# Introduction to Machine Learning  CS182

Lu Sun

School of Information Science and Technology

ShanghaiTech University

October 15, 2023

Today:
- Linear Methods for Regression I
  - Linear regression models
  - The Gauss-Markov theorem
  - Subsets selection

Readings:
- The Elements of Statistical Learning (ESL), Chapters 3
- Pattern Recognition and Machine Learning (PRML), Chapter 3

# Introduction

- A linear regression model assumes that,

（线性回归：假设 f 可以是线性方程）

$$f(x) = E(Y|X = x)$$

Regression function
$$\min_f \text{EPE}(f)$$

  - **linear** in the inputs $X_1, X_2, \ldots, X_p$.
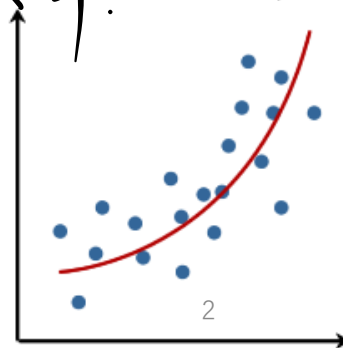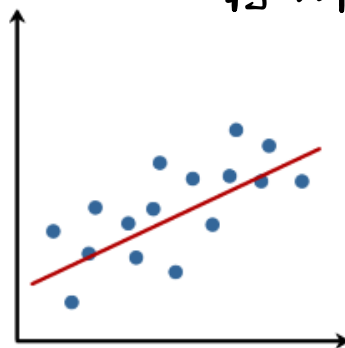
  - $p = 1 \ \rightarrow$ simple linear regression
  - $p > 1 \ \rightarrow$ multiple linear regression

- Suitable for the situations:
  - small number of training samples
  - low signal-to-noise ratio
  - sparse data

- Generalize to many nonlinear techniques.

（可以转换成线性. 将 Y 轴转掉）

Linear                    no linear relationship

2

# Linear Methods for Regression

--- Linear Regression Models

# Simple Linear Regression

- Training set: $(x_1, y_1), \ldots, (x_N, y_N)$  数据点
    - $x_i$: value of predictor $X$ (covariate, independent variable, feature,…)
    - $y_i$: value of response $Y$ (dependent variable, label,…)  标签
- We denote the regression function by

$$f(x) = \mathrm{E}(Y|X = x)$$

    - conditional expectation of $Y$ given $x$

- The linear regression model assumes a specific linear form

$$f(x) = \beta_0 + \beta x$$

    - usually thought of as an approximation to the truth

给预测、值 $\beta_0$ 不涉及 X，故能动天用.

为掉 $\beta$。为正则化打基础。

单时 $\hat{\beta}_0$, 加正则 — 手 于 低 化 后

# **Simple Linear Regression**

单变量时尽量从 $\hat{\beta}_0$ 入手.
因为如果要正则化 则 不可含 $\hat{\beta}_0$.

- Fitting the model by least squares

$$\hat{\beta}_0, \hat{\beta} = \boxed{\text{argmin}_{\beta_0, \beta}} \sum_{i=1}^{N} (y_i - \beta_0 - \beta x_i)^2$$

- Solutions are

求导为0

$$\hat{\beta} = \frac{\sum_{i=1}^{N}(x_i - \bar{x})(y_i - \boxed{\bar{y}})}{\sum_{i=1}^{N}(x_i - \boxed{\bar{x}})^2}$$

**Q**: How to get the solutions?

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}\bar{x}$$

sample mean:

$$\bar{x} = \frac{1}{N}\sum_{i=1}^{N} x_i$$

$$\bar{y} = \frac{1}{N}\sum_{i=1}^{N} y_i$$

相当于把均值移除, 过(0,0)点.

- $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}x_i$ are called the *fitted* or *predicted* values
然后同求导求 $\hat{\beta}$
- $r_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}x_i$ are called the *residuals*

# Multiple Linear Regression

发亢

假设：N个采样 相互独立，随机采样.

