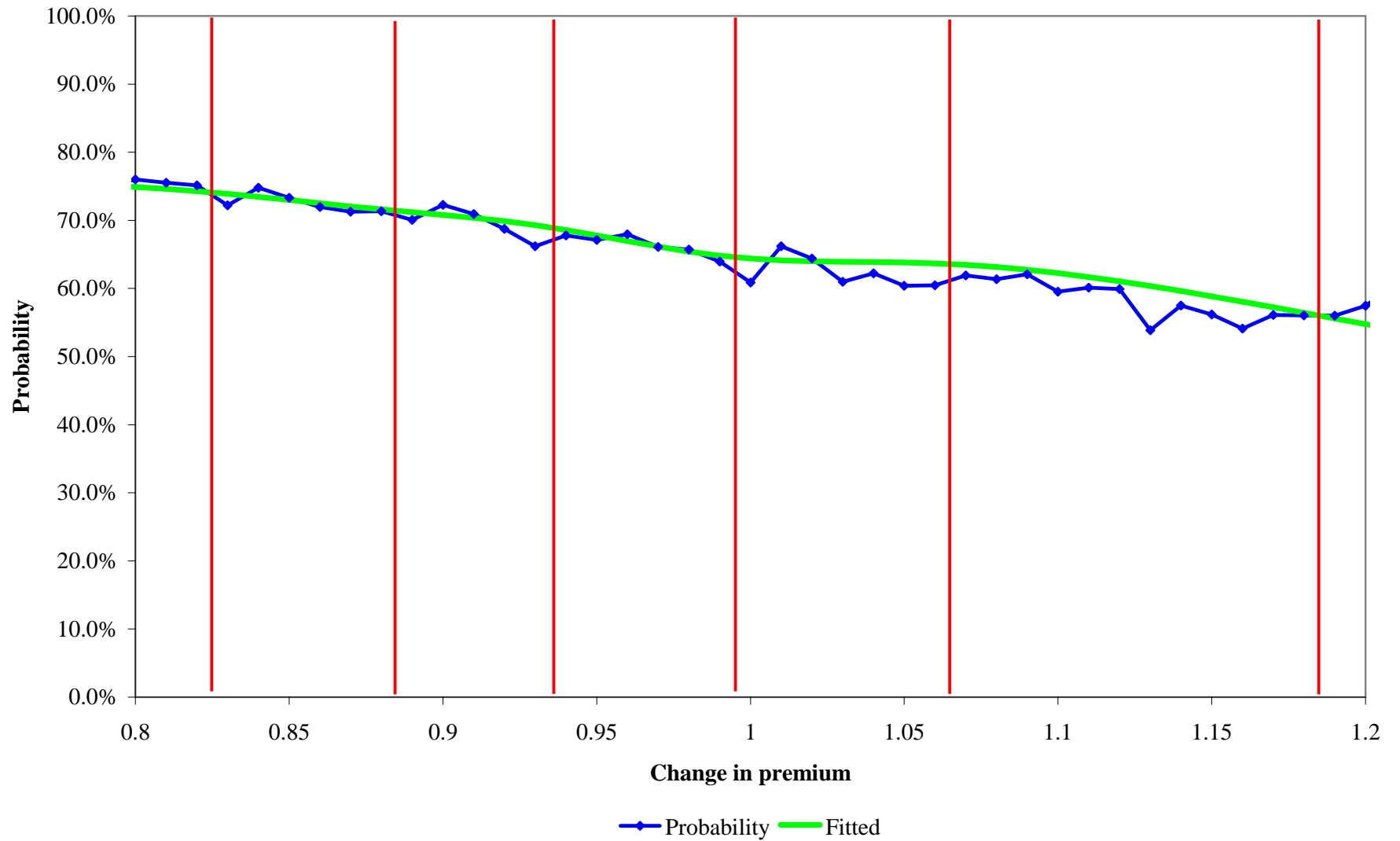


# X-Y plot

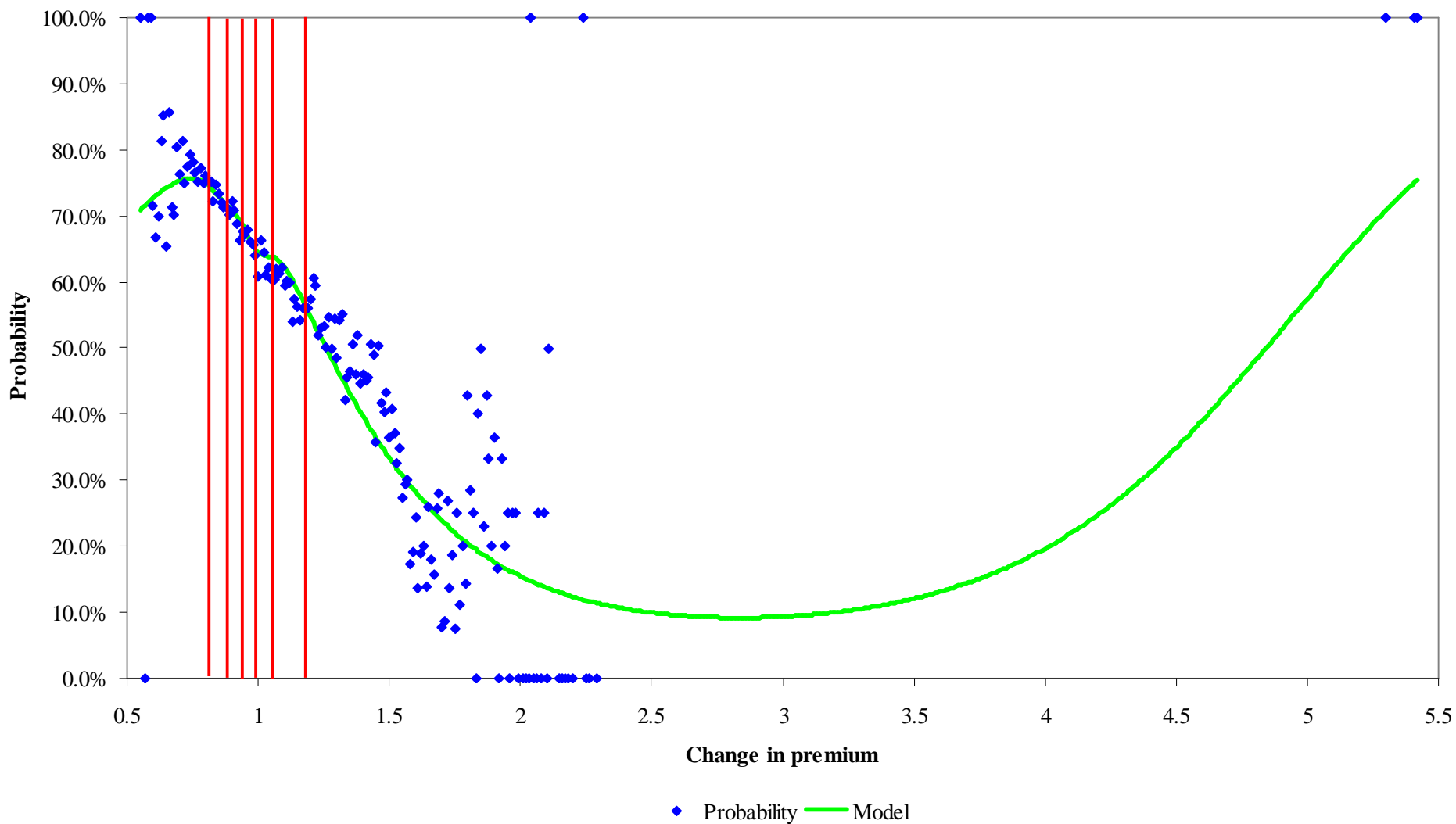


◆ Probability

# Model results



# Model results



# Agenda

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



# Interdependence of claim events

- Sometimes claim events are not individually independent
- A real life example:
  - health insurer recorded data such that individual medical 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 counts therefore appear in multiples per policy in the data
- This invalidates assumptions underlying the GLM framework

