

Bayes risk

$$r(f, \hat{\theta}) = \iint L(\theta, \hat{\theta}(x)) f(x, \theta) dx d\theta = \mathbb{E}_{\theta, X} [L(\theta, \hat{\theta}(X))]$$

$$r(f, \hat{\theta}) = \mathbb{E}_{\theta} \left[\mathbb{E}_{X|\theta} [L(\theta, \hat{\theta}(X))] \right] = \mathbb{E}_{\theta} [R(\theta, \hat{\theta})]$$

$$r(f, \hat{\theta}) = \mathbb{E}_X \left[\mathbb{E}_{\theta|X} [L(\theta, \hat{\theta}(X))] \right] = \mathbb{E}_X [r(\hat{\theta} | X)]$$

17.2 Admissibility

- $\hat{\theta}'$ dominates $\hat{\theta}$ if

$$\forall \theta : R(\theta, \hat{\theta}') \leq R(\theta, \hat{\theta})$$

$$\exists \theta : R(\theta, \hat{\theta}') < R(\theta, \hat{\theta})$$

- $\hat{\theta}$ is inadmissible if there is at least one other estimator $\hat{\theta}'$ that dominates it. Otherwise it is called admissible.

17.3 Bayes Rule

Bayes rule (or Bayes estimator)

- $r(f, \hat{\theta}) = \inf_{\tilde{\theta}} r(f, \tilde{\theta})$
- $\hat{\theta}(x) = \inf_{\theta} r(\theta | x) \forall x \implies r(f, \hat{\theta}) = \int r(\hat{\theta} | x) f(x) dx$

Theorems

- Squared error loss: posterior mean
- Absolute error loss: posterior median
- Zero-one loss: posterior mode

17.4 Minimax Rules

Maximum risk

$$\bar{R}(\hat{\theta}) = \sup_{\theta} R(\theta, \hat{\theta}) \quad \bar{R}(a) = \sup_{\theta} R(\theta, a)$$

Minimax rule

$$\sup_{\theta} R(\theta, \hat{\theta}) = \inf_{\tilde{\theta}} \bar{R}(\tilde{\theta}) = \inf_{\tilde{\theta}} \sup_{\theta} R(\theta, \tilde{\theta})$$

$$\hat{\theta} = \text{Bayes rule} \wedge \exists c : R(\theta, \hat{\theta}) = c$$

Least favorable prior

$$\hat{\theta}^f = \text{Bayes rule} \wedge R(\theta, \hat{\theta}^f) \leq r(f, \hat{\theta}^f) \forall \theta$$

18 Linear Regression

Definitions

- Response variable Y
- Covariate X (aka predictor variable or feature)

18.1 Simple Linear Regression

Model

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad \mathbb{E}[\epsilon_i | X_i] = 0, \quad \mathbb{V}[\epsilon_i | X_i] = \sigma^2$$

Fitted line

$$\hat{r}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$$

Predicted (fitted) values

$$\hat{Y}_i = \hat{r}(X_i)$$

Residuals

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)$$

Residual sums of squares (RSS)

$$\text{RSS}(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n \hat{\epsilon}_i^2$$

Least square estimates

$$\hat{\beta}^T = (\hat{\beta}_0, \hat{\beta}_1)^T : \min_{\hat{\beta}_0, \hat{\beta}_1} \text{RSS}$$

$$\hat{\beta}_0 = \bar{Y}_n - \hat{\beta}_1 \bar{X}_n$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n)}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} = \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2}$$

$$\mathbb{E}[\hat{\beta} | X^n] = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$$

$$\mathbb{V}[\hat{\beta} | X^n] = \frac{\sigma^2}{n s_X^2} \begin{pmatrix} n^{-1} \sum_{i=1}^n X_i^2 & -\bar{X}_n \\ -\bar{X}_n & 1 \end{pmatrix}$$

$$\hat{\text{se}}(\hat{\beta}_0) = \frac{\hat{\sigma}}{s_X \sqrt{n}} \sqrt{\frac{\sum_{i=1}^n X_i^2}{n}}$$

$$\hat{\text{se}}(\hat{\beta}_1) = \frac{\hat{\sigma}}{s_X \sqrt{n}}$$

where $s_X^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ and $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2$ (unbiased estimate).
Further properties:

