

# Linear Regression

## Linear Model

$$h(x) = \sum_{i=0}^d \theta_i x_i = \theta^T x$$

数据集  $D = \{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$ .  $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_d^{(i)})$

$\therefore \theta \in \mathbb{R}^d$ ,  $x \in \mathbb{R}^{n \times d}$ ,  $h(x) \in \mathbb{R}^n$ ,  $x_0 = 1$  (intercept term).

## Cost function

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2$$

求解方程其实有2种方法,

1° 线性方程, 多元方程组的直接求解方法.

2° 数值分析/优化的方法  $\Rightarrow$  先随机猜一个初始值, 然后不断迭代优化逼近正确值.

LMs algorithm (least mean square) (Gradient Descent)

根据梯度下降, 可以使用如下公式在每一步更新参数

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta), \quad j \in \{0, 1, \dots, d\}$$

1°  $\frac{\partial}{\partial \theta_j} J(\theta)$ :  $J(\theta)$  在  $\theta_j$  处的梯度, 此时  $J(\theta)$  沿梯度方向变化最快, 增长最快.

2°  $-\frac{\partial}{\partial \theta_j} J(\theta)$ : 负梯度方向, 此时  $J(\theta)$  下降最快, 从而得到最小的误差.

3°  $\alpha$ : learning rate, 在梯度方向上的步长.

假设只有一个训练样本, 对  $\theta_j$  求导:

$$\frac{d}{d\theta_j} J(\theta) = \frac{d}{d\theta_j} \frac{1}{2} (h(\theta) - y)^2$$

$$= \frac{1}{2} \cdot 2 (h(\theta) - y) \frac{d}{d\theta_j} (h(\theta) - y)$$

$$= (h(\theta) - y) \frac{d}{d\theta_j} \left( \sum_{i=0}^n \theta_j \cdot x_j - y \right)$$

$$= (h(\theta) - y) x_j$$

$$\Rightarrow \theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h(x^{(i)})) x_j^{(i)}$$

△ 对于全部的训练样本和各个维度的参数:

$$\theta := \theta + \alpha \sum_{i=1}^n [y^{(i)} - h(x^{(i)})] x^{(i)}$$

1° 每一次迭代, 更新量与误差成正比。

2° 上式称为 batch gradient descent  $\Rightarrow$  每次更新考虑全部的样本。

3° 由于  $J(\theta)$  是一个凸函数, 所以由梯度下降得到的解是全局最优解。

stochastic gradient descent

loop 1

for  $i = 1$  to  $n$ , 1

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)}, \text{ (for every } j \text{)}$$

1° SGD 可以尽可能逼近全局最优, 但无法保证必达到全局最优。

2° 通过逐步降低  $\alpha$  至 0, 可以到达全局最优。

# normal equations

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n [h_{\theta}(x^{(i)}) - y^{(i)}]^2$$

$$= \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$\therefore \nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (X\theta - y)^T (X\theta - y)$$

$$= \frac{1}{2} \nabla_{\theta} [(X\theta)^T X\theta - (X\theta)^T y - y^T X\theta + y^T y]$$

$$= \frac{1}{2} \nabla_{\theta} [\theta^T X^T X\theta - \theta^T X^T y - y^T X\theta]$$

$$= \frac{1}{2} \nabla_{\theta} [\theta^T X^T X\theta - 2y^T X\theta]$$

$$= \frac{1}{2} \nabla_{\theta} [\theta^T X^T X\theta - 2(X^T y)^T \theta]$$

$$= \frac{1}{2} (X^T X\theta + (X^T X)^T \theta - 2X^T y)$$

$$= X^T X\theta - X^T y = 0$$

$$\Rightarrow X^T X\theta = X^T y \Rightarrow \theta = (X^T X)^{-1} X^T y$$

## Probabilistic Interpretation (maximum likelihood)

假设  $y^{(i)}$  和  $x^{(i)}$  之间存在线性关系:

$$y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$$

若误差  $\varepsilon^{(i)}$  服从均值为0的正态分布,  $\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$ , 则:

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

则:  $y^{(i)} | x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)$ :

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

则有似然函数:

$$L(\theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

取对数似然:

$$\ell(\theta) = \log L(\theta) = \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

