# GLM III: Advanced Modeling Strategy

## 2005 CAS Seminar on Predictive Modeling

**Duncan Anderson MA FIA**

**Watson Wyatt Worldwide**

WWW.WATSONWYATT.COM
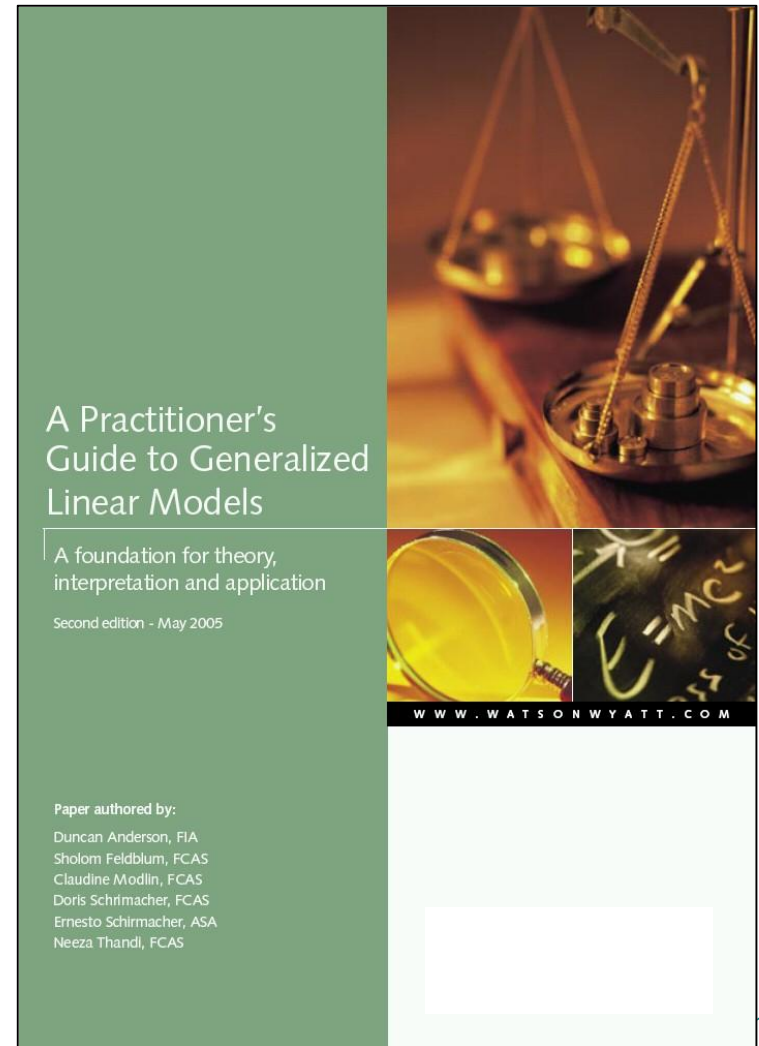
Watson Wyatt
*Worldwide*

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# "A Practitioner's Guide to GLMs"

- 2004 CAS Discussion Paper Program

- Discusses
  - testing the link function
  - the Tweedie distribution
  - aliasing / near aliasing
  - combining models across claim types
  - restricted models

- Copies available here

A Practitioner's
Guide to Generalized
Linear Models

A foundation for theory,
interpretation and application

Second edition - May 2005

WWW.WATSONWYATT.COM

Paper authored by:

Duncan Anderson, FIA
Sholom Feldblum, FCAS
Claudine Modlin, FCAS
Doris Schirmacher, FCAS
Ernesto Schirmacher, ASA
Neeza Thandi, FCAS

# Generalized linear models

$$E[Y_i] = \mu_i = g^{-1}(\Sigma X_{ij}\beta_j + \xi_i)$$

$$Var[Y_i] = \phi.V(\mu_i)/\omega_i$$

- Consider all factors simultaneously
- Provide statistical diagnostics
- Allow for nature of random process
- Robust and transparent
- Increasingly a global industry standard

# Generalized linear models

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

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

# Generalized linear models

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

Link function

Y-variate

Design matrix

Parameter estimates

Offset term

# Generalized linear models

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

Some function
(user defined)

Observed thing
(data)

Some matrix based
on data
(user defined)

Parameters
to be
estimated
(the answer!)

Known
effects

# Generalized linear models

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

Scale parameter

Prior weights

Variance function

- Usually assume exponential family, eg

- $\phi = \sigma^2$ (estimated), $V(x) = 1$ $\Rightarrow$ $\text{Var}[Y_i] = \sigma^2$   Normal

- $\phi = 1$  (specified),  $V(x) = x$ $\Rightarrow$ $\text{Var}[Y_i] = \mu_i$   Poisson

- $\phi = k$  (estimated), $V(x) = x^2$ $\Rightarrow$ $\text{Var}[Y_i] = k\mu_i^2$ Gamma

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# Model testing

- Use only those factors which are predictive
  - standard errors of parameter estimates
  - F tests / $\chi^2$ tests on deviances
  - stepwise approach (helpful if used with care)
  - consistency over time
  - human intuition

- Make sure the model is reasonable
  - variance function: residual plots
    (histograms / Q-Q / residual vs fitted value etc)
  - outliers: leverage / Cook's distance
  - link function: Box-Cox

# Box-Cox link function investigation

- GLM structure is

  $E[\underline{Y}] = \mu = g^{-1}(\mathbf{X}.\underline{\beta} + \underline{\xi})$     $Var[\underline{Y}] = \phi.V(\mu) / \underline{\omega}$

- Box Cox transforms defines

  $g(x) = ( x^{\lambda} - 1 ) / \lambda$ for $\lambda \neq 0$, $\ln(x)$ for $\lambda = 0$

- $\lambda = 1 \Rightarrow g(x) = x - 1 \Rightarrow$ additive (with base level shift)

- $\lambda \to 0 \Rightarrow g(x) \to \ln(x) \Rightarrow$ multiplicative (via maths)

- $\lambda = -1 \Rightarrow g(x) = 1 - 1/x \Rightarrow$ inverse (with base level shift)

- Try different values of $\lambda$ and measure goodness of fit to see which fits experience best

# Box-Cox link function investigation
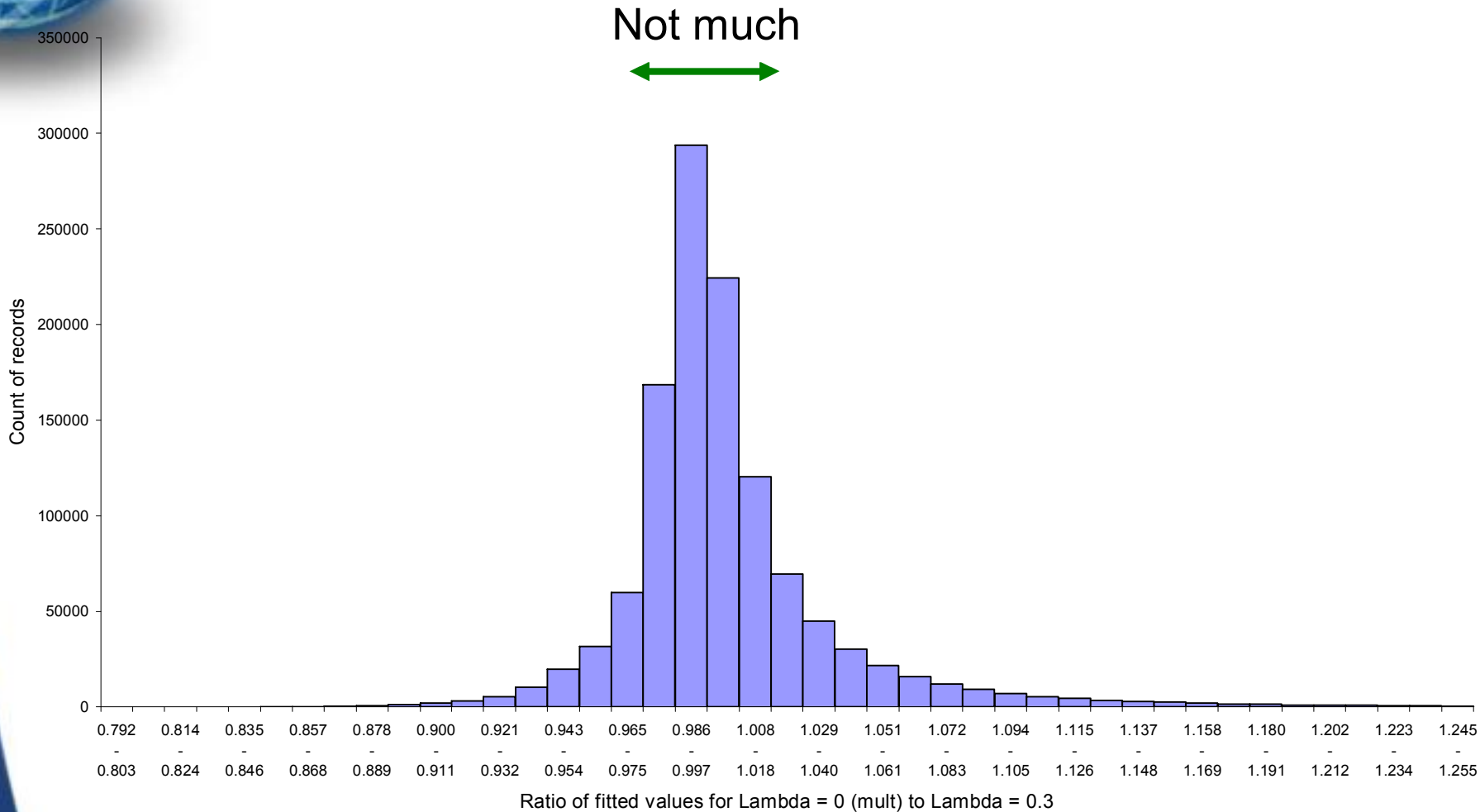## Auto third party property damage frequencies

# Box-Cox link function investigation
## Auto third party property damage average amounts

# Box-Cox link function investigation
## Comparing fitted values of different link functions



Not much

Count of records

350000
300000
250000
200000
150000
100000
50000
0

0.792 - 0.803 | 0.814 - 0.824 | 0.835 - 0.846 | 0.857 - 0.868 | 0.878 - 0.889 | 0.900 - 0.911 | 0.921 - 0.932 | 0.943 - 0.954 | 0.965 - 0.975 | 0.986 - 0.997 | 1.008 - 1.018 | 1.029 - 1.040 | 1.051 - 1.061 | 1.072 - 1.083 | 1.094 - 1.105 | 1.115 - 1.126 | 1.137 - 1.148 | 1.158 - 1.169 | 1.180 - 1.191 | 1.202 - 1.212 | 1.223 - 1.234 | 1.245 - 1.255

Ratio of fitted values for Lambda = 0 (mult) to Lambda = 0.3

# Agenda

- Introduction

- Testing the link function

- **The Tweedie distribution**

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# Standard approach

BI    Freq    x    Amt    = Cost 1

PD    Freq    x    Amt    = Cost 2

MED    Freq    x    Amt    = Cost 3

COL    Freq    x    Amt    = Cost 4

OTC    Freq    x    Amt    = Cost 5

# Tweedie distributions

- Incurred losses have a point mass at zero and then a continuous distribution

- Poisson and gamma not suited to this

- Tweedie distribution can have point mass at zero and parameters which can alter the shape to be like Poisson and gamma above zero

$$f_Y(y;\theta,\lambda,\alpha) = \sum_{n=1}^{\infty} \frac{\left\{ (\lambda\omega)^{1-\alpha} \kappa_\alpha(-1/y) \right\}^n}{\Gamma(-n\alpha)n! \, y} \cdot \exp\left\{ \lambda\omega[\theta_0 y - \kappa_\alpha(\theta_0)] \right\} \quad \text{for } y > 0$$

$$p(Y=0) = \exp\left\{ -\lambda\omega\kappa_\alpha(\theta_0) \right\}$$

# Tweedie distributions

$$\text{Tweedie: } \phi = k, \ V(x) = x^p \Rightarrow \text{Var}[\underline{Y}] = k\underline{\mu}^p$$

- p=1 corresponds to Poisson, p=2 to gamma

- Defines a valid distribution for p<0, 1<p<2, p>2

- Can be considered as Poisson/gamma process for 1<p<2

- Need to estimate both k and p when fitting models
  - often estimate a where p = (2-a)/(1-a)

- Typical values of p for insurance incurred claims around, or just under, 1.5

# Generalised linear models

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

Normal: $\phi = \sigma^2$, $V(x) = 1 \Rightarrow \mathrm{Var}[\underline{Y}] = \sigma^2 . \underline{1}$

Poisson: $\phi = 1$, $V(x) = x \Rightarrow \mathrm{Var}[\underline{Y}] = \underline{\mu}$

Gamma: $\phi = k$, $V(x) = x^2 \Rightarrow \mathrm{Var}[\underline{Y}] = k\underline{\mu}^2$

Tweedie: $\phi = k$, $V(x) = x^p \Rightarrow \mathrm{Var}[\underline{Y}] = k\underline{\mu}^p$

