

# Automating variable selection in distributional regression

---

Kevin Burke | University of Limerick

# Smooth information criterion (SIC)

---

Statistics and Computing (2023) 33:71

<https://doi.org/10.1007/s11222-023-10204-8>

ORIGINAL PAPER

## Variable selection using a smooth information criterion for distributional regression models

Meadhbh O'Neill<sup>1</sup>  · Kevin Burke<sup>1</sup> 

# Scientific question


---

Does  $X$  cause  $Y$ ?

Does  $X$  ~~cause~~  $Y$ ?  
relate to

# Scientific question

---

Does  $X$  ~~cause~~  $Y$ ?  
relate to  
  
linearly

## Scientific question

---

the mean of  
↓  
Does  $X$  ~~cause~~  $Y$ ?  
relate to  
↑  
linearly

# Mean regression (obsession?)

---

# Mean regression (obsession?)

---

- ▶  $Y = X^\top \beta + \sigma \varepsilon, \quad \varepsilon \sim N(0,1)$

- ▶  $\Rightarrow Y \sim N(X^\top \beta, \sigma^2)$

- ▶ Why should  $\sigma^2$  be constant?



# Mean regression (obsession?)

---

- $Y = X^\top \beta + \sigma \varepsilon, \quad \varepsilon \sim N(0,1)$

- $\Rightarrow Y \sim N(X^\top \beta, \sigma^2)$

- Why should  $\sigma^2$  be constant?



# Distributional regression

---

- Any/all distributional parameters could depend on  $X$

# Distributional regression

---

- Any/all distributional parameters could depend on  $X$
- $Y = \mu + \sigma\varepsilon, \quad \varepsilon \sim N(0,1)$
- $\mu = X^\top \beta, \quad \mathbf{\log \sigma^2 = X^\top \alpha}$
- $\Rightarrow Y \sim N(X^\top \beta, \mathbf{\exp(X^\top \alpha)})$

# Distributional regression

---

- ▶ Any/all distributional parameters could depend on  $X$
- ▶  $Y = \mu + \sigma\varepsilon, \quad \varepsilon \sim N(0,1)$
- ▶  $\mu = X^\top \beta, \quad \mathbf{\log \sigma^2 = X^\top \alpha}$
- ▶  $\Rightarrow Y \sim N(X^\top \beta, \mathbf{\exp(X^\top \alpha)})$
- ▶ Also known as “multi-parameter regression”

# Model space

---

- Consider  $X = (X_1, X_2, X_3)$

$\mu(X)$
$\emptyset$
$X_1$
$X_2$
$X_3$
$X_1, X_2$
$X_1, X_3$
$X_2, X_3$
$X_1, X_2, X_3$

# Model space

---

- Consider  $X = (X_1, X_2, X_3)$

$\mu(X)$

$\emptyset$

$X_1$

$X_2$

$X_3$

$X_1, X_2$

$X_1, X_3$

$X_2, X_3$

$X_1, X_2, X_3$

Mean regression

# Model space

---

- Consider  $X = (X_1, X_2, X_3)$

$\mu(X)$	$\sigma(X)$
$\emptyset$	$\emptyset$
$X_1$	$X_1$
$X_2$	$X_2$
$X_3$	$X_3$
$X_1, X_2$	$X_1, X_2$
$X_1, X_3$	$X_1, X_3$
$X_2, X_3$	$X_2, X_3$
$X_1, X_2, X_3$	$X_1, X_2, X_3$

Mean regression

Distributional regression (mean & variance)

# Model space

- Consider  $X = (X_1, X_2, X_3)$

$\mu(X)$	$\sigma(X)$
$\emptyset$	$\emptyset$
$X_1$	$X_1$
$X_2$	$X_2$
$X_3$	$X_3$
$X_1, X_2$	$X_1, X_2$
$X_1, X_3$	$X_1, X_3$
$X_2, X_3$	$X_2, X_3$
$X_1, X_2, X_3$	$X_1, X_2, X_3$

Mean regression

Distributional regression (mean & variance)

- Models

- $2^3 \times 2^3 = 2^{2 \times 3}$

- In general

- $2^{d \times p}$

- $d$  distributional parameters

- $p$  covariates



# How to select variables?

---

# How to select variables?

---

- $\text{BIC} = -2\ell + (\log n)(\|\beta\|_0 + \|\alpha\|_0)$

# How to select variables?

---

- $\text{BIC} = -2\ell + (\log n)(\|\beta\|_0 + \|\alpha\|_0)$
- $\ell_{\text{BIC}} = -\text{BIC}/2 = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$