## Mathematical implications

- Poisson model used for claim counts assumes

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

- Replacing every one claim with K claims gives

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



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

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



# Mathematical implications

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

## Generalized 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

## Estimating the scale parameter

- Deviance scale for Poisson

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

- Pearson scale

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

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

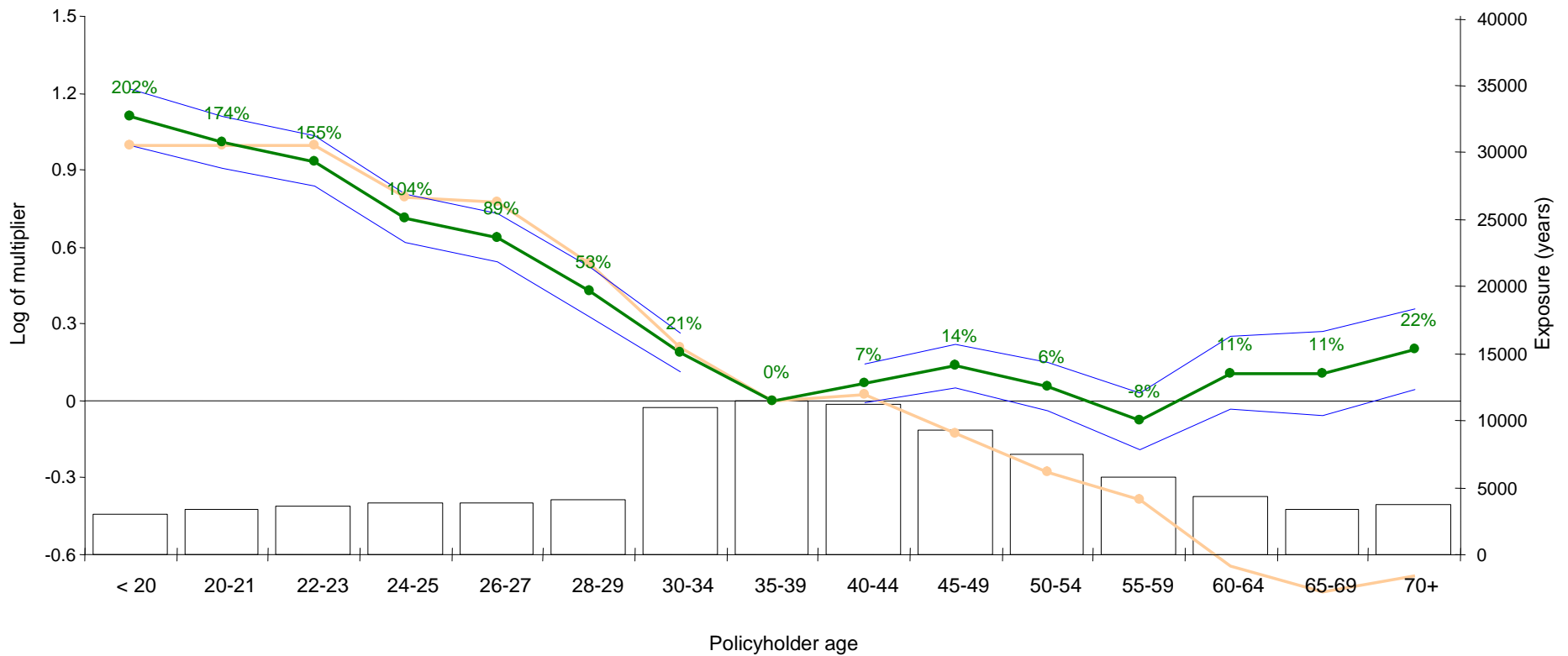
## Theoretical case study

- Fitted model to the third party property damage claim counts
- Data was adjusted by
  - multiplying exposure by 1,000
  - multiplying claims counts by 10
- Tried fitting models
  - Poisson
  - Over-dispersed Poisson (with Deviance scale)
  - Over-dispersed Poisson (with Pearson scale)