# Example 1: frequency

**Comparison of Tweedie model with traditional frequency/amounts approach**

Run 7 Model 2 - Frequency



Bonus Malus

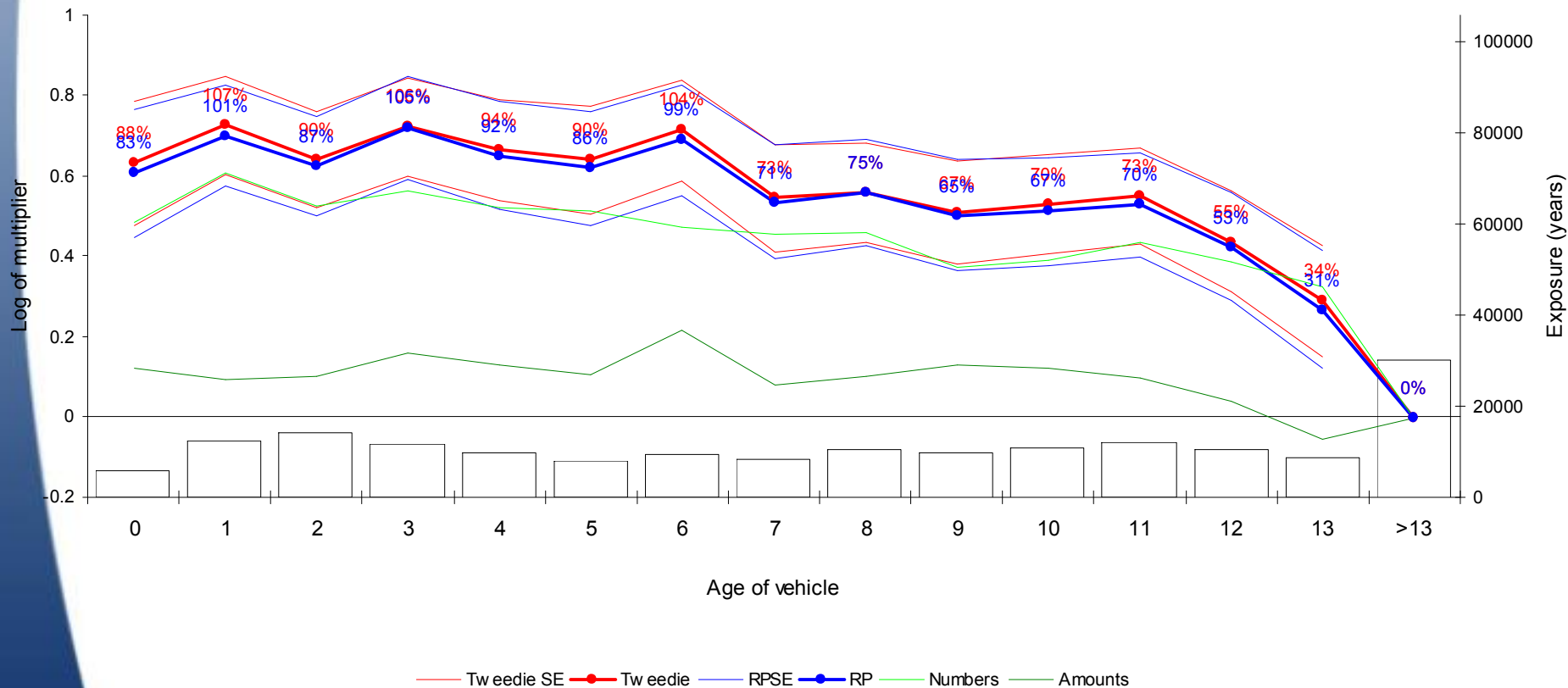Legend: Oneway relativities — Approx 95% confidence interval — Unsmoothed estimate — Smoothed estimate

P value = 0.0%

Rank 12/12

# Example 1: amounts

**Comparison of Tweedie model with traditional frequency/amounts approach**

Run 7 Model 6 - Amounts



EXCLUDED FACTOR

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

P value = 50.9%
Rank 4/12

# Example 1: traditional RP vs Tweedie

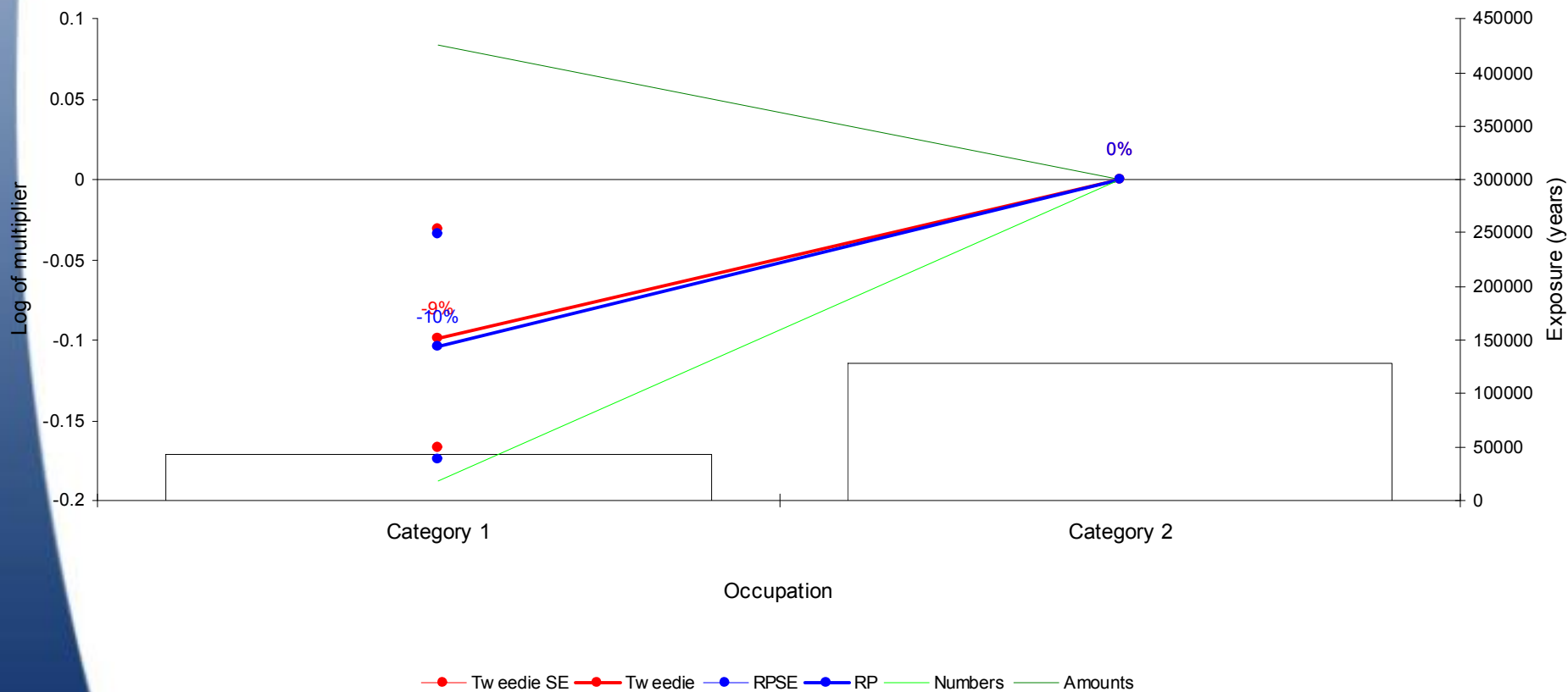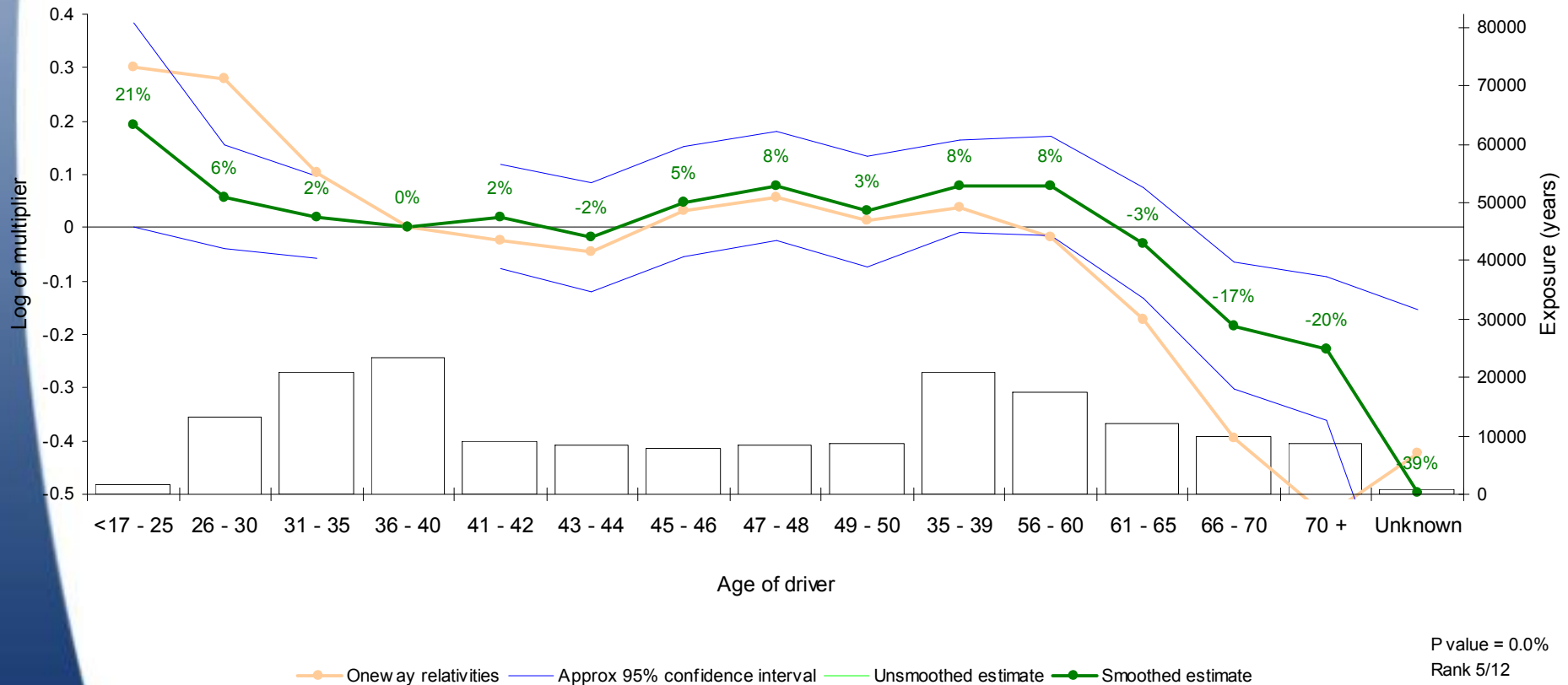**Comparison of Tweedie model with traditional frequency/amounts approach**

Run 11 Model 2 - Tweedie Models

# Example 2: frequency

## Comparison of Tweedie model with traditional frequency/amounts approach
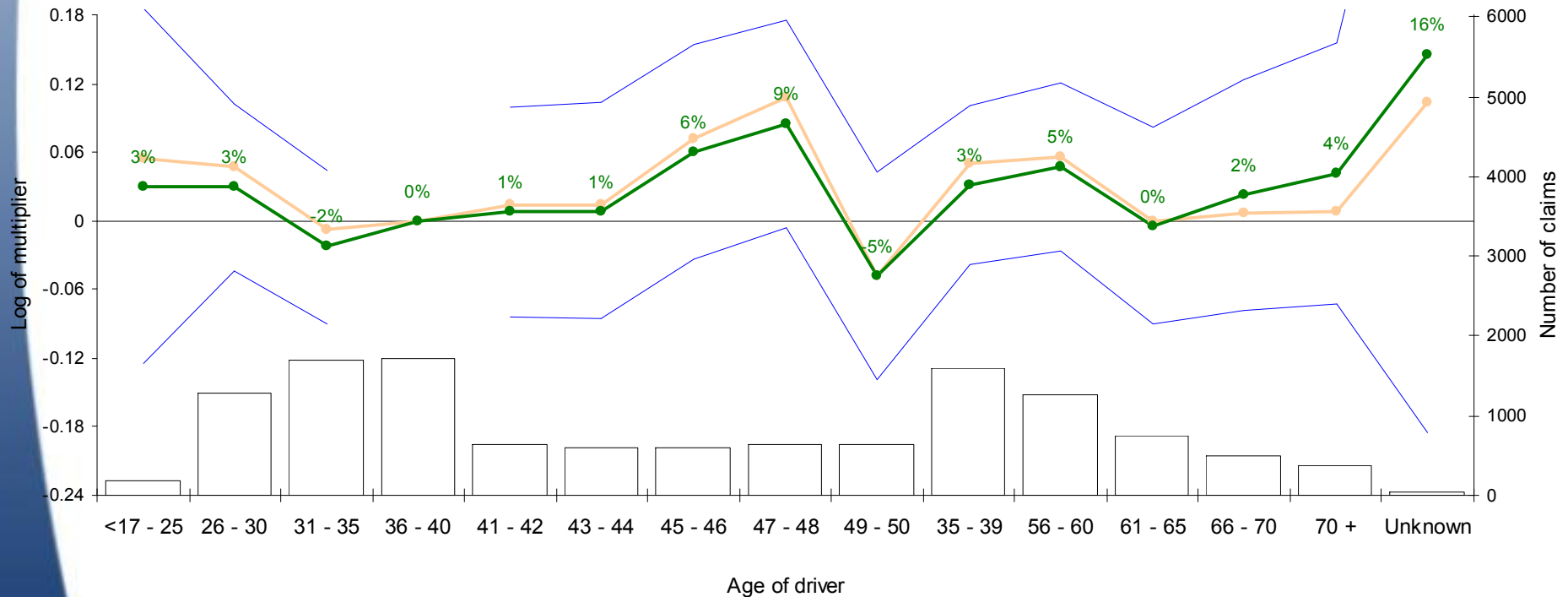Run 7 Model 1 - Frequency