# How to select variables?

---

- $\text{BIC} = -2\ell + (\log n)(\|\beta\|_0 + \|\alpha\|_0)$
- $\ell_{\text{BIC}} = -\text{BIC}/2 = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- $(\hat{\beta}, \hat{\alpha}) = \max_{\beta, \alpha} \ell_{\text{BIC}}$

# How to select variables?

---

- $\text{BIC} = -2\ell + (\log n)(\|\beta\|_0 + \|\alpha\|_0)$
- $\ell_{\text{BIC}} = -\text{BIC}/2 = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- $(\hat{\beta}, \hat{\alpha}) = \max_{\beta, \alpha} \ell_{\text{BIC}}$
- Discrete optimisation: fit all models

# How to select variables?

---

- $\text{BIC} = -2\ell + (\log n)(\|\beta\|_0 + \|\alpha\|_0)$
- $\ell_{\text{BIC}} = -\text{BIC}/2 = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- $(\hat{\beta}, \hat{\alpha}) = \max_{\beta, \alpha} \ell_{\text{BIC}}$
- Discrete optimisation: fit all models
- Heuristic: stepwise search

# Related problem

---

# Related problem

---

- Problem:  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$



# Related problem

---

- Problem:  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- Related:  $\ell_{\text{LASSO}} = \ell - \lambda(\|\beta\|_1 + \|\alpha\|_1)$

# Related problem

---

- Problem:  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- Related:  $\ell_{\text{LASSO}} = \ell - \lambda(\|\beta\|_1 + \|\alpha\|_1)$
- What about  $\lambda$ ?
  - Given  $\lambda$ , obtain  $(\hat{\beta}_{\text{LASSO}}, \hat{\alpha}_{\text{LASSO}}) = \arg \max \ell_{\text{LASSO}}$
  - Select  $\lambda$  using  $\ell_{\text{BIC}}$
  - Thus,  $\max \ell_{\text{BIC}}$  subject to solutions of the form  $(\hat{\beta}_{\text{LASSO}}, \hat{\alpha}_{\text{LASSO}})$

# Related problem

---

- Problem:  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- Related:  $\ell_{\text{LASSO}} = \ell - \lambda(\|\beta\|_1 + \|\alpha\|_1)$
- What about  $\lambda$ ?
  - Given  $\lambda$ , obtain  $(\hat{\beta}_{\text{LASSO}}, \hat{\alpha}_{\text{LASSO}}) = \arg \max \ell_{\text{LASSO}}$
  - Select  $\lambda$  using  $\ell_{\text{BIC}}$
  - Thus,  $\arg \max \ell_{\text{BIC}}$  subject to solutions of the form  $(\hat{\beta}_{\text{LASSO}}, \hat{\alpha}_{\text{LASSO}})$
- **Note that  $\|x\|_1$  is not differentiable**

# Recall the original problem

---

# Recall the original problem

---

- Maximise  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$

# Recall the original problem

---

- Maximise  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- Discrete optimisation challenging

# Recall the original problem

---

- Maximise  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- Discrete optimisation challenging
- Solve “related” LASSO-type problem

# Recall the original problem

---

- Maximise  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$
- Discrete optimisation challenging
- Solve “related” LASSO-type problem
- Not fully satisfying ...
  - $\lambda$  tuning required
  - $(\hat{\beta}, \hat{\alpha})$  not solution to original problem



# Idea: smooth information criterion (SIC)

---

# Idea: smooth information criterion (SIC)

---

- “Brilliant”

# Idea: smooth information criterion (SIC)

---

- “Brilliant”
- “Ingenuous”

# Idea: smooth information criterion (SIC)

---

- “Brilliant”
- “Ingenuous”
- “Remarkable”

# Idea: smooth information criterion (SIC)

---

■ “Brilliant”

■ “Ingenuous”

■ “Remarkable”

Robitzsch (2023)

a. L0 and Lp Loss Functions in Model-Robust Estimation of Structural Equation Models

b. Implementation aspects in regularized structural equation models

c. Implementation Aspects in Invariance

# Idea: smooth information criterion (SIC)

---

➤ “Brilliant”

➤ “Ingenious”

➤ “Remarkable”

Robitzsch (2023)

a. L0 and Lp Loss Functions in Model-Robust Estimation of Structural Equation Models

b. Implementation aspects in regularized structural equation models

c. Implementation Aspects in Invariance

➤  $\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0)$

# Idea: smooth information criterion (SIC)

---

➤ “Brilliant”

➤ “Ingenuous”

➤ “Remarkable”

Robitzsch (2023)

a. L0 and Lp Loss Functions in Model-Robust Estimation of Structural Equation Models

b. Implementation aspects in regularized structural equation models

c. Implementation Aspects in Invariance

$$\begin{aligned}\text{➤ } \ell_{\text{BIC}} &= \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0) \\ &= \ell - \{(\log n)/2\} \left( \sum_j |\beta_j|^0 + \sum_j |\alpha_j|^0 \right)\end{aligned}$$

