

DISCUSSION ON

***“Bayesian Smoothing and Feature Selection using Variational Automatic Relevance Determination”***

By Zihe Liu, Diptarka Saha and Feng Liang  
Presented by Feng Liang

O'Bayes Conference 2025  
Stavros Niarchos Foundation Cultural Center  
Athens, Greece

**Discussant: Xenia Miscouridou**

Dep of Mathematics and Statistics, University of Cyprus  
Dep of Mathematics and IX Centre, Imperial College London



Πανεπιστήμιο Κύπρου  
University of Cyprus

**Imperial College**  
London

# Problem Formulation

- Simultaneous smoothing and variable selection for additive models

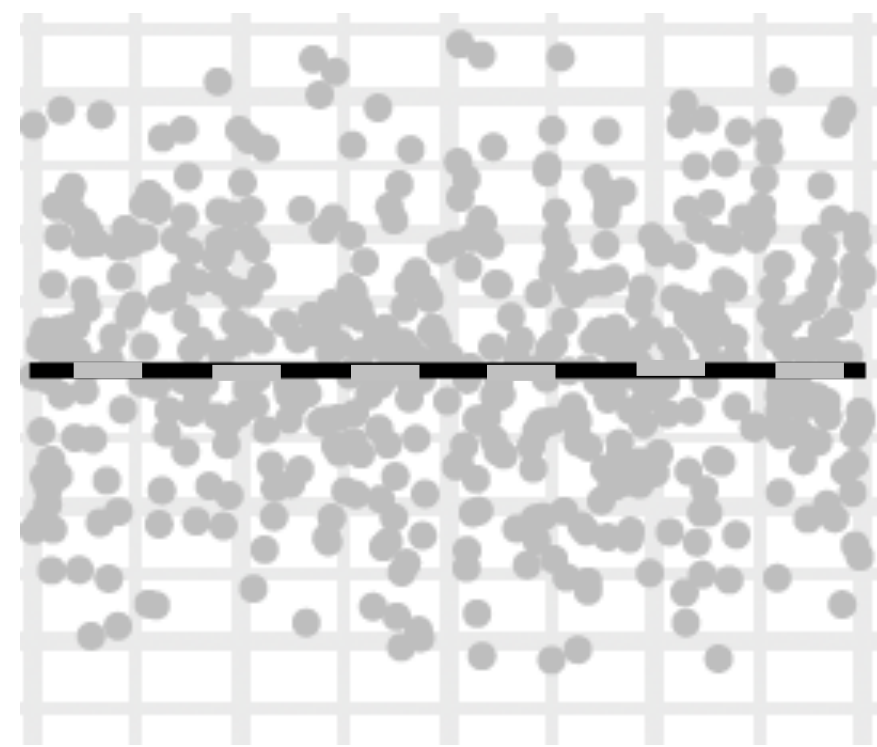
$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

# Problem Formulation

- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

## VARIABLE SELECTION



True  $f_j = 0$

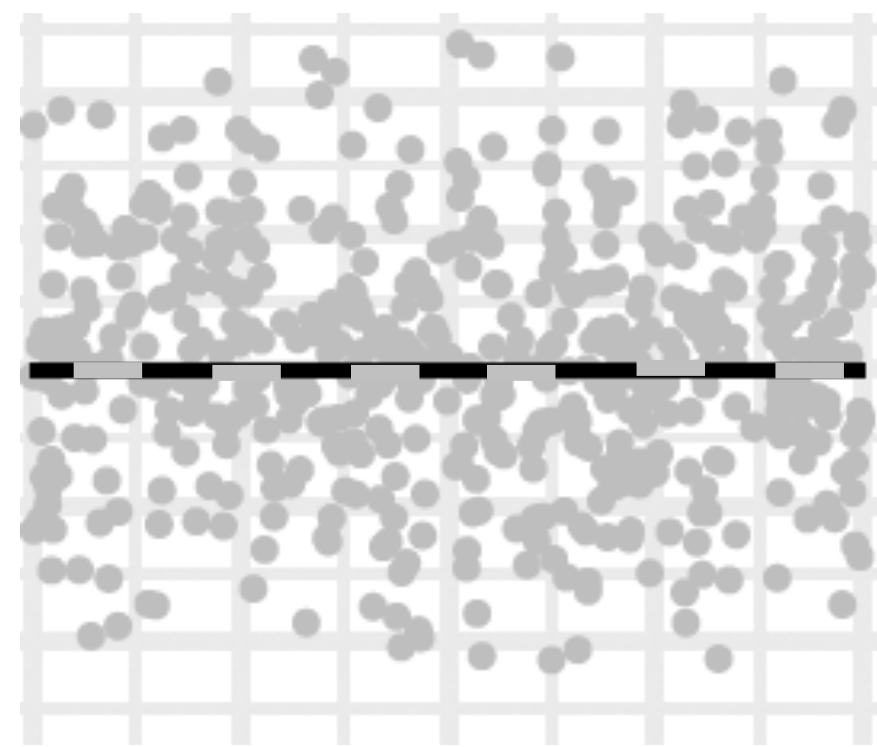
grey dots: data  
black line:  $\hat{f}_j$   
dotted line:  $f_j = 0$