P value = 0.0%

Rank 12/12

Legend: Oneway relativities — Approx 95% confidence interval — Unsmoothed estimate — Smoothed estimate

X-axis: Age of vehicle (0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, >13)

Left Y-axis: Log of multiplier

Right Y-axis: Exposure (years)

Smoothed estimate values: 62%, 84%, 69%, 75%, 68%, 67%, 61%, 58%, 58%, 45%, 48%, 54%, 47%, 38%, 0%

# Example 2: amounts

**Comparison of Tweedie model with traditional frequency/amounts approach**
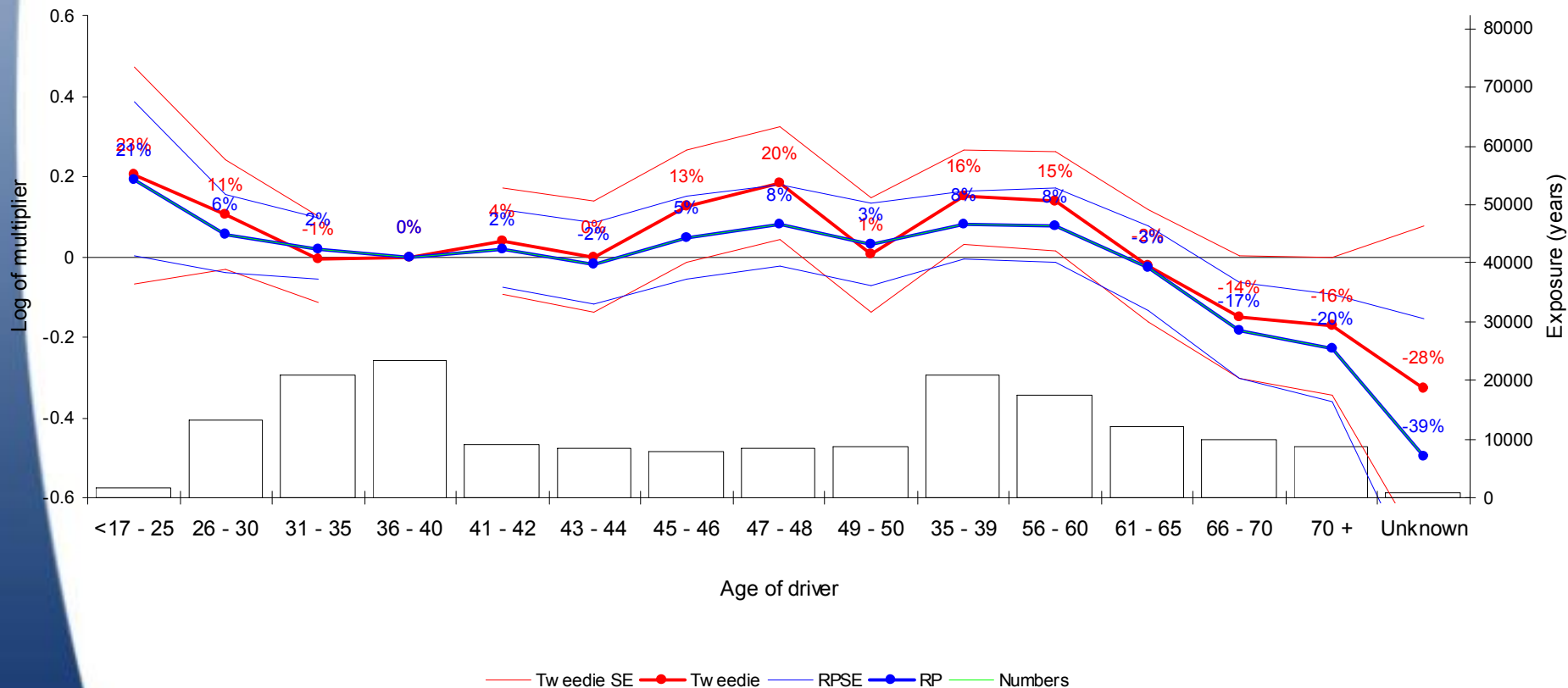
Run 7 Model 5 - Amounts



P value = 0.0%

Rank 5/7

# Example 2: traditional RP vs Tweedie

**Comparison of Tweedie model with traditional frequency/amounts approach**
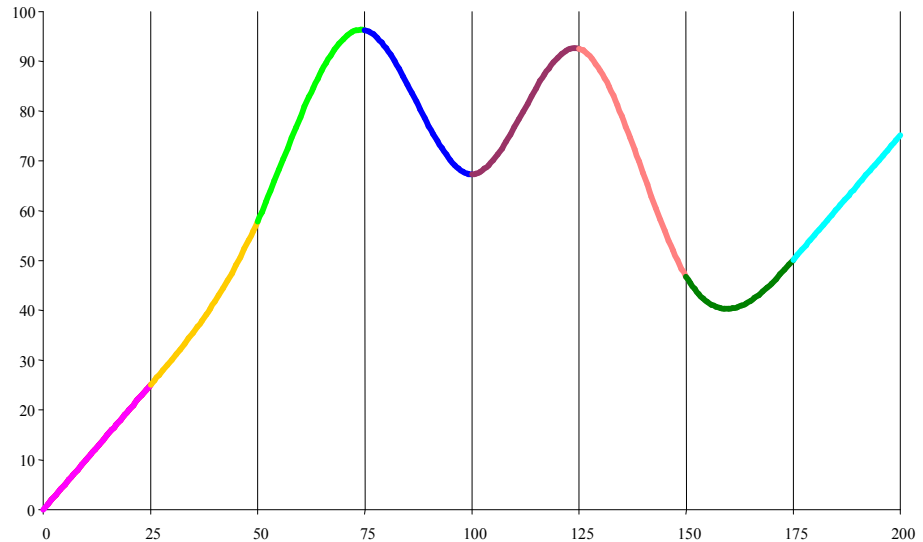
Run 11 Model 1 - Tweedie Models

# Example 3: traditional RP vs Tweedie

**Comparison of Tweedie model with traditional frequency/amounts approach**

Run 11 Model 1 - Tweedie Models

# Example 4: amounts

**Comparison of Tweedie model with traditional frequency/amounts approach**

Run 7 Model 5 - Amounts



EXCLUDED FACTOR

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

P value = 50.6%

Rank 4/9

# Example 4: traditional RP vs Tweedie



**Comparison of Tweedie model with traditional frequency/amounts approach**

Run 11 Model 1 - Tweedie Models

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# 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"

# Cubic splines

- Each polynomial is a cubic
    - $a + bx + cx^2 + dx^3$

- "Smoothness" at interior knots is defined as:
    - continuous
    - continuous first derivative
    - continuous second derivative

# Regression splines

- The position of the knots is specified by the user

- Standard GLMs can be used by careful definition of variates

- Pros
  - fits easily into existing structures
  - no complex re-sampling needed

- Cons
  - position of knots can effect final answer

# Smoothing splines

- One knot at each unique data value

- Additional curvature penalty prevents over fitting

- Curvature penalty selected by repeatedly sampling subsets and optimising generalised goodness of fit measure such as AIC

- Pros

  - allows data to guide final result

- Cons

  - 100s of knots required

  - optimisation process is time-consuming

  - difficult to produce new fitted values

# "Easy" regression splines

- Fit a cubic over the whole range

  - simply define $x$, $x^2$ and $x^3$ as variates and include in the model

- Fit additional cubic "correction" variates for each interval, defined as

  - 0 if $x < k_r$

  - $((x - k_{r+1})/(k_r - k_{r+1}))^3$ otherwise

# "Easy" regression splines

# "Easy" regression splines

# "Easy" regression splines

# "Easy" regression splines

- "Correction" variates get large quickly

- In practice GLM process can struggle with these large numbers

- Alternate basis is clearly desirable

# B-Splines

- Set of basis functions usually covering four segments (defined by five knots)

- Each function is itself a cubic spline

- Each basis function has the same shape, except for the three basis functions at each extreme which occupy fewer than four segments

