

Stat 577 “Statistical Learning Theory”

T2: Non-regularized regression methods

Xiongzhi Chen

Washington State University

Overview and Some Questions

Topics to cover

1. Regression with multivariate response
2. Regularized linear models
3. Other regularization techniques

Some questions: I

1. Let $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$. Let tr be the trace operator for matrices. Do you know that

$$\text{tr} \left(\mathbf{I} - \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \right) = n - (p + 1) \quad (1)$$

when $\text{rank}(\mathbf{X}) = p + 1$? Note $\mathbf{H} = \mathbf{X} \left(\mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T$ is often called the “hat matrix” in statistics.

2. What is your understanding of the term “degrees of freedom”? What is the degrees of freedom for

$$\hat{\sigma}^2 = \frac{1}{n - (p + 1)} \sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

where \hat{y}_i is the fitted value for y_i ?

Some questions: II

1. What is the purpose of regularization? What properties does an optimizer of the following regularization

$$R_\lambda(f) = \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda \int [f''(x)]^2 dx$$

have?

2. Can you explain what properties of f each of these controls: $\int [f''(x)]^2 dx$, $\int [f'(x)]^2 dx$ and $\int [f(x)]^2 dx$?
3. What is a general principle for variable/model selection?

Some questions: III

1. What is the key difficulty for Subset Selection?
2. Why cannot Backward Stepwise selection be applied when $N < p$, i.e., when the sample size is smaller than the number of predictors?
3. Among Best Subset selection, Forward-stepwise selection, and backward-stepwise selection, which produces a sequence of nested model?

Some questions: IV

1. Why in Figure 3.6 the mean-squared errors of the estimated coefficient are close for Best Subset, Forward Stepwise and Backward Stepwise selection? The text book claims that “Their performance is very similar, as is often the case.” Is there an example where their performances are not similar?
2. What are your comments about a variable selection method that includes a predictor if it is the most correlated with the current residual among all remaining predictors to be selected?
3. What are 3 principles of variable selection that were advocated by Dr. Jianqing Fan’s 2001 paper “Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties”? How to achieve these via a penalty?

Linear Models

Scalar response

- ▶ Predictor $X = (X_1, \dots, X_p)^T \in \mathbb{R}^p$ and response $Y \in \mathbb{R}$
- ▶ Coefficient vector $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T \in \mathbb{R}^{p+1}$ and population model

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon,$$

where $\varepsilon \in \mathbb{R}$ is the random error term

- ▶ Given observation vector $\mathbf{y} = (y_1, \dots, y_N)^T$ for Y and those \mathbf{x}_i for X , we define design matrix

$$\mathbf{X} = \begin{pmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_N^T \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{x}_1^T \\ \vdots & \vdots \\ 1 & \mathbf{x}_N^T \end{pmatrix} \in \mathbb{R}^{N \times (p+1)}$$

Scalar response

- ▶ Data version of model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, where $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^T \in \mathbb{R}^N$ and each $\varepsilon_i \sim \varepsilon$
- ▶ Minimizing “residual sums of squares (RSS)”

$$\text{RSS}(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

gives “least squares estimate (LSE)” of $\boldsymbol{\beta}$ as

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{y},$$

where \mathbf{A}^{-} is a generalized inverse of a matrix \mathbf{A}

- ▶ Least squares estimate of \mathbf{y} as

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-} \mathbf{X}^T \mathbf{y},$$

i.e., $\hat{\mathbf{y}}$ is the orthogonal projection of \mathbf{y} onto the column space of \mathbf{X}

- ▶ $\hat{\mathbf{y}}$ is unique but $\hat{\boldsymbol{\beta}}$ not always

Distributional properties of LSE

- ▶ Recall the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ and LSE $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T\mathbf{y}$
- ▶ Assume $\boldsymbol{\varepsilon} \sim \text{Normal}(0, \sigma^2\mathbf{I})$. Then we have unbiasedness

$$E[\hat{\boldsymbol{\beta}}] = E\left[(\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})\right] = \boldsymbol{\beta}$$

- ▶ Further, set $\mathbf{S} = \mathbf{X}^T\mathbf{X}$. Then $\hat{\boldsymbol{\beta}} - \boldsymbol{\beta} = \mathbf{S}^{-1}\mathbf{X}^T\boldsymbol{\varepsilon}$ and

$$\begin{aligned}\text{Var}(\hat{\boldsymbol{\beta}}) &= E\left\{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T\right\} = E\left\{\mathbf{S}^{-1}\mathbf{X}^T\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T(\mathbf{S}^{-1}\mathbf{X}^T)^T\right\} \\ &= (\mathbf{S}^{-1}\mathbf{X}^T) E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T) (\mathbf{S}^{-1}\mathbf{X}^T)^T \\ &= \sigma^2\mathbf{S}^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{S}^{-1})^T = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}\end{aligned}$$

Estimating variance of error term

- ▶ Recall the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, LSE $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ and $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = (\hat{y}_1, \dots, \hat{y}_N)^T$
- ▶ Assume $\boldsymbol{\varepsilon} \sim \text{Normal}(0, \sigma^2 \mathbf{I})$. Then an estimate of σ^2 is

$$\hat{\sigma}^2 = \frac{1}{N - p - 1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{\|\mathbf{y} - \hat{\mathbf{y}}\|^2}{N - p - 1},$$

and

$$(N - p - 1) \hat{\sigma}^2 \sim \sigma^2 \chi_{N-p-1}^2$$

- ▶ Further, $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ are statistically independent