- Given $X = (X_1, X_2, \ldots, X_p)^T$
- $E(Y|X)$ is (approximately) linear:

$$f(X) = \beta_0 + \sum_{j=1}^{p} X_j \beta_j$$

- Sources of the variable $X_j$
  - quantitative inputs
  - transformation
  - basis expansions
  - dummy coding
  - interaction
- Linear in the parameters $\beta$

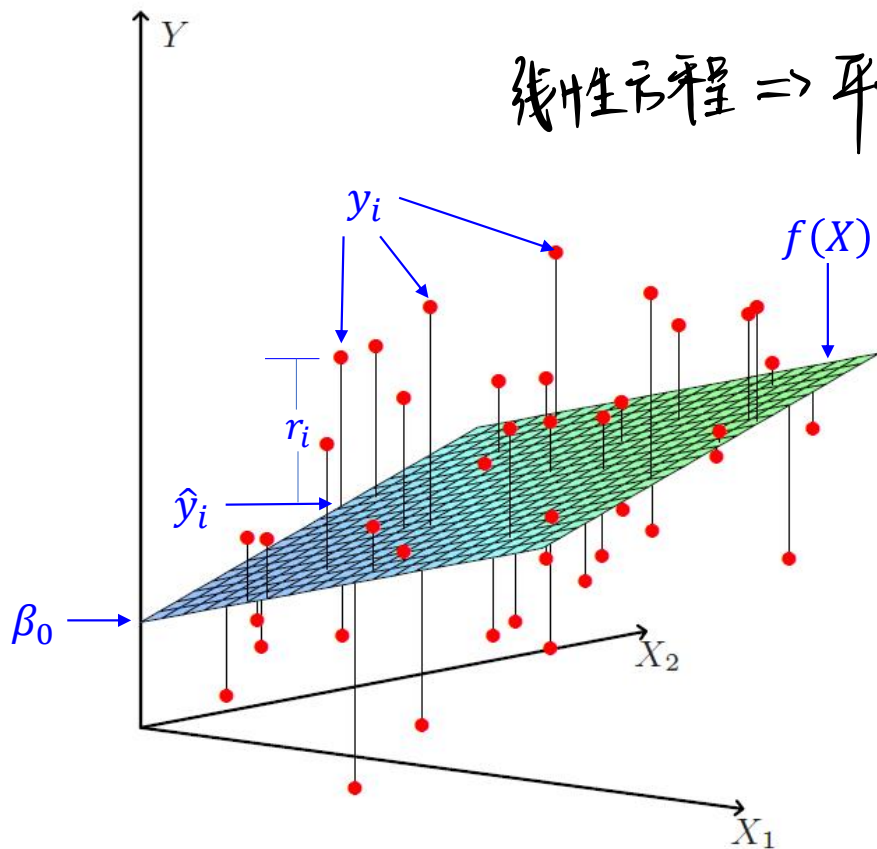- Training data $(x_1, y_1), \ldots, (x_N, y_N)$
- *Least squares*:

$$\text{RSS}(\beta) = \sum_{i=1}^{N}(y_i - f(x_i))^2$$

$$= \sum_{i=1}^{N}(y_i - \beta_0 - \sum_{j=1}^{p} x_{ij} \beta_j)^2$$

- It is reasonable once
  - Observations $(x_i, y_i)$ are randomly sampled from their population
  - Output $y_i$ is conditionally independent w.r.t. the inputs $x_i$
- No guarantee on the validity of model

# Multiple Linear Regression



线性方程 => 平面.

- Training data $(x_1, y_1), ..., (x_N, y_N)$
- *Least squares*:

$$\text{RSS}(\beta) = \sum_{i=1}^{N}(y_i - f(x_i))^2$$

$$= \sum_{i=1}^{N}(y_i - \beta_0 - \sum_{j=1}^{p} x_{ij}\beta_j)^2$$

- It is reasonable once
  - Observations $(x_i, y_i)$ are randomly sampled from their population
  - Output $y_i$ is conditionally independent w.r.t. the inputs $x_i$
- No guarantee on the validity of model