# The correct answer

## Fully worked example of the tutorial job

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



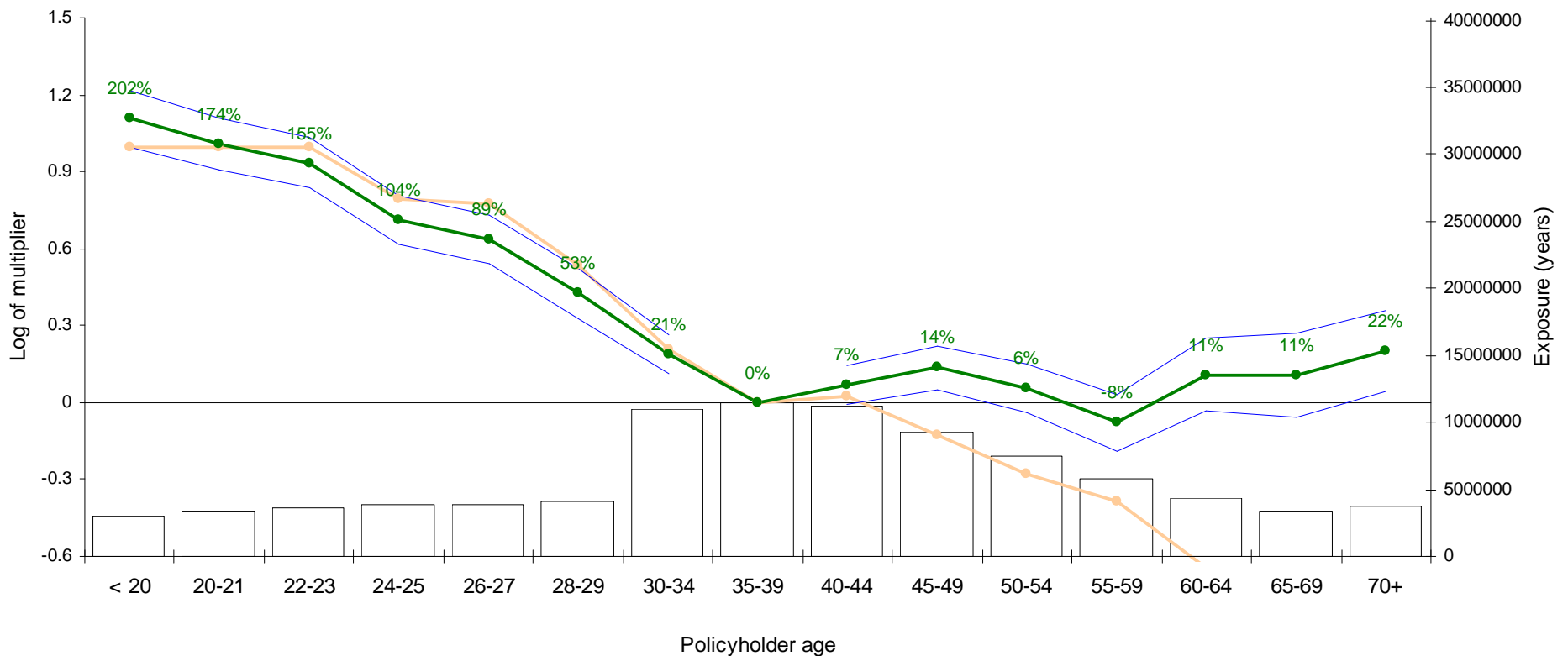
—○— Oneway relativities — Approx 95% confidence interval — Unsmoothed estimate —●— Smoothed estimate

P value = 0.0%  
Rank 10/11

# Exposure \* 1,000

## Fully worked example of the tutorial job

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



—●— Onew ay relativities 
 — Approx 95% confidence interval 
 —●— Unsmoothed estimate 
 —●— Smoothed estimate