Linear model: recap

- ▶ Recall the model

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon \text{ with } E(\varepsilon) = 0$$

- ▶ Equivalent model: $E(Y|X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$
- ▶ Caution: if $E(Y|X)$ is restricted in, e.g., $[0, 1]$ but $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T \in \mathbb{R}^{p+1}$ is allowed, then the above model is insensible
- ▶ In general, we model

$$g[E(Y|X)] = \beta_0 + \sum_{j=1}^p X_j \beta_j,$$

where g is a suitable “link function”

Derivation of LSE: I

► Recall

$$\begin{aligned} \text{RSS}(\boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 \\ &= \mathbf{y}^T \mathbf{y} - \underbrace{(\mathbf{X}\boldsymbol{\beta})^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\beta}}_{=2\mathbf{y}^T \mathbf{X}\boldsymbol{\beta}} + \boldsymbol{\beta}^T \underbrace{\mathbf{X}^T \mathbf{X} \boldsymbol{\beta}}_{=\mathbf{S}} \end{aligned} \quad (2)$$

- Vector calculus to derive $\hat{\boldsymbol{\beta}} \in \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^{p+1}} \text{RSS}(\boldsymbol{\beta})$
- Let $\mathbf{e}_j \in \mathbb{R}^{p+1}$ have its j th entry as 1 and the rest as 0. Set $\mathbf{a} = \mathbf{y}^T \mathbf{X} = (a_1, \dots, a_{p+1})$. Then

$$\partial_{\beta_j}(\mathbf{a}\boldsymbol{\beta}) = \lim_{t \rightarrow 0} \frac{\mathbf{a}(\boldsymbol{\beta} + t\mathbf{e}_j) - \mathbf{a}\boldsymbol{\beta}}{t} = \mathbf{a}\mathbf{e}_j = a_j,$$

and the gradient $\nabla_{\boldsymbol{\beta}}$, always a column vector, is

$$\nabla_{\boldsymbol{\beta}}(\mathbf{y}^T \mathbf{X}\boldsymbol{\beta}) = \nabla_{\boldsymbol{\beta}}(\mathbf{a}\boldsymbol{\beta}) = \mathbf{a}^T = \mathbf{X}^T \mathbf{y} \quad (3)$$

Derivation of LSE: II

- ▶ Obtain the gradient $\nabla_{\beta} (\beta^T \mathbf{S} \beta)$ as follows:

$$\begin{aligned}\partial_{\beta_j} (\beta^T \mathbf{S} \beta) &= \lim_{t \rightarrow 0} t^{-1} \left[(\beta + t \mathbf{e}_j)^T \mathbf{S} (\beta + t \mathbf{e}_j) - \beta^T \mathbf{S} \beta \right] \\ &= \lim_{t \rightarrow 0} t^{-1} \left[t \beta^T \mathbf{S} \mathbf{e}_j + t \mathbf{e}_j^T \mathbf{S} \beta + t^2 \mathbf{e}_j^T \mathbf{S} \mathbf{e}_j \right] \\ &= \lim_{t \rightarrow 0} t^{-1} \left[2t \mathbf{e}_j^T \mathbf{S} \beta + t^2 \mathbf{e}_j^T \mathbf{S} \mathbf{e}_j \right] \quad \text{since } \mathbf{S} = \mathbf{S}^T \\ &= 2 \mathbf{e}_j^T \mathbf{S} \beta = 2 (\mathbf{S} \beta)_j = \text{the } j\text{th entry of } \mathbf{S} \beta\end{aligned}$$

and

$$\nabla_{\beta} (\beta^T \mathbf{S} \beta) = 2 \mathbf{S} \beta = 2 \mathbf{X}^T \mathbf{X} \beta \quad (4)$$

- ▶ Note: $\mathbf{e}_j^T \mathbf{b}$ gives the j th entry of \mathbf{b} , and $\mathbf{e}_j^T \mathbf{A}$ the j th row of \mathbf{A} , and $\mathbf{A} \mathbf{e}_j$ the j th column of \mathbf{A}

Derivation of LSE: III

- ▶ Combining (2), (3) and (4) gives

$$\nabla_{\beta} \text{RSS}(\beta) = \nabla_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 = -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta$$

- ▶ The critical condition for $\hat{\beta} \in \operatorname{argmin}_{\beta \in \mathbb{R}^{p+1}} \text{RSS}(\beta)$ forces $\nabla_{\beta} \text{RSS}(\beta)|_{\hat{\beta}} = 0$. This gives the “normal equation” for $\hat{\beta}$ as

$$-2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X} \beta = 0 \Leftrightarrow \mathbf{X}^T \mathbf{X} \beta = \mathbf{X}^T \mathbf{y}, \quad (5)$$

where the symbol \Leftrightarrow denotes “equivalence” or “if and only if”