# Multiple Linear Regression

- Minimization of RSS($\beta$)
- Rewrite it by the vector form:

$$\text{RSS}(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$
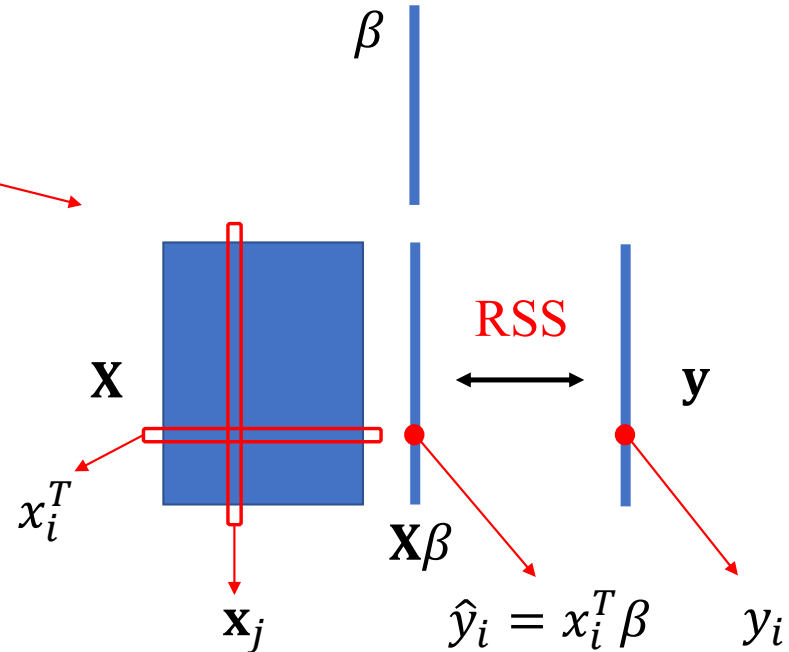
- Differentiating w.r.t. $\beta$

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

- Set the first derivative to zero

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0$$

- If $\mathbf{X}$ has full column rank,

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

$\beta$

$\mathbf{X}$

RSS

$\mathbf{y}$

$x_i^T$

$\mathbf{x}_j$

$\mathbf{X}\beta$

$\hat{y}_i = x_i^T\beta$

$y_i$

# Multiple Linear Regression

从加空间向Y的投影
结果为ŷ.

- Minimization of RSS($\beta$)
- Rewrite it by the vector form:

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

- Differentiating w.r.t. $\beta$

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

- Set the first derivative to zero

$$\boxed{\mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta) = 0}$$

- If $\mathbf{X}$ has full column rank,

$$\hat{\beta} = \underbrace{(\mathbf{X}^T\mathbf{X})}^{-1}\mathbf{X}^T\mathbf{y}$$

需要是满秩的（可逆）

- Prediction on a test sample $x_0$

$$\hat{f}(x_0) = (1\!:\!x_0)^T\hat{\beta}$$

- The fitted values at the training inputs

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{H}\mathbf{y}$$
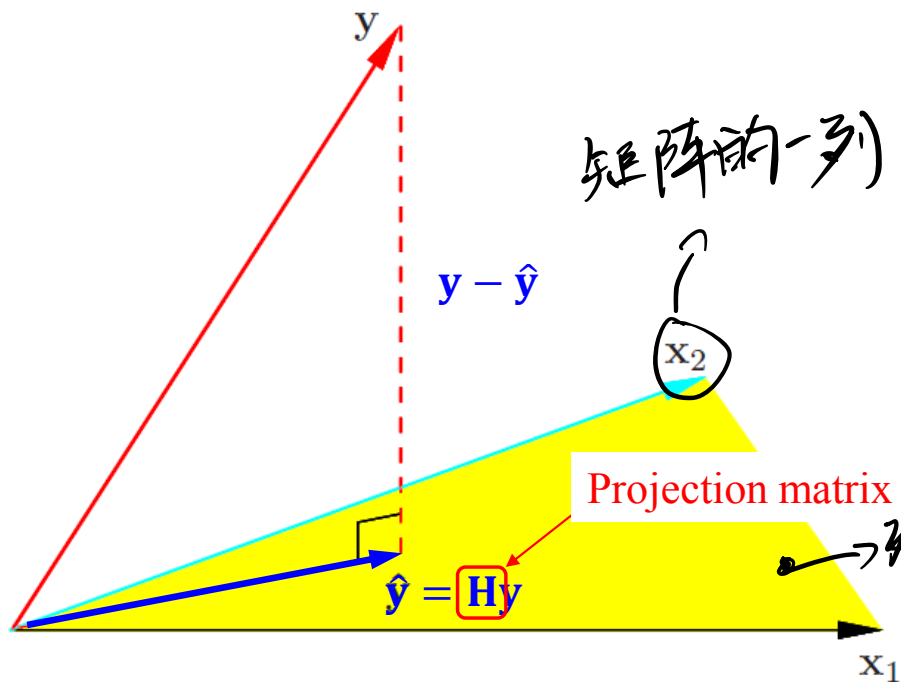
- The "hat" matrix $\mathbf{H}$
  - like a hat put on $\mathbf{y}$ 相当于是个投影.
- Geometrical interpretation
  - The optimal $\hat{\beta}$ makes the residual vector $\mathbf{y} - \hat{\mathbf{y}}$ orthogonal to the subspace spanned by the columns of $\mathbf{X}$

