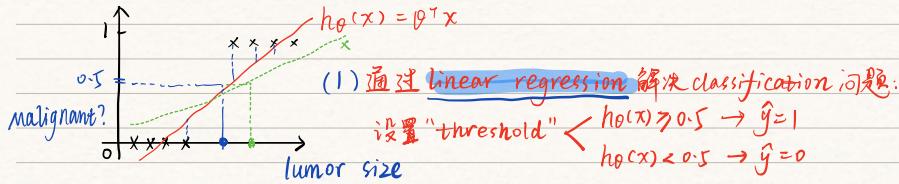


Classification ~ 逻辑回归 (Logistic Regression Model)

二分类 $y \in \{0, 1\}$ Binary classification problem



缺点: ① 线性回归结果很容易受影响 → 当不能线性可分时, 自然不适用线性回归
 ② $h_\theta(x)$ 可能 $> 1, < 0$ while $y=0$ or 1

(2) 通过 Logistic Regression 解决 (注: 名称虽为“回归”, 但实际上解决的是标签值为离散型的分类问题)

可以满足 $h_\theta(x) \in [0, 1]$.

假设 $h_\theta(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$

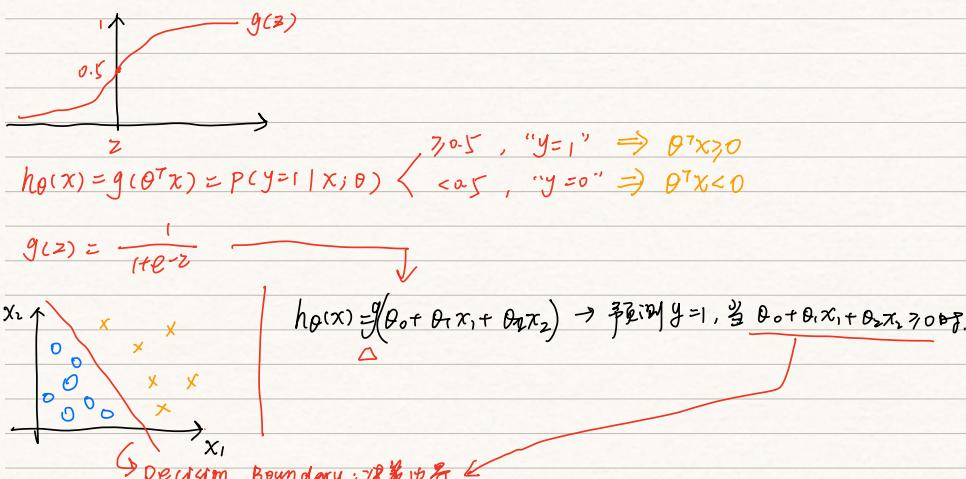
$g(z) : g(z) = \frac{1}{1 + e^{-z}}$

Sigmoid function / logistic function

$y = 0 \text{ or } 1$
 $P(y=1|x;\theta) + P(y=0|x;\theta) = 1$

输出介于 $[0, 1]$ 之间, 对此输出的解释是: 代表了患病的概率: $h_\theta(x) = P(y=1|x;\theta)$
 $h_\theta(x)$ give us the probability that our output, e.g.: $h_\theta(x) = 0.7$ 表示病人有 70% 的可能得病.

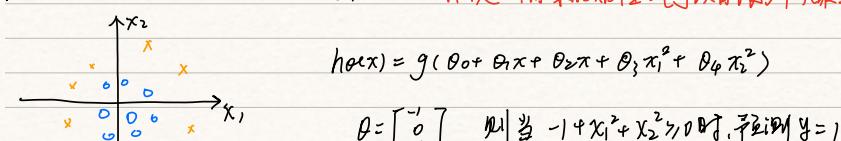
Decision Boundary



Non-linear decision boundaries:

(是假设 hypothesis 一个属性, 决定了参数即决定边界)

而不是训练集的属性! (可以有很多个决策边界).



x

x

0

decision Boundary

可能有更复杂的决策边界。

Cost Function

Loss function: 在整个训练集上的损失

对于线性回归，常用 $\text{MSE} = \frac{1}{m} \sum (y_i - \hat{y}_i)^2$ 作为 Cost function

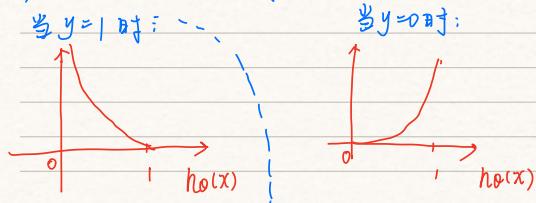
但如果 Logistic 回归也用这个作为 cost function，那么导致其非凸！(图示)