- ▶ So, LSE $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

Vector response

- ▶ Linear model with a scalar response: $Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon$ and $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- ▶ Linear model with vector response:
 - ▶ Let $Y = (Y_1, \dots, Y_m)$ and $X = (1, X_1, \dots, X_p)$
 - ▶ Let $\mathbf{B} \in \mathbb{R}^{(p+1) \times m}$ be the unknown coefficient matrix, and $\boldsymbol{\varepsilon} \in \mathbb{R}^m$ the random error vector
 - ▶ Models:

$$\left\{ \begin{array}{ll} \text{population version:} & Y = \mathbf{X}\mathbf{B} + \boldsymbol{\varepsilon} \\ \text{data version:} & \mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \\ \text{data version (rowwise):} & \mathbf{y}_i = \mathbf{x}_i \mathbf{B} + \boldsymbol{\varepsilon}_i \end{array} \right. ,$$

where the i th row \mathbf{y}_i of \mathbf{Y} contains the i th observation for Y , the i th row \mathbf{x}_i of \mathbf{X} the i th observation on X , and i th row $\boldsymbol{\varepsilon}_i$ of \mathbf{E} the i th realization of $\boldsymbol{\varepsilon}$ for $1 \leq i \leq N$

LSE of coefficient matrix

- ▶ Recall model $\mathbf{y}_i = \mathbf{x}_i \mathbf{B} + \varepsilon_i$ or $\mathbf{Y} = \mathbf{X} \mathbf{B} + \mathbf{E}$ and

$$\text{RSS}(\mathbf{B}) = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{x}_i \mathbf{B}\|^2$$

- ▶ Then minimizing $\text{RSS}(\mathbf{B})$ for $\mathbf{B} \in \mathbb{R}^{(p+1) \times m}$ gives the LSE

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^- \mathbf{X}^T \mathbf{Y},$$

where $(\mathbf{X}^T \mathbf{X})^-$ is a general inverse of $\mathbf{S} = \mathbf{X}^T \mathbf{X}$, and is \mathbf{S}^{-1} if \mathbf{S} is invertible

Derivation of LSE: overview

- ▶ Frobenius norm and trace operator
- ▶ Special matrices and basis
- ▶ Matrix calculus

Trace operator and Frobenius norm

Let $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times m}$.

- ▶ Define its “Frobenius norm” as

$$\|\mathbf{A}\| = \sqrt{\sum_{i=1}^n \sum_{j=1}^m a_{ij}^2}$$

- ▶ When $m = n$, i.e., \mathbf{A} is a square matrix, its “trace” is defined as

$$\text{trace}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$$

- ▶ For any two (compatible) matrices \mathbf{A} and \mathbf{C} , we have

$$\begin{cases} \text{trace}(\mathbf{A}) = \text{trace}(\mathbf{A}^T); \text{trace}(\mathbf{AC}) = \text{trace}(\mathbf{CA}); \\ \text{trace}(\mathbf{A} + \mathbf{C}) = \text{trace}(\mathbf{C} + \mathbf{A}) \end{cases}$$

since trace operates on the diagonal entries of a matrix

Frobenius norm

- ▶ Let $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times m}$ and write

$$\mathbf{A} = \underbrace{\begin{pmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{pmatrix}}_{\text{row layout}} = \underbrace{\begin{pmatrix} \tilde{\mathbf{a}}_1 & \cdots & \tilde{\mathbf{a}}_m \end{pmatrix}}_{\text{column layout}}$$

- ▶ Recall the Frobenius norm $\|\mathbf{A}\| = \sqrt{\sum_{i=1}^n \sum_{j=1}^m a_{ij}^2}$. Then

$$\|\mathbf{A}\| = \sqrt{\sum_{i=1}^n \|\mathbf{a}_i\|^2} = \sqrt{\sum_{i=1}^m \|\tilde{\mathbf{a}}_i\|^2}$$

- ▶ Also, $\|\mathbf{A}\| = \|\text{vec}(\mathbf{A})\|$, where

$$\text{vec}(\mathbf{A}) = \underbrace{\begin{pmatrix} \tilde{\mathbf{a}}_1^T \\ \vdots \\ \tilde{\mathbf{a}}_m^T \end{pmatrix}}_{\text{column concatenation}} \quad \text{or} \quad \text{vec}(\mathbf{A}) = \underbrace{(\mathbf{a}_1, \dots, \mathbf{a}_n)}_{\text{row concatenation}}$$

Trace operator and Frobenius norm

- ▶ Let $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times m}$ and write

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{pmatrix} = (\tilde{\mathbf{a}}_1 \quad \cdots \quad \tilde{\mathbf{a}}_m)$$

- ▶ Its Frobenius norm: $\|\mathbf{A}\| = \sqrt{\sum_{i=1}^n \|\mathbf{a}_i\|^2} = \sqrt{\sum_{j=1}^m \|\tilde{\mathbf{a}}_j\|^2}$
- ▶ Since $(\mathbf{A}\mathbf{A}^T)(i,i) = \mathbf{a}_i\mathbf{a}_i^T = \|\mathbf{a}_i\|^2$ and $(\mathbf{A}^T\mathbf{A})(j,j) = \|\tilde{\mathbf{a}}_j\|^2$, then

$$\begin{aligned} \|\mathbf{A}\|^2 &= \sum_{i=1}^n (\mathbf{A}\mathbf{A}^T)(i,i) = \text{trace}(\mathbf{A}\mathbf{A}^T) \\ &= \sum_{j=1}^m (\mathbf{A}^T\mathbf{A})(j,j) = \text{trace}(\mathbf{A}^T\mathbf{A}) \end{aligned}$$

