

LINEAR MODEL SELECTION

23/03/21

Standard linear model:

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$$

Many generalizations: glm, polynomial regression, spline models (Chapter 7)

Standard models are used very often as they are easily interpretable and have good predictive accuracy.

β is estimated via least-squares

LEAST-SQUARES

Data : $(y_i, x_{i1}, x_{i2}, \dots, x_{ip})$, $i=1, \dots, n$

In matrix form:

$$Y = X \beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$

$$\begin{pmatrix} \sigma^2 & & & \\ & \ddots & & \\ & & \ddots & \\ 0 & & & \ddots & \sigma^2 \end{pmatrix}$$

LS $\min_{\beta} \sum_{i=1}^n (y_i - x_i^t \beta)^2$ where $x_i = \begin{pmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix}$

$$= \min_{\beta} \underbrace{(Y - X \beta)^t}_{1 \times n} \underbrace{(Y - X \beta)}_{n \times 1}$$

$(a+b)(c+d)$

$$\begin{aligned}
 \min_{\beta} & (\underline{y} - \underline{x}\beta)^t (\underline{y} - \underline{x}\beta) \\
 & = \underline{y}^t \underline{y} - \underline{y}^t \underline{x} \beta - (\underline{x} \beta)^t \underline{y} + (\underline{x} \beta)^t (\underline{x} \beta) \\
 & = \underline{y}^t \underline{y} - \underline{y}^t \underline{x} \beta - \underline{\beta}^t \underline{x}^t \underline{y} + \underline{\beta}^t \underline{x}^t \underline{x} \beta \\
 & = \underline{y}^t \underline{y} - 2 \underline{\beta}^t \underline{x}^t \underline{y} + \underline{\beta}^t \underline{x}^t \underline{x} \beta
 \end{aligned}$$

$$(AB)^t = B^t A^t$$

$$\frac{d}{d\beta} (\underline{y} - \underline{x}\beta)^t (\underline{y} - \underline{x}\beta) = -2\underline{x}^t \underline{y} + 2\underline{x}^t \underline{x} \beta = 0$$

$$\begin{aligned}
 (\underline{x}^t \underline{x}) \beta &= \underline{x}^t \underline{y} \\
 \boxed{\hat{\beta}} &= (\underline{x}^t \underline{x})^{-1} \underline{x}^t \underline{y}
 \end{aligned}$$

p+1 equations
in p+1 parameters
(NORMAL EQUATIONS)

$\hat{\beta}$ PROPERTIES

$$\begin{aligned}
 ① E[\hat{\beta}] &= E[(\underline{x}^t \underline{x})^{-1} \underline{x}^t \underline{y}] = (\underline{x}^t \underline{x})^{-1} \underline{x}^t E[\underline{y}] \\
 &= \underbrace{(\underline{x}^t \underline{x})^{-1} \underline{x}^t}_{I} \underline{x} \beta = \beta \quad \text{UNBIASED}
 \end{aligned}$$

$$\begin{aligned}
 ② \text{Var}(\hat{\beta}) &= \text{Var}((\underline{x}^t \underline{x})^{-1} \underline{x}^t \underline{y}) = (\underline{x}^t \underline{x})^{-1} \underline{x}^t \text{Var}(\underline{y}) \underline{x} (\underline{x}^t \underline{x})^{-1} \\
 &\text{covariance matrix} \\
 &\quad \text{Var}(A\underline{y}) \\
 &\quad = A \text{Var}(\underline{y}) A^t
 \end{aligned}$$

$$\begin{aligned}
 &= (\underline{x}^t \underline{x})^{-1} \underline{x}^t \sigma^2 I \underline{x} (\underline{x}^t \underline{x})^{-1} \\
 &= \sigma^2 \underbrace{(\underline{x}^t \underline{x})^{-1} \underline{x}^t \underline{x}}_I (\underline{x}^t \underline{x})^{-1} = \sigma^2 (\underline{x}^t \underline{x})^{-1}
 \end{aligned}$$

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 (\underline{x}^t \underline{x})_{jj})$$

$$\left(\begin{array}{c} \text{var}(\hat{\beta}_0) \ \text{cov}(\hat{\beta}_0, \hat{\beta}_1) \dots \text{cov}(\hat{\beta}_0, \hat{\beta}_p) \\ \vdots \ \vdots \ \ddots \ \vdots \\ \text{cov}(\hat{\beta}_p, \hat{\beta}_0) \ \text{var}(\hat{\beta}_p) \end{array} \right)$$

BIAS-VARIANCE TRADE-OFF FOR STANDARD LINEAR MODELS

$$\begin{aligned} E[(y_0 - \hat{f}(x_0))^2] &= E_{x_0, \epsilon}[(y_0 - \underline{x}_0^t \hat{\beta})^2] \\ &= [f(x_0) - E[\hat{f}(x_0)]]^2 + \underbrace{\text{Var}(\epsilon)}_{\sigma^2 \text{ irreducible error}} + \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{VARIANCE}} \end{aligned}$$

BIAS $x_0^t \beta - E(x_0^t \hat{\beta}) = x_0^t \beta - x_0^t E(\hat{\beta}) = x_0^t \beta - x_0^t \beta = 0$

ZERO
BIAS

VARIANCE $\begin{aligned} \text{Var}(x_0^t \hat{\beta}) &= x_0^t \text{Var}(\hat{\beta}) x_0 = \sigma^2 x_0^t (x^t x)^{-1} x_0 \\ &= \sigma^2 \text{Trace}(x_0^t (x^t x)^{-1} x_0) \\ &= \sigma^2 \text{Trace}((x^t x)^{-1} \underbrace{x_0 x_0^t}_{p \times p}) \end{aligned}$