# Problem Formulation

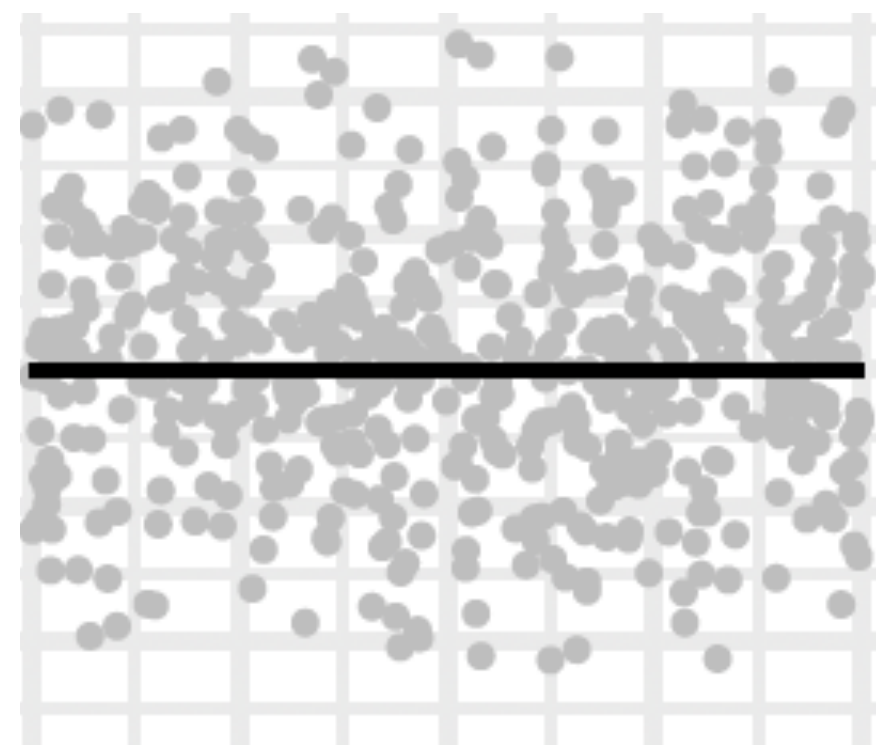
- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