↓

⇒ 最大化对数似然, 等价于最小化均方误差, 且最优的  $\theta$  值与  $\sigma^2$  的取值无关。

$$\frac{\partial \ell(\theta)}{\partial \theta} = 0 - \frac{1}{\sigma^2} \cdot \frac{1}{2} \cdot 2 \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) \cdot (-x^{(i)}) = 0$$

$$\Rightarrow \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)}) \cdot x^{(i)} = 0 \Rightarrow X^T Y = (X^T X) \cdot \theta$$

$$\Rightarrow \theta = (X^T X)^{-1} X^T Y$$



## Locally weighted linear regression

损失函数调整为:

$$J(\theta) = \sum_{i=1}^n w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

→ 权的权重

$w^{(i)}$  通常取值为:

$$w^{(i)} = \exp\left[-\frac{(x^{(i)} - x)^2}{2\tau^2}\right]$$

→ 预测值.  
→ 高斯核函数.  
→ 带宽 (bandwidth)

1° 当  $x^{(i)}$  与  $x$  的距离  $|x^{(i)} - x|$  很大时,  $w^{(i)}$  很小, 此时对  $J(\theta)$  影响很小, 不必关注  $y^{(i)} - \theta^T x^{(i)}$ .

2° 当  $x^{(i)}$  与  $x$  的距离  $|x^{(i)} - x|$  很小时,  $w^{(i)} = 1$ , 此时对  $J(\theta)$  影响最大, 需尽量减小误差.

对  $J(\theta)$  求导得:

$$J(\theta) = W(y - X\theta)^T (y - X\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta} W(y - X\theta)^T (y - X\theta)$$

$$= \frac{\partial}{\partial \theta} [W(y^T - \theta^T X^T) (y - X\theta)]$$

$$= \frac{\partial}{\partial \theta} [(Wy^T - W\theta^T X^T) (y - X\theta)]$$

$$= \frac{\partial}{\partial \theta} (Wy^T y - Wy^T X\theta - W\theta^T X^T y + W\theta^T X^T X\theta)$$

$$= (0 - Wy^T X - W y^T X + W\theta^T X^T X + W\theta^T X^T X)$$

$$= -2Wy^T X + 2W\theta^T X^T X = 0$$

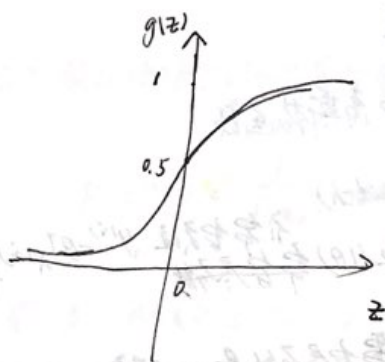
$$\Rightarrow W\theta^T X^T X = Wy^T X \Rightarrow W\theta^T = Wy^T X (X^T X)^{-1} \Rightarrow \theta W^T = Wy^T X (X^T X)^{-1}$$

$$\Rightarrow \theta = Wy^T X (X^T X)^{-1} (W^T)^{-1} = Wy^T X (W^T X^T X)^{-1}$$

# Logistic Regression $\Rightarrow$ binary classification $\{0, 1\}$

$$h(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}, \quad g(z) = \frac{1}{1 + e^{-z}}.$$

$\hookrightarrow$  sigmoid function.



Sigmoid 函数在 0 附近以斜率很大，所以将一些临界值映射到 0 或 1，避免分类错误。

$$\begin{aligned} g(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\ &= \frac{-1}{(1 + e^{-z})^2} \cdot (e^{-z})(-1) \\ &= \frac{e^{-z}}{(1 + e^{-z})^2} \end{aligned}$$

$$\begin{aligned} &= \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}}\right) \\ &= g(z) \cdot (1 - g(z)). \end{aligned}$$

## Probabilistic Interpret.

假设  $p(y=1|x;\theta) = h(x)$ .

$$p(y=0|x;\theta) = 1 - h(x)$$

$$\therefore p(y|x;\theta) = (h(x))^y (1 - h(x))^{1-y} \Rightarrow \text{伯努利分布.}$$

$$\therefore L(\theta) = \prod_{i=1}^n p(y^{(i)}|x^{(i)};\theta) = \prod_{i=1}^n (h(x^{(i)}))^{y^{(i)}} (1 - h(x^{(i)}))^{1-y^{(i)}}.$$

$$\therefore \ell(\theta) = \log L(\theta) = \sum_{i=1}^n y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log (1 - h(x^{(i)}))$$