# B-Splines
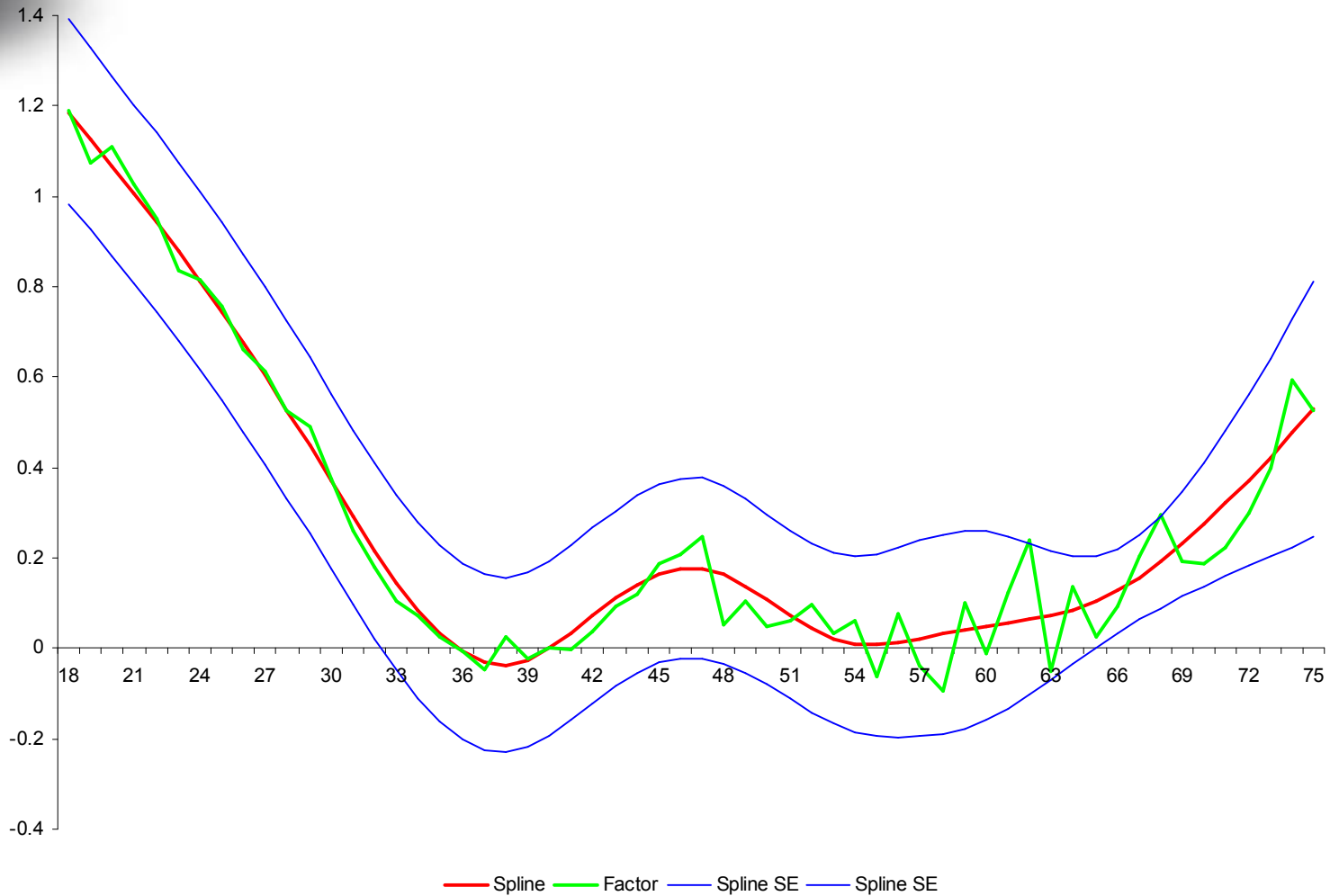
# B-Splines

# B-Splines

# B-Splines
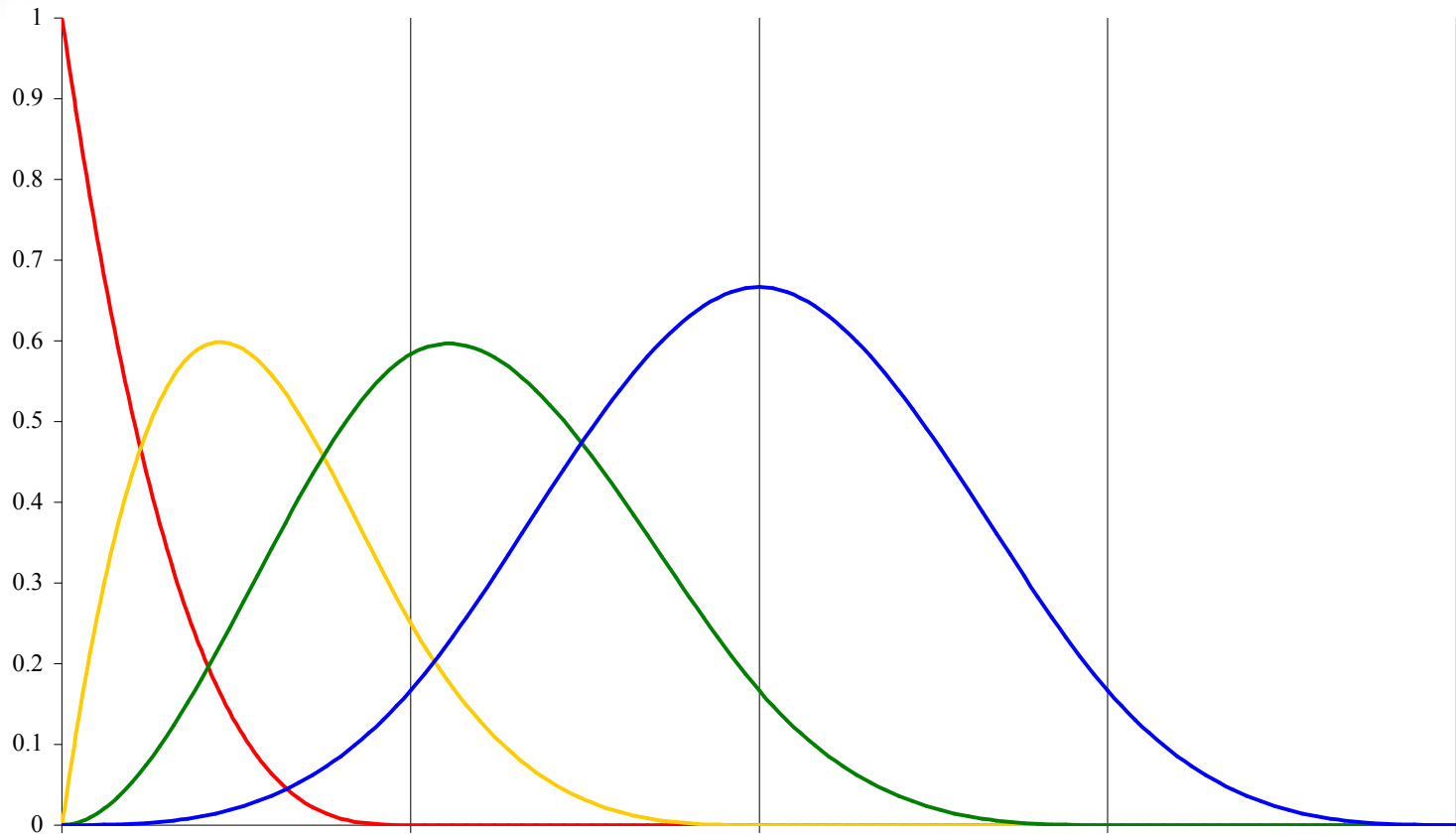
# Example



Factor

# Example
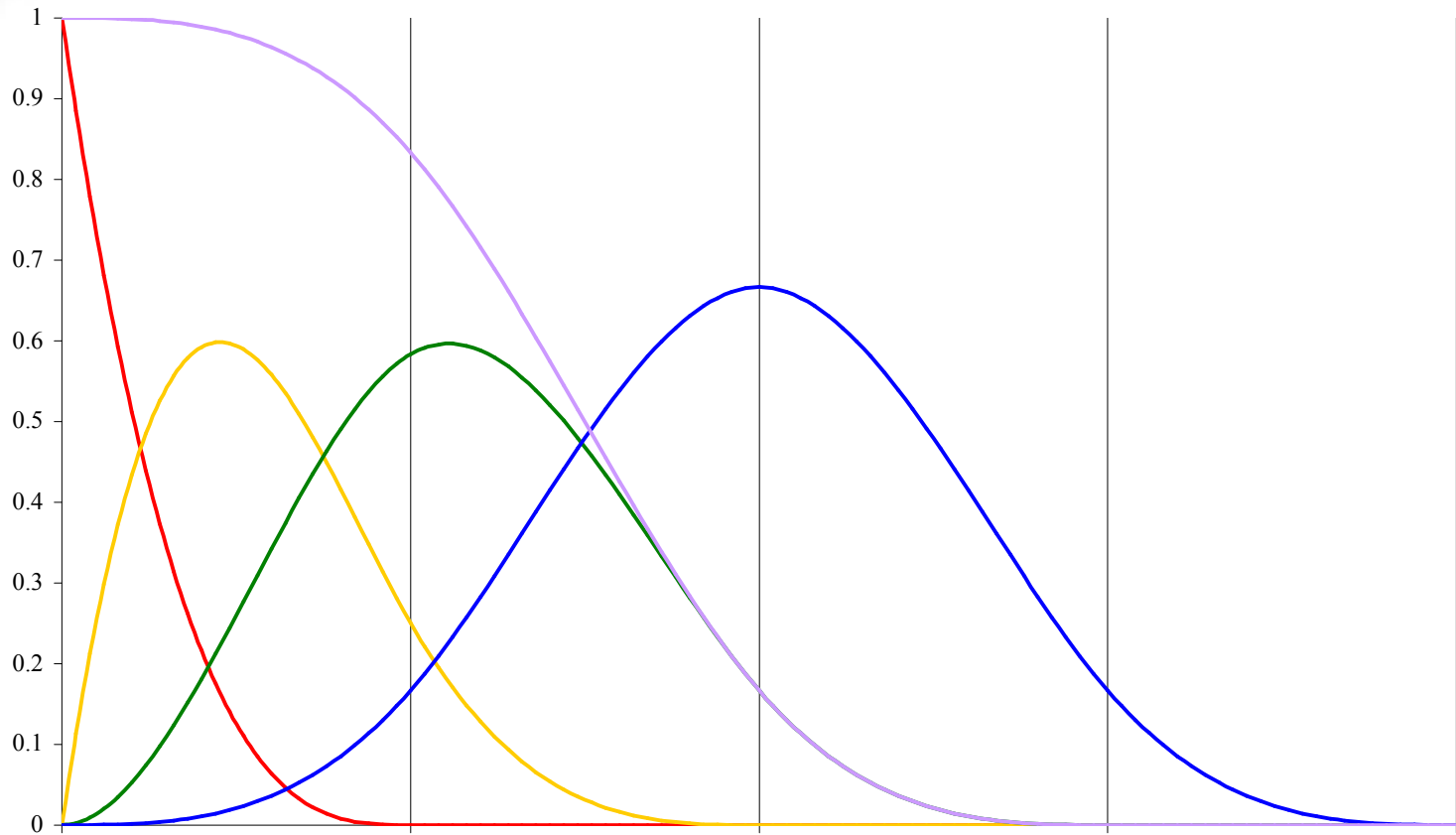
# Example

# Example

# Extrapolation

- Can combine the three boundary basis functions to achieve linear extrapolation

- Sum of three functions tends to 1 at the exterior knot, and continues as 1 when extrapolated

- Can combine the last two functions to give function that tends to a straight line at the exterior knot, and can be extrapolated as such
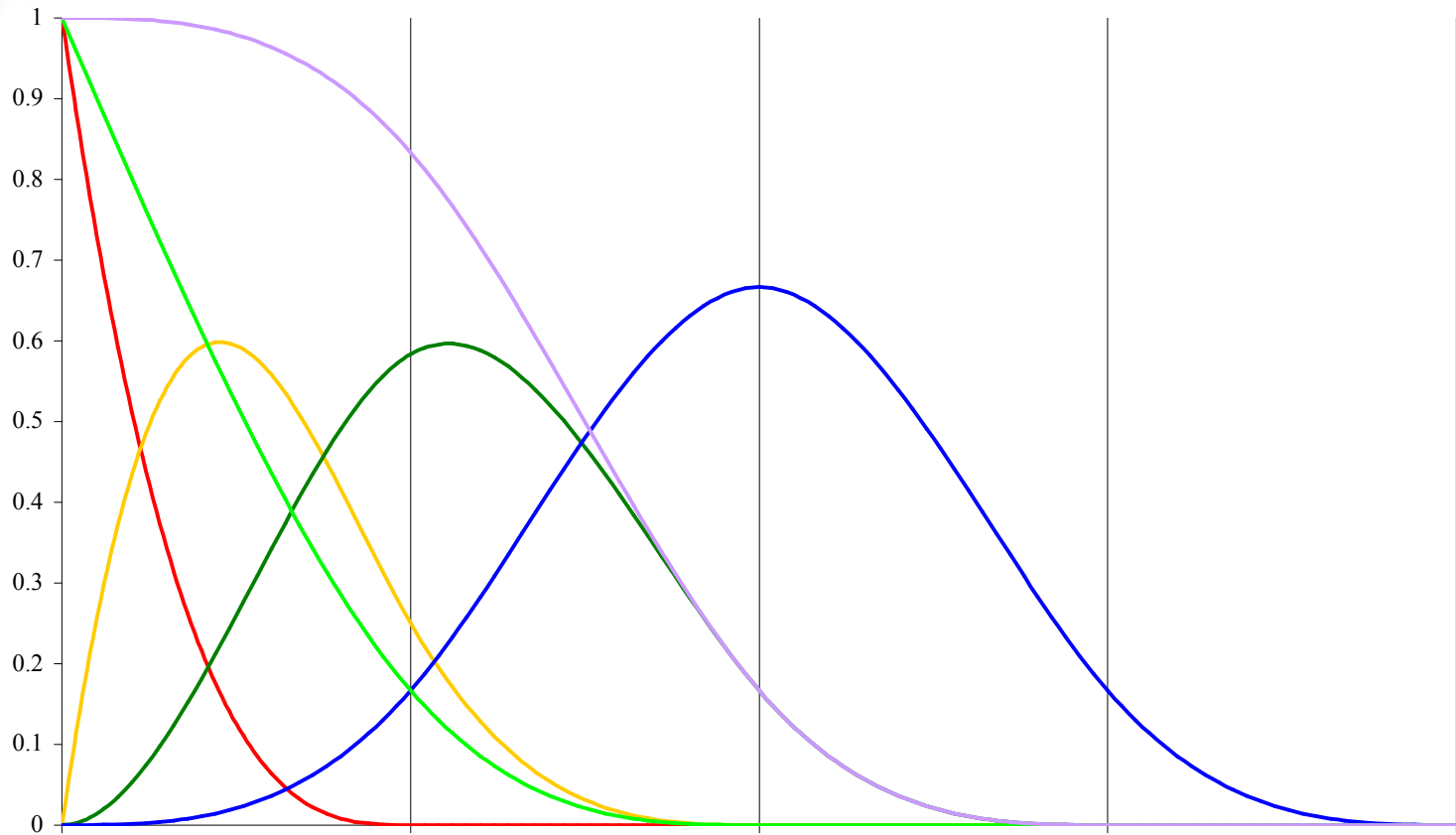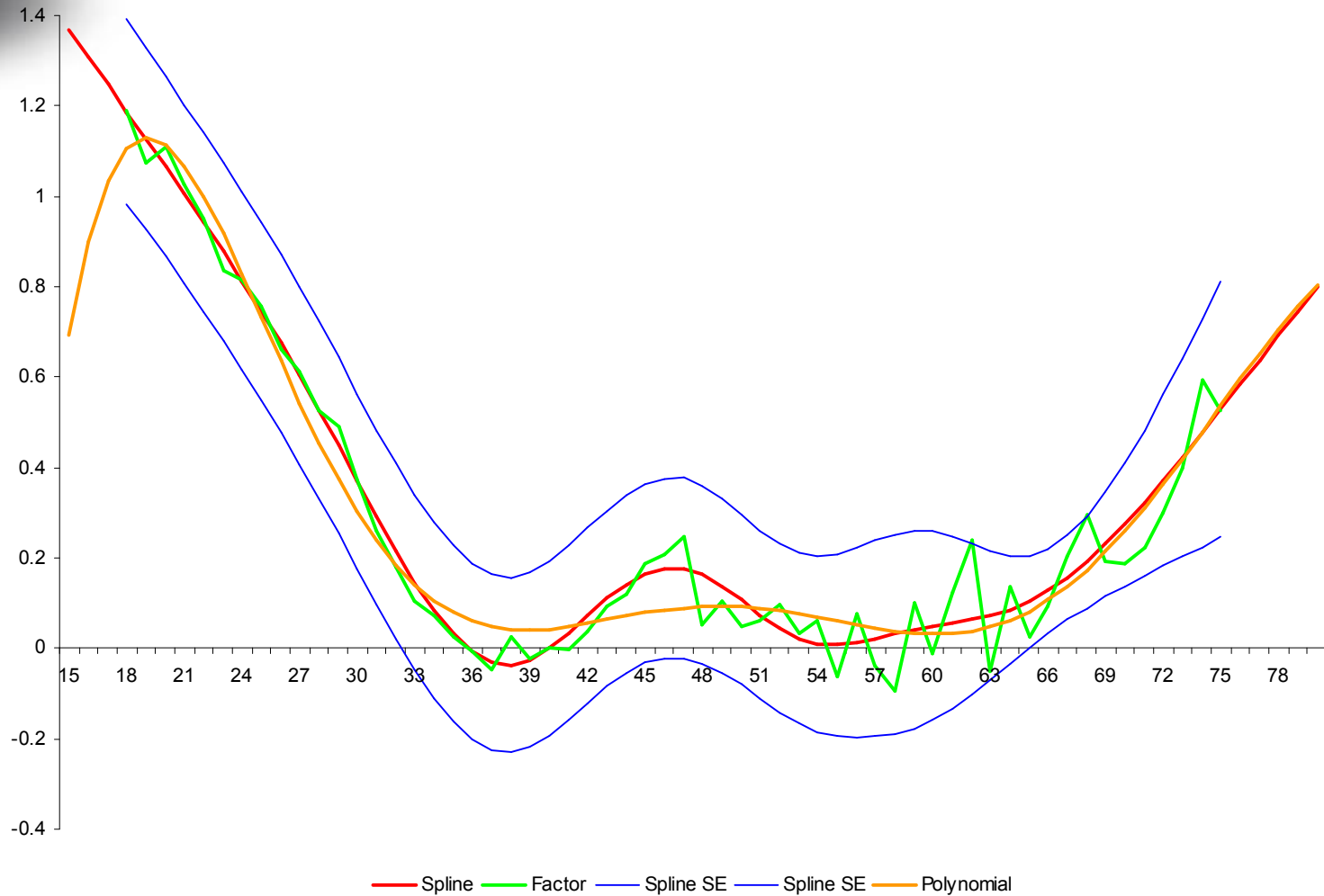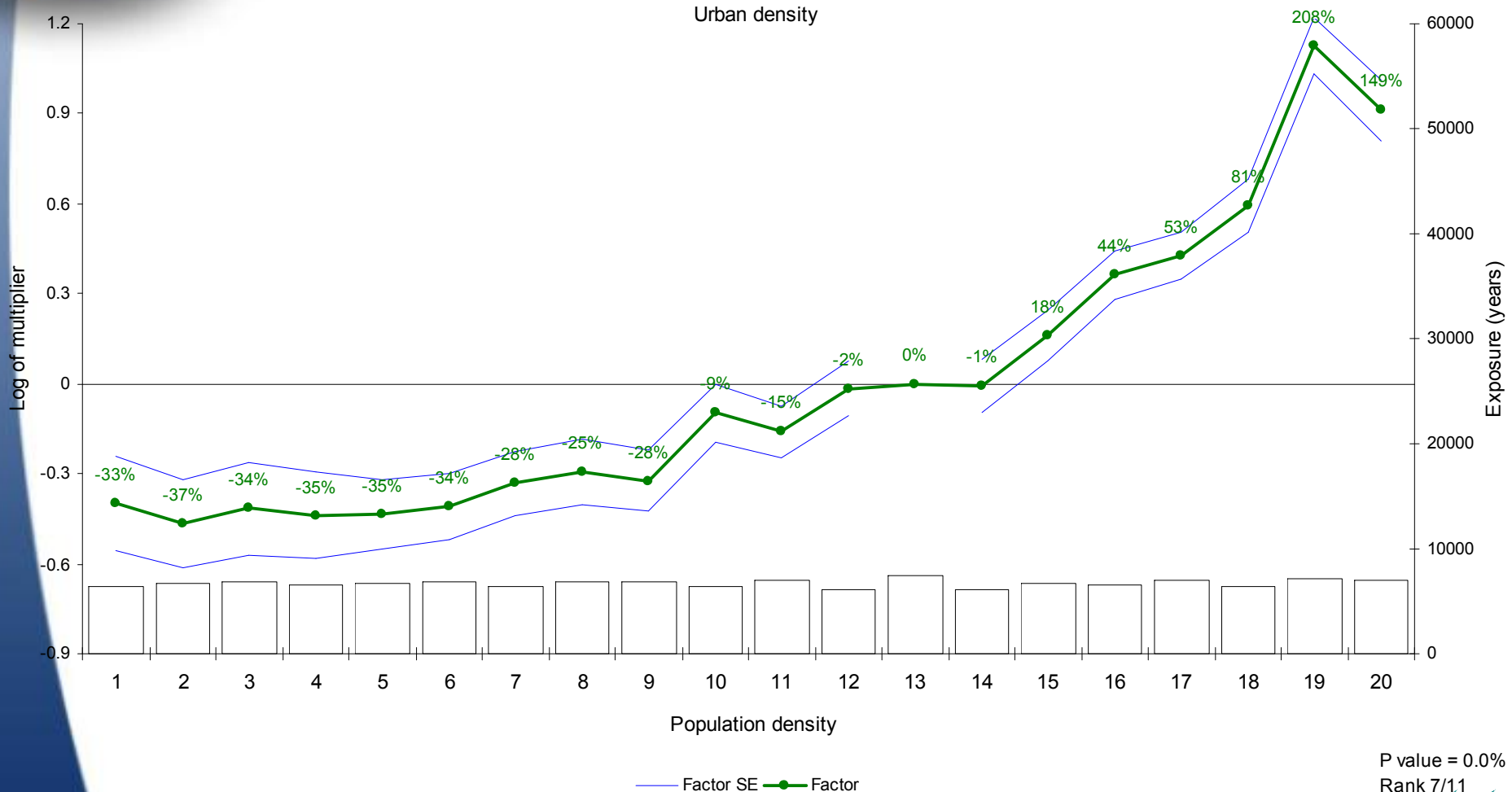
# B-Splines

