

Logistic Regression

基础概念

Logistic Regression 虽然被称为回归，但实际上是分类模型，并常用于二分类。**Logistic Regression** 因其简单、可并行化、可解释性强深受工业界喜爱。

Logistic 回归的本质是：假设数据服从这个分布，然后使用极大似然估计做参数的估计。

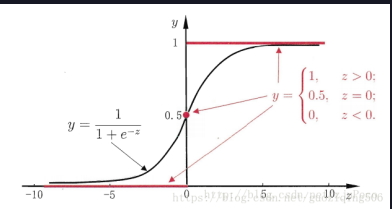
阶跃函数Sigmoid Function

直接依靠拟合曲线的函数值是不能得到类标号的，还需要一种理想的“阶跃函数”，将函数值按照正负性分别映射为**0**、**1**类标号。

Sigmoid函数：

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

Sigmoid函数的图像如下图所示，当**z>0**时，**Sigmoid**函数大于**0.5**；当**z<0**时，**Sigmoid**函数小于**0.5**。所以，我们可以将拟合曲线的函数值带入**Sigmoid**函数，观察**φ(z)**与**0.5**的大小确定其类标号。

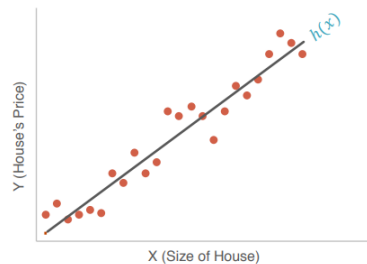


Sigmoid函数还有一个好处，那就是因为其取值在**0**、**1**之间。所以可以看做是测试元组属于类**1**的后验概率，即**p(y=1|X)**。其实这一点从图像也可以看出来：**z**的值越大，表明元组的空间位置距离分类面越远，他就越可能属于类**1**，所以图中**z**越大，函数值也就越接近**1**；同理，**z**越小，表明元组越不可能属于类**1**。

Linear regression vs Logistic regression

Linear regression hypothesis $h_{\theta}(x)$

- Univariate linear regression (Simple linear regression)



$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- Multivariate linear regression

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \sum_{i=0}^n \theta_i x_i$$



$$h_{\theta}(x) = \theta^T x \quad // \text{rewrite the equation in vectorized form}$$

Logistic regression hypothesis $h_{\theta}(x)$

In a binary classification we use a different hypothesis.

We predict the probability that a given example belongs to the "1" class versus the probability that it belongs to the "0" class

- 1 (positive) : e.g. malignant tumor
- 0 (negative) : e.g. non-malignant tumor

We want $0 \leq h_{\theta}(x) \leq 1$

So we need to find a function for our hypothesis so that the output is bounded (0,1). Sigmoid function can help us.

$$h_{\theta}(x) = \theta^T x$$



$$h_{\theta}(x) = \sigma(\theta^T x)$$



$$h_{\theta}(x) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$$

Use sigmoid function to "transform" the linear regression equation $h_{\theta}(x) = \theta^T x$

Sigmoid/Logistic function: $\sigma(z) = \frac{1}{1 + \exp(-z)}$

代价函数cost function

$$p(y|X; W) = \phi(z)^y (1 - \phi(z))^{1-y}$$

在参数**W**下，元组类标号为**y**的后验概率

上述公式表达的含义是在参数**W**下，元组类标号为**y**的后验概率。假设现在已经得到了一个抽样样本，那么联合概率的大小就可以反映模型的代价：联合概率越大，说明模型的学习结果与真实情况越接近；联合概率越小，说明模型的学习结果与真实情况越背离。而对于这个联合概率，我们可以通过计算参数的最大似然估计的那一套方法来确定使得联合概率最大的参数**W**，此时的**W**就是我们要选的最佳参数，它使得联合概率最大（即代价函数最小）。

Minimize the cost function $J(\theta)$ to find the best choice of θ

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

As the cost function is convex, so Gradient Descent can be used to find the global minimum (if the learning rate is not too large and you wait long enough).

Partial derivative of $J(\theta)$

The partial derivative of $J(\theta)$ as given on the left side with respect to θ_j is

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

Gradient Descent for logistic regression

As we want to minimize the cost function $J(\theta)$, so

repeat until convergence {

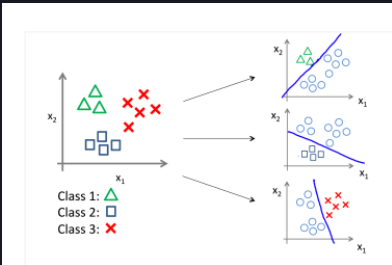
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

// (Simultaneously update θ_0 for every $j=0, 1, \dots, n$)

}

多分类问题

用one vs all



LR优点：

- 1.直接对分类的可能性建模，无需事先假设数据分布，避免了假设分布不准确带来的问题
- 2.不仅预测出类别，还可得到近似概率预测
- 3.对率函数是任意阶可导凸函数，有很好得数学性质，很多数值优化算法可直接用于求取最优解
- 4.容易使用和解释，计算代价低
- 5.LR对时间和内存需求上相当高效
- 6.可应用于分布式数据，并且还有在线算法实现，用较小资源处理较大数据
- 7.对数据中小噪声鲁棒性很好，并且不会受到轻微多重共线性影响
- 8.因为结果是概率，可用作排序模型

LR缺点：

- 1.容易欠拟合，分类精度不高
- 2.特征有缺失或特征空间很大时效果不好