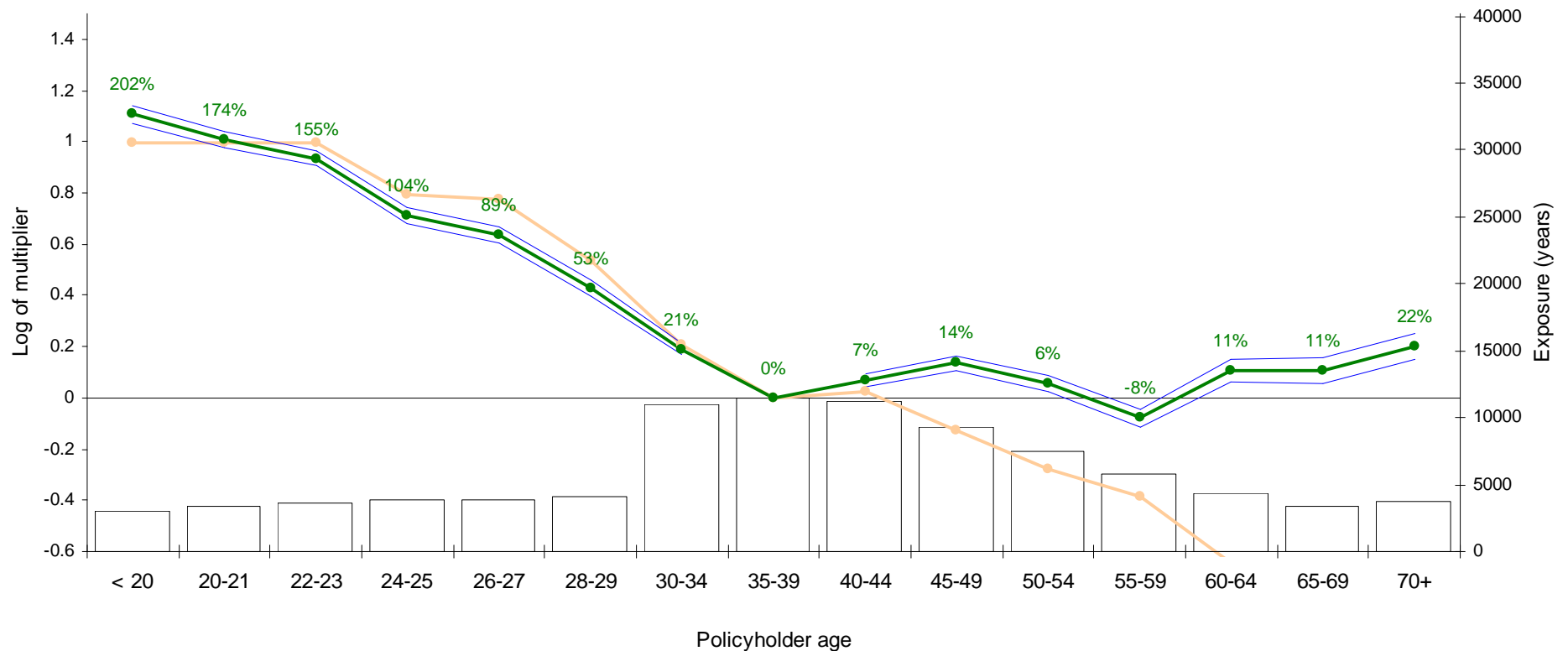
P value = 0.0%  
Rank 10/11



# Counts \* 10

## Fully worked example of the tutorial job

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



