

$$E[(Y - \hat{f}(x))^2 | X = x] = [f(x) - \hat{f}(x)]^2 + \text{Var}(\epsilon)$$

$$\text{MSE} = 1/n \sum (y_i - \hat{f}(x))^2 = E[(Y - \hat{f}(x))^2]$$

$$E[(y_0 - \hat{f}(x_0))^2] = \text{Var}(\hat{f}(x_0)) + (f(x_0) - E(\hat{f}(x_0)))^2 + \text{Var}(\epsilon)$$

$$\text{Classification Error (Loss)} = 1/n \sum I[y_i \neq \hat{C}(x_i)]$$

$$\text{RSS} = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \text{ -- Matrix form } (Y - X\beta)^T (Y - X\beta) = \epsilon^T \epsilon$$

$$\text{TSS} = \sum (y_i - \bar{y})^2; R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}}; \text{RSE} = \sqrt{\frac{\text{RSS}}{n-p}}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}; \beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum x_i (y_i - \bar{y})}{\sum x_i (x_i - \bar{x})}$$

$$\text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right], \text{SE}(\hat{\beta}_1)^2 = \sigma^2 \left[ \frac{1}{\sum (x_i - \bar{x})^2} \right], \text{ where } \sigma^2 = \text{Var}(\epsilon);$$

$$95\% \text{ Confidence Interval for } \beta_1 = [\hat{\beta}_1 - 2 * \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 * \text{SE}(\hat{\beta}_1)]$$

$$\hat{\sigma}^2 = \frac{\text{RSS}}{n-2}; \hat{\beta} = (X^T X)^{-1} X^T y \leftarrow \text{OLS Estimate}$$

$$\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2 \leftarrow \text{Matrix Form} \mid \text{Non-Matrix Form} \rightarrow \text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$

$$\beta_{\text{px1}} = \text{argmin}_{\beta} (y_i - x\beta)^T (y_i - x\beta); \text{ T-statistic } (\beta_1) = \frac{\hat{\beta}_1}{\text{SE}(\hat{\beta}_1)}$$

$$\text{Correlation Coefficient: } r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}; \text{ F-statistic} = \frac{(\text{TSS} - \text{RSS})/p}{\text{RSS}/(n-p-1)}$$

$$\text{Logistic Regression } p(x) = \frac{e^{(\beta_0 + \beta_1 x)}}{1 + e^{(\beta_0 + \beta_1 x)}}$$

$$\text{If } Y_i = \beta_0 + \beta_1 x_i + \epsilon_i \text{ and } \epsilon_i \sim N(0, \sigma^2)$$

$$\ell = \prod_{i=1}^p \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\epsilon_i - 0)^2}{2\sigma^2}} = \prod_{i=1}^p \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}}$$

$$\text{Maximum Likelihood Estimate} = \text{Min RSS}$$

$$\text{Bayes Theorem: } \Pr(Y=k|X=x) = (\Pr(X=x|Y=k) \cdot \Pr(Y=k)) / \Pr(X=x)$$

$$\text{Bayes Classifier: } p_k(x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)} \text{ and } \log(p_k(x)) = \log(\pi_k) + \log(f_k(x)) + \text{const}$$

$$\hat{\delta}(x) = x \cdot \frac{\hat{\mu}}{\hat{\sigma}^2} - \frac{\hat{\mu}^2}{2\hat{\sigma}^2} + \log(\hat{\pi}) \leftarrow \text{LDA Single-variate Discriminant}$$

$$\delta_k = x^T \Sigma^{-1} \mu_k - \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k \leftarrow \text{LDA Multivariate Discriminant}$$

$$\delta_k = -\frac{1}{2} x^T \Sigma_k^{-1} x + x^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma_k^{-1} \mu_k - \frac{1}{2} \log |\Sigma_k| + \log \pi_k \leftarrow$$