# Multiple Linear Regression

y

**y**

$\mathbf{y} - \hat{\mathbf{y}}$

矩阵的一列)

$\mathbf{x}_2$

Projection matrix

→列空间).

$\hat{\mathbf{y}} = \boxed{\mathbf{H}}\mathbf{y}$

$\mathbf{x}_1$

- Prediction on a test sample $x_0$

$$\hat{f}(x_0) = (1 : x_0)^T \hat{\beta}$$

- The fitted values at the training inputs

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{H}\mathbf{y}$$

- The "hat" matrix $\mathbf{H}$
  - like a hat put on $\mathbf{y}$

- Geometrical interpretation
  - The optimal $\hat{\beta}$ makes the residual vector $\mathbf{y} - \hat{\mathbf{y}}$ orthogonal to the subspace spanned by the columns of $\mathbf{X}$

$\mathbf{X} = (\mathbf{x_1}, \dots, \mathbf{x_p})$, where $\mathbf{x}_j = (x_{1j}, \dots, x_{Nj})^T \in \mathbb{R}^N$

# Multiple Linear Regression

奇异性生

On the singularity of $\mathbf{X}^T\mathbf{X}$

- *Fat* data matrix $\mathbf{X}$   X^TX是p×p的矩阵
  - singular   非奇导需要满秩,即 rank(X^TX)=P
- *Square* data matrix $\mathbf{X}$
  - probably singular
  - nonsingular if rank($\mathbf{X}$) = $p$
- *Skinny* data matrix $\mathbf{X}$
  - probably nonsingular
  - singular if rank($\mathbf{X}$) < $p$

列可能是奇异的

需要有 rank(X)=p

$p$

| | |
|---|---|
| X | **Fat** ($N < p$) |

rank($\mathbf{X}$) ≤ $N$ < $p$

| | |
|---|---|
| X | **Square** ($N = p$) |

rank($\mathbf{X}$) ≤ $N, p$

| | |
|---|---|
| X | **Skinny** ($N > p$) |

rank($\mathbf{X}$) ≤ $p$ < $N$

The solution $\hat{\beta}$ is unique once $\mathbf{X}^T\mathbf{X}$ is nonsingular (rank($\mathbf{X}$) = $p$)

非满秩: 有 余信息.

信号: 维度高, 样本少

解决: 特征选择(降维, 去掉

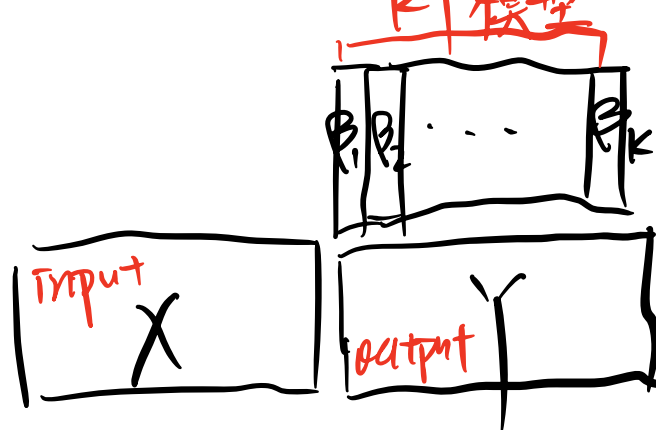正则化 : $\hat{\beta} = (X^TX + \lambda I)^{-1} X^T y$

# Multiple Linear Regression

- Rank deficient **X**
  - coding qualitative inputs
    - ➢ redundancy in columns of X
  - image and signal analysis
    - ➢ more features ($p > N$)
- Two ways to overcome it
  - feature selection (dimension reduction)
  - regularization

$p$

$N$ | X | **Fat** $(N < p)$ | $\text{rank}(\mathbf{X}) \le N < p$

$N$ | X | **Square** $(N = p)$ | $\text{rank}(\mathbf{X}) \le N, p$

$N$ | X | **Skinny** $(N > p)$ | $\text{rank}(\mathbf{X}) \le p < N$

**Multiple Output Regression***

多输出.