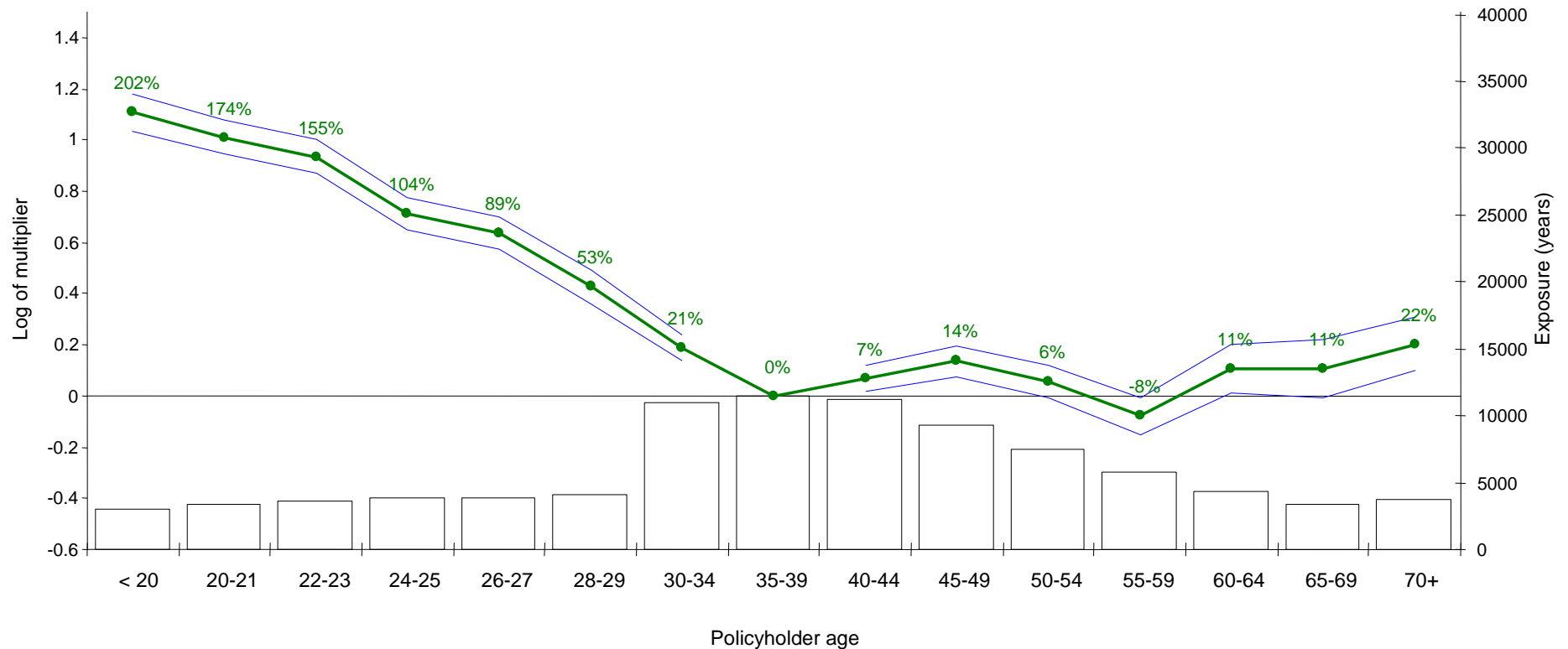
—●— Onew ay relativities 
 — Approx 95% confidence interval 
 —●— Unsmoothed estimate 
 —●— Smoothed estimate