- ▶ Namely, $\|\mathbf{A}\| = \sqrt{\text{trace}(\mathbf{A}^T\mathbf{A})} = \sqrt{\text{trace}(\mathbf{A}\mathbf{A}^T)}$

Frobenius norm and inner product

- ▶ Let $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times m}$ and recall its Frobenius norm

$$\|\mathbf{A}\| = \sqrt{\sum_{i=1}^n \sum_{j=1}^m a_{ij}^2} = \sqrt{\text{trace}(\mathbf{A}^T \mathbf{A})} = \sqrt{\text{trace}(\mathbf{A} \mathbf{A}^T)}$$

- ▶ For $\mathbf{C} \in \mathbb{R}^{n \times m}$, we can define “inner product”

$$\langle \mathbf{A}, \mathbf{C} \rangle = \text{trace}(\mathbf{A}^T \mathbf{C}) = \text{trace}(\mathbf{C}^T \mathbf{A}) \quad (6)$$

- ▶ When \mathbf{A} and \mathbf{C} are vectors in \mathbb{R}^s , $\langle \mathbf{A}, \mathbf{C} \rangle$ is just the Euclidean inner product
- ▶ With (6), we can endow a geometric structure on $\mathbb{R}^{n \times m}$, where we can discuss angle, length, etc and where essentially we regard \mathbf{A} as $\text{vec}(\mathbf{A})$

Special matrices

- ▶ Let $\mathbf{P}_{ij} \in \mathbb{R}^{n \times m}$ have only its (i, j) th entry as 1 and the rest as 0
- ▶ The set $\{\mathbf{P}_{ij}\}_{1 \leq i \leq n, 1 \leq j \leq m}$ form a basis for $\mathbb{R}^{n \times m}$, i.e., any $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times m}$ is a linear combination of the \mathbf{P}_{ij} 's. In fact,

$$\mathbf{A} = \sum_{i=1}^n \sum_{j=1}^m a_{ij} \mathbf{P}_{ij} \Leftrightarrow \mathbf{A} = (a_{ij})$$

- ▶ Note $\mathbf{P}_{ij}^T = \mathbf{P}_{ji}$. Let $\mathbf{P}_{ij}^T \mathbf{A} \in \mathbb{R}^{m \times m}$. Then

$$\mathbf{P}_{ij}^T \mathbf{A} = \begin{pmatrix} 0 \\ a_{i1}, a_{i2}, \dots, a_{im} \\ 0 \end{pmatrix},$$

i.e., the j th row of $\mathbf{P}_{ij}^T \mathbf{A}$ is the i th row of $\mathbf{A} = (a_{ij})$, and all other rows of $\mathbf{P}_{ij}^T \mathbf{A}$ are zero

Matrix calculus

- ▶ Recall $\mathbf{P}_{ij} \in \mathbb{R}^{n \times m}$ has only its (i, j) th entry as 1 and the rest as 0, and pick any $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times m}$
- ▶ Let $g : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$ and define

$$\frac{\partial g(\mathbf{A})}{\partial \mathbf{A}} = \left(\frac{\partial g(\mathbf{A})}{\partial a_{ij}} \right) \in \mathbb{R}^{n \times m},$$

where the partial derivative

$$\frac{\partial}{\partial a_{ij}} g(\mathbf{A}) = \lim_{t \rightarrow 0} \frac{1}{t} [g(\mathbf{A} + t\mathbf{P}_{ij}) - g(\mathbf{A})]$$

- ▶ An example g is

$$\begin{aligned} g(\mathbf{B}) &= \text{RSS}(\mathbf{B}) = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{x}_i \mathbf{B}\|^2 = \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|^2 \\ &= \text{trace} \left[(\mathbf{Y} - \mathbf{X}\mathbf{B})^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right] \end{aligned}$$

Derivation of LSE: strategy

- ▶ Recall $\mathbf{B} = (b_{ij})$ and

$$\text{RSS}(\mathbf{B}) = g(\mathbf{B}) = \text{trace} \left[(\mathbf{Y} - \mathbf{XB})^T (\mathbf{Y} - \mathbf{XB}) \right]$$

and

$$\frac{\partial g(\mathbf{B})}{\partial \mathbf{B}} = \left(\frac{\partial g(\mathbf{B})}{\partial b_{ij}} \right) = \left(\lim_{t \rightarrow 0} \frac{1}{t} [g(\mathbf{B} + t\mathbf{P}_{ij}) - g(\mathbf{B})] \right)$$