- Multiple outputs $Y_1, Y_2, \ldots, Y_K$
- Assume a linear model for each output

$$Y_k = \beta_{0k} + \sum_{j=1}^{p} X_j \beta_{jk} + \varepsilon_k = f_k(X) + \varepsilon_k$$

现在变成矩阵:

- In matrix notation

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$$

有K个标签，K种模型
独立处理K个模型
放入Y中

where $\mathbf{X} \in \mathbb{R}^{N \times (p+1)}$, $\mathbf{B} \in \mathbb{R}^{(p+1) \times K}$ and $\mathbf{E} \in \mathbb{R}^{N \times K}$.

- A generalization of the univariate loss function

$$\text{RSS}(\mathbf{B}) = \sum_{k=1}^{K} \sum_{i=1}^{N} \left( y_{ik} - f_k(x_i) \right)^2 = \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2$$

F-范数（矩阵的F-范数）

对角线之和（迹）

For an arbitrary matrix $\mathbf{A}$, the Frobenius-norm is defined by $\|\mathbf{A}\|_F^2 = \text{Tr}(\mathbf{A}^T \mathbf{A}) = \sum_{ij} a_{ij}^2$.

# Multiple Output Regression*

- Our problem:
$$\widehat{\mathbf{B}} = \text{argmin}_{\mathbf{B}} \text{RSS}(\mathbf{B}) = \text{argmin}_{\mathbf{B}} \|\mathbf{Y} - \mathbf{XB}\|_F^2$$

- A quadratic function with global minimum 是一个点

- Rewrite RSS($\mathbf{B}$) as follows    Matrix trace

$$\text{RSS}(\mathbf{B}) = \boxed{\text{Tr}}\big((\mathbf{Y} - \mathbf{XB})^T(\mathbf{Y} - \mathbf{XB})\big)$$

$\mathbb{R}^{K \times K}$

$$= \text{Tr}(\mathbf{Y}^T\mathbf{Y} - \boxed{\mathbf{Y}^T\mathbf{XB} - \mathbf{B}^T\mathbf{X}^T\mathbf{Y}} + \mathbf{B}^T\mathbf{X}^T\mathbf{XB})$$

$$= \text{Tr}(\mathbf{Y}^T\mathbf{Y}) - 2\text{Tr}(\mathbf{B}^T\mathbf{X}^T\mathbf{Y}) + \text{Tr}(\mathbf{B}^T\mathbf{X}^T\mathbf{XB})$$

- Differentiating w.r.t. $\mathbf{B}$
$$\frac{\partial \text{RSS}(\mathbf{B})}{\partial \mathbf{B}} = -2\mathbf{X}^T\mathbf{Y} + 2\mathbf{X}^T\mathbf{XB}$$

k个X

- If $\mathbf{X}^T\mathbf{X}$ is nonsingular, $\widehat{\mathbf{B}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$ ⟶ $\hat{\beta}_k = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}_k, \forall k$

Multiple outputs do not affect one another's least squares estimates.

# Linear Methods for Regression

--- The Gauss-Markov Theorem

# The Gauss-Markov Theorem

但实际上，现实中无偏估计不多

对线性无偏估计中，最小二乘法有最小方差

The least squares estimator has the lowest sampling variance within the class of linear unbiased estimators.