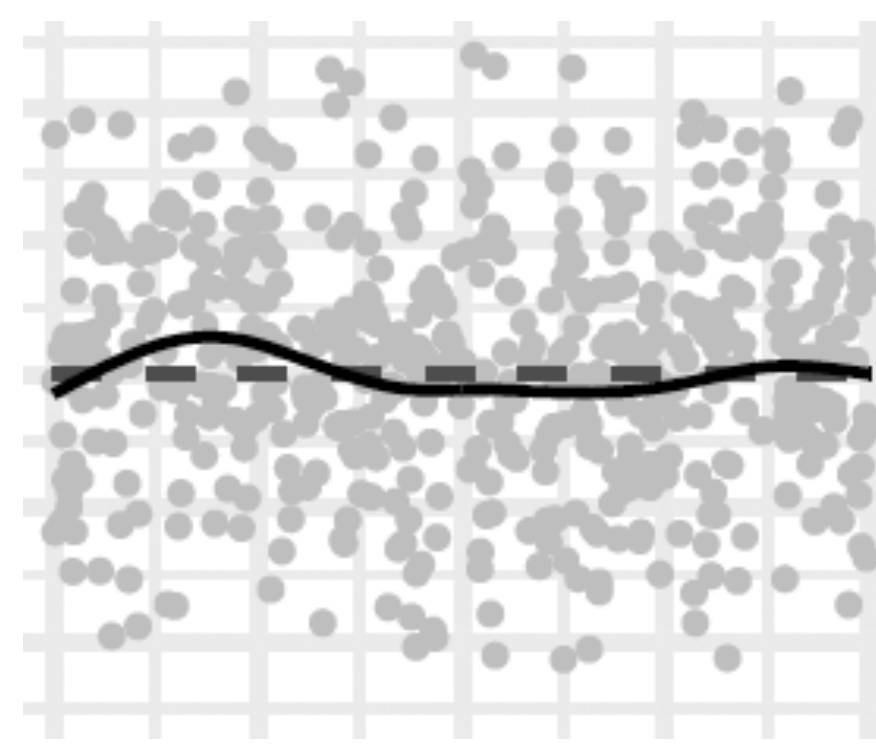
## VARIABLE SELECTION



True  $f_j = 0$



Good estimation



Bad estimation

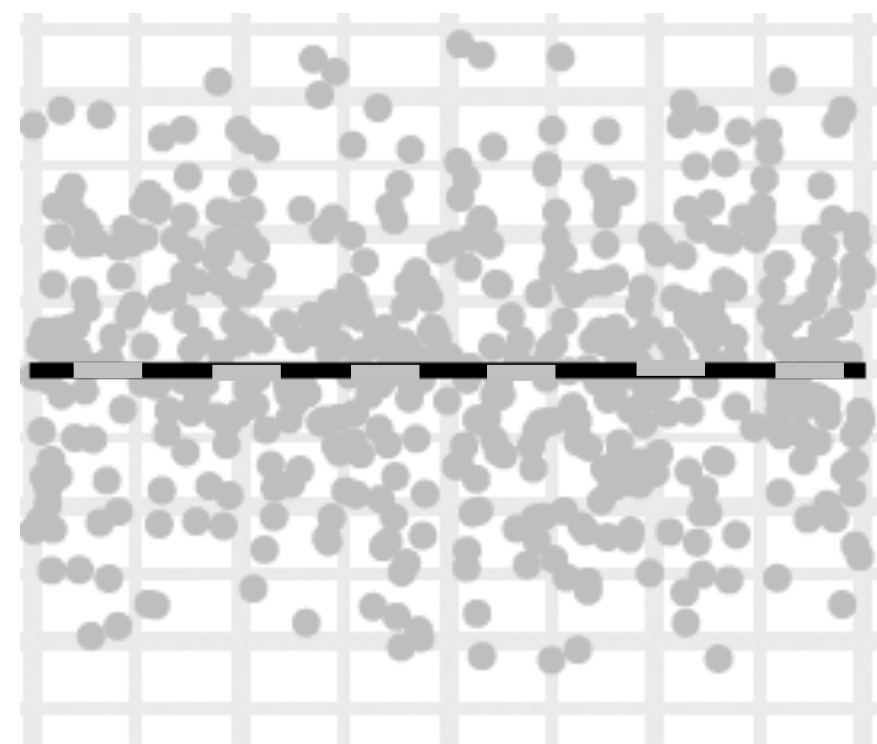
grey dots: data  
black line:  $\hat{f}_j$   
dotted line:  $f_j = 0$

# Problem Formulation

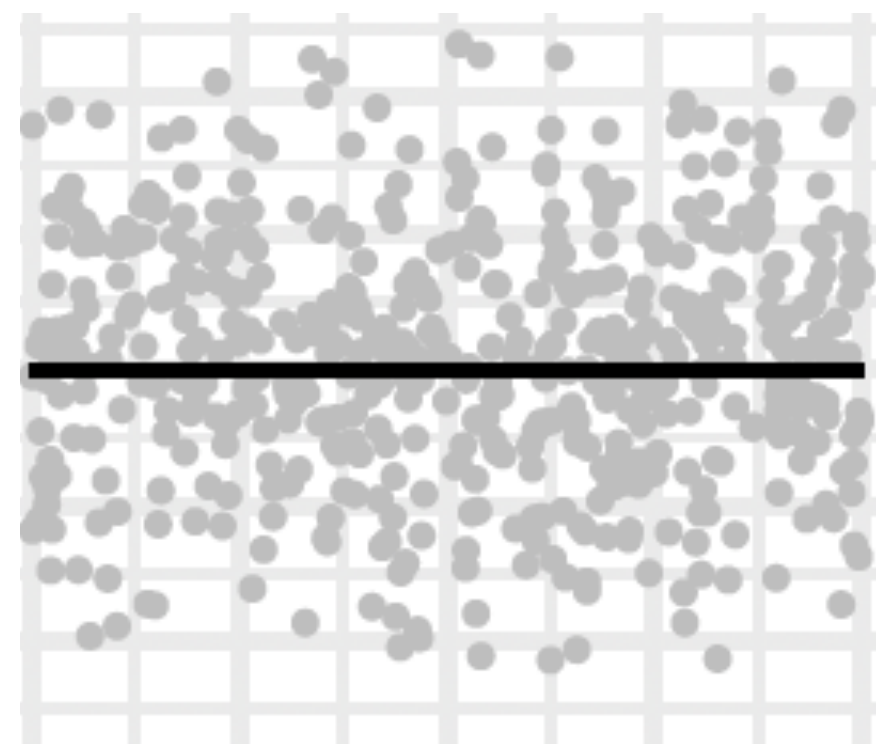
- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