假设有  
对一个样本, 对  $\theta$  求导:

$$\frac{\partial}{\partial \theta_j} l(\theta) = \frac{\partial}{\partial \theta_j} [y \log h_{\theta}(x) + (1-y) \log (1-h_{\theta}(x))]$$

$$= \left[ y \frac{1}{h_{\theta}(x)} \cdot h_{\theta}(x) \cdot (1-h_{\theta}(x)) + (1-y) \frac{1}{1-h_{\theta}(x)} \cdot (-1) \cdot h_{\theta}(x) (1-h_{\theta}(x)) \right] \cdot \frac{\partial \theta^T x}{\partial \theta_j}$$

$$= [y(1-h_{\theta}(x)) + (y-1)h_{\theta}(x)] \cdot \frac{\partial \theta^T x}{\partial \theta_j}$$

$$= (y - h_{\theta}(x)) \cdot x_j$$

所以, 对于随机梯度的下降有:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Newton's method:

$$\theta := \theta - \frac{f'(\theta)}{f''(\theta)}, \quad f: \mathbb{R} \rightarrow \mathbb{R}.$$

~~maximum  $f(\theta)$   $\Rightarrow$  find  $\theta$  which makes  $f'(\theta) = 0$ .~~

maximum  $l(\theta)$  equals find  $\theta^*$  makes  $l'(\theta^*) = 0$ .

$$\theta := \theta - \frac{l'(\theta)}{l''(\theta)}$$

考虑向量表示:

$$\theta := \theta - H^{-1} \nabla_{\theta} l(\theta), \quad \text{其中 } H_{ij} = \frac{\partial^2 l(\theta)}{\partial \theta_i \partial \theta_j}.$$

Newton-Raphson method.

当年预估值应用于最大化  
逻辑回归的似然函数时,  
得到的函数是 Fisher  
Scoring.

# Generalized Linear Models

## exponential family

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

$\downarrow$  natural parameter.       $\downarrow$  log partition function.       $\nearrow$  sufficient statistic.

## Bernoulli distribution

$$p(y; \phi) = \phi^y (1-\phi)^{1-y}$$

$$= \exp(y \log \phi + (1-y) \log(1-\phi))$$

$$= \exp(y \log \frac{\phi}{1-\phi} + \log(1-\phi))$$

$$\begin{matrix} b(y) & T(y) & \eta & a(\eta) \end{matrix}$$

$$\Rightarrow \begin{cases} b(y) = 1 \\ T(y) = y \\ a(\eta) = -\log(1-\phi) = \log(1+e^\eta) \end{cases}$$

$$\eta = \log \frac{\phi}{1-\phi}$$

$$\therefore \frac{\phi}{1-\phi} = e^\eta$$

$$\therefore \phi = (1-\phi)e^\eta$$

$$\therefore \phi(1+e^\eta) = e^\eta$$

$$\therefore \phi = \frac{e^\eta}{1+e^\eta} = \frac{1}{1+e^{-\eta}}$$

## Gaussian distribution

假设  $\sigma^2 = 1$ , 则有:

$$p(y; \eta) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(y-\mu)^2\right]$$

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \exp\left(\eta y - \frac{1}{2}\eta^2\right)$$

$$\begin{matrix} b(y) & T(y) & \eta & a(\eta) \end{matrix}$$

$$\therefore \begin{cases} b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \\ \eta = \mu \\ T(y) = y \\ a(\eta) = \frac{1}{2}\mu^2 = \frac{1}{2}\eta^2 \end{cases}$$



## 构建 GLMs

线性模型基于以下 3 个假设：

1. 给定  $x$  和  $\theta$ ,  $y$  的分布服从指数家族分布:  $y|x; \theta \sim \text{Exponential Family}(\eta)$
2. 给定  $x$ , GLM 的目标是使得学习到的假设满足  $h(x) = E[y|x]$
3.  $\eta$  和  $x$  呈线性相关:  $\eta = \theta^T x$

### 线性回归 (Least Square) (Gaussian Distribution)

$$h(x) = E[y|x; \theta] \Rightarrow \text{上面的假设 (2)}$$

$$= \mu$$

$\Rightarrow y|x; \theta$  服从高斯分布.

$$= \eta$$

$\Rightarrow$  高斯分布属于指数家族 ( $\eta = \mu$ )

$$= \theta^T x$$

$\Rightarrow y|x; \theta \sim \mathcal{N}(\theta^T x, \sigma^2)$  / 上面的假设 (3)

### 逻辑回归 (Logistic Regression) (Bernoulli Distribution)

$$h(x) = E[y|x; \theta] \Rightarrow \text{上面的假设 (2)}$$

$$= \phi$$

$\Rightarrow y|x; \theta$  服从伯努利分布.

$$= \frac{1}{1 + e^{-\eta}}$$

$\Rightarrow$  伯努利分布属于指数家族.

