

P13 Background

→ The results are output directly

Frequentists: Statistical Machine Learning \xrightarrow{eg} Linear Regression ① Linearity ② Global ③ Data isn't processed

Bayesians : Probabilistic Graphical Model Create a new model by changing ①②③

① Attribute Nonlinearity: Feature Conversion (Polynomial Regression e.g. $w_i x_i^2$)
 ② Global Nonlinearity: Linear Classification (The results are output to a nonlinear activation function)

③ Coefficient nonlinearity: Neural Networks (Coefficient are not fixed) \Rightarrow ~~$x \times x$~~ \because w initial is different, the result may be different.

② Segmenting the data: Regression Splines (样条回归: 数据分段并拟合), Decision Tree

③ Use after data processing : PCA

Linear Regression $\xrightarrow[\text{Dimensionality Reduction}]{\text{Activation Functions}}$ Linear Classification $\Rightarrow \begin{cases} y = f(w^T x + b) & y \in \{0, 1\} \\ f: w^T x + b \rightarrow y & f^{-1}: y \rightarrow w^T x + b \end{cases}$

Drop to one dimension, set the threshold. Greater than the threshold is 1, less than the threshold is 0.

P14 Preception (Only binary classification is possible for one preception)

↳ Error Driver $D: \{\text{Misclassified samples}\}$ Model: $y = \text{sign}(w^T x)$ $\text{sign}(a) = \begin{cases} 1 & a \geq 0 \\ -1 & a < 0 \end{cases}$

Loss Function: $\mathcal{L}(a) = \sum_{i=1}^N I\{y_i w^T x_i \leq 0\}$ $I(a) = \begin{cases} 1 & a = \text{True} \\ 0 & a = \text{False} \end{cases}$ $y_i w^T x_i \leq 0$ indicates that the sample was misclassified

Not Differentiable (2) $L(w) = \sum_{x_i \in D} -y_i w^T x_i$

P15 Linear Discriminant Analysis (LDA) (Multiclass Classification)

$$x_{c_1} = \{x_i \mid y_i = +1\}, \quad x_{c_2} = \{x_i \mid y_i = -1\}$$

Coordinate after projection: $Z = W^T X$; $\bar{Z} = \frac{1}{N} \sum_{i=1}^N W^T X_i$; $S_Z = \frac{1}{N} \sum_{i=1}^N (Z_i - \bar{Z})(Z_i - \bar{Z})^T = \frac{1}{N} \sum_{i=1}^N (W^T X_i - \bar{Z})(W^T X_i - \bar{Z})^T$

Let category C.G. Algorithm idea: Small sample distance within the same category, large sample distance between different categories

Within the same category: S_1, S_2 Between different categories: $(\bar{z}_1 - \bar{z}_2)^2$ Loss Function: $J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{S_1 + S_2}$ $\hat{w} = \arg \max_w J(w)$

$$(\bar{x}_1 - \bar{x}_2)^2 = \left[w^T \left(\frac{1}{N_1} \sum_{i=1}^{N_1} x_i - \frac{1}{N_2} \sum_{j=1}^{N_2} x_j \right) \right]^2 = w^T (\bar{x}_{c1} - \bar{x}_{c2}) (\bar{x}_{c1} - \bar{x}_{c2})^T w$$
$$S_1 = \frac{1}{N_1} \sum_{i=1}^{N_1} (w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j) (w^T x_i - \frac{1}{N_1} \sum_{j=1}^{N_1} w^T x_j)^T = w^T \frac{N_1}{N_1} \frac{1}{N_1} (x_i - \bar{x}_{c1}) (x_i - \bar{x}_{c1})^T w = w^T \zeta_{c1} w$$
$$J(w) = \frac{w^T (\bar{x}_{c1} - \bar{x}_{c2}) (\bar{x}_{c1} - \bar{x}_{c2})^T w}{w^T (S_{c1} + S_{c2}) w}$$

P16 LDA Model Solution

$$J(w) = \frac{(\bar{z}_1 - \bar{z}_2)^2}{s_1 + s_2} = \frac{w^T(\bar{x}_0 - \bar{x}_2)(\bar{x}_0 - \bar{x}_2)^T w}{w^T(s_{c1} + s_{c2})w} \quad \text{in } S_b = (\bar{x}_{c1} - \bar{x}_{c2})(\bar{x}_{c1} - \bar{x}_{c2})^T \text{ between-class variance}$$
$$S_{w, \text{PXD}} = S_{c1} + S_{c2} \text{ with -class variance}$$
$$J(w) = \frac{w^T S_b w}{w^T S_w w} \quad \frac{\partial (J(w))}{\partial w} = \frac{S_b w S_w^T S_w w - w^T S_b w S_w w}{(w^T S_w w)^2} = 0 \quad \therefore S_b w \underbrace{w^T S_w w}_{=1} = \underbrace{w^T S_b w}_{=1} S_w w$$