- ▶ The critical condition for

$$\hat{\mathbf{B}} \in \underset{\mathbf{B} \in \mathbb{R}^{(p+1) \times m}}{\text{argmin}} \text{RSS}(\boldsymbol{\beta}) \implies \left. \frac{\partial g(\mathbf{B})}{\partial \mathbf{B}} \right|_{\hat{\mathbf{B}}} = 0,$$

which induces the normal equation for $\hat{\mathbf{B}}$ as

$$-2\mathbf{X}^T (\mathbf{Y} - \mathbf{XB}) = 0 \Leftrightarrow \mathbf{X}^T \mathbf{Y} = \mathbf{X}^T \mathbf{XB} \Rightarrow \hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- ▶ Step 1: compute difference as

$$\begin{aligned} & g(\mathbf{B} + t\mathbf{P}_{ij}) - g(\mathbf{B}) \\ &= \text{trace} \left\{ [\mathbf{Y} - \mathbf{X}(\mathbf{B} + t\mathbf{P}_{ij})]^T [\mathbf{Y} - \mathbf{X}(\mathbf{B} + t\mathbf{P}_{ij})] \right\} \\ &\quad - \text{trace} \left[(\mathbf{Y} - \mathbf{X}\mathbf{B})^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right] \\ &= \text{trace} \left\{ [(\mathbf{Y} - \mathbf{X}\mathbf{B}) - t\mathbf{X}\mathbf{P}_{ij}]^T [(\mathbf{Y} - \mathbf{X}\mathbf{B}) - t\mathbf{X}\mathbf{P}_{ij}] \right\} \\ &\quad - \text{trace} \left[(\mathbf{Y} - \mathbf{X}\mathbf{B})^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right] \\ &= \text{trace} \left\{ -t(\mathbf{Y} - \mathbf{X}\mathbf{B})^T \mathbf{X}\mathbf{P}_{ij} - t\mathbf{P}_{ij}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) + t^2 \mathbf{P}_{ij}^T \mathbf{X}^T \mathbf{X} \mathbf{P}_{ij} \right\} \end{aligned}$$

(note linear terms and quadratic term in t)

Derivation of LSE: II

- ▶ Step 2: compute partial derivative as

$$\begin{aligned}\frac{\partial}{\partial b_{ij}} g(\mathbf{B}) &= \lim_{t \rightarrow 0} \frac{1}{t} [g(\mathbf{B} + t\mathbf{P}_{ij}) - g(\mathbf{B})] \\ &= \lim_{t \rightarrow 0} \frac{1}{t} \text{trace} \left\{ -t(\mathbf{Y} - \mathbf{XB})^T \mathbf{XP}_{ij} - t\mathbf{P}_{ij}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{XB}) \right\} \\ &\quad + \lim_{t \rightarrow 0} \frac{1}{t} \text{trace} \left\{ t^2 \mathbf{P}_{ij}^T \mathbf{X}^T \mathbf{XP}_{ij} \right\} \\ &= \text{trace} \left\{ -(\mathbf{Y} - \mathbf{XB})^T \mathbf{XP}_{ij} - \mathbf{P}_{ij}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{XB}) \right\} \\ &= -2 \text{trace} \left\{ \mathbf{P}_{ij}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{XB}) \right\},\end{aligned}$$

since $\text{trace}(\mathbf{C}) = \text{trace}(\mathbf{C}^T)$ with $\mathbf{C} = \mathbf{P}_{ij}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{XB})$

Derivation of LSE: III

- ▶ Step 3: obtain matrix derivative. Set

$$\mathbf{A} = \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) = (a_{ij}) \quad \text{and} \quad \mathbf{Q} = \mathbf{P}_{ij}^T \mathbf{A}.$$

Then

$$\mathbf{Q} = \mathbf{P}_{ij}^T \mathbf{A} = \begin{pmatrix} 0 \\ (a_{i1}, \dots, a_{ij}, \dots, a_{im}) \\ 0 \end{pmatrix} \quad (j\text{th row}),$$

i.e., the j th row of \mathbf{Q} is the i th row of \mathbf{A} and the rest are 0

- ▶ So,

$$\frac{\partial g(\mathbf{B})}{\partial b_{ij}} = -2 \text{trace} \left\{ \mathbf{P}_{ij}^T \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right\} = -2a_{ij}$$

i.e., $\frac{\partial}{\partial b_{ij}} g(\mathbf{B})$ is the (i, j) th entry of \mathbf{A}

- ▶ By the definition of $\nabla_{\mathbf{B}}$ as $\nabla_{\mathbf{B}} = \left(\frac{\partial}{\partial b_{ij}} \right)$, we see

$$\nabla_{\mathbf{B}} g(\mathbf{B}) = -2\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{B})$$

- ▶ Obtain normal equation: since

$$\nabla_{\mathbf{B}} g(\mathbf{B}) = -2\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{B}),$$

setting $\nabla_{\mathbf{B}} g(\mathbf{B}) = 0$ gives

$$-2\mathbf{X}^T (\mathbf{Y} - \mathbf{X}\mathbf{B}) = 0 \Leftrightarrow \mathbf{X}^T \mathbf{Y} = \mathbf{X}^T \mathbf{X}\mathbf{B}$$

- ▶ When $\mathbf{X}^T \mathbf{X}$ is invertible, the unique solution is

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \tag{7}$$

Dependence adjustment: settings and LSE

- ▶ Recall model:

$$\begin{cases} \text{data version:} & \mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \\ \text{data version (rowwise):} & \mathbf{y}_i = \mathbf{x}_i\mathbf{B} + \varepsilon_i \end{cases}$$