# B-Splines

# B-Splines

# Example

# Further example


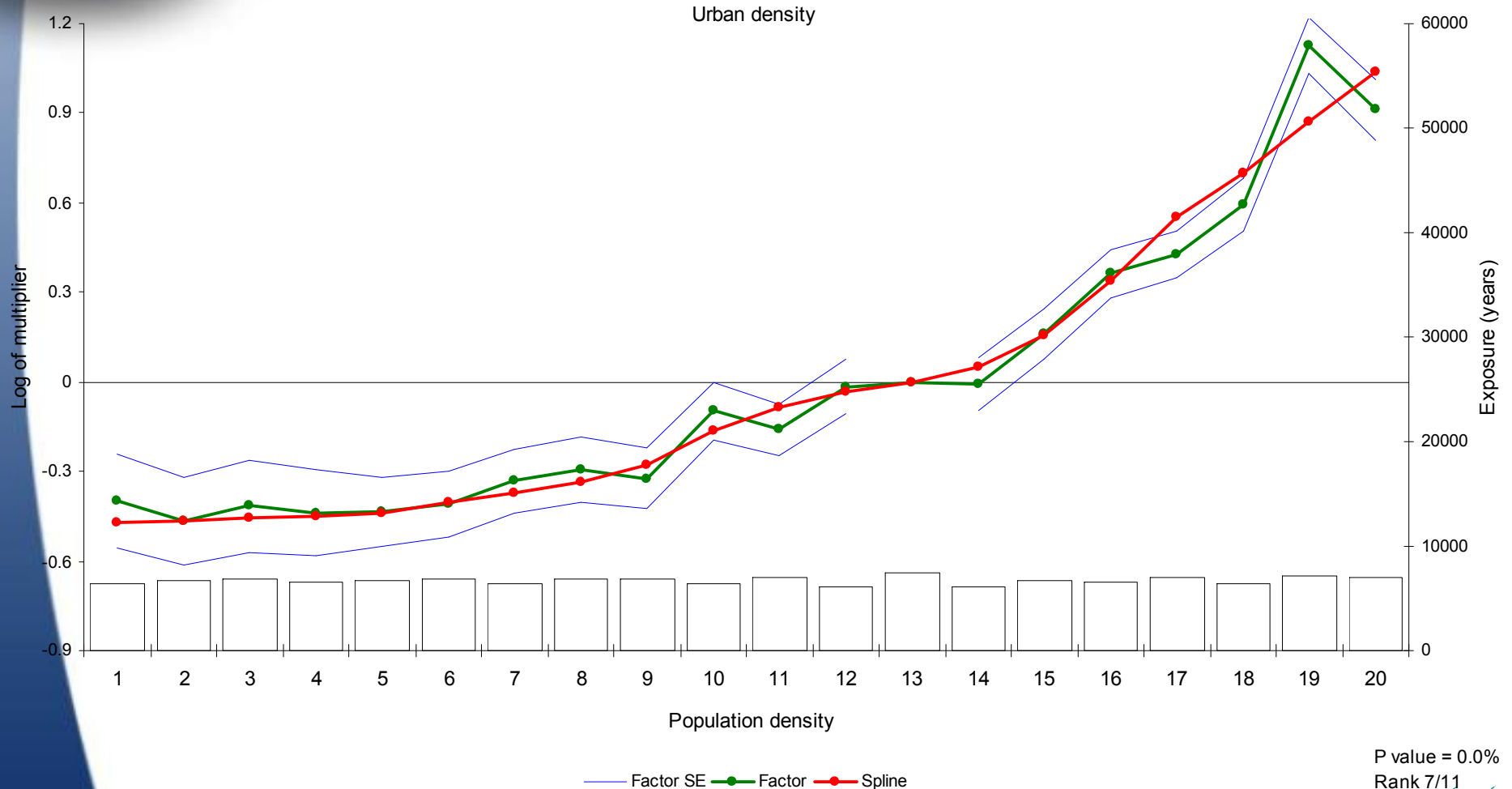
**Comparison of factor with spline**

Urban density

# Further example
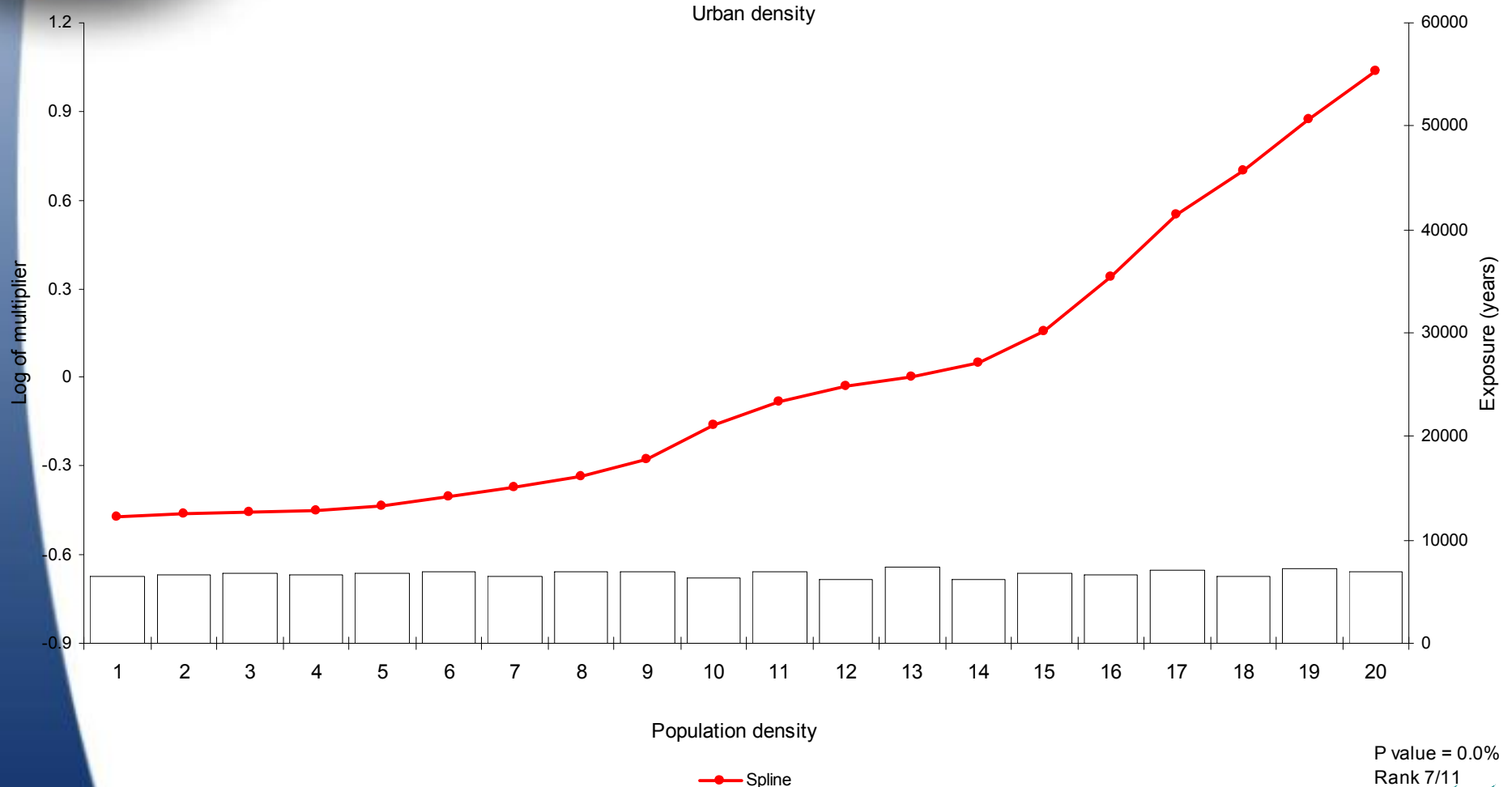


**Comparison of factor with spline**
Urban density

# Further example

## Comparison of factor with spline

Urban density



P value = 0.0%
Rank 7/11

# Splines

- Practical way of modeling continuous variables

- Often better than polynomials

- Increases complexity, therefore best used
  - when it is important that rates vary continuously with a variable
  - when modeling elasticity to be used in price optimization analyses