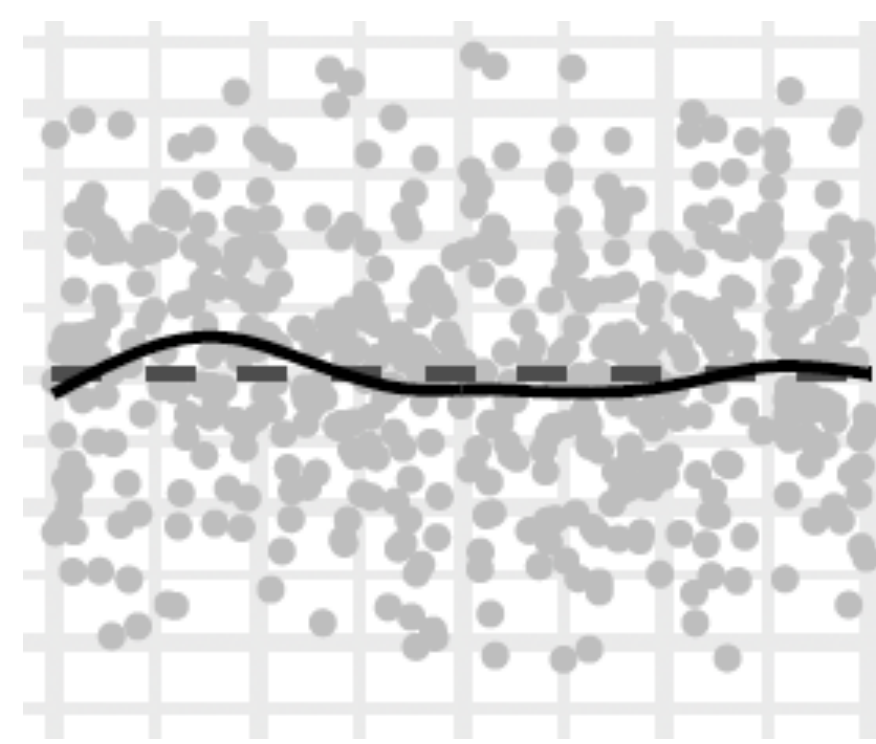
## VARIABLE SELECTION



True  $f_j = 0$



Good estimation



Bad estimation

grey dots: data  
black line:  $\hat{f}_j$   
dotted line:  $f_j = 0$

## METHODS

- gLasso-type penalty: COSSO, SPAM, GAMSEL
- BAYESIAN: spike and slab

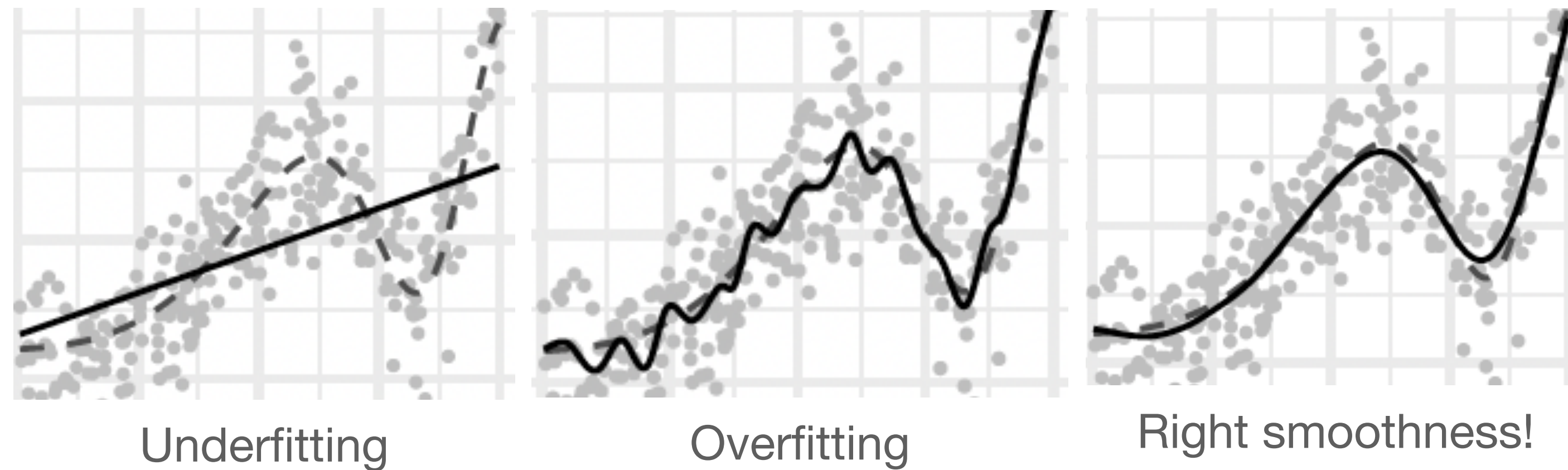


# Problem Formulation

- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

## SMOOTHING



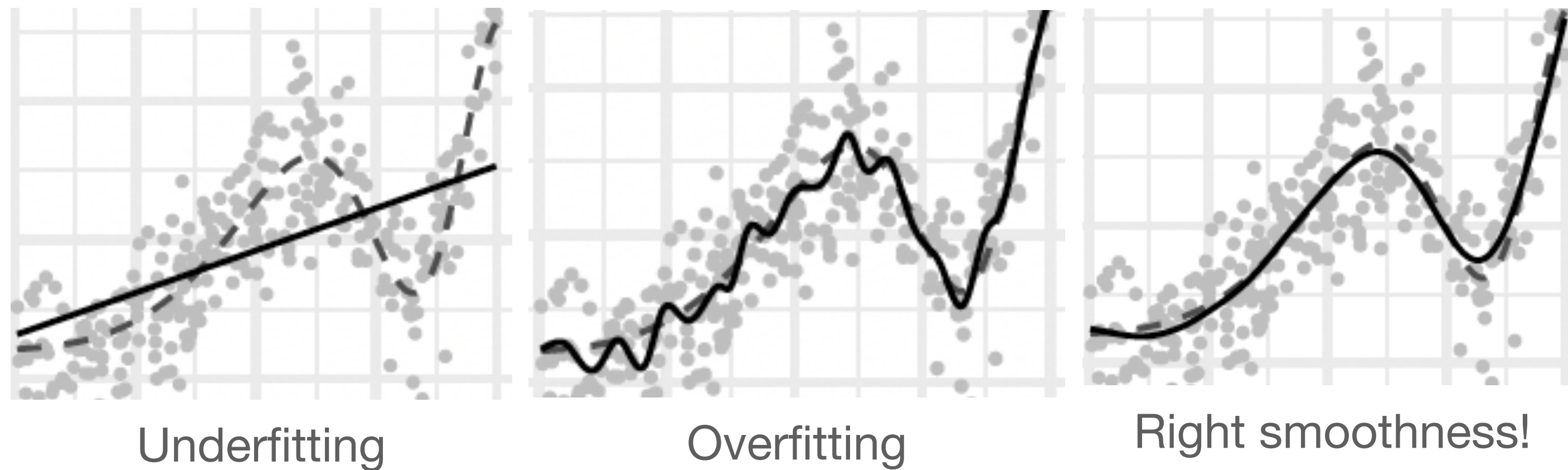
grey dots: data  
black line:  $\hat{f}_j$   
dotted line:  $f_j \neq 0$

# Problem Formulation

- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

## SMOOTHING



grey dots: data  
black line:  $\hat{f}_j$   
dotted line:  $f_j \neq 0$

## METHODS

- Smoothing splines with ridge type penalty
- BAYESIAN: Normal prior on the coefficients

# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models  $f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$



# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models  $f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$

## The proposed method

✓ performs both smoothing and variable selection

# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models  $f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$

## The proposed method

- ✓ performs both smoothing and variable selection
- ✓ can classify a feature's contribution as linear, non linear or zero

# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

## The proposed method

- ✓ performs both smoothing and variable selection
- ✓ can classify a feature's contribution as linear, non linear or zero
- ✓ can achieve exact sparsity

# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models

$$f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$$

## The proposed method

- ✓ performs both smoothing and variable selection
- ✓ can classify a feature's contribution as linear, non linear or zero
- ✓ can achieve exact sparsity
- ✓ is tuned with a single hyper parameter

# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models  $f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$

## The proposed method

- ✓ performs both smoothing and variable selection
- ✓ can classify a feature's contribution as linear, non linear or zero
- ✓ can achieve exact sparsity
- ✓ is tuned with a single hyper parameter
- ✓ is efficient



# Variational Automatic Relevance Determination

- Simultaneous smoothing and variable selection for additive models  $f(x_1, \dots, x_p) = f_1(x_1) + \dots + f_p(x_p)$

## The proposed method

- ✓ performs both smoothing and variable selection
- ✓ can classify a feature's contribution as linear, non linear or zero
- ✓ can achieve exact sparsity
- ✓ is tuned with a single hyper parameter
- ✓ is efficient
- ✓ outperforms other methods in accuracy in estimation and selection