# Idea: smooth information criterion (SIC)

■ “Brilliant”

■ “Ingenuous”

■ “Remarkable”

Robitzsch (2023)

a. L0 and Lp Loss Functions in Model-Robust Estimation of Structural Equation Models

b. Implementation aspects in regularized structural equation models

c. Implementation Aspects in Invariance

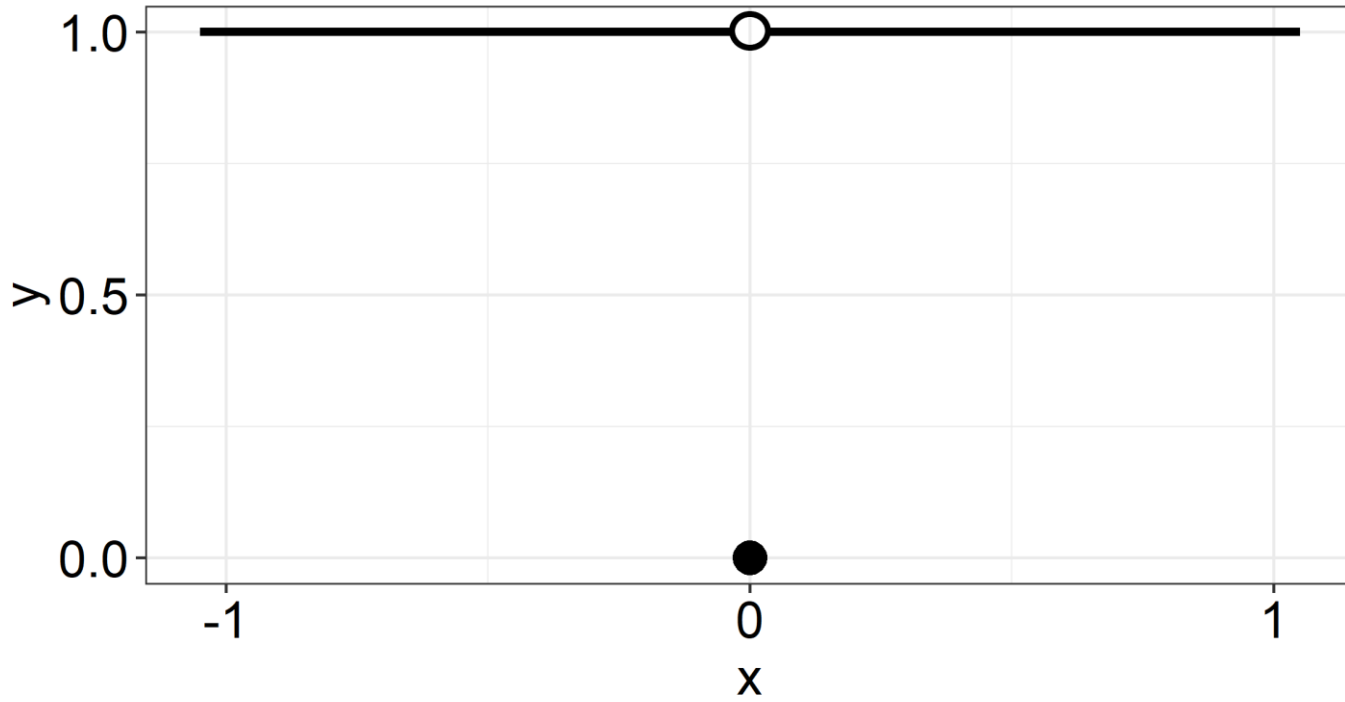
$$\begin{aligned}\ell_{\text{BIC}} &= \ell - \{(\log n)/2\} (\|\beta\|_0 + \|\alpha\|_0) \\ &= \ell - \{(\log n)/2\} \left( \sum_j |\beta_j|^0 + \sum_j |\alpha_j|^0 \right)\end{aligned}$$

$$\ell_{\text{SIC}} = \ell - \{(\log n)/2\} \left( \sum_j \phi_\epsilon(\beta_j) + \sum_j \phi_\epsilon(\alpha_j) \right)$$



