

# The Elements of Statistical Learning

## Chapter 2.

predictors

responses

(categorical) factors

ordered categorical var: chapter 4

components (predictor  $j$ )  $X_j$

$i$ th observation  $X_i$

prediction of the output  $Y \in \mathbb{R}$ :  $\hat{Y} \in \mathbb{R}$

### 2.3 Two simple approaches to Prediction:

#### 2.3.1 Least Squares

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^p X_j \hat{\beta}_j$$

$\uparrow$  bias in machine learning

(an inner product)  $\hat{Y} = X^T \hat{\beta}$

$$X^T = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}_{n \times p}$$

where  $X_i$  is  $1 \times p$

Note:  $X$  is a column vector ( $p \times n$ )? column vector:  $p \times 1$  ( $p$  columns 1 row)

In general,  $Y$  can be a  $K$ -vector, the  $\beta$  would be  $p \times K$  matrix of coefficients.

Solution: minimize the residual sum of squares:

$$\min_{\beta} \text{RSS}(\beta) = \sum_{i=1}^N (y_i - X_i^T \beta)^2$$

quadratic function: minimum always exists but may not be unique.

$$\text{RSS}(\beta) = (y - X\beta)^T (y - X\beta) \quad X (N \times p)$$

$$\Rightarrow \text{Differentiating } \text{RSS}(\beta) \text{ w.r.t } \beta: X^T (y - X\beta) = 0 \quad X^T y = X^T X \beta \Rightarrow \underline{\beta = (X^T X)^{-1} X^T y} \quad \star$$

Fitted value of  $i$ th input:  $\hat{y}_i = X_i^T \beta$

## 2.3.2 Nearest-Neighbor Methods

fitted  $\hat{f}$ :  $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$  average or sometimes: majority vote (categorical response)

## 2.3.3 comparison

① Least squares:  $\rightarrow$  Gaussian distr. assumption ~~\*~~  
Smoothness & stability 平滑且稳定 low-variance & high-bias

② KNN: not rely on any stringent assumptions, can adapt to any situation high-var & low-bias

## 2.4 Statistical Decision Theory

EPE: expected (squared) prediction error

① Quantitative response ( $Y$ )

$$EPE(f) = E(Y - f(X))^2 = \int [y - f(x)]^2 P_r(dx, dy) = \int_x \int_y (y - f(x))^2 P_r(y|x) P_r(x)$$

Conditioning on input  $X$ : (joint density  $P_r(X, Y) = P_r(Y|X) \cdot P_r(X)$ )

Therefore  $EPE(f) = E_x E_{Y|X}([Y - f(x)]^2 | X)$   $\leftarrow$

To minimize EPE:  $f(x) = \arg\min_c E_{Y|X}([Y - c]^2 | X = x)$

solution:  $f(x) = E(Y|X=x)$  the conditional expectation

KNN:  $\hat{f}(x) = \text{Ave}(y_i | x_i \in N_k(x))$

$f(x) \approx X^T \beta$  plug into  $EPE(f) = E(Y - X^T \beta)^2$

differentiate w.r.t  $\beta \Rightarrow \beta = [E(X^T X)]^{-1} E(XY)$

② Qualitative response ( $G$ )  $\rightarrow k$  levels

$$EPE = E[L(G, \hat{G}(X))] \quad P_r(G, X)$$

$$EPE = E_x \sum_{k=1}^K L[G_k, \hat{G}(X)] P_r(G_k | X)$$

Loss function:  $K \times K$  matrix  $L$  where  $L(k, l)$  is the price paid for classifying an observation belonging to class  $G_k$  to class  $G_l$ .

To minimize EPE:

$$\hat{G}(x) = \arg\min_{g \in G} \sum_{k=1}^K L(G_k, g) P_r(G_k | X = x)$$

With 0-1 loss function:

When  $g = G_k$

$$\hat{G}(x) = \underset{g \in G}{\operatorname{argmin}} (1 - \Pr(g|X=x))$$

or simply  $\hat{G}(x) = G_k$  if  $\Pr(G_k|X=x) = \max_{g \in G} \Pr(g|X=x)$  Bayes classifier

Example: if  $G$  only has 2 levels and be denoted as  $Y=0$  or  $1$ .

Then  $\hat{f}(x) = E(Y|x) = \Pr(G = G_1|x)$  if  $G_1$  corresponds to  $Y=1$ .

## 2.5 Local Methods in High Dimensions

### \* Curse of dimensionality

Example:  $Y = f(x) = e^{-8\|x\|^2}$   $\rightarrow$  true relationship  
predict  $y_0$  at the test point  $x_0 = 0$ .

Training set  $T$ .

