

Machine Learning - Week 3

赵燕

目录

1	Classification and Representation	2
1.1	Classification	2
1.2	Hypothesis Representation	3
1.3	Decision Boundary	6
2	Logistic Regression Model	8
2.1	Cost Function	8
2.2	Simplified Cost Function and Gradient Descent	11
2.3	Advanced Optimization	12
3	Multiclass Classification	15
3.1	Multiclass Classification:One-vs-all	15
4	Solving the Problem of Overfitting	15
4.1	The Problem of Overfitting	15
4.2	Cost Function	15
4.3	Regularized Linear Regression	15
4.4	Regularized Logistic Regression	15

1 Classification and Representation

1.1 Classification

To attempt classification, one method is to use linear regression and map all predictions greater than 0.5 as a 1 and all less than 0.5 as a 0. However, this method doesn't work well because classification is not actually a linear function.

在分类问题中，需要预测的变量 y 是离散的值，引出要学习的逻辑回归算法（Logistic Regression），这是目前最流行使用的一种学习算法。

分类问题举例：

- (1) 判断一封电子邮件是否是垃圾邮件；
- (2) 判断一次金融交易是否是欺诈；
- (3) 判断肿瘤是 良性还是恶性；

Classification

- Email: Spam / Not Spam?
- Online Transactions: Fraudulent (Yes / No)?
- Tumor: Malignant / Benign ?

图 1: 分类问题举例

从二元的问题开始讨论：

将因变量（dependent variable）可能属于两个类分别称为负向类（negative class）和正向类（positive class），则因变量 $y \in \{0, 1\}$ ，其中0表示负向类，1表示正向类。

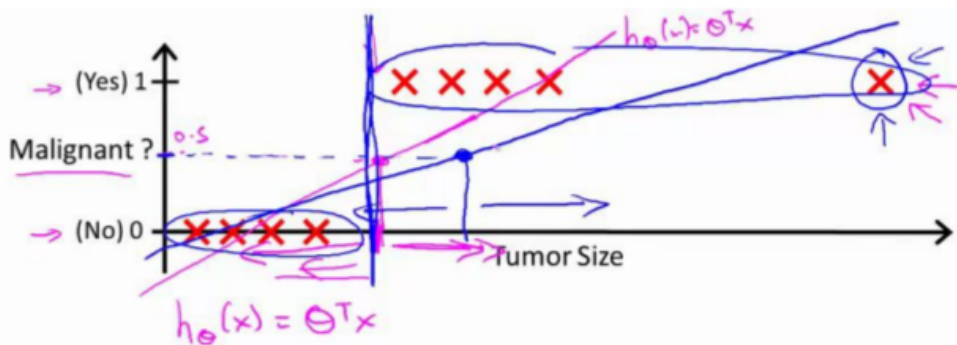


图 2: 图示

The classification problem is just like the regression problem, except that the values we now want to predict take on only a small number of discrete values. For now, we will focus on the binary classification problem in which y can take on only two values, 0 and 1. (Most of what we say here will also generalize to the multiple-class case.) For instance, if we are trying to build a spam classifier for email, then $x(i)$ may be some features of a piece of email, and y may be 1 if it is a piece of spam mail, and 0 otherwise. Hence, $y \in \{0, 1\}$. 0 is also called the negative class, and 1 the positive class, and they are sometimes also denoted by the symbols “-” and “+.” Given $x(i)$, the corresponding $y(i)$ is also called the label for the training example.

如果我们要用线性回归算法来解决一个分类问题,对于分类, y 取值为 0 或者 1,但如果你使用的是线性回归,那么假设函数的输出值可能远大于 1,或者远小于 0,即使所有训练样本的标签 y 都等于 0 或 1。尽管我们知道标签应该取值 0 或者 1,但是如果算法得到的值远大于 1 或者远小于 0 的话,就会感觉很奇怪。所以我们在接下来的要研究的算法就叫做逻辑回归算法,这个算法的性质是:它的输出值永远在 0 到 1 之间。

Classification: $y = 0 \text{ or } 1$

$h_{\theta}(x)$ can be > 1 or < 0

Logistic Regression: $0 \leq h_{\theta}(x) \leq 1$

图 3: 逻辑回归算法

逻辑回归算法是分类算法，我们将它作为分类算法使用，有时候可能因为这个算法的名字中出现了“回归”使你感到困惑，但逻辑回归算法实际上是一种分类算法，它适用于标签 y 取值离散的情况下，如：1，0，0，1。