# Smooth $L_0$ norm

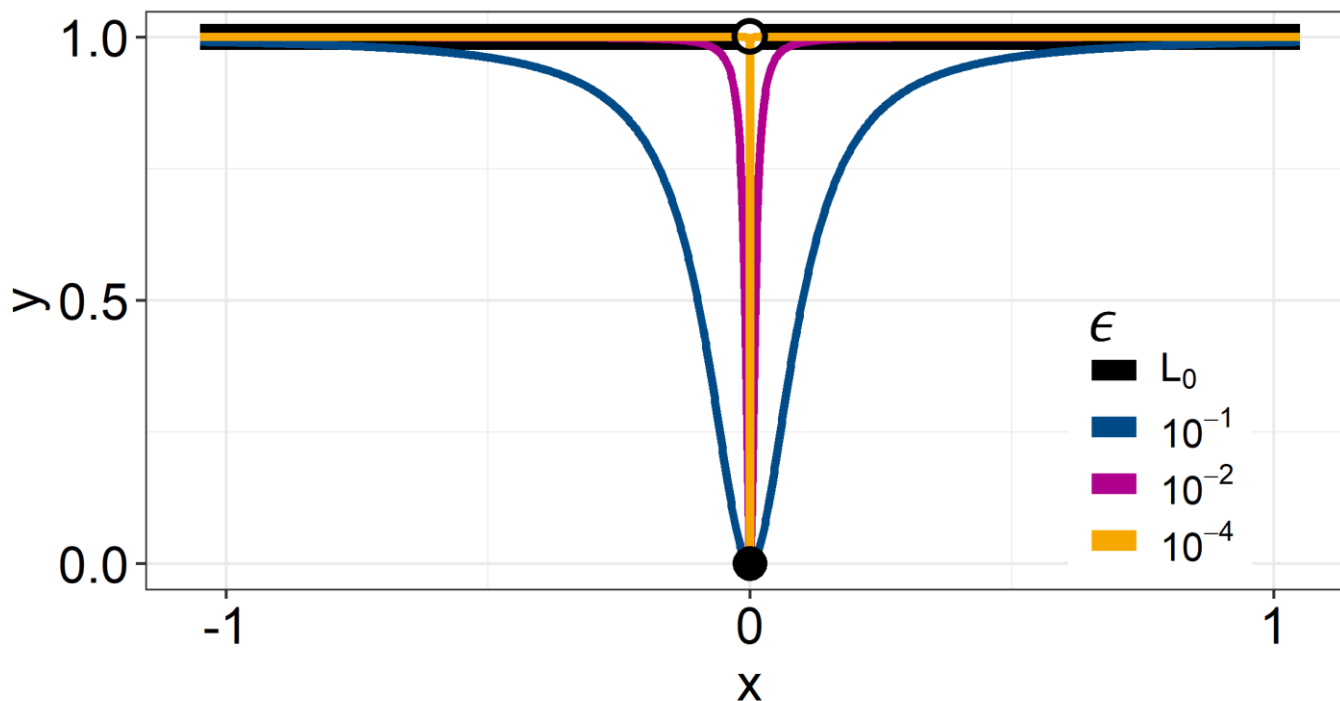
---



# Smooth $L_0$ norm

$$\phi_\epsilon(x) = \frac{x^2}{x^2 + \epsilon^2}$$

- Differentiable for  $\epsilon > 0$
- $\lim_{\epsilon \rightarrow 0} \phi_\epsilon(x) = |x|^0$



# $\epsilon$ -telescoping

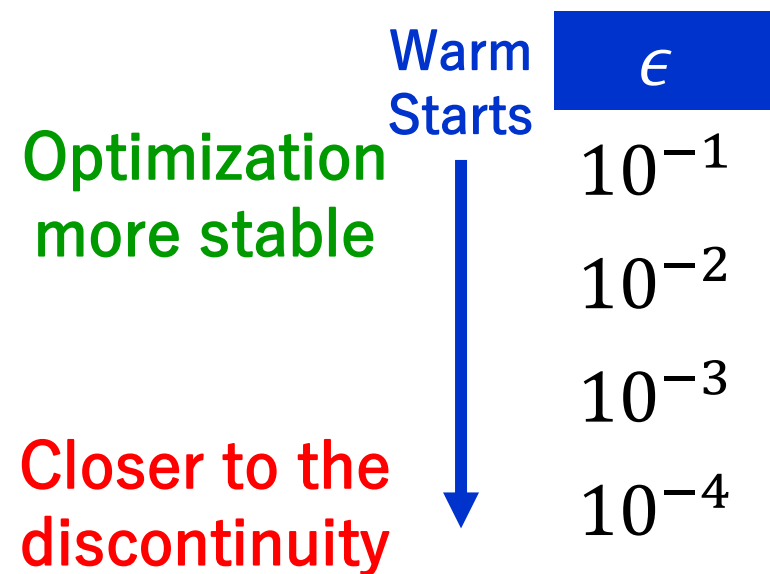
---

- $\lim_{\epsilon \rightarrow 0} \ell_{\text{SIC}} = \ell_{\text{BIC}}$
- Maximising  $\ell_{\text{SIC}}$  equivalent to  $\ell_{\text{BIC}}$  for  $\epsilon \approx 0$

