1. 线性回归. → 回归任务

学习 $f(z_i) = wx + b$ 使得 $f(x_i) \simeq y_i$     用均方误差最小化求求 $w$ 和 $b$

$\dfrac{\partial E(w,b)}{\partial w} = 2(w \sum x_i^2 - \sum(y_i - b)x_i)$     $\Rightarrow$  $w = \dfrac{\sum y_i(x_i - \bar{x})}{\sum x_i^2 - \frac{1}{m}(\sum x_i)^2}$     closed-form 的解     更一般来说 $\dfrac{\partial E}{\partial \hat{w}} = 2X^T(X\hat{w} - y)$

$\dfrac{\partial E(w,b)}{\partial b} = 2(mb - \sum(y_i - wx_i))$     $b = \dfrac{1}{m} \sum(y_i - wx_i)$     $\hat{w}^* = (X^TX)^{-1}X^Ty$

例:

$y = g^{-1}(w^Tx + b)$ 广义线性模型     $\ln y = w^Tx + b$

$g(y) = w^Tx + b$ ↑     $y = e^{w^Tx + b}$

2. logistic regression

$\ln \dfrac{y}{1-y} = w^Tx + b$     $\dfrac{y}{1-y}$ 称为 odds   $y$ 表示视样本 $x$ 为正样本的几率

用极大似然法估计参数

优点: 无需假设数据分布
可得到预测概率
sigmoid $\Rightarrow$ 任意阶可导

3. Linear Discriminant Analysis LDA 线性判别分析     → 监督降维技术

将样例投影到直线上 使同类样例投影尽可能接近

def.
两类样本的中心在直线上的投影 $w^T\mu_0$, $w^T\mu_1$; 两类样本协方差 $w^T\Sigma_0 w$, $w^T\Sigma_1 w$

欲使同类样本接近 异类样本远离  $w^T\Sigma_0 w + w^T\Sigma_1 w \downarrow$  $\|w^T\mu_0 - w^T\mu_1\|$ 应尽可能大 $\Rightarrow$ $J = \dfrac{\|w^T\mu_0 - w^T\mu_1\|_2^2}{w^T\Sigma_0 w + w^T\Sigma_1 w}$

类内散度矩阵 $S_w = \Sigma_0 + \Sigma_1 = \sum\limits_{x \in X_0}(x-\mu_0)(x-\mu_0)^T + \sum\limits_{x \in X_1}(x-\mu_1)(x-\mu_1)^T$

$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T$

$= \dfrac{w^T(\mu_0 - \mu_1)(\mu_0 - \mu_1)^Tw}{w^T(\Sigma_0 + \Sigma_1)w}$

$= \dfrac{w^TS_b w}{w^TS_w w}$ $\Rightarrow$ LDA 最大化目标

$\Rightarrow$ $\min -w^TS_b w$
s.t. $w^TS_w w = 1$

$\Rightarrow$ $S_b w = \lambda S_w w = \lambda(\mu_0 - \mu_1)$
$w = S_b^{-1}(\mu_0 - \mu_1)$

多分类学习 $\begin{cases} OvO & N(N-1)/2 \text{ 两两配对} \\ OvR \\ MvM \end{cases}$ 特例

类别不均衡 $\begin{cases} \text{理论解决方法: } \dfrac{y}{1-y} > 1 \text{ 视为正例} \Rightarrow \dfrac{y}{1-y} > \dfrac{m^+}{m^-} \text{ 视为正例. 再缩放} \quad \text{代价敏感学习} \\ \text{实际操作方法: } \begin{cases} \text{欠采样: 去除反例} \quad \text{EasyEnsemble} \\ \text{过采样: 增加正例} \quad \text{SMOTE} \\ \text{阈值移动 threshold-moving} \end{cases} \end{cases}$