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# Standard approach

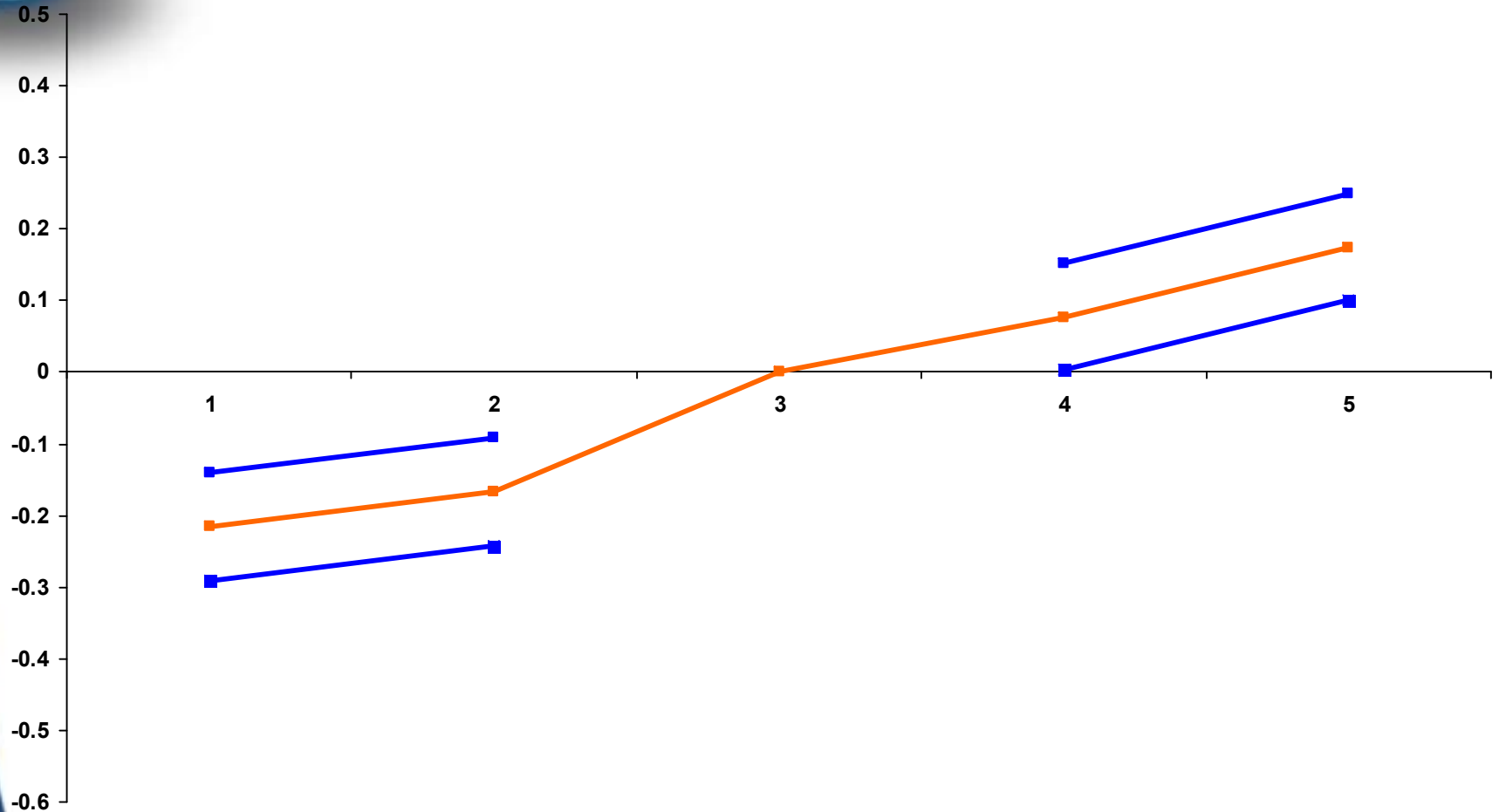| | Freq | | Amt | |
|---|---|---|---|---|
| BI | Freq | x | Amt | = Cost 1 |
| PD | Freq | x | Amt | = Cost 2 |
| MED | Freq | x | Amt | = Cost 3 |
| COL | Freq | x | Amt | = Cost 4 |
| OTC | Freq | x | Amt | = Cost 5 |

# Binomial reference models

# Offset reference model

# Offset reference model

# Offset reference model

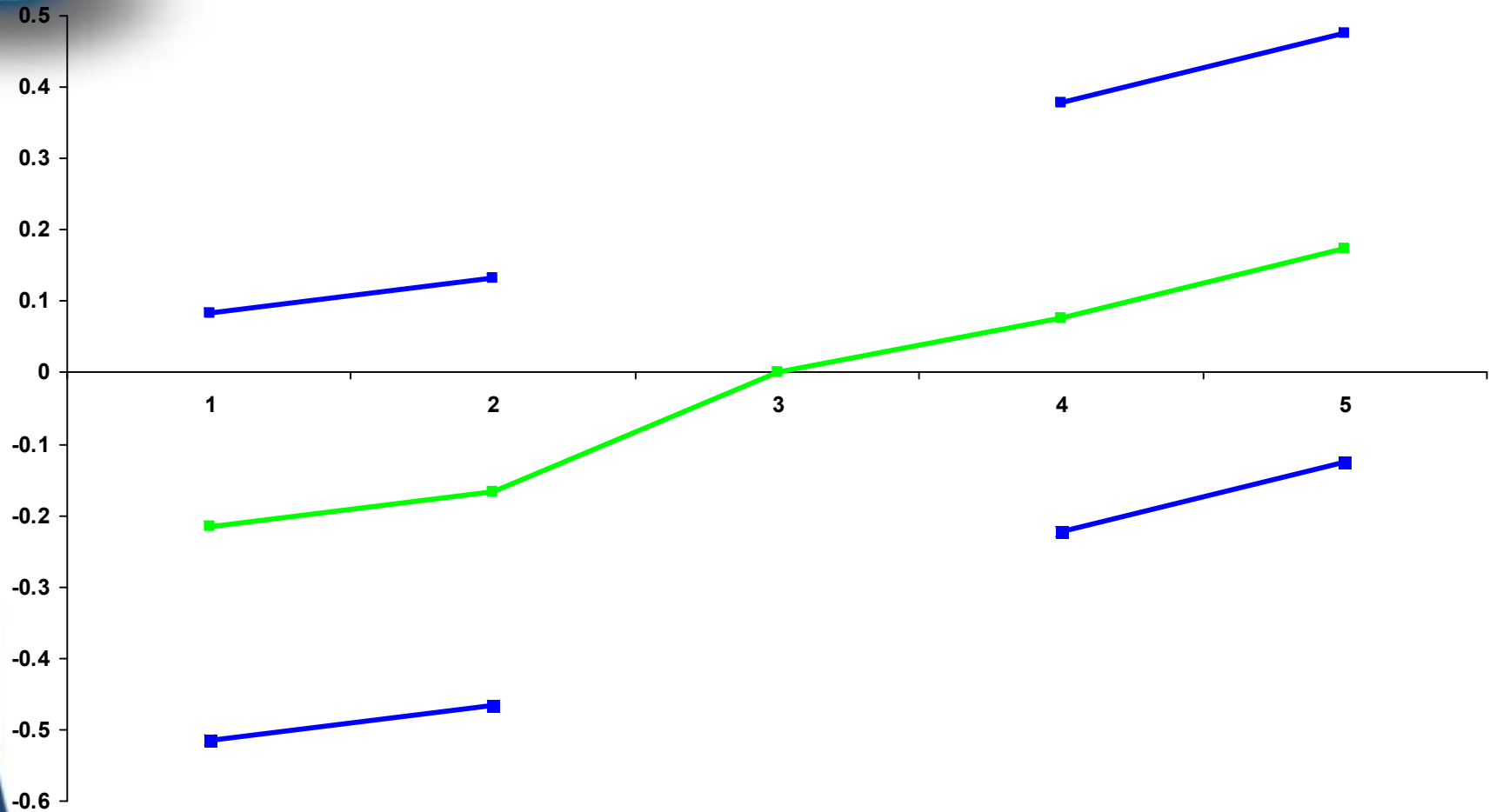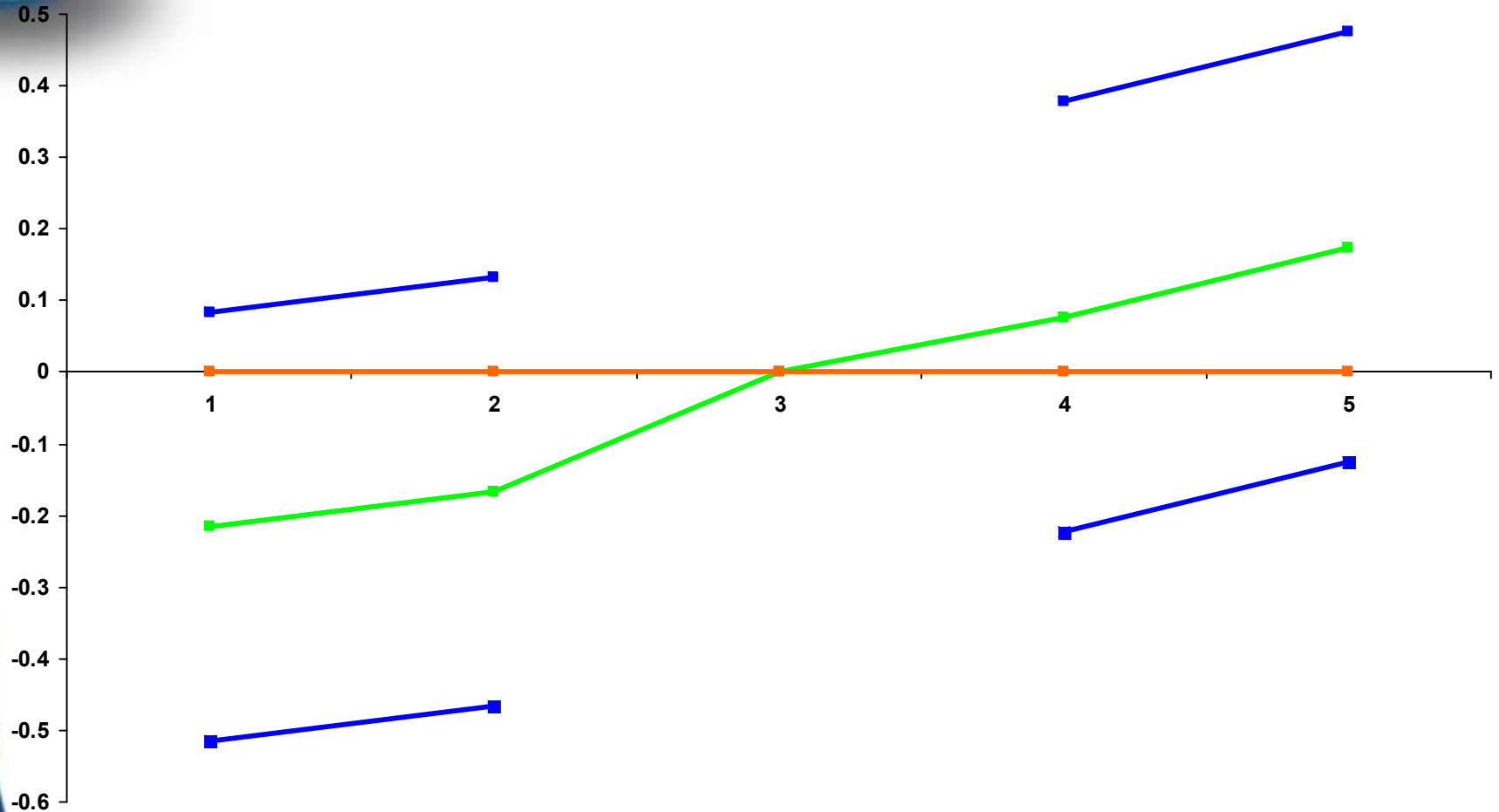# Offset reference model

# Offset reference model

# Offset reference model

# Offset reference model

# Offset reference model

# Testing the reference model approach

(1) Fit to BI claims on all data - the "correct answer"

100% of large company

10%

Random sample to emulate small company

(2) Model BI claims with standard approach

(3) Model BI claims referencing PD experience on this small sample

# Example of reference model method working



Green: BI 100% ("correct answer")

Blue: Standard BI GLM on 10%

Purple: Referencing PD on 10%

Standard BI GLM factor rejected

27%
19%
14%
14%
0%

Log of multiplier

Exposure (years)

Level 1    Level 2    Level 3

Factor X

— Approx 2 s.e. from estimate - Full model   —●— Unsmoothed estimate - Full model   —●— PD model

# Example of reference model method working



Green: BI 100% ("correct answer")

Blue: Standard BI GLM on 10%

Purple: Referencing PD on 10%

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# Aliasing and "near aliasing"

- Aliasing
  - the removal of unwanted redundant parameters

- Intrinsic aliasing
  - occurs by the design of the model

- Extrinsic aliasing
  - occurs "accidentally" as a result of the data

# Example

- Suppose we wanted a model of the form:

$$\mu = \alpha + \beta_1 \text{ if } \underline{age} < 30$$

$$+ \beta_2 \text{ if } \underline{age}\ 30 - 40$$

$$+ \beta_3 \text{ if } \underline{age} > 40$$

$$+ \gamma_1 \text{ if } \underline{sex}\ male$$

$$+ \gamma_2 \text{ if } \underline{sex}\ female$$

# Form of X.β in this case

|  | Age | | | Sex | |
|---|---|---|---|---|---|
|  | <30 | 30-40 | >40 | M | F |
| 1 | 1 | 0 1 0 | | 1 0 | |
| 2 | 1 | 1 0 0 | | 1 0 | |
| 3 | 1 | 1 0 0 | | 0 1 | |
| 4 | 1 | 0 0 1 | | 1 0 | |
| 5 | 1 | 0 1 0 | | 0 1 | |
| | ............................... | | | | |
| | ............................... | | | | |

$$\begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \gamma_1 \\ \gamma_2 \end{pmatrix}$$