# Setup

# Setup

$n$  responses  $\mathbf{y} = (y_1, \dots, y_n)^T$   
 $p$  predictors  $\mathbf{x}_j = \{(x_{j1}, \dots, x_{nj})^T\}_{j=1}^p$   
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$  iid  $N(0, \sigma^2)$

$$\mathbf{y} = \sum_{j=1}^p f_j(\mathbf{x}_j) + \boldsymbol{\epsilon}$$

where

Each  $f_j$  is represented through a basis expansion of linear and nonlinear terms

$$f_j(x) = \beta_0 x + \sum_{k=1}^{d_j} \beta_{jk} h_{jk}(x)$$

# Setup

$n$  responses  $\mathbf{y} = (y_1, \dots, y_n)^T$   
 $p$  predictors  $\mathbf{x}_j = \{(x_{j1}, \dots, x_{nj})^T\}_{j=1}^p$   
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$  iid  $N(0, \sigma^2)$

$$\mathbf{y} = \sum_{j=1}^p f_j(\mathbf{x}_j) + \boldsymbol{\epsilon}$$

where

Each  $f_j$  is represented through a basis expansion of linear and nonlinear terms

$$f_j(x) = \beta_0 x + \sum_{k=1}^{d_j} \beta_{jk} h_{jk}(x)$$

We then get  
a matrix  
representation

$$\mathbf{y} = \sum_{j=1}^{2p} \mathbf{Z}_j \boldsymbol{\beta}_j + \boldsymbol{\epsilon}$$

# Setup

$n$  responses  $\mathbf{y} = (y_1, \dots, y_n)^T$   
 $p$  predictors  $\mathbf{x}_j = \{(x_{j1}, \dots, x_{nj})^T\}_{j=1}^p$   
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T \text{ iid } N(0, \sigma^2)$

$$\mathbf{y} = \sum_{j=1}^p f_j(\mathbf{x}_j) + \boldsymbol{\epsilon}$$

where

Each  $f_j$  is represented through a basis expansion of linear and nonlinear terms

$$f_j(x) = \beta_0 x + \sum_{k=1}^{d_j} \beta_{jk} h_{jk}(x)$$

We then get  
a matrix  
representation

$$\mathbf{y} = \underbrace{Z_1 \boldsymbol{\beta}_1 + \dots + Z_p \boldsymbol{\beta}_p}_{p \text{ nonlinear terms}} + \underbrace{Z_{p+1} \boldsymbol{\beta}_{p+1} + \dots + Z_{2p} \boldsymbol{\beta}_{2p}}_{p \text{ linear terms}} + \boldsymbol{\epsilon}$$

$p$  nonlinear terms

$p$  linear terms



# Setup

$n$  responses  $\mathbf{y} = (y_1, \dots, y_n)^T$   
 $p$  predictors  $\mathbf{x}_j = \{(x_{j1j}, \dots, x_{nj})^T\}_{j=1}^p$   
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$  iid  $N(0, \sigma^2)$

$$\mathbf{y} = \sum_{j=1}^p f_j(\mathbf{x}_j) + \boldsymbol{\epsilon}$$

where

Each  $f_j$  is represented through a basis expansion of linear and nonlinear terms

$$f_j(x) = \beta_0 x + \sum_{k=1}^{d_j} \beta_{jk} h_{jk}(x)$$

Estimating the  $\beta_j$  can classify linear, non linear and zero contributions

$$\mathbf{y} = \underbrace{Z_1 \boldsymbol{\beta}_1 + \dots + Z_p \boldsymbol{\beta}_p}_{p \text{ nonlinear terms}} + \underbrace{Z_{p+1} \boldsymbol{\beta}_{p+1} + \dots + Z_{2p} \boldsymbol{\beta}_{2p}}_{p \text{ linear terms}} + \boldsymbol{\epsilon}$$

$p$  nonlinear terms

$p$  linear terms

# Setup

$n$  responses  $\mathbf{y} = (y_1, \dots, y_n)^T$   
 $p$  predictors  $\mathbf{x}_j = \{(x_{j1}, \dots, x_{nj})^T\}_{j=1}^p$   
 $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T \text{ iid } N(0, \sigma^2)$

$$\mathbf{y} = \sum_{j=1}^p f_j(\mathbf{x}_j) + \epsilon$$

where

Each  $f_j$  is represented through a basis expansion of linear and nonlinear terms

$$f_j(x) = \beta_0 x + \sum_{k=1}^{d_j} \beta_{jk} h_{jk}(x)$$

Estimating the  $\beta_j$  can classify linear, non linear and zero contributions

$$\mathbf{y} = \underbrace{Z_1 \beta_1 + \dots + Z_p \beta_p}_{p \text{ nonlinear terms}} + \underbrace{Z_{p+1} \beta_{p+1} + \dots + Z_{2p} \beta_{2p}}_{p \text{ linear terms}} + \epsilon$$

$p$  nonlinear terms

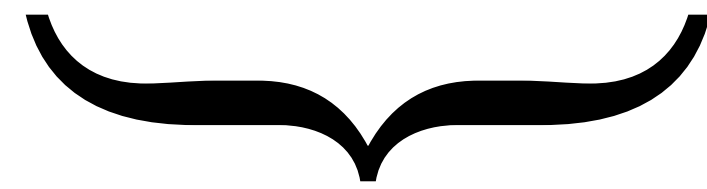
$p$  linear terms

**Q1:**

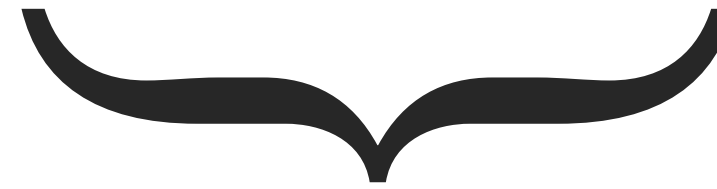
*What about large  $p$ ?  
In relation to feasibility, efficiency, and theoretical results*