# $\epsilon$ -telescoping

---

- $\lim_{\epsilon \rightarrow 0} \ell_{\text{SIC}} = \ell_{\text{BIC}}$
- Maximising  $\ell_{\text{SIC}}$  equivalent to  $\ell_{\text{BIC}}$  for  $\epsilon \approx 0$



# $\epsilon$ -telescoping

- ▶  $\lim_{\epsilon \rightarrow 0} \ell_{\text{SIC}} = \ell_{\text{BIC}}$
- ▶ Maximising  $\ell_{\text{SIC}}$  equivalent to  $\ell_{\text{BIC}}$  for  $\epsilon \approx 0$

		true zero	true nonzero
		$\beta_1 = 0$	$\beta_2 = 1$
Optimization more stable	Warm Starts		
	$10^{-1}$	-0.0249883798	1.008
	$10^{-2}$	-0.0005789613	1.009
	$10^{-3}$	-0.0000058339	1.009
Closer to the discontinuity	$10^{-4}$	<b>-0.0000000006</b>	1.009

# smoothic package

---





## Prostate data example

- Level of prostate-specific antigen (PSA)
- $p = 8$ , various clinical measures
- $n = 97$



## Prostate data example

- Level of prostate-specific antigen (PSA)
- $p = 8$ , various clinical measures
- $n = 97$

```
# Smooth Information Criterion -----  
library(smoothic)  
fit <- smoothic(formula = lpsa ~ .,  
                 data = pcancer,  
                 family = "normal",  
                 model = "mpr")    # or model = "spr"
```