*Proof*: suppose $\tilde{\beta} = \mathbf{C}\mathbf{y}$ is a linear estimator of $\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$,

where $\mathbf{C} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T + \mathbf{D}$, and $\mathbf{D} \in \mathbb{R}^{p \times N}$ is a non-zero matrix

噪声 ⇒ 噪声的方差

$$\begin{aligned}
\mathrm{E}[\tilde{\beta}] &= \mathrm{E}[Cy] \\
&= \mathrm{E}[((X'X)^{-1}X' + D)(X\beta + \varepsilon)] \\
&= ((X'X)^{-1}X' + D)X\beta + ((X'X)^{-1}X' + D)\,\mathrm{E}[\varepsilon] \\
&= ((X'X)^{-1}X' + D)X\beta \\
&= (X'X)^{-1}X'X\beta + DX\beta \\
&= (I_p + DX)\beta.
\end{aligned}$$

$E[\varepsilon] = 0$

假设无偏

If and only if $\mathbf{DX} = 0$, $\tilde{\beta}$ is unbiased.

$$\begin{aligned}
\mathrm{Var}(\tilde{\beta}) &= \mathrm{Var}(Cy) \\
&= C\,\mathrm{Var}(y)\,C' \\
&= \sigma^2 CC' \\
&= \sigma^2 ((X'X)^{-1}X' + D)(X(X'X)^{-1} + D') \\
&= \sigma^2 ((X'X)^{-1}X'X(X'X)^{-1} + (X'X)^{-1}X'D' + DX(X'X)^{-1} + DD') \\
&= \sigma^2 (X'X)^{-1} + \sigma^2 (X'X)^{-1}(DX)' + \sigma^2 DX(X'X)^{-1} + \sigma^2 DD' \\
&= \sigma^2 (X'X)^{-1} + \sigma^2 DD' \\
&= \mathrm{Var}(\hat{\beta}) + \sigma^2 DD'
\end{aligned}$$

$\mathrm{Var}(\mathbf{y}) = E[\mathbf{y} - E[\mathbf{y}]]^2 = \mathrm{Var}(\varepsilon)$

$\mathbf{DX} = 0$

$\mathrm{Var}(\hat{\beta}) = \sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$

Positive semidefinite

# The Gauss-Markov Theorem

> *The least squares estimator has the lowest sampling variance within the class of linear unbiased estimators.*

*Proof*: suppose $\tilde{\beta} = \mathbf{C}\mathbf{y}$ is a linear estimator of $\hat{\beta} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$,

where $\mathbf{C} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T + \mathbf{D}$, and $\mathbf{D} \in \mathbb{R}^{p \times N}$ is a non-zero matrix

Given an arbitrary test point $x_0$, we have

$$\mathrm{Var}(\tilde{y}_0) = \mathrm{Var}(x_0^T\tilde{\beta})$$
$$= x_0^T\mathrm{Var}(\tilde{\beta})x_0$$
$$= x_0^T\mathrm{Var}(\hat{\beta})x_0 + \sigma^2 x_0^T\mathbf{D}\mathbf{D}^T x_0$$
$$= \mathrm{Var}(\hat{y}_0) + \sigma^2\underline{x_0^T\mathbf{D}\mathbf{D}^T x_0} > 0$$

$$\mathrm{Var}(\tilde{\beta}) = \mathrm{Var}(Cy)$$
$$= C\,\mathrm{Var}(y)C'$$
$$= \sigma^2 CC'$$
$$= \sigma^2\left((X'X)^{-1}X' + D\right)\left(X(X'X)^{-1} + D'\right)$$
$$= \sigma^2\left((X'X)^{-1}X'X(X'X)^{-1} + (X'X)^{-1}X'D' + DX(X'X)^{-1} + DD'\right)$$
$$= \sigma^2(X'X)^{-1} + \sigma^2(X'X)^{-1}(DX)' + \sigma^2 DX(X'X)^{-1} + \sigma^2 DD'$$
$$= \sigma^2(X'X)^{-1} + \sigma^2 DD'$$
$$= \mathrm{Var}\left(\hat{\beta}\right) + \sigma^2 DD'$$

# The Gauss-Markov Theorem

*The least squares estimator has the lowest sampling variance within the class of linear unbiased estimators.*

Remarks

- Among the unbiased linear methods, least squares has the lowest MSE
  - MSE = Var + Bias$^2$
- A biased methods probably has lower MSE
  - Var-Bias trade-off
  - A small increase in Bias might gives rise to a large reduction in Var ← Model selection

# Linear Methods for Regression

--- Subset Selection

# Introduction

局限性:

Two limitations of least squares
- prediction accuracy 准确相对不高（相对线性回归）
  - low bias and high variance
    → sacrifice a little bias to reduce the variance
- interpretation 解释性不足（无法解释哪种特征影响大）
  - hard to interpret a large number of input features
    → find a subset of features exhibiting strong effects

模型选择.