- Consistency: $\hat{\beta}_0 \xrightarrow{P} \beta_0$ and $\hat{\beta}_1 \xrightarrow{P} \beta_1$

- Asymptotic normality:

$$\frac{\widehat{\beta}_0 - \beta_0}{\widehat{\text{se}}(\widehat{\beta}_0)} \xrightarrow{D} \mathcal{N}(0, 1) \quad \text{and} \quad \frac{\widehat{\beta}_1 - \beta_1}{\widehat{\text{se}}(\widehat{\beta}_1)} \xrightarrow{D} \mathcal{N}(0, 1)$$

- Approximate $1 - \alpha$ confidence intervals for β_0 and β_1 :

$$\widehat{\beta}_0 \pm z_{\alpha/2} \widehat{\text{se}}(\widehat{\beta}_0) \quad \text{and} \quad \widehat{\beta}_1 \pm z_{\alpha/2} \widehat{\text{se}}(\widehat{\beta}_1)$$

- Wald test for $H_0 : \beta_1 = 0$ vs. $H_1 : \beta_1 \neq 0$: reject H_0 if $|W| > z_{\alpha/2}$ where $W = \widehat{\beta}_1 / \widehat{\text{se}}(\widehat{\beta}_1)$.

R^2

$$R^2 = \frac{\sum_{i=1}^n (\widehat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n \widehat{\epsilon}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

Likelihood

$$\mathcal{L} = \prod_{i=1}^n f(X_i, Y_i) = \prod_{i=1}^n f_X(X_i) \times \prod_{i=1}^n f_{Y|X}(Y_i | X_i) = \mathcal{L}_1 \times \mathcal{L}_2$$

$$\mathcal{L}_1 = \prod_{i=1}^n f_X(X_i)$$

$$\mathcal{L}_2 = \prod_{i=1}^n f_{Y|X}(Y_i | X_i) \propto \sigma^{-n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_i (Y_i - (\beta_0 + \beta_1 X_i))^2 \right\}$$

Under the assumption of Normality, the least squares estimator is also the MLE but the least squares variance estimator is not the MLE.

$$\widehat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \widehat{\epsilon}_i^2$$

18.2 Prediction

Observe $X = x_*$ of the covariate and want to predict their outcome Y_* .

$$\widehat{Y}_* = \widehat{\beta}_0 + \widehat{\beta}_1 x_*$$

$$\mathbb{V}[\widehat{Y}_*] = \mathbb{V}[\widehat{\beta}_0] + x_*^2 \mathbb{V}[\widehat{\beta}_1] + 2x_* \text{Cov}[\widehat{\beta}_0, \widehat{\beta}_1]$$

Prediction interval

$$\widehat{\xi}_n^2 = \widehat{\sigma}^2 \left(\frac{\sum_{i=1}^n (X_i - X_*)^2}{n \sum_i (X_i - \bar{X})^2} + 1 \right)$$

$$\widehat{Y}_* \pm z_{\alpha/2} \widehat{\xi}_n$$

18.3 Multiple Regression

$$Y = X\beta + \epsilon$$

where

$$X = \begin{pmatrix} X_{11} & \cdots & X_{1k} \\ \vdots & \ddots & \vdots \\ X_{n1} & \cdots & X_{nk} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad \epsilon = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Likelihood

$$\mathcal{L}(\mu, \Sigma) = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \text{RSS} \right\}$$

$$\text{RSS} = (y - X\beta)^T (y - X\beta) = \|Y - X\beta\|^2 = \sum_{i=1}^n (Y_i - x_i^T \beta)^2$$

If the $(k \times k)$ matrix $X^T X$ is invertible,

$$\widehat{\beta} = (X^T X)^{-1} X^T Y$$

$$\mathbb{V}[\widehat{\beta} | X^n] = \sigma^2 (X^T X)^{-1}$$

$$\widehat{\beta} \approx \mathcal{N}(\beta, \sigma^2 (X^T X)^{-1})$$

Estimate regression function

$$\widehat{r}(x) = \sum_{j=1}^k \widehat{\beta}_j x_j$$

Unbiased estimate for σ^2

$$\widehat{\sigma}^2 = \frac{1}{n-k} \sum_{i=1}^n \widehat{\epsilon}_i^2 \quad \widehat{\epsilon} = X\widehat{\beta} - Y$$

