

# Linear Model

## 1. Linear Regression

### 1.1 Simple Linear Regression

- *model formulation*

$$\hat{y}_i = wx_i + b$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

$$(w^*, b^*) = \arg \min_{(w, b)} J(w, b)$$

- *parameter estimation*

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{2}{m} \sum_i (\hat{y}_i - y_i) \frac{\partial}{\partial w} (\hat{y}_i - y_i) = \frac{2}{m} \sum_i (\hat{y}_i - y_i) x_i \\ &= \frac{2}{m} \left( w \sum_i x_i^2 - \sum_i (y_i - b) x_i \right) \\ \frac{\partial J}{\partial b} &= \frac{2}{m} \sum_i (\hat{y}_i - y_i) \frac{\partial}{\partial b} (\hat{y}_i - y_i) = \frac{2}{m} \sum_i (\hat{y}_i - y_i) \\ &= \frac{2}{m} \left( mb - \sum_i (y_i - wx_i) \right) \end{aligned}$$

- *closed-form solution*

$$\begin{aligned} w^* &= \frac{\sum y_i (x_i - \bar{x})}{\sum x_i^2 - \frac{1}{m} (\sum x_i)^2} \\ b^* &= \frac{1}{m} \sum_{i=1}^m (y_i - wx_i) \end{aligned}$$

### 1.2 Multiple Linear Regression

- *model formulation*

$$\mathbf{w} = (\mathbf{w}; b) \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_m \\ 1 & 1 & \cdots & 1 \end{pmatrix}^T = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix}$$

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} = \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{pmatrix} \quad \hat{\mathbf{y}}^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)}$$

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{m} \sum_{i=1}^m \left( \hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)} \right)^2 = \frac{1}{m} \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \\ &= \frac{1}{m} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad \mapsto [MSE_{train}] \end{aligned}$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$$

度量模型性能的一种方法是计算在测试集上的均方误差MSE, 为了减小MSE, 一种直观方式是最小化训练集上的均方误差

## • *parameter estimation*

- Normal Equation

$$\begin{aligned} \nabla_{\mathbf{w}} J(\mathbf{w}) = 0 &\Rightarrow \nabla_{\mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0 \\ &\Rightarrow \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

- Gradient Descent

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{w}_j} &= \frac{\partial}{\partial \mathbf{w}_j} \left( \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)} \right) \mathbf{x}_j^{(i)} \end{aligned}$$

## 1.3 Probabilistic interpretation for cost function

$\mathbf{y}^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} + \epsilon^{(i)}$ , assume  $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$  and are distributed IID

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

for given  $\mathbf{x}^{(i)}$  and  $\boldsymbol{\theta}$

$$p(y_i|x_i;\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\mathbf{y}^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right)$$

$$L(\boldsymbol{\theta}) = L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = p(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) = \prod_{i=1}^m p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

似然函数(likelihood function),  $L(\boldsymbol{\theta})$ 表示在概率密度函数的参数是 $\boldsymbol{\theta}$ 时, 得到这组样本的概率

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \log L(\boldsymbol{\theta})$$

$$\ell(\boldsymbol{\theta}) = \log L(\boldsymbol{\theta}) = \sum_{i=1}^m \log p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \boldsymbol{\theta})$$

$$= \sum \left( \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{(\mathbf{y}^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2}{2\sigma^2} \right)$$

$$= m \ln \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^m (\hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)})^2$$

$$\arg \max \ell(\boldsymbol{\theta}) \Rightarrow \arg \min \|\hat{\mathbf{y}} - \mathbf{y}\|^2$$

最大化关于 $w$ 的对数似然和最小化均方误差会得到相同的参数估计

## 1.4 Generalized Linear Models

For a monotonic and differentiable function  $g(\cdot)$

$$y = g^{-1}(\mathbf{w}^T \mathbf{x} + b)$$

These broader family of models Generalized Linear Models, 其中函数 $g(\cdot)$ 称为“联系函数”(link function)