We use model selection to overcome the limitations
- variable subset selection, shrinkage, dimension reduction.

子集选择    正则化    维度降低

- not restricted to linear models

# Subset Selection

- Best-subset selection 复杂度极高

  - For each $s \in \{0,1,...,p\}$, find the subset in size of $s$ that gives lowest $\text{RSS}(\beta) = \left\| \mathbf{y} - \mathbf{x}^{(s)}\beta \right\|_2^2$

P>40时 不可使用
这种方法

$\binom{4}{2} = 6$

只是选了2个特征
还可选 1个/3个/4个 特征去训练.

| $p = 4$ <br> $s = 2$ | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $\mathbf{X}^{(s)}$ |
|---|---|---|---|---|---|
| Model 1 | √ | √ | × | × | $(\mathbf{x}_1, \mathbf{x}_2)$ |
| Model 2 | √ | × | √ | × | $(\mathbf{x}_1, \mathbf{x}_3)$ |
| Model 3 | √ | × | × | √ | $(\mathbf{x}_1, \mathbf{x}_4)$ |
| Model 4 | × | √ | √ | × | $(\mathbf{x}_2, \mathbf{x}_3)$ |
| Model 5 | × | √ | × | √ | $(\mathbf{x}_2, \mathbf{x}_4)$ |
| Model 6 | × | × | √ | √ | $(\mathbf{x}_3, \mathbf{x}_4)$ |

# Subset Selection

- Best-subset selection
  - For each $s \in \{0, 1, \dots, p\}$, find the subset in size of $s$ that gives lowest
    $$\mathrm{RSS}(\beta) = \left\| \mathbf{y} - \mathbf{X}^{(s)}\beta \right\|_2^2$$

- Example
  - prostate cancer example ($p = 8$) $\leftarrow$ 40
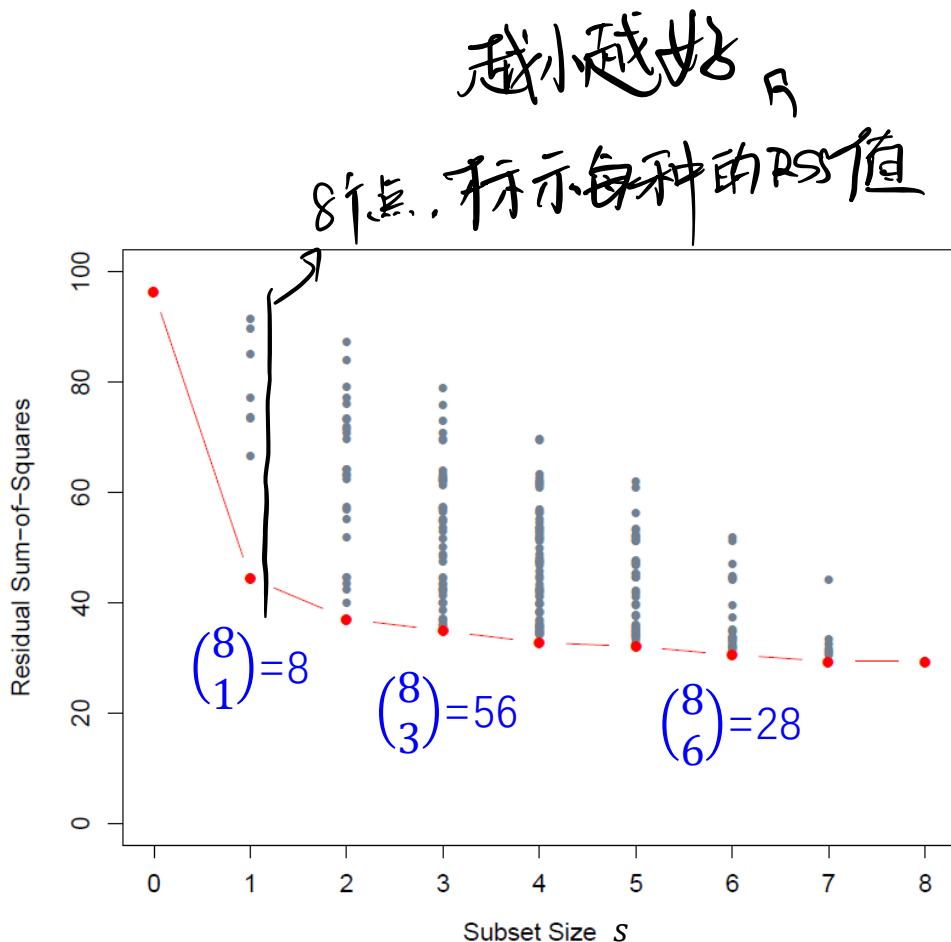  - the red lower bound denotes the models eligible for selection
  - the red lower bound keeps decreasing (s = 8?)
  - *cross-validation* to estimate prediction error and select $s$