# smoothic package

---

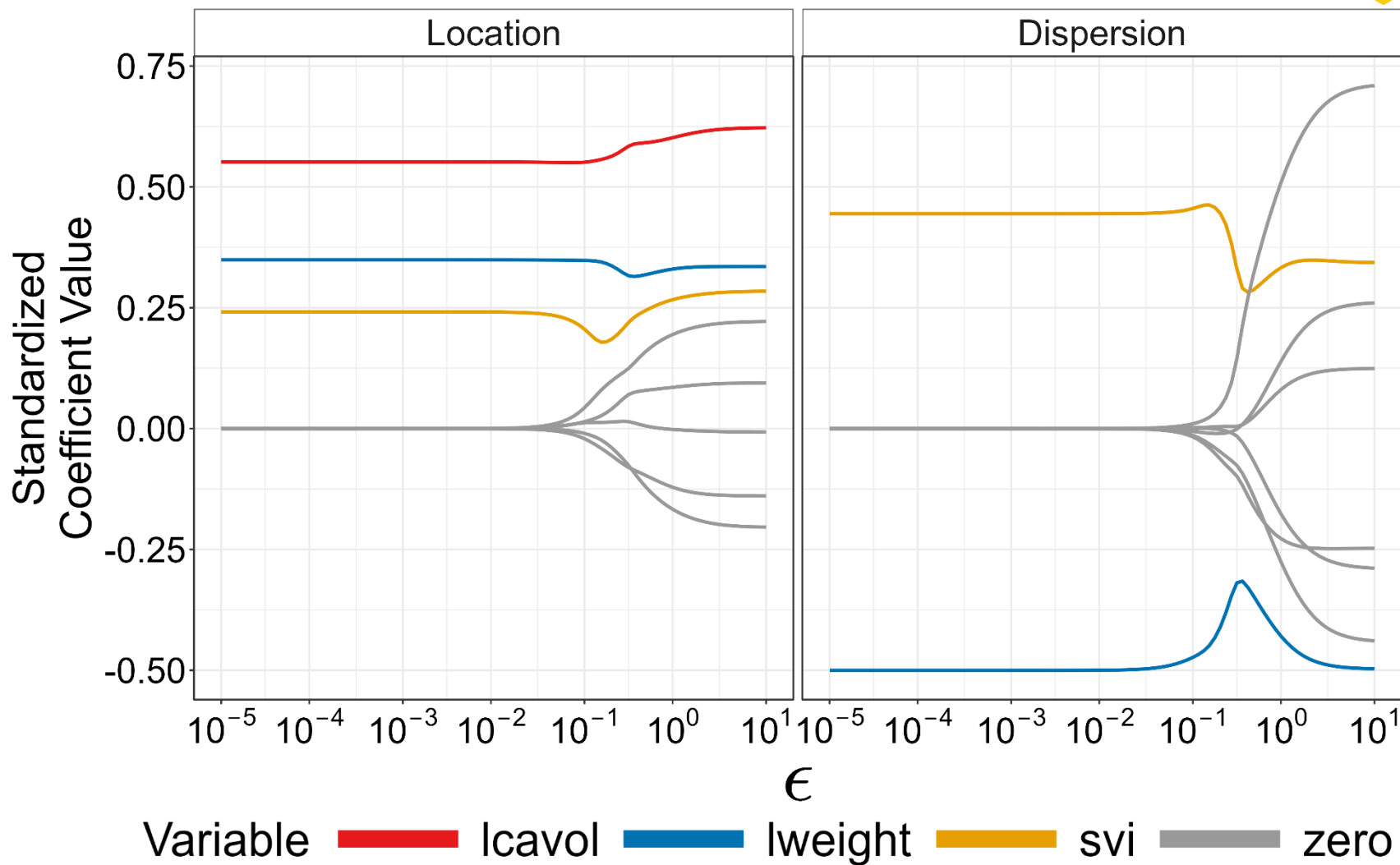
$\epsilon$ -telescope





`plot_paths()`

$\epsilon$ -telescope



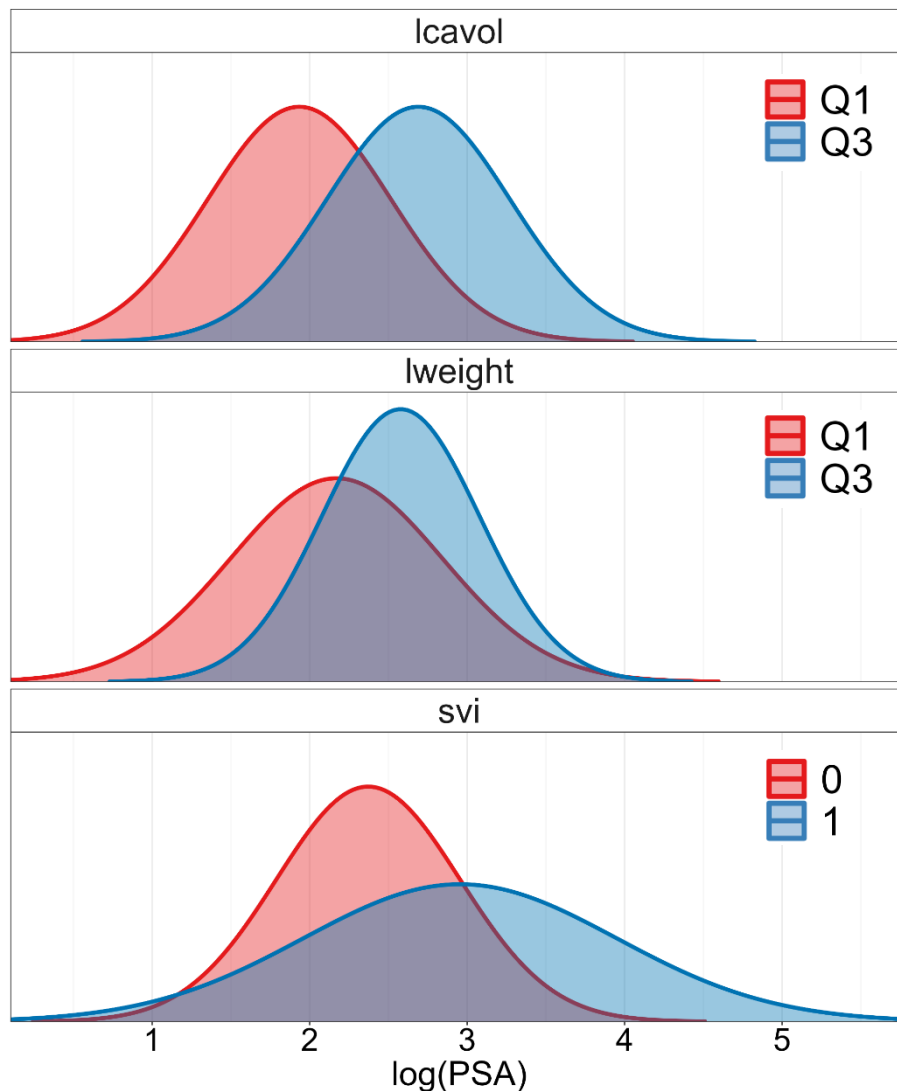


## Conditional Density Curves



`plot_effects()`

## Conditional Density Curves



- `lcavol` | log(cancer volume)
- `lweight` | log(prostate weight)
- `svi` | presence of seminal vesicle invasion