- ▶ Assume: $\varepsilon_i \sim G(0, \Sigma)$ with $|\Sigma| \neq 0$ and $(\mathbf{X}^T\mathbf{X})^{-1}$ exists
- ▶ Then $\hat{\mathbf{B}} = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T\mathbf{Y}$ minimizes

$$h(\mathbf{B}) = \sum_{i=1}^N \underbrace{(\mathbf{y}_i - \mathbf{x}_i\mathbf{B}) \Sigma^{-1} (\mathbf{y}_i - \mathbf{x}_i\mathbf{B})^T}_{=r_i(\Sigma) = \|r_i(\Sigma)\|^2} = \left\| (\mathbf{Y} - \mathbf{X}\mathbf{B}) \Sigma^{-1/2} \right\|^2,$$

- ▶ Note: residual vector $\mathbf{r}_i = \mathbf{y}_i - \mathbf{x}_i\mathbf{B}$ is adjusted by $\Sigma^{-1/2}$ into

$$\mathbf{r}_i(\Sigma) = (\mathbf{y}_i - \mathbf{x}_i\mathbf{B}) \Sigma^{-1/2}$$

LSE based on adjusted residuals

- ▶ “Spectral decomposition”: $\Sigma = \mathbf{U}\mathbf{D}\mathbf{U}^T$, where $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ and \mathbf{D} is a diagonal matrix. Then $\Sigma^{-1/2} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{U}^T$
- ▶ Let $\tilde{\mathbf{y}}_i = \mathbf{y}_i\Sigma^{-1/2}$ and $\tilde{\mathbf{Y}} = \mathbf{Y}\Sigma^{-1/2}$. Then the covariance of $\tilde{\mathbf{y}}_i$ is \mathbf{I}
- ▶ Set $\tilde{\mathbf{B}} = \mathbf{B}\Sigma^{-1/2}$. Then

$$\begin{aligned}h(\mathbf{B}) &= \sum_{i=1}^N (\mathbf{y}_i - \mathbf{x}_i\mathbf{B}) \Sigma^{-1} (\mathbf{y}_i - \mathbf{x}_i\mathbf{B})^T \\ &\iff h(\tilde{\mathbf{B}}) = \sum_{i=1}^N (\tilde{\mathbf{y}}_i - \mathbf{x}_i\tilde{\mathbf{B}}) (\tilde{\mathbf{y}}_i - \mathbf{x}_i\tilde{\mathbf{B}})^T\end{aligned}$$

- ▶ The minimizer of $h(\tilde{\mathbf{B}})$ is

$$\mathbf{B}^* = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T\tilde{\mathbf{Y}} = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T\mathbf{Y}\Sigma^{-1/2},$$

and $\tilde{\mathbf{B}} = \mathbf{B}\Sigma^{-1/2}$ gives $\hat{\mathbf{B}} = (\mathbf{X}^T\mathbf{X})^{-1} \mathbf{X}^T\mathbf{Y}$ as the minimizer of $h(\mathbf{B})$

LSE fails to incorporate dependence adjustment

- ▶ $\hat{\mathbf{B}} = (\mathbf{X}\mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ is the minimizer of the “weighted RSS”

$$h(\mathbf{B}) = \sum_{i=1}^N (\mathbf{y}_i - \mathbf{x}_i \mathbf{B}) \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \mathbf{x}_i \mathbf{B})^T = \left\| (\mathbf{Y} - \mathbf{X}\mathbf{B}) \boldsymbol{\Sigma}^{-1/2} \right\|^2,$$

where $\mathbf{r}_i = \mathbf{y}_i - \mathbf{x}_i \mathbf{B}$ is adjusted by $\boldsymbol{\Sigma}^{-1/2}$ (encoding dependence)

- ▶ However, $\hat{\mathbf{B}} = (\mathbf{X}\mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ is also the minimizer of

$$g(\mathbf{B}) = \sum_{i=1}^N (\mathbf{y}_i - \mathbf{x}_i \mathbf{B}) (\mathbf{y}_i - \mathbf{x}_i \mathbf{B})^T = \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|^2$$

- ▶ Namely, for a linear model, least squares method cannot take into account the dependence structure of random error vector

Weighted least squares: settings and LSE

- ▶ Model with $\tilde{\boldsymbol{\varepsilon}}_j \sim G(0, \boldsymbol{\Sigma})$,

$$\begin{cases} \text{data version:} & \mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E} \\ \text{data version (columnwise):} & \tilde{\mathbf{y}}_j = \mathbf{X}\tilde{\mathbf{b}}_j + \tilde{\boldsymbol{\varepsilon}}_j \end{cases}$$

- ▶ Consider the weighted RSS

$$u(\mathbf{B}) = \sum_{j=1}^N \left[\boldsymbol{\Sigma}^{-1/2} (\tilde{\mathbf{y}}_j - \mathbf{X}\tilde{\mathbf{b}}_j) \right]^T \left[\boldsymbol{\Sigma}^{-1/2} (\tilde{\mathbf{y}}_j - \mathbf{X}\tilde{\mathbf{b}}_j) \right]$$

Then we have the “weighted LSE”

$$\hat{\mathbf{B}} = \left(\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{Y}$$