- Typically intractable for $p > 40$

地里可接上：

越小越好 R

8个点，表示每种的RSS值



$\binom{8}{1} = 8$   $\binom{8}{3} = 56$   $\binom{8}{6} = 28$

All the subset models for the prostate cancer example.

复杂度：$\beta = (X'X)^{-1}X'y$

$O(N^3)$  $O(p^3)$  $\to (p\times p) \times (p\times1) = N+1$

可整体看.

$\to O(p^3 + Np^2)$

# Forward- and Backward-Stepwise Selection

贪心：局部最优

- Forward-stepwise
  - starts with intercept
  - sequentially adds the best predictor
- Greedy algorithm
  - sub-optimal
- Advantages
  - Computational
    - even $p \gg N$
  - Statistical
    - constrained search 有约束的搜索
    - lower variance, more bias

# Forward- and Backward-Stepwise Selection

- Forward-stepwise
  - starts with intercept
  - sequentially adds the best predictor
- Greedy algorithm
  - sub-optimal
- Advantages
  - Computational
    - even $p \gg N$
  - Statistical
    - constrained search
    - lower variance, more bias

- Backward-stepwise
  - starts with the full model
  - sequentially deletes the worst predictor
- Greedy algorithm
- Only useful when $N > p$
  - linear regression

- Smart stepwise
  - group of variables
  - add or drop whole groups at a time

扔掉

使用限制,

N<P时是奇异的.

# Forward- and Backward-Stepwise Selection



$E\|\hat{\beta}(k) - \beta\|^2$

Best Subset
Forward Stepwise
Backward Stepwise
Forward Stagewise

Performance becomes stable once $p \geq 10$

Subset Size $s$

MSE on $\hat{\beta}$

- Example
  - $Y = X^T\beta + \varepsilon$
  - $N = 300, p = 31$
  - only 10 variables are effective
  - similar performance

以10个做为验证

# K-Fold Cross-Validation

若极少,可以使用 LT (...) (极端情况)

从训练集中抽由一部分做为 validation
用其它部分训练,用这个部分做 validation.
为防止 train 太少使得 N<P (Fat) K-Fold Cross-Validation

- Each has a complexity parameter $\lambda$   成本:训练太多次.(时间成本)
  - the subset size in subset selection
  - the neighborhood size in $k$-NN
  - The coefficient of regularization
- $K$-fold cross validation   分成 K 等份
  - divide the training data into $K$ roughly equal parts ($K = 5$ or $10$)
  - for $k = 1, ..., K$,   用一份做验证.剩下的训练
    - fit the model with $K - 1$ parts
    - compute the error $E_k$ on the rest part
  - The $K$-fold cross validation error

一共 K 份都分开 所以跑 K 次训练

$$E(\lambda) = \frac{1}{K} \sum_{k=1}^{K} E_k(\lambda)$$

Repeat this for many values of $\lambda$, and choose the best value that makes $E(\lambda)$ lowest.



Training set

Training folds    Test fold

1st iteration    $\Rightarrow E_1$

2nd iteration    $\Rightarrow E_2$

3rd iteration    $\Rightarrow E_3$

...

10th iteration    $\Rightarrow E_{10}$

$E = \frac{1}{10} \sum_{i=1}^{10} E_i$

validation error ← $E$

validation set

best $\lambda$

26