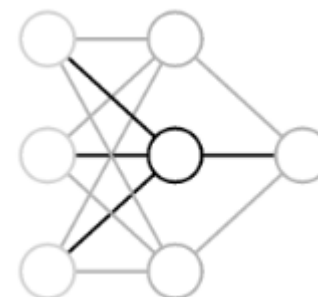
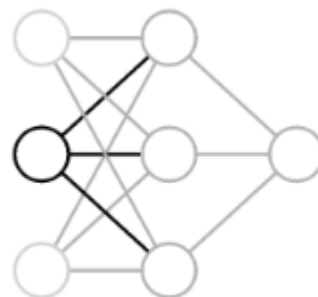
# Neural network selection

- Developing  $\ell_{\text{SIC}}$  in neural networks

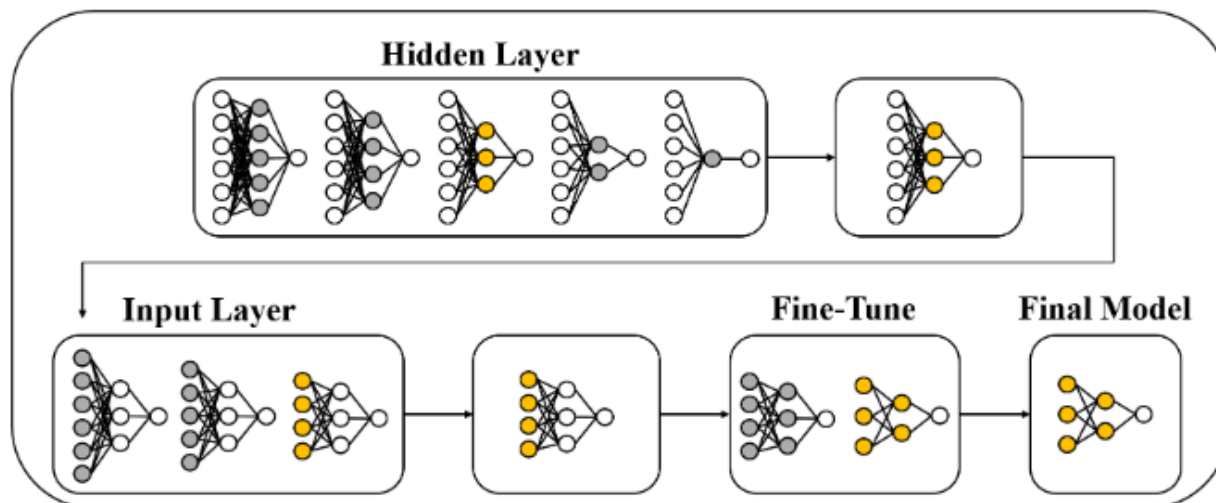
Input neurons

Hidden neurons

- Penalty types



- Traditional stepwise procedures



# selectnn package



```
nn <- selectnn(charges ~ ., data = insurance,  
               Q = 8, n_init = 5)  
summary(nn)
```

Number of input nodes: 3

Number of hidden nodes: 1

Value: 235.005

Inputs:

Covariate	Selected	Delta.BIC
lcavol	Yes	33.685
lweight	Yes	10.325
svi	Yes	6.794
age	No	
lbph	No	
lcp	No	
gleason	No	
pgg45	No	

... also see  
[interpretnn](#)



# Summary

---

- Distributional regression more flexible than standard mean regression
- Variable selection more challenging but the smooth information criterion is efficient
- Developing extensions to neural networks
- O'Neill & Burke (2023) Variable selection using a smooth information criterion for distributional regression models. Stats & Computing.
- Also see: [kevinburke.ie](https://kevinburke.ie) and [arxiv.org/a/burke\\_k\\_1](https://arxiv.org/a/burke_k_1)