## 2. Logistic Regression

### • model formulation

$$\hat{y} = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

$$\rightarrow \ln \frac{y}{1-y} = \mathbf{w}^T \mathbf{x} + b$$

$y/(1-y)$ : 几率(odds)反映了 $x$ 为正例的相对可能性,  $\ln y/(1-y)$ : 对数几率(log odds亦称logit)

LR模型本质上在回归真值的对数几率，即将GLM模型中的联系函数设置为logistic function

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$Cost(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$J(w) = -\frac{1}{m} \sum_{i=1}^m \left[ \mathbf{y}^{(i)} \log \hat{\mathbf{y}}^{(i)} + (1 - \mathbf{y}^{(i)}) \log(1 - \hat{\mathbf{y}}^{(i)}) \right]$$

在这里不使用均方误差代价函数是由于均方误差代价函数在这里不是凸函数，存在多个局部极小值点

## • *parameter estimation*

- Gradient Descent

$$\frac{\partial J}{\partial \mathbf{w}_j} = \frac{1}{m} \sum_{i=1}^m \left( \hat{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)} \right) \mathbf{x}_j^{(i)}$$

对logistic function求导得  $g' = g(1-g)$

$$\frac{d}{dz} g(z) = \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-z}} \cdot \frac{e^{-1} + 1 - 1}{1 + e^{-z}} = \frac{1}{1 + e^{-z}} \cdot \left( 1 - \frac{1}{1 + e^{-z}} \right) = g(z)(1 - g(z))$$

for one training example, assume  $z = \mathbf{w}^T \mathbf{x}$  in this,  $\mathbf{w} = (\mathbf{w}; b)$  and  $\mathbf{x} = (\mathbf{x}; 1)$

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{w}_j} &= \frac{\partial J}{\partial g(z)} \cdot \frac{\partial g(z)}{\partial z} \cdot \frac{\partial z}{\partial \mathbf{w}_j} \\ &= - \left( y \frac{1}{g(z)} + (1 - y) \frac{1}{1 - g(z)} \cdot (-1) \right) \frac{\partial g(z)}{\partial z} \cdot \frac{\partial z}{\partial \mathbf{w}_j} \\ &= - \left( y \frac{1}{g(z)} - (1 - y) \frac{1}{1 - g(z)} \right) g(z)(1 - g(z)) \cdot x_j \\ &= - (y(1 - g(z)) - (1 - y)g(z)) x_j \\ &= -(y - g(z)) x_j \\ &= (\hat{y} - y) x_j \end{aligned}$$

## • *probabilistic interpretation*

Assume  $y$  satisfies Bernoulli distribution

$$\begin{aligned} P(y = 1 | \mathbf{x}) &= \hat{y} \\ P(y = 0 | \mathbf{x}) &= 1 - \hat{y} \end{aligned}$$

$$\Rightarrow p(y|x; w) = \hat{y}^y (1 - \hat{y})^{1-y} \quad (1)$$

将 $y=1$ 带入(1)式得 $p(y|x)=h(x)$ , 将 $y=0$ 带入(1)式得 $p(y|x)=1-h(x)$

$$L(\mathbf{w}) = p(\mathbf{y}|\mathbf{X}; \mathbf{w}) = \prod_{i=1}^m p(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}; \mathbf{w})$$

$$\ell(\mathbf{w}) = \log L(\mathbf{w}) = \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

$$\Rightarrow J(\mathbf{w}) = -\ell(\mathbf{w}) \quad \arg \min J(\mathbf{w}) \Rightarrow \arg \max \ell(\mathbf{w})$$

### 3. LDA

### PLA 与样本线性组合

### 4. Reference

[1] Andrew Ng, Machine Learning, <https://www.coursera.org/learn/machine-learning/>

[2] 周志华, 《机器学习》, 清华大学出版社

[3] Ian Goodfellow and Yoshua Bengio and Aaron Courville, "Deep Learning"