P value = 0.0%  
Rank 7/11

# Over-dispersed Poisson (deviance)

## Fully worked example of the tutorial job

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



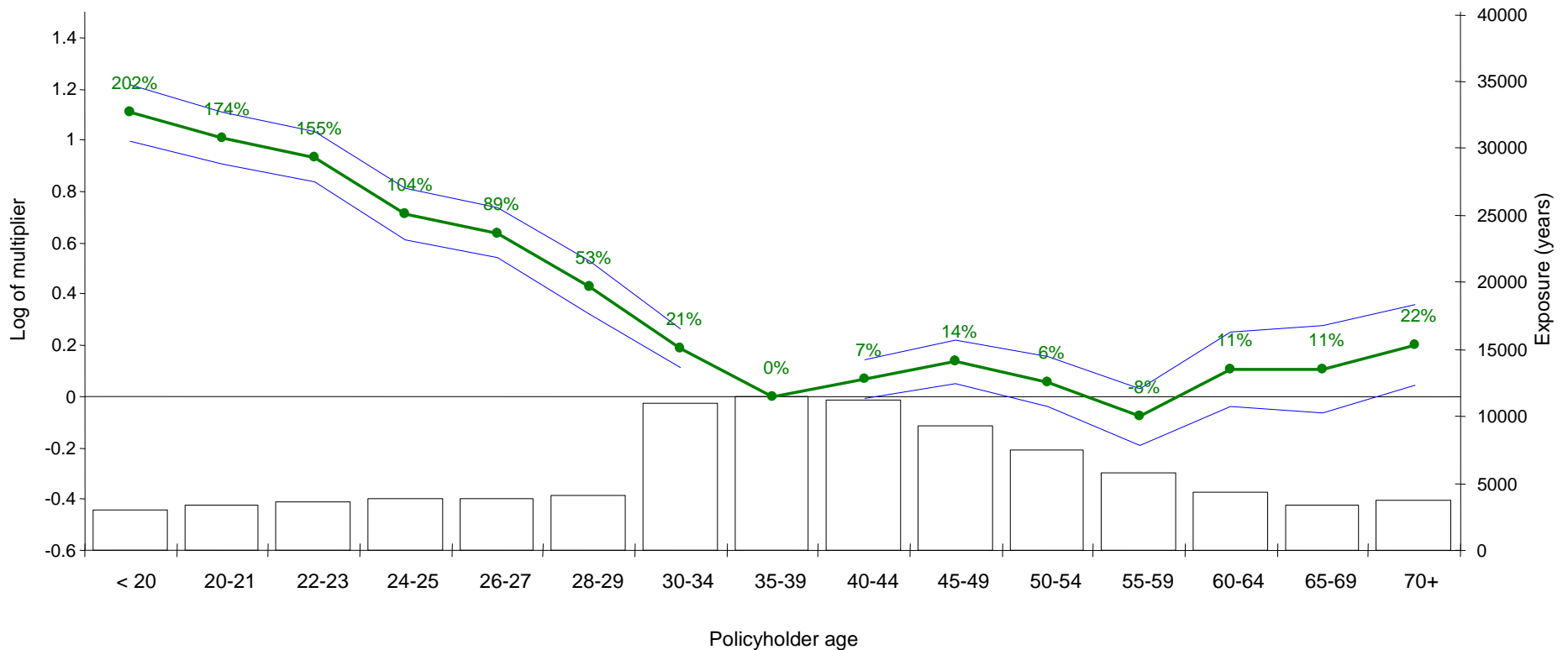
— Approx 95% confidence interval — Unsmoothed estimate — Smoothed estimate