# Example

- Suppose we wanted a model of the form:

$$\underline{\mu} = \alpha + \beta_1 \text{ if } \underline{age} < 30$$

$$+ \beta_2 \text{ if } \underline{age}\ 30 - 40$$

"Base levels"    $$+ \beta_3 \text{ if } \underline{age} > 40$$

$$+ \gamma_1 \text{ if } \underline{sex}\ male$$

$$+ \gamma_2 \text{ if } \underline{sex}\ female$$

# X.β having adjusted for base levels

$$
\begin{array}{c}
\text{Age} \qquad \text{Sex} \\
\begin{array}{ccccc}
<30 & 30\text{-}40 & >40 & M & F
\end{array}
\end{array}
$$

$$
\begin{array}{l}
1 \\
2 \\
3 \\
4 \\
5
\end{array}
\begin{pmatrix}
1 & 0 & 0 & 1 & 0 \\
1 & 1 & 0 & 1 & 0 \\
1 & 1 & 0 & 0 & 1 \\
1 & 0 & 1 & 1 & 0 \\
1 & 0 & 0 & 0 & 1 \\
\cdots\cdots\cdots\cdots\cdots \\
\cdots\cdots\cdots\cdots\cdots
\end{pmatrix}
\begin{pmatrix}
\alpha \\
\beta_1 \\
\beta_2 \\
\beta_3 \\
\gamma_1 \\
\gamma_2
\end{pmatrix} .
$$

# X.β having adjusted for base levels

$$
\begin{array}{c}
\phantom{1} \\
1 \\
2 \\
3 \\
4 \\
5 \\
\phantom{1}
\end{array}
\begin{pmatrix}
& \text{Age} & & \text{Sex} \\
<30 & >40 & & F \\
1 & 0 & 0 & 0 \\
1 & 1 & 0 & 0 \\
1 & 1 & 0 & 1 \\
1 & 0 & 1 & 0 \\
1 & 0 & 0 & 1 \\
\cdots\cdots\cdots\cdots \\
\cdots\cdots\cdots\cdots
\end{pmatrix}
.
\begin{pmatrix}
\alpha \\
\beta_1 \\
\beta_3 \\
\gamma_2
\end{pmatrix}
$$

# Intrinsic aliasing

## Example job

Run 16 Model 3 - Small interaction - Third party material damage, Numbers

# Extrinsic aliasing

- If a perfect correlation exists, one factor can alias levels of another

- Eg if doors declared first:

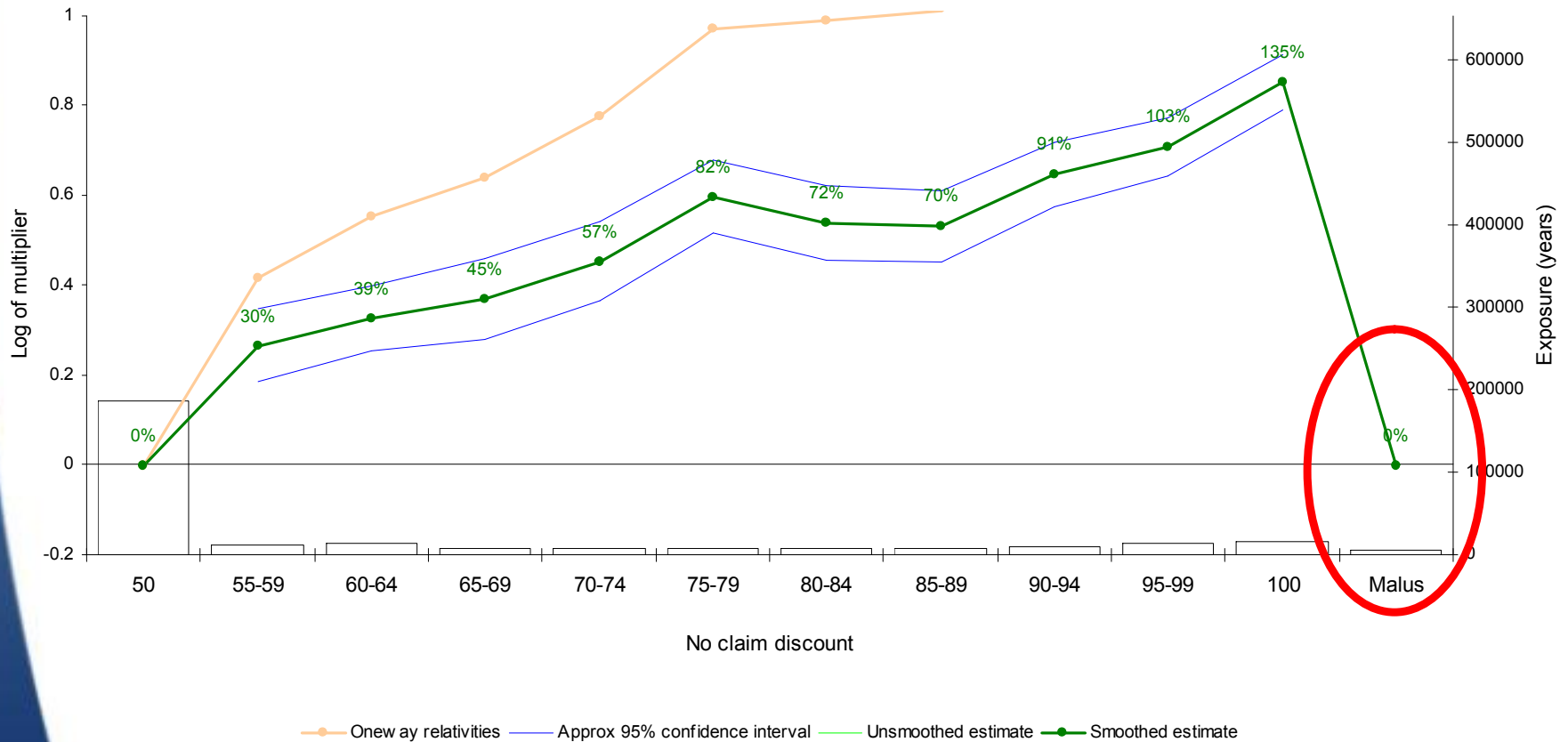| Exposure: # Doors → Color ↓ | 2 | 3 | 4 (Selected base) | 5 | Unknown |
|---|---|---|---|---|---|
| Red (Selected base) | 13,234 | 12,343 | 13,432 | 13,432 | 0 |
| Green | 4,543 | 4,543 | 13,243 | 2,345 | 0 |
| Blue | 6,544 | 5,443 | 15,654 | 4,565 | 0 |
| Black | 4,643 | 1,235 | 14,565 | 4,545 | 0 |
| Unknown (Further aliasing) | 0 | 0 | 0 | 0 | 3,242 |

- This is the only reason the order of declaration can matter (fitted values are unaffected)

# Extrinsic aliasing

## Example job

Run 16 Model 3 - Small interaction - Third party material damage, Numbers

# "Near aliasing"

- If two factors are almost perfectly, but not quite aliased, convergence problems can result and/or results can become hard to interpret

| Exposure: # Doors → | 2 | 3 | Selected base 4 | 5 | Unknown |
|---|---|---|---|---|---|
| Color ↓ | | | | | |
| Selected base Red | 13,234 | 12,343 | 13,432 | 13,432 | 0 |
| Green | 4,543 | 4,543 | 13,243 | 2,345 | 0 |
| Blue | 6,544 | 5,443 | 15,654 | 4,565 | 0 |
| Black | 4,643 | 1,235 | 14,565 | 4,545 | 2 |
| Unknown | 0 | 0 | 0 | 0 | 3,242 |

- Eg if the 2 black, unknown doors policies had no claims, GLM would try to estimate a very large negative number for unknown doors, and a very large positive number for unknown color

# "Near aliasing" - solution

1. Spot it

2. Fix the data!

| Exposure: # Doors → Colour ↓ | 2 | 3 | 4 | 5 | Unknown |
|---|---|---|---|---|---|
| Red | 13,234 | 12,343 | 13,432 | 13,432 | 0 |
| Green | 4,543 | 4,543 | 13,243 | 2,345 | 0 |
| Blue | 6,544 | 5,443 | 15,654 | 4,565 | 0 |
| Black | 4,643 | 1,235 | 14,565 | 4,545 | 2 |
| Unknown | 0 | 0 | 0 | 0 | 3,242 |

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# Combining claim elements - I

BI $\quad$ Freq $\times$ Amt $= $ Cost 1

PD $\quad$ Freq $\times$ Amt $= $ Cost 2

MED $\quad$ Freq $\times$ Amt $= $ Cost 3

COL $\quad$ Freq $\times$ Amt $= $ Cost 4

OTC $\quad$ Freq $\times$ Amt $= $ Cost 5

- Multiply factors for frequencies and amounts

- Calculate risk premium as sum of claim elements

# Combining claim elements - II

BI    Freq    x    Amt    = Cost 1

PD    Freq    x    Amt    = Cost 2

MED    Freq    x    Amt    = Cost 3

COL    Freq    x    Amt    = Cost 4

OTC    Freq    x    Amt    = Cost 5

- Consider current exposure

- Calculate expected frequency and amount for each claim type for each record

- Combine to give expected total cost of claims for each record

- Fit model to this expected value

# Calculation of risk premium

|  |  | TPPD Numbers | TPPD Amounts | TPBI Numbers | TPBI Amounts |
|---|---|---|---|---|---|
| Intercept |  | 32% | £1000 | 12% | £4860 |
| Sex | Male | 1.000 | 1.000 | 1.000 | 1.000 |
|  | Female | 0.750 | 1.200 | 0.667 | 0.900 |
| Area | Town | 1.000 | 1.000 | 1.000 | 1.000 |
|  | Country | 1.250 | 0.700 | 0.750 | 0.833 |

| Policy | Sex | Area | WWNUM1 | WWAMT1 | WWNUM2 | WWAMT2 | WWCC1 | WWCC2 | WWRSKPRM |
|---|---|---|---|---|---|---|---|---|---|
| … | … | … | … | … | … | … | … | … | … |
| 82155654 | M | T | 32% | 1000 | 12% | 4860 | 320 | 583.20 | 903.20 |
| 82168746 | F | T | 24% | 1200 | 8% | 4374 | 288 | 349.92 | 637.92 |
| 82179481 | M | C | 40% | 700 | 9% | 4050 | 280 | 364.50 | 644.50 |
| 82186845 | F | C | 30% | 840 | 6% | 3645 | 252 | 218.70 | 470.70 |
| … | … | … | … | … | … | … | … | … | … |

# Risk premium standard errors

- Risk premium model standard errors are small owing to the smoothness of the expected value

- It is possible to approximate standard error of risk premium parameter estimates based on standard errors of parameter estimates in underlying models

- Care needed in interpreting such approximations since they do not reflect model error, eg deciding to exclude a marginal factor
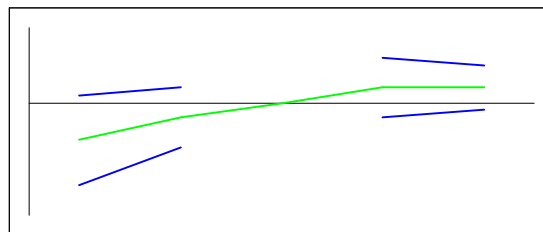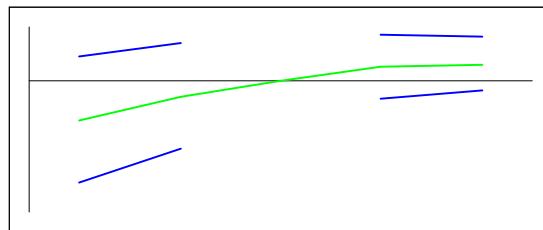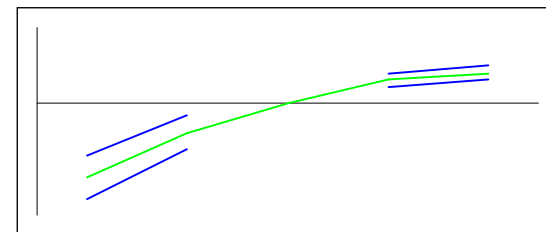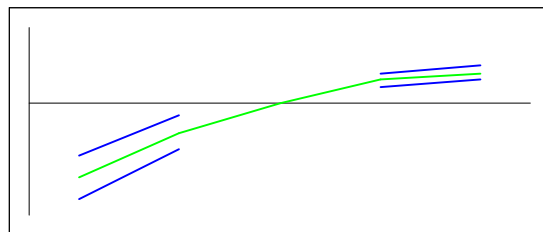
# Risk premium standard errors - failings
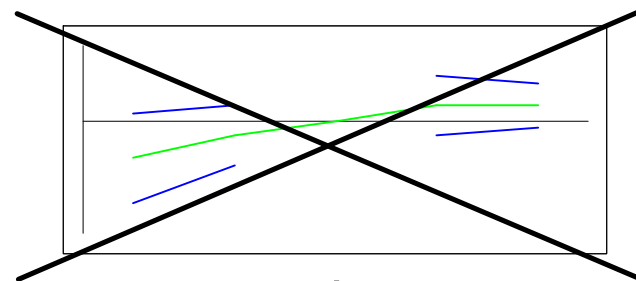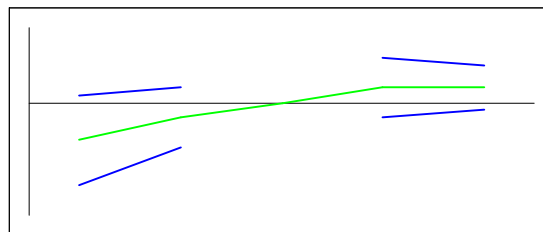
Numbers