# 39<sup>th</sup> IWSM

## Limerick, Ireland

### 13<sup>th</sup> – 18<sup>th</sup>, July 2025



MATHEMATICS  
APPLICATIONS CONSORTIUM  
FOR SCIENCE & INDUSTRY



UNIVERSITY OF  
**LIMERICK**  
OLLSCOIL LUIMNIGH



Brendan Murphy | Ireland

Ruth King | UK

Sonja Greven | Germany

Daniele Durante | Italy

Cynthia Rudin | US

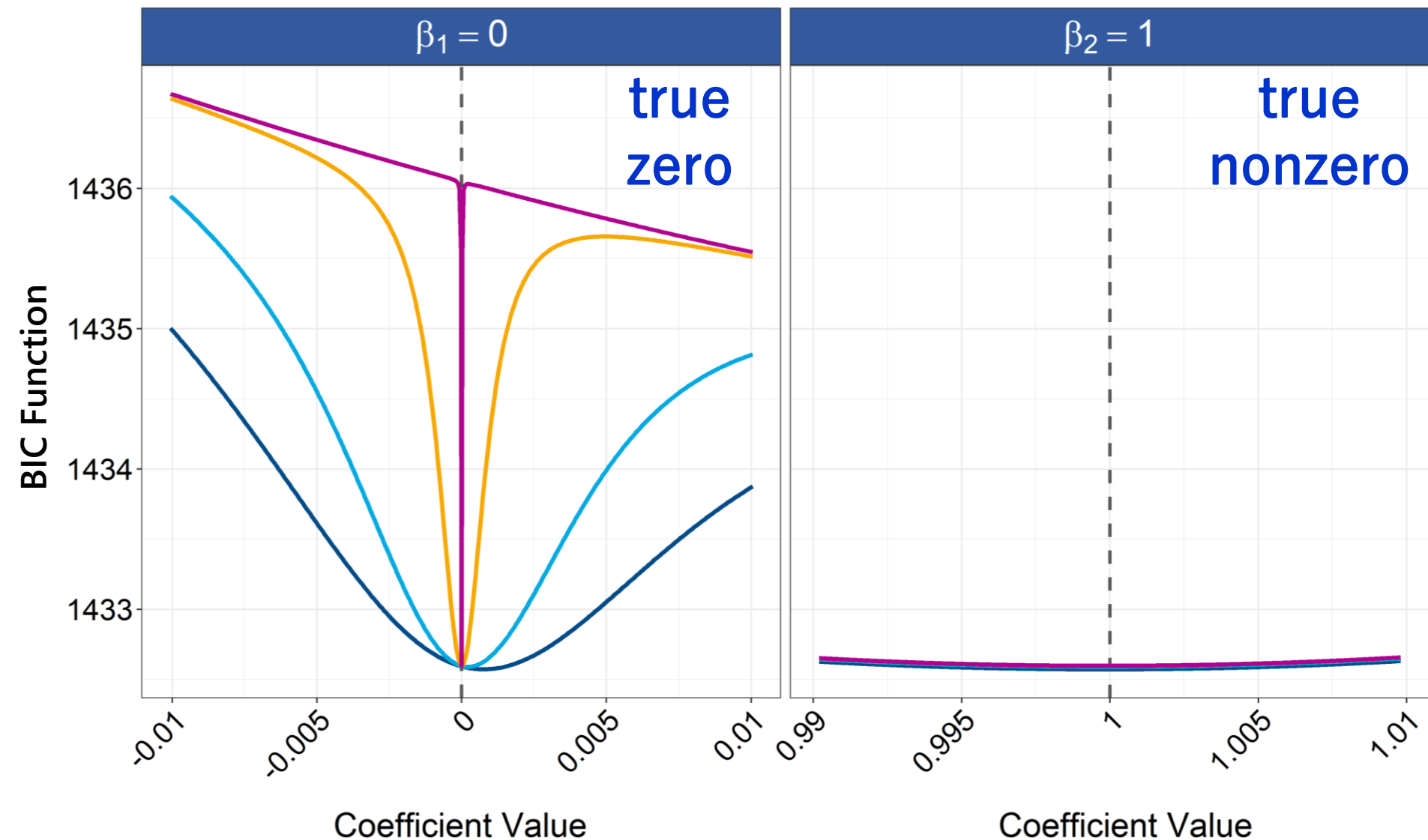




# $\epsilon$ -telescoping

$\epsilon$

0.00001 0.001 0.005 0.01



# Further details

- Separate tuning parameters

$$\ell_{\text{SAL}} = \ell - \left( \lambda_{\beta} \sum_j w_{\beta_j} a_{\epsilon}(\beta_j) + \lambda_{\alpha} \sum_j w_{\alpha_j} a_{\epsilon}(\alpha_j) \right)$$

- Standard errors

$$\text{cov}(\hat{\theta})$$

$$= \left\{ -\nabla_{\theta} \nabla_{\theta}^{\top} \ell_{\text{SAL}}(\hat{\theta}) \right\}^{-1} \left\{ -\nabla_{\theta} \nabla_{\theta}^{\top} \ell(\hat{\theta}) \right\} \left\{ -\nabla_{\theta} \nabla_{\theta}^{\top} \ell_{\text{SAL}}(\hat{\theta}) \right\}^{-1}$$

- Effective degrees of freedom

$$\ell_{\text{BIC}} = \ell - \{(\log n)/2\} (\text{edf})$$

$$\text{edf} = \text{trace} \left[ \left\{ -\nabla_{\theta} \nabla_{\theta}^{\top} \ell_{\text{SAL}}(\hat{\theta}) \right\}^{-1} \left\{ -\nabla_{\theta} \nabla_{\theta}^{\top} \ell(\hat{\theta}) \right\} \right]$$