QDA

$$Pr(y = k | x = x) = \frac{e^{\delta_k(x)}}{\sum_{\ell} e^{\delta_{\ell}(x)}}$$

$$\text{Maximum Log Likelihood QDA} = \max \sum \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{(-\frac{(x_i - \mu_{y_i})^2}{2\sigma^2})} \right) = \max$$

$$\sum_{i=1}^N \left( -\log \sigma - \frac{1}{2\sigma^2} (x_i - \mu_{y_i})^2 + \text{consts} \right)$$

$$X = \begin{bmatrix} a & b & c & d \end{bmatrix}; \det(X) = ad - bc; X^{-1} = (1/\det X) \begin{bmatrix} d & -b & -c & a \end{bmatrix}$$

↑Flexibility → ↑Var + ↓Bias

Confusion Matrix

**True**

**Pred**                      Correct                      Type 1(False +)

                                 Type 2(False -)                      Correct

FPR: FP / (TN + FP); (type 1)

TPR: TP / (TP + FN)

$$C(x) : j, \text{ if } p_j(x) = \max\{p_1(x), \dots, p_k(x)\}, \text{ where } p_k(x) = \Pr(Y=k | X=x)$$

Multivariate Gaussian PDF:

$$f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\text{Logistic regression model: } \log\left(\frac{p(x)}{1-p(x)}\right) = \beta_0 + \beta_1 X$$

$$\text{Log reg. PDF: } p(x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

$$\text{Log Likelihood Log reg.: } \log \ell(x) = \log \left( \prod_{i=1}^N (p(x_i))^{y_i} (1 - p(x_i))^{1-y_i} \right)$$

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y); \text{Var}(aX) = a^2 \text{Var}(X); \text{Cov}(aX, bY) = ab\text{Cov}(X, Y);$$

$$\text{TSS} = \sqrt{RSE/(n-2)/(1-R^2)}; \text{LOOCV: CV} = 1/n \sum \text{MSE}$$

$$\text{Speeding up LOOCV: CV} = 1/n \sum ((y_i - \hat{y}_i)/(1 - h_i))^2$$

$$h_i = 1/n + ((x_i - \bar{x})^2 / \sum (x_j - \bar{x})^2)$$

K-fold cv:  $\sum n_k/n$  MSE

classification cv:  $\sum n_k/n$  Err

$$\text{Err} = \sum I(y_i \neq \hat{y}_i)/n_k$$

$$C_p = \frac{1}{n} (\text{RSS} + 2d\hat{\sigma}^2), \quad \text{AIC} = -2 \log L + 2 \cdot d$$

L is likelihood function

$$\text{BIC} = \frac{1}{n} (\text{RSS} + \log(n)d\hat{\sigma}^2).$$

$$\text{Maximum Log Likelihood} = \max \sum \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu_{yi})^2}{2\sigma^2}} \right)$$

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

$$\text{Adjusted } R^2 = 1 - \frac{\text{RSS}/(n - d - 1)}{\text{TSS}/(n - 1)}.$$

► **Ridge regression** (parameter  $\lambda$ ),  $\ell_2$  penalty

$$\begin{aligned} \min_{\beta} \text{RSS}(\beta) + \lambda \sum_j \beta_j^2 = \\ \min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_j \beta_j^2 \end{aligned}$$

**Constrained Ridge**

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \|Y - X^T \beta\|^2 + \lambda \|\beta\|^2,$$

Another way to formulate the ridge regression is

$$\begin{aligned} \hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \|Y - X^T \beta\|^2 \\ \text{subject to } \|\beta\|^2 \leq t, \end{aligned}$$

For ridge regression:  $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}} \quad \leq \text{Standardized predictors}$$

$$\begin{aligned} \min_{\beta} \text{RSS}(\beta) + \lambda \sum_j |\beta_j| = \\ \min_{\beta} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_j |\beta_j| \end{aligned} \quad \leq \text{Lasso}$$

$$\underset{\beta}{\text{minimize}} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq s$$