# Setup

$$y = Z_1\beta_1 + \dots + Z_p\beta_p + Z_{p+1}\beta_{p+1} + \dots + Z_{2p}\beta_{2p} + \epsilon$$



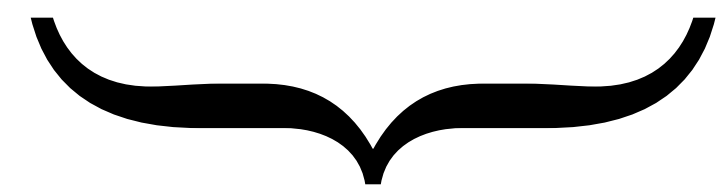
p nonlinear terms



p linear terms

# Group ARD Prior

$$y = Z_1\beta_1 + \dots + Z_p\beta_p + Z_{p+1}\beta_{p+1} + \dots + Z_{2p}\beta_{2p} + \epsilon$$



p nonlinear terms



p linear terms

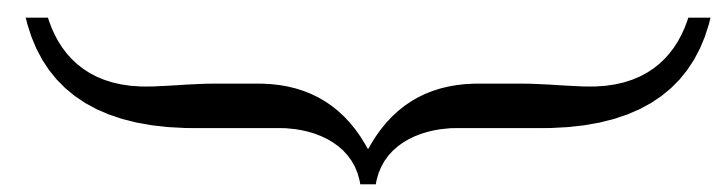
Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(0, r_j^2 I_{d_j})$$

smoothness and  
sparsity parameter

# Variational Inference

$$y = Z_1\beta_1 + \dots + Z_p\beta_p + Z_{p+1}\beta_{p+1} + \dots + Z_{2p}\beta_{2p} + \epsilon$$



p nonlinear terms



p linear terms

Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(0, r_j^2 I_{d_j})$$

smoothness and  
sparsity parameter

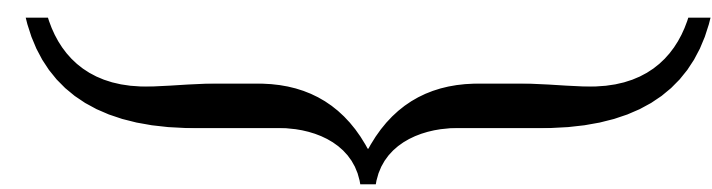
Variational  $q(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(\mu_j, \Phi_j)$$



# Variational Inference

$$y = Z_1\beta_1 + \dots + Z_p\beta_p + Z_{p+1}\beta_{p+1} + \dots + Z_{2p}\beta_{2p} + \epsilon$$



p nonlinear terms



p linear terms

Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(0, r_j^2 I_{d_j})$$

Q2:

*Implications of a different covariance matrix?*

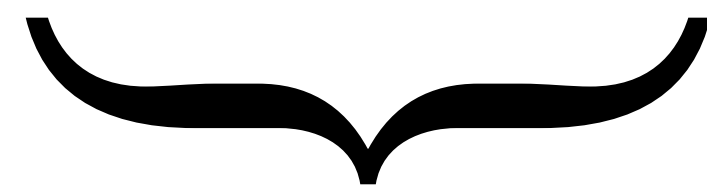
smoothness and  
sparsity parameter

Variational  $q(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(\mu_j, \Phi_j)$$

# $\alpha$ -Variational ELBO

$$y = Z_1 \beta_1 + \dots + Z_p \beta_p + Z_{p+1} \beta_{p+1} + \dots + Z_{2p} \beta_{2p} + \epsilon$$