因此，我们采用另一种 cost function

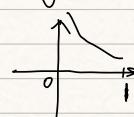
$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)), & \text{if } y=1 \\ -\log(1-h_{\theta}(x)), & \text{if } y=0 \end{cases}$$

$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$

其图像为：



若预测为 1, 则 cost=0
若预测为 0, 则 cost=∞



采用了这样的 cost function 后， $J(\theta)$ 就是凸的了（可以用 GD 得到最优解）
(没有 MSE 这样的名称)。

PS:

另一个角度解释为何要用这样的 cost function: MLE 的角度

$$P(y=1 | x; \theta) = h_{\theta}(x)$$

$$P(y=0 | x; \theta) = 1 - h_{\theta}(x)$$

$$\text{则: } P(y | x; \theta) = h_{\theta}(x)^y \cdot [1 - h_{\theta}(x)]^{1-y} \quad \leftarrow \begin{array}{l} y=0 \\ y=1 \end{array} \text{ 分别是各自的式子。通过指教巧妙地把两项函数写成一行。}$$

$$\text{显然: } L(\theta) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}; \theta)^{y^{(i)}} \cdot P(y^{(i)} = 0 | x^{(i)}; \theta)^{1-y^{(i)}}$$

$$\text{反对数: } -\ln L(\theta) = \sum_{i=1}^n \left[y^{(i)} \ln h_{\theta}(x^{(i)}) + (1-y^{(i)}) \ln [1 - h_{\theta}(x^{(i)})] \right] = J(\theta)$$

$$\text{即: } J(h_{\theta}(x), y | \theta) = \begin{cases} -\ln(h_{\theta}(x)), & y=1 \\ -\ln(1-h_{\theta}(x)), & y=0 \end{cases} \text{ 一致。}$$

化简 cost function 与 Gradient Descent.

回顾 Logistic 回归:

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta x}}$$

$$\text{Cost}(x, y) = -y \ln(h_{\theta}(x)) - (1-y) \ln(1-h_{\theta}(x)) \quad \leftarrow \begin{array}{l} y=0 \text{ 两个角度来看, 与上面是} \\ y=1 \text{ 一样的} \end{array}$$

$$\text{整个数据集上: } J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{entire cost function} = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))]$$

向量形式:

$$h = g(X\theta)$$

$$J(\theta) = \frac{1}{m} (-y^T \log(h) - (1-y^T) \log(1-h))$$

目标: $\min J(\theta)$

Repeat: {

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta) \quad | \text{ simultaneously update all } \theta_j$$

Repeat: { $\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$ 和线性回归的迭代式一样 }

$$\text{向量形式: } \theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \bar{y})$$

$$(g(\boxed{\square}) - \boxed{\square})$$

其他的优化方式 (除了梯度下降外): Conjugate gradient, BFGS, L-BFGS

\backslash Classification = ~ 多分类 one vs all (Multi-class Classification)

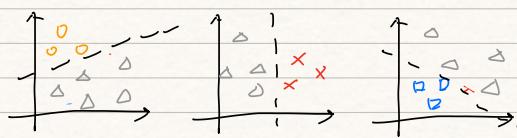
当 $y \in \{0, 1\}$ 二分类

$y \in \{0, 1, 2, \dots, n\}$ 多分类

我们可以将问题视作 $n+1$ 个二分类问题:

训练 $(n+1)$ 个二分类器

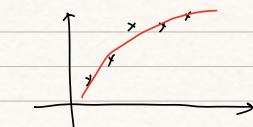
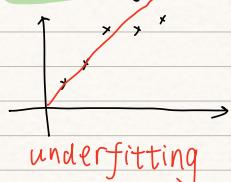
$(n+1)$ 条



$$h_\theta^{(0)}(x) = P\{y=0|x; \theta\} \quad h_\theta^{(1)}(x) = P\{y=1|x; \theta\} \quad h_\theta^{(2)}(x) = P\{y=2|x; \theta\}$$

$$\text{Prediction} = \max_i \{h_\theta^{(i)}(x)\}$$

过拟合与欠拟合:



1. 欠拟合 / : 模型假设对现有数据拟合得很差, 通常是因为模型太简单, 利用了较少的特征.

/ 高偏差 high bias

过拟合 / : 模型对现有数据拟合得很好, 但泛化能力差, 不能很好地预测新数据; 复杂; 多

| 高方差: high variance

如何解决过拟合? (两条原则): 目标: (降低模型复杂程度?)