Derivation of weighted LSE

- ▶ “Spectral decomposition”: $\Sigma = \mathbf{U}\mathbf{D}\mathbf{U}^T$, where $\mathbf{U}\mathbf{U}^T = \mathbf{I}$ and \mathbf{D} is a diagonal matrix. Then $\Sigma^{-1/2} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{U}^T$ and

$$u(\mathbf{B}) = \text{trace} \left\{ \left[\Sigma^{-1/2} (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right]^T \left[\Sigma^{-1/2} (\mathbf{Y} - \mathbf{X}\mathbf{B}) \right] \right\}$$

- ▶ Set $\tilde{\mathbf{X}} = \Sigma^{-1/2}\mathbf{X}$ and $\tilde{\mathbf{Y}} = \Sigma^{-1/2}\mathbf{Y}$. Then

$$u(\mathbf{B}) = \text{trace} \left\{ (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\mathbf{B})^T (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}}\mathbf{B}) \right\}$$

- ▶ The minimizer $\hat{\mathbf{B}}$ for $u(\mathbf{B})$, $\mathbf{B} \in \mathbb{R}^{(p+1) \times m}$ is

$$\begin{aligned} \hat{\mathbf{B}} &= \left[\left(\Sigma^{-1/2}\mathbf{X} \right)^T \Sigma^{-1/2}\mathbf{X} \right]^{-1} \left(\Sigma^{-1/2}\mathbf{X} \right)^T \Sigma^{-1/2}\mathbf{Y} \\ &= \left(\mathbf{X}^T \Sigma^{-1}\mathbf{X} \right)^{-1} \mathbf{X}^T \Sigma^{-1}\mathbf{Y} \end{aligned}$$

Comparisons: LSE and weighted LSE

- ▶ For model $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$, i.e.,

rowwise: $\mathbf{y}_i = \mathbf{x}_i\mathbf{B} + \varepsilon_i$ with $\varepsilon_i \sim G(0, \Sigma)$,

the LSE $\hat{\mathbf{B}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$ is the minimizer of both

$$\begin{cases} g(\mathbf{B}) = \sum_{i=1}^N (\mathbf{y}_i - \mathbf{x}_i\mathbf{B})(\mathbf{y}_i - \mathbf{x}_i\mathbf{B})^T \\ h(\mathbf{B}) = \sum_{i=1}^N \left[(\mathbf{y}_i - \mathbf{x}_i\mathbf{B}) \Sigma^{-1/2} \right] \left[(\mathbf{y}_i - \mathbf{x}_i\mathbf{B}) \Sigma^{-1/2} \right]^T \end{cases}$$

- ▶ For model $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$, i.e.,

columnwise: $\tilde{\mathbf{y}}_j = \mathbf{X}\tilde{\mathbf{b}}_j + \tilde{\varepsilon}_j$ with $\tilde{\varepsilon}_j \sim G(0, \Sigma)$,

the LSE $\hat{\mathbf{B}} = (\mathbf{X}^T\Sigma^{-1}\mathbf{X})^{-1}\mathbf{X}^T\Sigma^{-1}\mathbf{Y}$ is the minimizer of

$$u(\mathbf{B}) = \sum_{j=1}^N \left[\Sigma^{-1/2} (\tilde{\mathbf{y}}_j - \mathbf{X}\tilde{\mathbf{b}}_j) \right]^T \left[\Sigma^{-1/2} (\tilde{\mathbf{y}}_j - \mathbf{X}\tilde{\mathbf{b}}_j) \right]$$

Traditional variable/model selection

Linear model with scalar response

- ▶ Predictor $X = (X_1, \dots, X_p)^T \in \mathbb{R}^p$ and response $Y \in \mathbb{R}$
- ▶ Coefficient vector $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T \in \mathbb{R}^{p+1}$ and population model

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon, \quad (8)$$

where $\varepsilon \in \mathbb{R}$ is the random error term $E(\varepsilon) = 0$ and $\text{var}(\varepsilon) = \sigma^2$

- ▶ Data version: set $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)^T \in \mathbb{R}^N$ and each $\varepsilon_i \sim \varepsilon$, then

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- ▶ Variable/model selection for (8) is the same as retaining X_j whose $\beta_j \neq 0$

Motivations for variable selection

- ▶ “Simple model with good predictive performance” is desired
- ▶ There is need for variable/model selection
 - ▶ when sample size $N \geq p + 1$ to obtain simpler models
 - ▶ when $N < p + 1$ to obtain more stable models
 - ▶ when structural assumptions are reasonable to obtain more efficient models
- ▶ Usually variable/model selection is performed in
 - ▶ “Non-interpolation regime” for, i.e., small and non-zero training error, for “small models” such as generalized additive models
 - ▶ “Interpolation regime”, i.e., zero training error, for “big models” such as neural networks
- ▶ Variable/model selection requires criteria

Selection criteria

- ▶ Always keep β_0 in $Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon$
- ▶ Let d be the number of predictors retained in a model and $\hat{\sigma}^2$ an estimate of σ^2
- ▶ Criterion on this model:

Mallow's C_p :	$C_p = \frac{1}{n} (\text{RSS} + 2d\hat{\sigma}^2)$
Akaike information criterion (AIC):	$\text{AIC} \simeq \frac{1}{n\hat{\sigma}^2} (\text{RSS} + 2d\hat{\sigma}^2)$
Bayesian information criterion (BIC):	$\text{BIC} \simeq \frac{1}{n\hat{\sigma}^2} (\text{RSS} + d\hat{\sigma}^2 \log n)$
Adjusted R-square:	$R_{\text{Adj}}^2 = 1 - \frac{\text{RSS}/(n-d-1)}{\text{TSS}/(n-1)}$
Cross-validation:	no explicit formula