p nonlinear terms



p linear terms

Prior  $p(\beta)$

$$\left( \beta_j \right)_{j=1}^{2p} \stackrel{ind}{\sim} N(0, r_j^2 I_{d_j})$$

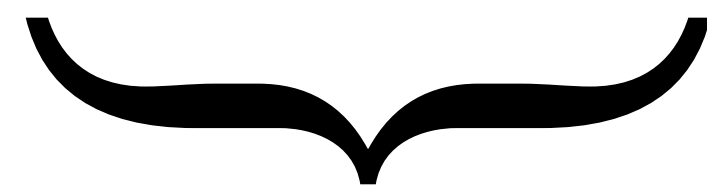
Variational  $q(\beta)$

$$\left( \beta_j \right)_{j=1}^{2p} \stackrel{ind}{\sim} N(\mu_j, \Phi_j)$$

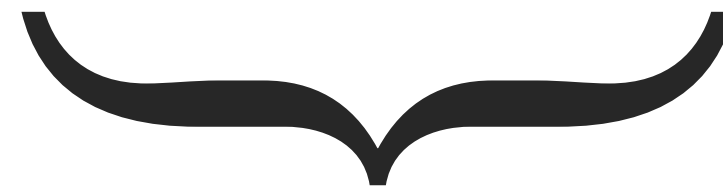
$$L(\mu, \Phi, r^2) = - E_q \left[ \log p \left( y | \beta_1, \dots, \beta_{2p} \right) \right] + \tilde{\alpha} KL(q || p)$$

# $\alpha$ -Variational ELBO

$$y = Z_1 \beta_1 + \dots + Z_p \beta_p + Z_{p+1} \beta_{p+1} + \dots + Z_{2p} \beta_{2p} + \epsilon$$



p nonlinear terms



p linear terms

Prior  $p(\beta)$

$$\left( \beta_j \right)_{j=1}^{2p} \stackrel{ind}{\sim} N(0, r_j^2 I_{d_j})$$

Variational  $q(\beta)$

$$\left( \beta_j \right)_{j=1}^{2p} \stackrel{ind}{\sim} N(\mu_j, \Phi_j)$$

$$L(\mu, \Phi, r^2) = -E_q \left[ \log p \left( y | \beta_1, \dots, \beta_{2p} \right) \right] + \tilde{\alpha} KL(q || p)$$

Q3:

Other forms of modifications on KL or on other divergences?

# Coordinate Descent

Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(0, r_j^2 I_{d_j})$$

Variational  $q(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{ind}{\sim} N(\mu_j, \Phi_j)$$

$$L(\mu, \Phi, r^2) = -E_q \left[ \log p \left( y | \beta_1, \dots, \beta_{2p} \right) \right] + \tilde{\alpha} KL(q || p)$$

# Coordinate Descent

Learn  $\left(\Phi_j, \mu_j, r_j^2\right)_{j=1}^{2p}$

- Can simplify into a univariate problem on  $r_j^2$
- Grid search to find the optimal values

Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{\text{ind}}{\sim} N(0, r_j^2 I_{d_j})$$

Variational  $q(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{\text{ind}}{\sim} N(\mu_j, \Phi_j)$$

$$L(\mu, \Phi, r^2) = -E_q \left[ \log p \left( y | \beta_1, \dots, \beta_{2p} \right) \right] + \tilde{\alpha} KL(q || p)$$



# Coordinate Descent

Learn  $\left(\Phi_j, \mu_j, r_j^2\right)_{j=1}^{2p}$

- Can simplify into a univariate problem on  $r_j^2$
- Grid search to find the optimal values

Q4:

*Other way than grid search? In relation to the feasibility in high dimensions.*

Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{\text{ind}}{\sim} N(0, r_j^2 I_{d_j})$$

Variational  $q(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{\text{ind}}{\sim} N(\mu_j, \Phi_j)$$

$$L(\mu, \Phi, r^2) = -E_q \left[ \log p \left( y | \beta_1, \dots, \beta_{2p} \right) \right] + \tilde{\alpha} KL(q || p)$$

# Coordinate Descent

Learn  $\left(\Phi_j, \mu_j, r_j^2\right)_{j=1}^{2p}$

- Can simplify into a univariate problem on  $r_j^2$
- Grid search to find the optimal values

Q4:

*Other way than grid search? In relation to the feasibility in high dimensions.*

Prior  $p(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{\text{ind}}{\sim} N(0, r_j^2 I_{d_j})$$

Variational  $q(\beta)$

$$\left(\beta_j\right)_{j=1}^{2p} \stackrel{\text{ind}}{\sim} N(\mu_j, \Phi_j)$$

Q5:

*How practical is tuning  $\alpha$ ?*

$$L(\mu, \Phi, r^2) = -E_q \left[ \log p \left( y | \beta_1, \dots, \beta_{2p} \right) \right] + \tilde{\alpha} KL(q || p)$$

Note:  $\alpha = \tilde{\alpha} \sigma^2$

# Lastly

Experiments with synthetic and real-world datasets demonstrating

- effectiveness of VARD in feature selection and individual smoothing
- capacity to differentiate nonlinear, linear, and zero functions
- estimation accuracy
- competing performance to other methods

# Lastly

Experiments with synthetic and real-world datasets demonstrating

- effectiveness of VARD in feature selection and individual smoothing
- capacity to differentiate nonlinear, linear, and zero functions
- estimation accuracy
- competing performance to other methods

**Q6:**

*Where does VARD stand in comparison to deep learning methods? How can it compare?*

Congratulations to the authors!

Thank you for your attention!