P value = 0.0%  
Rank 7/11

# Over-dispersed Poisson (Pearson)

## Fully worked example of the tutorial job

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



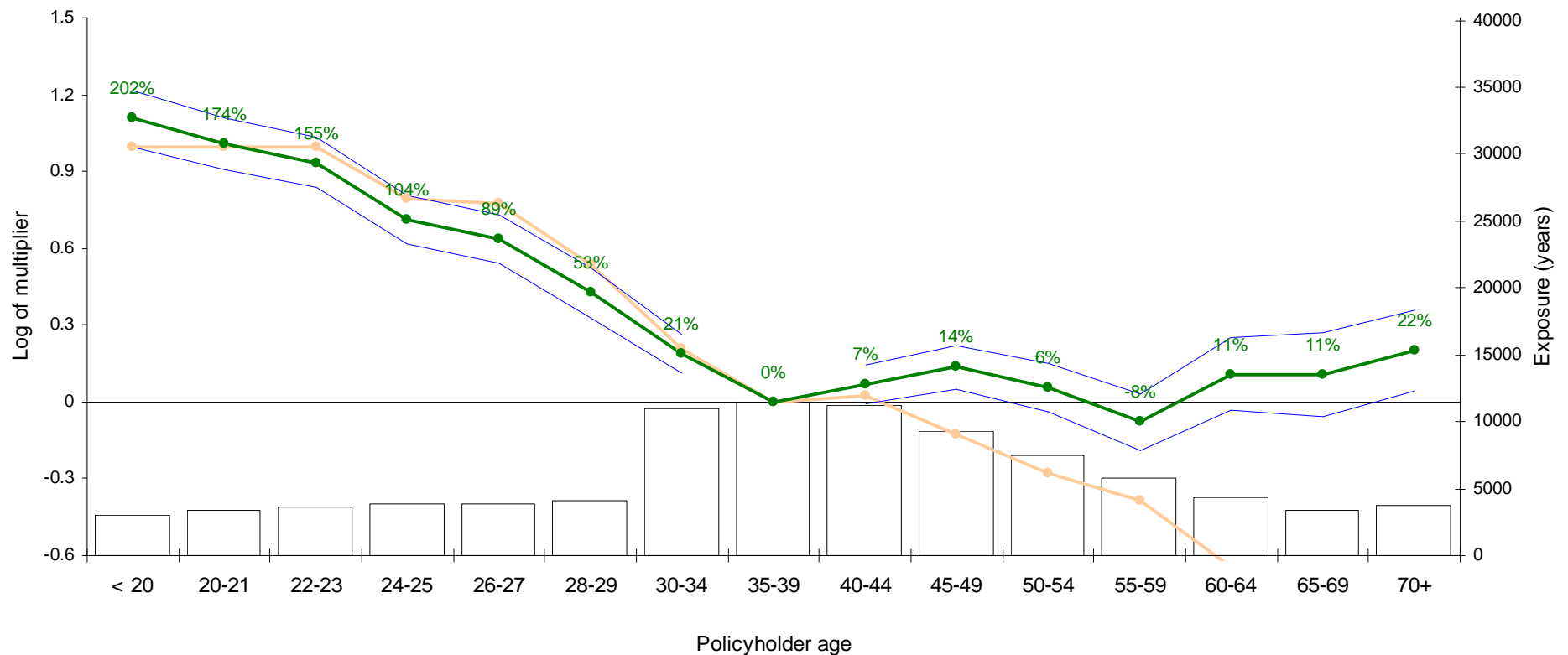
— Approx 95% confidence interval — Unsmoothed estimate — Smoothed estimate

P value = 0.0%  
Rank 7/11

# The correct answer

## Fully worked example of the tutorial job

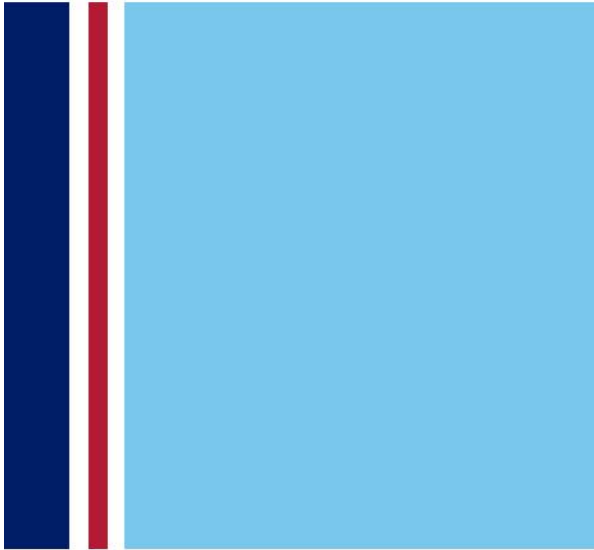
Run 8 Model 1 - Final models with analysis - TPPD numbers



—○— Onew ay relativities — Approx 95% confidence interval — Unsmoothed estimate —●— Smoothed estimate

P value = 0.0%  
Rank 10/11

watsonwyatt.com



## Practical Issues in Model Design

### **CAS Special Interest Seminar on Predictive Modeling**

Claudine Modlin, FCAS, MAAA  
October 11, 2007