MSE for estimating  $f(0)$ :  
 $\downarrow$  true value  
 $\nearrow$  estimator

$$\begin{aligned} \text{MSE}(x_0) &= E_T [f(x_0) - \hat{y}_0]^2 = E_T [f(x_0) - E_T(\hat{y}_0) + E_T(\hat{y}_0) - \hat{y}_0]^2 \\ &= E_T [\hat{y}_0 - E_T(\hat{y}_0)]^2 + (E_T[\hat{y}_0] - f(x_0))^2 \\ &= \text{Var}_T(\hat{y}_0) + \text{Bias}^2 \hat{y}_0 \end{aligned}$$

$$\begin{aligned} (2.27) \quad \text{EPE}(x_0) &= E_{y_0|x_0} E_T (y_0 - \hat{y}_0)^2 \\ &= E_{y_0|x_0} E_T (x_0^T \beta) \end{aligned}$$



2.6

$T = \{x_i, y_i\} \quad i=1, 2, \dots, N$        $T$  : training set

① linear basis expansions:  $f_{\theta}(x) = \sum_{k=1}^K h_k(x) \theta_k$

### 3.4 Shrinkage Methods

#### 3.4.1 Ridge

$$\hat{\beta}^{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \quad \lambda > 0$$

exclude  $\beta_0$  from the penalty term

normally standardize the inputs before solving the above formula

In matrix form:  $RSS(\lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$

Then  $\hat{\beta}^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$

$$-2X^T(y - X\beta) + 2\lambda\beta = 0$$

$$-X^T y + X^T X \beta + \lambda \beta = 0$$

$$(X^T X + \lambda I) \beta = X^T y$$

For orthonormal inputs,  $\hat{\beta}^{\text{ridge}} = \frac{\hat{\beta}}{1 + \lambda}$

$$df(\lambda) = \operatorname{tr}[X(X^T X + \lambda I)^{-1} X^T] = \operatorname{tr}(H\lambda) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + \lambda}$$

Exercise 3.6 prior dist.  $\beta_j \stackrel{\text{i.i.d.}}{\sim} N(0, \tau^2)$

$$f(\beta|y) = \frac{f(\beta, y)}{f(y)} = \frac{f(y|\beta)f(\beta)}{\int_{\beta} f(y|\beta)f(\beta) d\beta} \propto f(y|\beta)f(\beta)$$

Constant

$f(\beta|y)$  is a function of  $f(y|\beta)f(\beta)$  or like follow the same (scale different) distribution

$$\int_{\beta|y} f(\beta|y) d\beta = 1$$

Thin SVD of  $X$  matrix.

$$X = UDV^T$$

$\downarrow \quad \downarrow \quad \downarrow$   
 $N \times p \quad p \times p \quad p \times p$

only calculated the  $p$  column vectors which correspond to  $p$  row vectors of  $V^T$

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

$$U D V^T (V D^T U^T U D V^T)^{-1} V D^T U^T y$$

$$U D V^T (V^T)^{-1} D^{-1} (D^T)^{-1} V^{-1} V D^T U^T y$$

1      1      1

Then ridge  $\hat{\beta}^{\text{ridge}} = X(X^T X + \lambda I)^{-1} X^T y$

$$\begin{aligned}
 &= X(V D^2 V^T + \lambda V V^T)^{-1} V D U^T y \\
 &= X(V(D^2 + \lambda I) V^T)^{-1} V D U^T y \\
 &= \cancel{X} (V^T)^{-1} (D^2 + \lambda I)^{-1} \underbrace{V^{-1} V}_{I} D U^T y \\
 &= X(V^T)^{-1} (D^2 + \lambda I)^{-1} D U^T y \\
 &= U D V^T (V^T)^{-1} (D^2 + \lambda I)^{-1} D U^T y \\
 &= U D (D^2 + \lambda I)^{-1} D U^T y \\
 &= \sum_{j=1}^p \underbrace{\mu_j}_{\substack{\downarrow \\ \text{column of } U}} \frac{d_j^2}{d_j^2 + \lambda} \mu_j^T y
 \end{aligned}$$

sample covariance matrix:  $S = X^T X / N$

Where  $X^T X = V D^2 V^T \rightarrow$  eigen decomposition

first principle component direction  $v_1$ :  $z_1 = X v_1$  has the largest <sup>sample</sup> variance.

$$\begin{aligned}
 \text{Var}(z_1) &= \text{Var}(X v_1) \\
 &= v_1^T \text{Var}(X) v_1 \\
 &= v_1^T \frac{X^T X}{N} v_1 \\
 &= v_1^T \frac{V D^2 V^T}{N} v_1 = \frac{d_1^2}{N}
 \end{aligned}$$

Where  $v_1^T V = [1 \ 0 \ 0 \ \dots]$

$$\text{Var}(z_1) = \frac{[1 \ 0 \ 0 \ \dots] \begin{bmatrix} d_1^2 & & 0 \\ & d_2^2 & \\ 0 & & \ddots \\ & & & d_p^2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}}{N} = \frac{d_1^2}{N}$$

### 3.4.2 Lasso

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^P x_{ij} \beta_j)^2 \quad \text{subject to } \sum_{j=1}^P |\beta_j| \leq t$$

equals

$$\hat{\beta}^{\text{lasso}} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^P x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^P |\beta_j| \right\}$$

↓  
nonlinear in the  $y_i$  -  
no closed form expression as in ridge

eg:  $t_0 = \sum_{j=1}^P |\beta_j|$  if  $t = \frac{t_0}{2}$ , LSE coefficients are shrunk by about 50% on average.

## Appendix. - Terminology

### Chapter 2.

hyperplane (超平面) = a subspace whose dimension is one less than that of its ambient space. <sup>环绕空间</sup>

An affine set (仿射集) affine space (仿射空间)  
codimension:

Bivariate Gaussian (/normal) distribution