Amounts

Risk premium

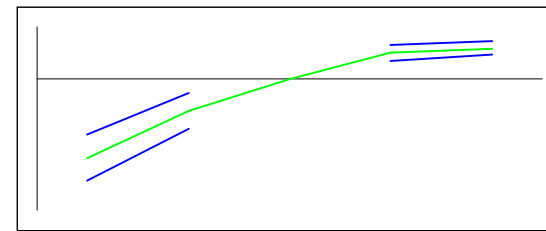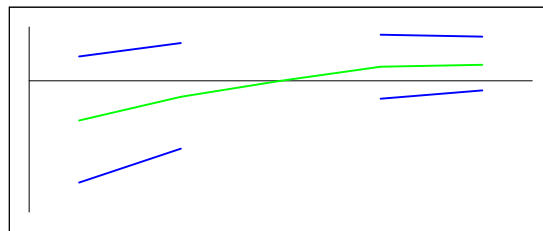# Risk premium standard errors - failings

Numbers

Amounts

Risk premium

# Agenda

- Introduction

- Testing the link function

- The Tweedie distribution

- Splines

- Reference models

- Aliasing / near aliasing

- Combining models across claim types

- Restricted models

# Restricted models

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

Offset

- Offset term used for known effects, eg exposure in a numbers model

- Can also be used to constrain model (eg claim free years / payment frequency / amount of cover)

- Other factors adjusted to compensate

# Restricted models

$$
\begin{array}{c}
\phantom{} \\
1 \\
2 \\
3 \\
4 \\
5 \\
\phantom{} \\
\phantom{}
\end{array}
\left(
\begin{array}{cccc}
 & \text{Age} & & \text{Sex} \\
 & \text{<30} & \text{>40} & \text{F} \\
1 & 0 & 0 & 0 \\
1 & 1 & 0 & 0 \\
1 & 1 & 0 & 1 \\
1 & 0 & 1 & 0 \\
1 & 0 & 0 & 1 \\
\cdots\cdots\cdots\cdots \\
\cdots\cdots\cdots\cdots
\end{array}
\right)
\cdot
\left(
\begin{array}{c}
\alpha \\
\beta_1 \\
\beta_3 \\
\gamma_2
\end{array}
\right)
$$

# Restricted models

$$E[\underline{Y}] = \underline{\mu} = g^{-1}(\, \mathbf{X}.\underline{\beta}\,)$$

|   | Age | | Sex |
|---|-----|---|-----|
|   | <30 | >40 | F |
| 1 | 1 | 0 | 0 | 0 |
| 2 | 1 | 1 | 0 | 0 |
| 3 | 1 | 1 | 0 | 1 |
| 4 | 1 | 0 | 1 | 0 |
| 5 | 1 | 0 | 0 | 1 |

$$\begin{pmatrix} \alpha \\ \beta_1 \\ \beta_3 \\ \gamma_2 \end{pmatrix}$$

# Restricted models

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

Age

|   | <30 | >40 |
|---|-----|-----|
| 1 | 1 | 0 | 0 |
| 2 | 1 | 1 | 0 |
| 3 | 1 | 1 | 0 |
| 4 | 1 | 0 | 1 |
| 5 | 1 | 0 | 0 |

$$\begin{pmatrix} \alpha \\ \beta_1 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 0.1 \\ 0 \\ 0.1 \\ \ldots \\ \ldots \end{pmatrix}$$

# Restricted models

$$
\begin{array}{c}
\phantom{x} \\
\phantom{x} \\
1 \\
2 \\
3 \\
4 \\
5 \\
\phantom{x}
\end{array}
\begin{pmatrix}
 & \text{Age} & & \text{Sex} \\
 & \text{<30} & \text{>40} & \text{F} \\
1 & 0 & 0 & 0 \\
1 & 1 & 0 & 0 \\
1 & 1 & 0 & 1 \\
1 & 0 & 1 & 0 \\
1 & 0 & 0 & 1 \\
\multicolumn{4}{c}{\dotfill} \\
\multicolumn{4}{c}{\dotfill}
\end{pmatrix}
\cdot
\begin{pmatrix}
\alpha \\
\beta_1 \\
\beta_3 \\
\mathbf{0.1}
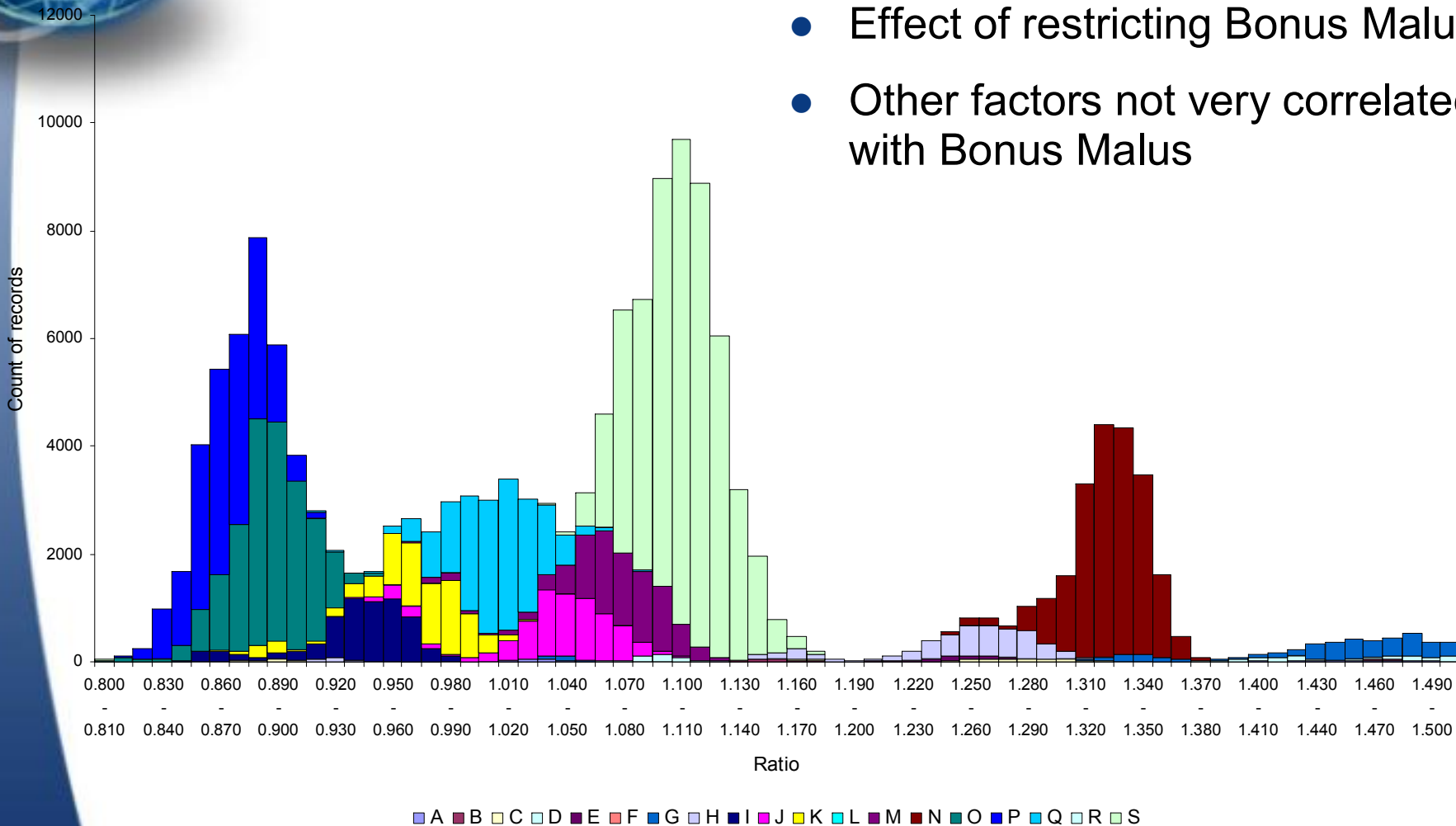\end{pmatrix}
$$

# Restricted models
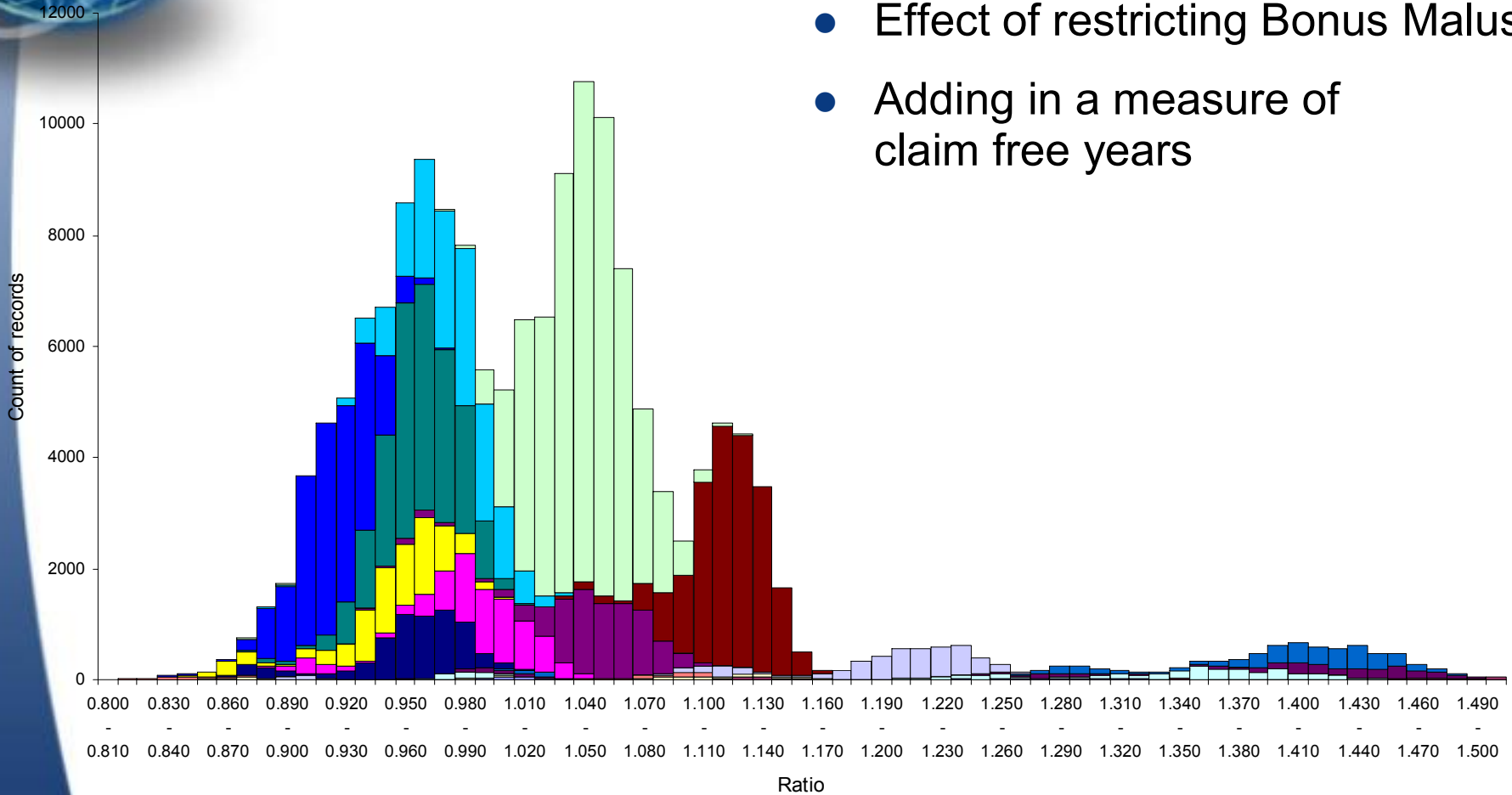


Restriction

# Restricted models

# Testing the effectiveness of restrictions



- Effect of restricting Bonus Malus

- Other factors not very correlated with Bonus Malus

# Testing the effectiveness of restrictions

- Effect of restricting Bonus Malus

- Adding in a measure of claim free years

# Restrictions

- Only use to "get around" restrictions

- A commercial smoothing is a commercial smoothing

- Apply at risk premium stage

# GLM III: Advanced Modeling Strategy

**2005 CAS Seminar on Predictive Modeling**

**Duncan Anderson MA FIA**

**Watson Wyatt Worldwide**

WWW.WATSONWYATT.COM

Watson Wyatt
*Worldwide*