(1) 减少特征数量

- 手动选择特征
- 使用模型选择算法

(2) 正则化

- 保留所有特征, 但限制参数大小 θ ;

Regularization 的 Cost Function

使用正则化, 在不减少特征数且保留原来的数据的情况下, 通过减小项的权重, 来克服过拟合 (项的权重很小 $= 0.0001$, 视为没有, 则模型简单). 在 cost function 里加上这些权重, 且提高其比重, 则可以实现.

1. 对某些项进行正则:

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h^{(i)}(x) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2 \quad h(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

2. 对所有项进行正则:

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h^{(i)}(x) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \quad \begin{array}{l} \text{正则化参数} \rightarrow \text{若太小, 则欠拟合} \\ \text{平方项方便求梯度} \rightarrow \text{若太大, 仍然过拟合.} \end{array}$$

(ps: θ_0 项不进行正则)

Regularized Linear Regression

应用正则化到线性回归:

(1) 方式一: 梯度下降。因为修改了损失函数, 所以在用梯度下降法更新参数时会对参数加以限制。
特别地, 我们分离出 θ_0 , 因为我们不对 θ_0 进行惩罚。

Repeat {

$$\theta_0 := \theta_0 - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \quad \theta_j \text{ 是当前更新的参数.}$$

$$\theta_j := \theta_j - \alpha \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j \in \{1, 2, \dots, n\}$$

$$\theta_j := (1 - \alpha \frac{\lambda}{m}) \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

合并

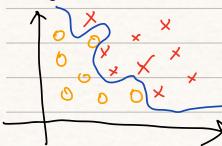
$\alpha > 0, m > 0, \theta_j | (1 - \alpha \frac{\lambda}{m}) < 1$, 起到缩小 θ_j 的目的。

(2) 方式二: 通过 Normal Equation:

$$\Theta = (X^T X + \lambda \cdot L)^{-1} X^T y$$

$$L = \begin{bmatrix} 0 & & \\ & 1 & \\ & & I_m \end{bmatrix} \text{。注意若 } X^T X \text{ 不可逆时, 加上 } L \text{ 可使其变为可逆。}$$

Regularized Logistic Function



逻辑回归的 Cost Function:

$$J(\theta) = -\left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right]$$

$$\text{添加正则项: } + \frac{1}{2m} \sum_{j=1}^n \theta_j^2$$

Gradient descent:

$$\text{Repeat } f \quad \theta_0 := \theta_0 + \sum_{j=1}^m$$

附注: Logistic 回归 梯度下降推导

预备知识: sigmoid 函数求导:

$$g(z) = \frac{1}{1+e^{-z}} \quad ; \quad g'(z) = \frac{-e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{(1+e^{-z}) - 1}{1+e^{-z}}$$

求梯度推导:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \cdot \frac{1}{h_\theta(x^{(i)})} \cdot \frac{\partial h_\theta(x^{(i)})}{\partial \theta_j} + (1-y^{(i)}) \cdot \frac{(1)}{1-h_\theta(x^{(i)})} \cdot \frac{\partial h_\theta(x^{(i)})}{\partial \theta_j} \right] \\ &= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \cdot \frac{1}{h_\theta(x^{(i)})} - (1-y^{(i)}) \cdot \frac{1}{1-h_\theta(x^{(i)})} \right] \frac{\partial h_\theta(x^{(i)})}{\partial \theta_j} \end{aligned}$$

$$\frac{\partial h_\theta(x^{(i)})}{\partial \theta_j} = h_\theta(x^{(i)}) \cdot [1-h_\theta(x^{(i)})] \cdot [\theta^T x^{(i)}]'$$

其中 $\theta^T x^{(i)} = \theta_0 + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} \dots \therefore [\theta^T x^{(i)}] \text{ 关于 } \theta_j \text{ 的系数是 } x_j^{(i)}, \text{ 为一常数。}$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \cdot \frac{1}{h_\theta(x^{(i)})} - (1-y^{(i)}) \cdot \frac{1}{1-h_\theta(x^{(i)})} \right] \cdot h_\theta(x^{(i)}) \cdot [1-h_\theta(x^{(i)})] \cdot x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} [1-h_\theta(x^{(i)})] - (1-y^{(i)}) \cdot h_\theta(x^{(i)}) \right] \cdot x_j^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - h_\theta(x^{(i)})] x_j^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

梯度下降法更新参数: $\theta_j := \theta_j - \alpha \cdot \frac{\partial h_\theta(x^{(i)})}{\partial \theta_j}$

$$\theta_j = \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)}$$