

Modern Statistical Methods

Rajen Shih

Chapters

1. Kernel machines
2. The Lasso and extensions
3. Graphical modelling and causal inference
4. Multiple testing and high-dimensional inference

Books

"Elements of statistical learning"
"Statistics for high-dimensional data"
"Statistical learning with sparsity"

Webpage

- Notes
- Feedback form

Introduction

Ordinary Least Squares (OLS)

Data $(Y_i, x_i) \in \mathbb{R} \times \mathbb{R}^p$, $i = 1, \dots, n$
 \uparrow \nwarrow
 response predictors

Linear model assumes $Y = X\beta^0 + \varepsilon$ where $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix}$ $E\varepsilon = 0$
 $\text{Var}(\varepsilon) = \sigma^2 I$
Can estimate β^0 by OLS provided X has full col rank (so $X^T X$ invertible).

$$\hat{\beta}^{OLS} = \arg \min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2 = (X^T X)^{-1} X^T Y$$

$$E \hat{\beta}^{OLS} = E (X^T X)^{-1} X^T (X \beta^0 + \epsilon) = \beta^0 \text{ unbiased}$$

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

Gauss-Markov theorem: $\hat{\beta}^{OLS}$ is the best linear unbiased estimator of β^0

Any other unbiased estimator $\tilde{\beta} = AY$

$$\text{Var}(\tilde{\beta}) = \sigma^2 (X^T X)^{-1} \text{ positive semi-definite}$$

Maximum Likelihood

Specify a density for Y , $f(y; \theta)$, θ is an unknown vector of parameters in \mathbb{R}^d

log-likelihood $l(\theta) = \log(f(Y; \theta))$

Maximum likelihood estimation maximizes $l(\theta)$

- Can maximize $\hat{\theta}$.

$$\text{Fisher information matrix } I(\theta) = -E \left(\frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\theta) \right)$$

Cramér-Rao lower bound if $\tilde{\theta}$ is an unbiased estimator, then

$$\text{Var}(\tilde{\theta}) - I^{-1}(\theta) \text{ is +ve semi-definite}$$

MLEs - asymptotically ($n \rightarrow \infty$) unbiased, achieves the C-R bound, normally distributed

Chapter 1 Kernel machines

Consider linear model $Y = X\beta^0 + \epsilon$ $E\epsilon = 0$ $\text{Var}(\epsilon) = \sigma^2 I$

For obtained estimator $\tilde{\beta}$ should stuff the mean-squared error matrix

$$\begin{aligned} E(\tilde{\beta} - \beta^0)(\tilde{\beta} - \beta^0)^T &= E\{(\tilde{\beta} - E\tilde{\beta} + E\tilde{\beta} - \beta^0)(\tilde{\beta} - E\tilde{\beta} + E\tilde{\beta} - \beta^0)^T\} \\ &= \text{Var}(\tilde{\beta}) + \underbrace{E(\tilde{\beta} - \beta^0) \{E(\tilde{\beta} - \beta^0)\}^T}_{\text{squared bias}} \end{aligned}$$

1.1. Ridge regression

Idea: Reduce variance of OLS by shrinking estimated coefficients towards 0

Ridge regression solves

$$(\hat{\mu}_\lambda^R, \hat{\beta}_\lambda^R) = \arg \min_{(\mu, \beta) \in \mathbb{R} \times \mathbb{R}^p} \{ \|Y - \mu \mathbf{1} - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \}$$

where $\mathbf{1}$ is an n -vector of 1's

$\lambda \geq 0$ is a tuning parameter

Intercept term is not penalised, so that if response is $Y' = Y + c\mathbf{1}$ then

$$\hat{\mu}_\lambda^R(Y') = \hat{\mu}_\lambda^R(Y) + c$$

It is also standard practice to centre each col of X and scale them to have equal l_2 -norms.

Then have $\hat{\mu}_\lambda^R = \bar{Y} = \frac{1}{n} \sum Y_i$ so by replacing Y with $Y - \bar{Y}\mathbf{1}$ can omit intercept from the objective.

$$\text{Hence } \hat{\beta}_\lambda^R = (X^T X + \lambda I)^{-1} X^T Y$$

Theorem 1 Suppose $\text{rank}(X) = p$. For λ suff. small (depending on β^0 and σ^2)

$$E(\hat{\beta}^{\text{OLS}} - \beta^0)(\hat{\beta}^{\text{OLS}} - \beta^0)^T - E(\hat{\beta}_\lambda^R - \beta^0)(\hat{\beta}_\lambda^R - \beta^0)^T \text{ is positive definite.}$$

- EOSL p 61-68
241-249 cross-validation