Trace of a Matrix

$$A = \begin{pmatrix} & & & \\ & & & \\ & & & \\ & & & \end{pmatrix}$$

sum of diagonal elements

$$\begin{aligned} \text{Tr}(ABC) &= \\ &= \text{Tr}(BCA) \end{aligned}$$

$$\begin{aligned} \text{cov}(x, y) &= E((x - E(x))(y - E(y))) \\ E(x) &= E(y) = 0 \\ \text{cov}(x, y) &= \frac{\sum_{i=1}^n x_i y_i}{n} \end{aligned}$$

$$\begin{pmatrix} x_{01}^2 & x_{01}x_{02} & \dots & x_{01}x_{0p} \end{pmatrix}$$

sample covariance

$$= \sigma^2 \text{Trace} \left((\mathbf{x}^t \mathbf{x})^t \mathbf{x}_0 \mathbf{x}_0^t \right)$$

$$= \sigma^2 \text{Trace} \left(\underbrace{\left(\frac{\mathbf{x}^t \mathbf{x}}{n} \cdot n \right)^{-1}}_{\mathbb{I}} \underbrace{\mathbf{x}_0 \mathbf{x}_0^t}_{\mathbb{I}} \right)$$

$$= \sigma^2 \text{Trace} \left(\frac{\mathbb{I}}{n} \right)$$



$$= \sigma^2 \frac{p}{n}$$

$$E \left[(y_0 - \mathbf{x}_0^t \hat{\beta})^2 \right] = \sigma^2 + \sigma^2 \frac{p}{n}$$

Assuming that predictors are standardized (x_0 mean and unit variance) and uncorrelated, then $\mathbf{x}_0 \mathbf{x}_0^t = \mathbb{I}$

Similarly $\frac{\mathbf{x}^t \mathbf{x}}{n}$ is the sample covariance of (x_1, \dots, x_p) from data, so $\frac{\mathbf{x}^t \mathbf{x}}{n} \approx \mathbb{I}_p$

VARIANCE OF STANDARD LINEAR MODEL

It depends on $\left(\frac{p}{n}\right)$ We can distinguish 3 cases:

① $n \gg p$

→ bias = 0, variance very small, so LS will do very well in these cases

② n not much larger than p

→ variability in LS fit → overfitting, poor predictions

③ $n < p$ $X^t X$ is not invertible

$$n \approx 100$$

$$p \approx 3000$$

→ no unique LS estimate, infinite variance

There are 3 approaches to improve predictive accuracy of a standard linear model:

① Limit p ~ VARIABLE SELECTION

→ Reduces variance

→ Increases interpretability by identifying the most important predictors and removing redundant ones

② Shrinkage methods : Replace LS by alternative methods so that bias is increased (a bit) and variance is reduced
→ next lecture

③ Dimensionality Reduction Project the p predictors into a M -dimensional space with $M < p$ and use the projections in a LS model
→ end of course

SUBSET SELECTION

Identify a subset of important predictors out of X_1, \dots, X_p

Two procedures: best subset and stepwise selection

BEST SUBSET

Consider all possible combinations of the p predictors, fit a model by LS to each of those and choose best one according to a model selection criterion.

Algorithm ① Let \mathcal{M}_0 denote the null model (only intercept)
→ each observation is predicted to \bar{y}

② For $K=1, \dots, p$:

(a) Fit all $\binom{p}{K}$ models with exactly K predictors

(b) Pick the \uparrow best among those and call this \mathcal{M}_K
via RSS, R^2

③ Select a single best model out of $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$ using some model selection criterion: AIC, BIC, C_p , adj R^2 , CV

HOW MANY MODELS?

p predictors $\rightarrow 2^p$ possible models

$p=2 \rightarrow \mathcal{M}_0$ only intercept (null model)

X_1

X_2

$\mathcal{M}_2 = X_1, X_2$ (full model)

4 models

P

2

3

8

16

20

100

2^P ↓ too many

4

8

256

65536

1048576

1.26765×10^{30}

$1 + \frac{p(p+1)}{2}$

37

211

5051

Much less!

STEPWISE METHODS

3 versions : FORWARD, BACKWARD, STEPWISE (hybrid)

FORWARD

- ① Let \mathcal{M}_0 be the null model
- ② For $K = 0, \dots, p-1$:
 - (a) Fit all $p-K$ models that augment \mathcal{M}_K with one predictor
 - (b) choose best among the $p-K$ models and call \mathcal{M}_{K+1}
- ③ Select best model among $\mathcal{M}_0, \mathcal{M}_1, \dots, \mathcal{M}_p$ according to a model selection criterion

How many models now?

Null model $\rightarrow 1$

$K=0 \rightarrow p$ models

$K=1 \rightarrow p-1$ models

:

$K=p-1 \rightarrow 1$ model

$$\begin{aligned} \text{In total : } 1 + \sum_{K=0}^{p-1} (p-K) &= 1 + \sum_{K=1}^p (p-K+1) = 1 + p^2 + p - \sum_{K=1}^p K \\ &= 1 + p^2 + p - \frac{p(p+1)}{2} = 1 + \frac{p(p+1)}{2} \end{aligned}$$

BACKWARD

- Starting from the full model (if you can!)
- Remove the least important predictor each time.

The main issue with forward or backward procedures is that once a predictor is added/removed from the model, then it remains so for all the subsequent steps. The problem is particularly pronounced for correlated predictors.

J₁: X₁
J₂: X₂, X₃ Hybrid approaches can be used here : STEPWISE

At each step of a forward procedure, one variable is added and any of the existing variables can be removed if it does not provide an improvement to the model fit.