Constrained Lasso =>

$$\text{Polynomial} \Rightarrow Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_p X^p$$

$$\text{Step function} \Rightarrow Y = \beta_0 + \beta_1 I(X < c_1) + \beta_2 I(c_1 < X < c_2) + \dots + \beta_p I(c_p < X)$$

$$\text{Cubic Spline} \Rightarrow Y_i = \beta_0 + \beta_1 X_i + \beta_2 X_i^2 + \beta_3 X_i^3 + \beta_4 (X_i - \xi_1)_+ + \dots + \beta_{k+3} (X_i - \xi_k)_+$$

$$\underset{g \in \mathcal{S}}{\text{minimize}} \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g''(t)^2 dt$$

Smoothing Spline =>

$$\text{Gini Index} \Rightarrow G = \sum_{k=1}^K \hat{P}_{mk} (1 - \hat{P}_{mk})$$

$$\text{Cross-Entropy} \Rightarrow D = - \sum_{k=1}^K \hat{P}_{mk} \log_2(\hat{P}_{mk})$$

$$\text{If } \text{Var}(\alpha X + (1 - \alpha)Y) \text{ then minimizing this is } \alpha = \frac{\sigma_Y^2 - \sigma_{XY}}{\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}}$$

SVM (Hard-Margin)

SVM (Soft-Margin)

$$\underset{\beta_0, \beta_1, \dots, \beta_p, M}{\text{maximize}} \quad M$$

$$\text{subject to} \quad \sum_{j=1}^p \beta_j^2 = 1,$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n.$$

$$\underset{\beta_0, \beta_1, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M}{\text{maximize}} \quad M$$

$$\text{subject to} \quad \sum_{j=1}^p \beta_j^2 = 1,$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i),$$

$$\epsilon_i \geq 0, \quad \sum_{i=1}^n \epsilon_i \leq C, \quad \epsilon_1, \dots, \epsilon_n \text{ are slack variables}$$

$$f(x) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p = \hat{\beta}_0 + \sum_{i=1}^n \hat{\alpha}_i \langle x, x_i \rangle$$

Linear support vector classifier =>

$$f(x) = \hat{\beta}_0 + \sum_{i=1}^n \hat{\alpha}_i k(x, x_i)$$

Kernel version ^^ =>

$$K(x_i, x_{i'}) = \left( 1 + \sum_{j=1}^p x_{ij} x_{i'j} \right)^d$$

Polynomial Kernel =>

$$K(x_i, x_{i'}) = \exp(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2).$$

Radial Kernel =>

$$\underset{\beta_0, \beta_1, \dots, \beta_p}{\text{minimize}} \left\{ \sum_{i=1}^n \max[0, 1 - y_i f(x_i)] + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

SVM Hinge Loss =>

PCA and Eigen: *maximize*  $\phi$  with  $\phi_1^T \Sigma \phi_1$  *subject to*  $\phi_1^T \phi_1 = 1$  which is equivalent to  $\phi_1^T \Sigma \phi_1 - \lambda(\phi_1^T \phi_1 - 1)$  and then taking gradient and setting to zero ends up being  $\lambda$

PCA:  $Z_1 = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{p1}X_p$  for single component

$$\underset{\phi_{11}, \dots, \phi_{p1}}{\text{maximize}} \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{j1} x_{ij} \right)^2 \text{ subject to } \sum_{j=1}^p \phi_{j1}^2 = 1.$$

$$Z_{n \times 1} \Phi_1^T = x_{n \times 2}^* = \text{denoised data}$$

$$\hat{\Sigma} = \frac{X^T X}{n}$$

$$\sum_{j=1}^p \text{Var}(x_j) = \sum_{m=1}^M \text{Var}(Z_m)$$

K-means:

$$\min C_1 \dots C_K \sum_{k=1}^K WCV(C_k) \text{ where } WCV(C_k) = \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^n (x_{ij} - x_{i'j})^2 = \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j=1}^n (x_{ij} - \bar{x}_{kj})^2$$

As K increases, cost function decreases