Focus mainly on the direction of w , not the size

$$W = \frac{W^T S_W W}{W^T S_W W} S_W^{-1} S_B \cdot W = \frac{W^T S_W W}{W^T S_W W} S_W^{-1} (\bar{X}_C - \bar{X}_A) (\bar{X}_C - \bar{X}_A)^T W \quad 1 \times p \times p \times 1 \in \mathbb{R} \quad \rightarrow \quad S_W^{-1} S_B W = \frac{W^T S_B W}{W^T S_W W} W, \therefore W \text{ is eigenvector of } S_W^{-1} S_B$$
$$w \propto S_w^{-1}(\bar{x}_1 - \bar{x}_2) \text{ if } S_w^{-1} \text{ is diagonal matrix, isotropic (各向同性)} \quad S_w^{-1} \propto I \Rightarrow w \propto (\bar{x}_1 - \bar{x}_2)$$

P17 Logistic Regression

Discriminative models find $P(y|x)$ directly $\Rightarrow \hat{y} = \arg \max_{y \in \{0,1\}} P(y|x)$ which can maximum $P(y|x)$

Sigmoid function: $\sigma(z) = \frac{1}{1+e^{-z}} \Rightarrow \ln \frac{\sigma(z)}{1-\sigma(z)} = z$, $\frac{\sigma(z)}{1-\sigma(z)}$ is odds (几率)

$$\text{Let } P(y=1|x) = \sigma(z) \left\{ \begin{aligned} p_1 &= P(y=1|x) = \sigma(w^T x) = \frac{1}{1+e^{-w^T x}}, y=1 \\ p_0 &= P(y=0|x) = 1 - \sigma(w^T x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}, y=0 \end{aligned} \right\} \rightarrow P(y|x) = p_1^y p_0^{(1-y)}$$
$$\text{MLE: } \hat{\omega} = \arg \max_{\omega} P(Y|x) = \arg \max_{\omega} \log \prod_{i=1}^N P(y_i|x_i) = \arg \max_{\omega} \sum_{i=1}^N (\log P(y_i|x_i))$$

$$= \arg \max_{\omega} \sum_{i=1}^N [y_i \cdot (\log p_i) + (1-y_i) \cdot (\log \bar{p}_i)] \quad \text{cross Entropy}$$

max MLE \Rightarrow min Loss Function (min Cross Entropy)

P18 Gaussian Discriminant Analysis

Generative models \Rightarrow Compare $P(y=0|x)$ and $P(y=1|x)$ by $p(x|y)p(y)$ without finding $P(y|x) \Rightarrow \hat{y} = \arg \max_{y \in \{0,1\}} P(x|y)P(y)$

$$y \sim \text{Bernoulli}(\phi) \Rightarrow y = \begin{cases} \phi & y=1 \\ 1-\phi & y=0 \end{cases} \Rightarrow \hat{\phi} \Rightarrow y = \phi^y (1-\phi)^{(1-y)}$$
$$\left. \begin{aligned} x|y=1 &\sim N(\mu_1, \Sigma) \\ x|y=0 &\sim N(\mu_2, \Sigma) \end{aligned} \right\} \Rightarrow P(x|y) = N(\mu_1, \Sigma)^y \cdot N(\mu_2, \Sigma)^{(1-y)}$$
$$\begin{aligned} \text{log-likelihood: } L(\theta) &= \log_{\frac{1}{2}} P(y_i | x_i) = \sum_{i=1}^N [\log P(x_i | y_i) + \log P(y_i)] \\ &= \sum_{i=1}^N [\log \pi(\mathcal{H}_i, \Sigma) + \log \pi(\mathcal{H}_i, \Sigma)^{(1-y_i)} + \log \phi(y_i) (1-\phi)^{(1-y_i)}] \end{aligned}$$
$$\theta = (\mu_1, \mu_2, \Sigma, \phi) \quad \hat{\theta} = \arg \max_{\theta} l(\theta)$$

P19 GDA Model Solution

$$\textcircled{1} \text{ Find } \phi: \frac{\partial \mathcal{L}(\theta)}{\partial \phi} = \frac{\partial}{\partial \phi} \sum_{i=1}^N [y_i \log \phi + (1-y_i) \log(1-\phi)] = \sum_{i=1}^N [y_i \cdot \frac{1}{\phi} - (1-y_i) \frac{1}{1-\phi}] = 0$$
$$\therefore \sum_{i=1}^N [y_i(1-\phi) - (1-y_i)\phi] = 0 \Rightarrow \sum_{i=1}^N (y_i - \phi) = 0 \Rightarrow \phi = \frac{1}{N} \sum_{i=1}^N y_i$$
$$\therefore y=1: N_1 \quad y=0: N_2 \quad N_1 + N_2 = N \quad \therefore \phi = \frac{N_1}{N}$$