Best-subset selection (BSS)

- ▶ Always keep β_0 in $Y = \beta_0 + \sum_{j=1}^p X_j\beta_j + \varepsilon$
- ▶ There are 2^p submodels to check

Algorithm 6.1 *Best subset selection*

1. Let \mathcal{M}_0 denote the *null model*, which contains no predictors. This model simply predicts the sample mean for each observation.
 2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here *best* is defined as having the smallest RSS, or equivalently largest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Forward-stepwise selection

- ▶ Always keep β_0 in $Y = \beta_0 + \sum_{j=1}^p X_j\beta_j + \varepsilon$
- ▶ There are $1 + \sum_{k=0}^{p-1} (p - k) = 1 + 2^{-1}p(p + 1)$ submodels to check

Algorithm 6.2 *Forward stepwise selection*

1. Let \mathcal{M}_0 denote the *null* model, which contains no predictors.
 2. For $k = 0, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the *best* among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Backward-stepwise selection

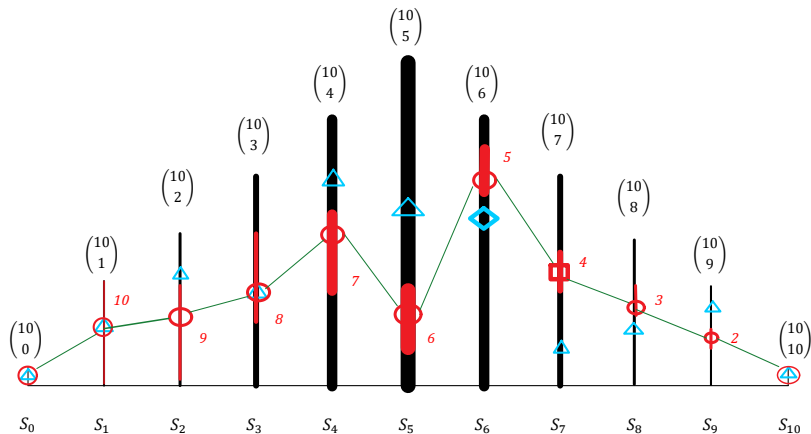
- ▶ Always keep β_0 in $Y = \beta_0 + \sum_{j=1}^p X_j\beta_j + \varepsilon$
- ▶ There are $1 + \sum_{k=0}^{p-1} (p - k) = 1 + 2^{-1}p(p + 1)$ submodels to check

Algorithm 6.3 *Backward stepwise selection*

1. Let \mathcal{M}_p denote the *full* model, which contains all p predictors.
 2. For $k = p, p - 1, \dots, 1$:
 - (a) Consider all k models that contain all but one of the predictors in \mathcal{M}_k , for a total of $k - 1$ predictors.
 - (b) Choose the *best* among these k models, and call it \mathcal{M}_{k-1} . Here *best* is defined as having smallest RSS or highest R^2 .
 3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using cross-validated prediction error, C_p (AIC), BIC, or adjusted R^2 .
-

Trajectory of selected models

- ▶ “Triangle”: trajectory for BSS with “Diamond” as best model
- ▶ “Circle”: trajectory for FSS with “Rectangle” as best model



An illustration

- Always keep β_0 in $Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon$

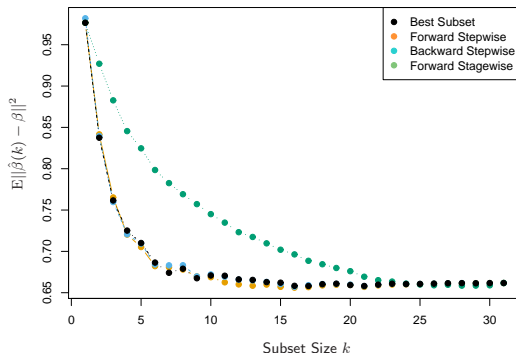


FIGURE 3.6. Comparison of four subset-selection techniques on a simulated linear regression problem $Y = X^T \beta + \varepsilon$. There are $N = 300$ observations on $p = 31$ standard Gaussian variables, with pairwise correlations all equal to 0.85. For 10 of the variables, the coefficients are drawn at random from a $N(0, 0.4)$ distribution; the rest are zero. The noise $\varepsilon \sim N(0, 6.25)$, resulting in a signal-to-noise ratio of 0.64. Results are averaged over 50 simulations. Shown is the mean-squared error of the estimated coefficient $\hat{\beta}(k)$ at each step from the true β .

Recent progresses on BSS

- ▶ Best subset selection (BSS) is NP-hard but Mixed Integer Optimization (MIO) can help
- ▶ “Best subset selection via a modern optimization lens” by Dimitris Bertsimas, Angela King and Rahul Mazumder, *The Annals of Statistics*, 2016
- ▶ “Sparse high-dimensional regression: Exact scalable algorithms and phase transitions” by Dimitris Bertsimas and Bart Van Parys, *The Annals of Statistics*, 2020