MLE

$$\widehat{\mu} = \bar{X} \quad \widehat{\sigma}^2 = \frac{n-k}{n} \sigma^2$$

$1 - \alpha$ Confidence interval

$$\widehat{\beta}_j \pm z_{\alpha/2} \widehat{\text{se}}(\widehat{\beta}_j)$$

18.4 Model Selection

Consider predicting a new observation Y^* for covariates X^* and let $S \subset J$ denote a subset of the covariates in the model, where $|S| = k$ and $|J| = n$.
Issues

- Underfitting: too few covariates yields high bias
- Overfitting: too many covariates yields high variance

Procedure

1. Assign a score to each model
2. Search through all models to find the one with the highest score

Hypothesis testing

$$H_0 : \beta_j = 0 \text{ vs. } H_1 : \beta_j \neq 0 \quad \forall j \in J$$

Mean squared prediction error (MSPE)

$$\text{MSPE} = \mathbb{E} \left[(\hat{Y}(S) - Y^*)^2 \right]$$

Prediction risk

$$R(S) = \sum_{i=1}^n \text{MSPE}_i = \sum_{i=1}^n \mathbb{E} \left[(\hat{Y}_i(S) - Y_i^*)^2 \right]$$

Training error

$$\hat{R}_{tr}(S) = \sum_{i=1}^n (\hat{Y}_i(S) - Y_i)^2$$

R^2

$$R^2(S) = 1 - \frac{\text{RSS}(S)}{\text{TSS}} = 1 - \frac{\hat{R}_{tr}(S)}{\text{TSS}} = 1 - \frac{\sum_{i=1}^n (\hat{Y}_i(S) - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

The training error is a downward-biased estimate of the prediction risk.

$$\mathbb{E} \left[\hat{R}_{tr}(S) \right] < R(S)$$

$$\text{bias}(\hat{R}_{tr}(S)) = \mathbb{E} \left[\hat{R}_{tr}(S) \right] - R(S) = -2 \sum_{i=1}^n \text{Cov} \left[\hat{Y}_i, Y_i \right]$$

Adjusted R^2

$$R^2(S) = 1 - \frac{n-1}{n-k} \frac{\text{RSS}}{\text{TSS}}$$

MALLOW'S C_p statistic

$$\hat{R}(S) = \hat{R}_{tr}(S) + 2k\hat{\sigma}^2 = \text{lack of fit} + \text{complexity penalty}$$

AKAIKE Information Criterion (AIC)

$$\text{AIC}(S) = \ell_n(\hat{\beta}_S, \hat{\sigma}_S^2) - k$$

Bayesian Information Criterion (BIC)

$$\text{BIC}(S) = \ell_n(\hat{\beta}_S, \hat{\sigma}_S^2) - \frac{k}{2} \log n$$

Validation and training

$$\hat{R}_V(S) = \sum_{i=1}^m (\hat{Y}_i^*(S) - Y_i^*)^2 \quad m = |\{\text{validation data}\}|, \text{ often } \frac{n}{4} \text{ or } \frac{n}{2}$$

Leave-one-out cross-validation

$$\hat{R}_{CV}(S) = \sum_{i=1}^n (Y_i - \hat{Y}_{(i)})^2 = \sum_{i=1}^n \left(\frac{Y_i - \hat{Y}_i(S)}{1 - U_{ii}(S)} \right)^2$$

$$U(S) = X_S(X_S^T X_S)^{-1} X_S \text{ ("hat matrix")}$$

19 Non-parametric Function Estimation

19.1 Density Estimation

Estimate $f(x)$, where $f(x) = \mathbb{P}[X \in A] = \int_A f(x) dx$.

Integrated square error (ISE)

$$L(f, \hat{f}_n) = \int \left(f(x) - \hat{f}_n(x) \right)^2 dx = J(h) + \int f^2(x) dx$$

Frequentist risk

$$R(f, \hat{f}_n) = \mathbb{E} \left[L(f, \hat{f}_n) \right] = \int b^2(x) dx + \int v(x) dx$$

$$b(x) = \mathbb{E} \left[\hat{f}_n(x) \right] - f(x)$$

$$v(x) = \mathbb{V} \left[\hat{f}_n(x) \right]$$