② Find μ_1 : $\sum_{i=1}^N \log N(\mu_1, \Sigma)^{-1} = \sum_{i=1}^N y_i \left(\log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right\} \right)$

$$\therefore \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\therefore y=1: N_1, y=0: N_2, N_1+N_2=N \therefore \phi = \frac{N_1}{N}$$

$$\textcircled{1} \text{ Find } \mu_1: \frac{1}{N} \sum_{i=1}^N \log N(\mu_1, \Sigma) y_i = \frac{1}{N} \sum_{i=1}^N y_i \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right\}$$

$$\hat{\mu}_1 = \arg \max_{\mu_1} \frac{1}{N} \sum_{i=1}^N y_i \left[-\frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right] \therefore \frac{\partial f(\mu_1)}{\partial \mu_1} = 0$$

$$\therefore \frac{1}{N} \sum_{i=1}^N y_i \Sigma^{-1} (x_i - \mu_1) = 0 \Rightarrow \sum_{i=1}^N y_i \mu_1 = \sum_{i=1}^N y_i x_i \Rightarrow \mu_1 = \frac{\sum_{i=1}^N y_i x_i}{\sum_{i=1}^N y_i} = \frac{\sum_{i=1}^N y_i x_i}{N_1}$$

P20 GDA Model Solution ②

$$\text{Find } \Sigma: \frac{1}{N} \sum_{i=1}^N y_i \log N(\mu_i, \Sigma) = \sum_{x_i \in C_1} \log \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right\}$$

$$= \sum_{x_i \in C_1} \left[\log(2\pi)^{-\frac{p}{2}} - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right]$$

$$= -\frac{1}{2} N_1 \log |\Sigma| - \frac{1}{2} \sum_{i=1}^{N_1} (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) + c \quad \text{tr}(c) = c \text{ 矩阵的迹}$$

$$\text{Formula: } \frac{\partial \text{tr}(AB)}{\partial A} = B^T; \frac{\partial |A|}{\partial A} = |A| A^{-1}; \text{tr}(ABC) = \text{tr}(CAB)$$

$$\therefore \text{Original formular} = -\frac{1}{2} N_1 \log |\Sigma| - \frac{1}{2} \text{tr} \left[\sum_{i=1}^{N_1} (x_i - \mu_1)(x_i - \mu_1)^T \Sigma^{-1} \right] + c$$

$$\therefore \frac{1}{N_1} \sum_{i=1}^{N_1} (x_i - \mu_1)(x_i - \mu_1)^T = S_1, \therefore \text{Original formular} = -\frac{1}{2} N_1 \log |\Sigma| - \frac{1}{2} N_1 \text{tr}(S_1 \Sigma^{-1}) + c$$

$$\therefore \hat{\Sigma} = \arg \max_{\Sigma} -\frac{1}{2} (N_1 + N_2) \log |\Sigma| - \frac{1}{2} N_1 \text{tr}(S_1 \Sigma^{-1}) - \frac{1}{2} N_2 \text{tr}(S_2 \Sigma^{-1}) + c$$

$$\frac{\partial f(\Sigma)}{\partial \Sigma} = 0 \therefore \frac{\partial \log |\Sigma|}{\partial \Sigma} = \frac{1}{|\Sigma|} S^{-1} \quad \frac{\partial \text{tr}(S_1 \Sigma^{-1})}{\partial \Sigma} = S_1^{-1} \Sigma^{-2}$$

$$\therefore -\frac{1}{2} N \Sigma^{-1} + \frac{1}{2} N_1 S_1 \Sigma^{-2} + \frac{1}{2} N_2 S_2 \Sigma^{-2} = 0 \Rightarrow N \Sigma = N_1 S_1 + N_2 S_2 = 0$$

$$\therefore \hat{\Sigma} = \frac{1}{N} (N_1 S_1 + N_2 S_2)$$

P21 Naive Bayes Classifier

Algorithm idea: Naive Bayesian Hypothesis \Rightarrow Conditional Independence Hypothesis

\Rightarrow The simplest probability graphical model (directed graph)

$$\begin{array}{c} \textcircled{y} \\ \swarrow \quad \searrow \\ \textcircled{x_1} \quad \textcircled{x_2} \\ \downarrow \quad \downarrow \\ \textcircled{x_p} \end{array} \quad \begin{array}{l} x_i \perp x_j | y \quad (i \neq j) \\ P(x|y) = \prod_{i=1}^p P(x_i|y) \end{array} \quad X \in \mathbb{R}^p$$

Single Experiment Multiple Experiments

Binomial classification Bernoulli (伯努利分布) Binomial (二项式分布)

Multiclassification Categorical (分类分布) Multinomial (多项式分布)

$$P(y = \Delta | X) = \frac{P(X|y = \Delta)P(y = \Delta)}{P(X)} \propto P(X|y = \Delta)P(y = \Delta)$$

The purpose is to compare the size of $P(y|X)$ of different categories and make a classification.

Without caring about the value of $P(y|X)$