$$= \frac{1}{1 + e^{-\theta^T x}}$$

$\Rightarrow$  上面的假设 (3)

# Softmax Regression

$k$  分类问题:  $y \in \{1, 2, \dots, k\}$ , 使用  $\phi_i$  表示  $y$  取  $k$  类别的概率。

实际上, 可以使用  $k-1$  个自由参数拟合模型, 因为所有概率之和为 1。

$$T(y) \in \mathbb{R}^{k-1}$$

$$T_1(y) = \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}, \dots, T_{k-1}(y) = \begin{bmatrix} 0 \\ 0 \\ \vdots \end{bmatrix}, T_k(y) = \begin{bmatrix} 0 \\ 0 \\ \vdots \end{bmatrix}$$

使用多项式建模:  $y|x; \phi \sim \text{Mult}(\phi_1, \dots, \phi_k)$

$$\therefore E[T_i(y)] = P(y=i) = \phi_i$$

证明多项式属于指数家族。

$$\therefore P(y; \phi) = \phi_1^{I(y=1)} \phi_2^{I(y=2)} \dots \phi_k^{I(y=k)}$$

$$= \phi_1^{T_1(y)} \phi_2^{T_2(y)} \dots \phi_k^{T_k(y)}$$

$$= \phi_1^{T_1(y)} \phi_2^{T_2(y)} \dots \phi_{k-1}^{T_{k-1}(y)} \phi_k^{1 - \sum_{i=1}^{k-1} T_i(y)}$$

$$= \exp \left[ T_1(y) \log \phi_1 + T_2(y) \log \phi_2 + \dots + T_{k-1}(y) \log \phi_{k-1} + \left(1 - \sum_{i=1}^{k-1} T_i(y)\right) \log \phi_k \right]$$

$$= \exp \left[ T_1(y) \log \frac{\phi_1}{\phi_k} + T_2(y) \log \frac{\phi_2}{\phi_k} + \dots + T_{k-1}(y) \log \frac{\phi_{k-1}}{\phi_k} + \log \phi_k \right]$$

$$\therefore \begin{cases} b(y) = 1 \\ \eta = \begin{bmatrix} \log(\phi_1 / \phi_k) \\ \log(\phi_2 / \phi_k) \\ \vdots \\ \log(\phi_{k-1} / \phi_k) \end{bmatrix} \end{cases}$$

$$T(y) = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^{k-1}$$

$$a(\eta) = -\log \phi_k$$

设  $\eta_i$  表示  $\eta$  中的第  $i$  个元素, 则  $\eta_i = \log \frac{\phi_i}{\phi_k}$ , 且  $\eta_k = \log \frac{\phi_k}{\phi_k} = 0$

$$\therefore e^{\eta_i} = \frac{\phi_i}{\phi_k}$$

$$\therefore \phi_i = \phi_k e^{\eta_i}$$

$$\therefore \sum_{i=1}^k \phi_i = \sum_{i=1}^k \phi_k e^{\eta_i} = 1$$

$$\therefore \phi_k = \frac{1}{\sum_{i=1}^k e^{\eta_i}} \Rightarrow \phi_i = \phi_k e^{\eta_i} = \frac{e^{\eta_i}}{\sum_{i=1}^k e^{\eta_i}}$$

$$\therefore \frac{\phi_i}{\phi_k} = \frac{\phi_i}{\phi_k} =$$

$$\therefore p(y=i | x; \theta) = \phi_i = \frac{e^{\eta_i}}{\sum_{i=1}^k e^{\eta_i}} = \frac{e^{\theta_i^T x}}{\sum_{i=1}^k e^{\theta_i^T x}}$$

$$h_{\theta}(x) = E [T(y) | x; \theta]$$

$$= E \left[ \begin{matrix} I(y=1) \\ \vdots \\ I(y=k-1) \end{matrix} \middle| x; \theta \right]$$

$$= \begin{bmatrix} \phi_1 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

$$= \begin{bmatrix} e^{\theta_1^T x} / \sum_{i=1}^k e^{\theta_i^T x} \\ \vdots \\ e^{\theta_{k-1}^T x} / \sum_{i=1}^k e^{\theta_i^T x} \end{bmatrix}$$

$\Rightarrow$  假设 (2) 对  $T(y) | x; \theta$  为均值

$\Rightarrow T(y) | x; \theta$  服从多项式分布.

$\hookrightarrow$  多项式分布属于指数家族.

且假设 (3)  $\eta_i = \theta_i^T x$

$\Rightarrow$  softmax regression.