1.2 Hypothesis Representation

假设函数表达式：

我们希望分类器的输出在0和1之间，因此，要想出一个满足某个性质的假设函数，这个性质是它的预测值要在0和1之间。

回顾肿瘤分类问题：可以用线性回归的方法求出一条适合数据的一条直线：

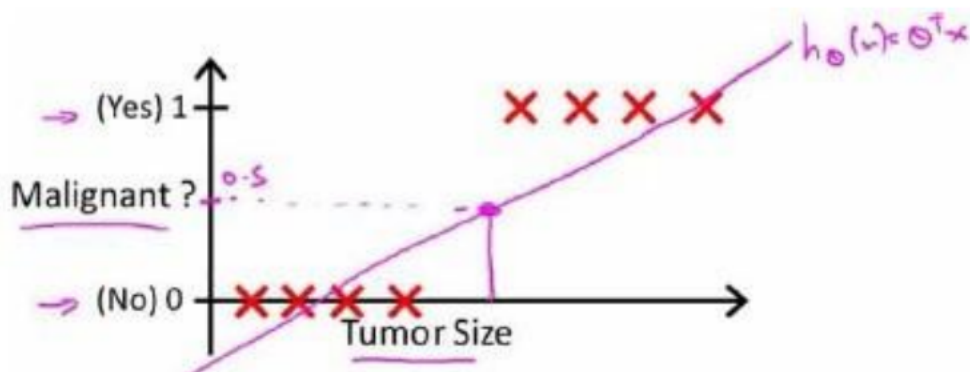


图 4: 肿瘤分类线性回归

根据线性回归我们只能预测连续的值，然而对于分类问题，我们要输出0或1，我们可以预测：（1）当 h_θ 大于等于 0.5 时,预测 $y=1$ ；（2）当 h_θ 小于 0.5 时,预测 $y=0$ 对于上图所示的数据,这样的线性模型似乎能很好地完成分类任务。假使我们又观测到一个非常大尺寸的恶性肿瘤,将其作为实例加入到我们的训练集中来,这将使得我们获得一条新的直线。

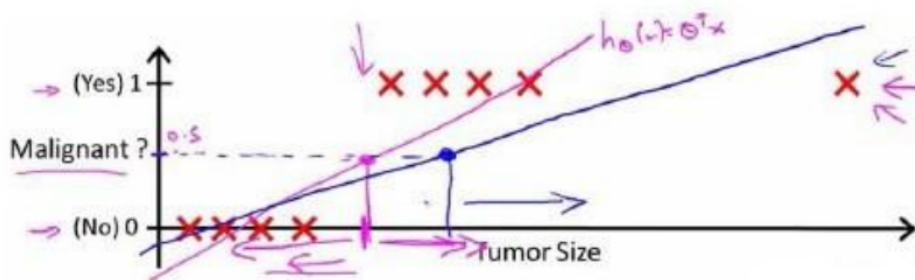


图 5: 肿瘤问题线性回归超范围

这时，再使用 0.5 作为阈值来预测肿瘤是良性还是恶性便不合适了。可以看出线性回归模型，因为其预测的值可以超越 $[0,1]$ 的范围，并不适合解决这样的问题。

We could approach the classification problem ignoring the fact that y is discrete-valued, and use our old linear regression algorithm to try to predict y given x . However, it is easy to construct examples where this method performs very poorly. Intuitively, it also doesn't make sense for $h_\theta(x)$ to take values larger than 1 or smaller than 0 when we know that $y \in [0, 1]$. To fix this, let's change the form for our hypotheses $h_\theta(x)$ to satisfy $0 \leq h_\theta(x) \leq 1$. This is accomplished by plugging $\theta^T x$ into the Logistic Function.

引入一个新的模型，逻辑回归，该模型的输出变量范围始终在 0 和 1 之间。

逻辑回归模型的假设是：

$$h_\theta(x) = g(\theta^T x) \quad (1)$$

其中：

x :代表特征向量；

g :代表逻辑函数(logistic function)，是一个常用的逻辑函数为 S 形函数(Sigmoid function)，公式为：

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

该函数的图像为:

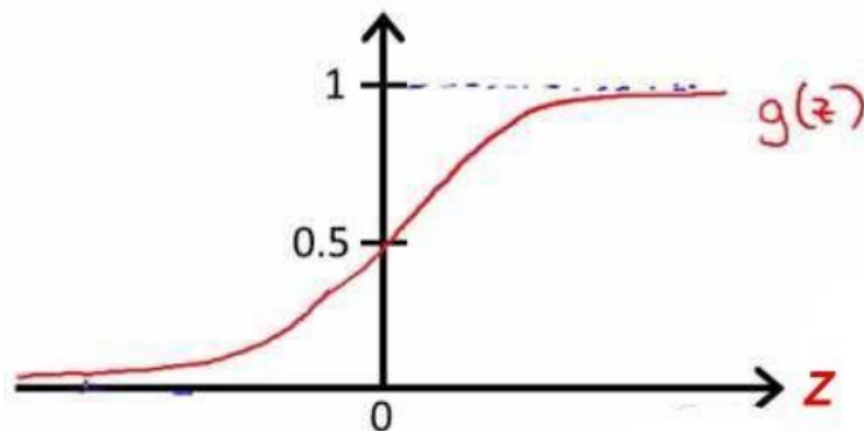


图 6: 逻辑函数 (S函数) 图像

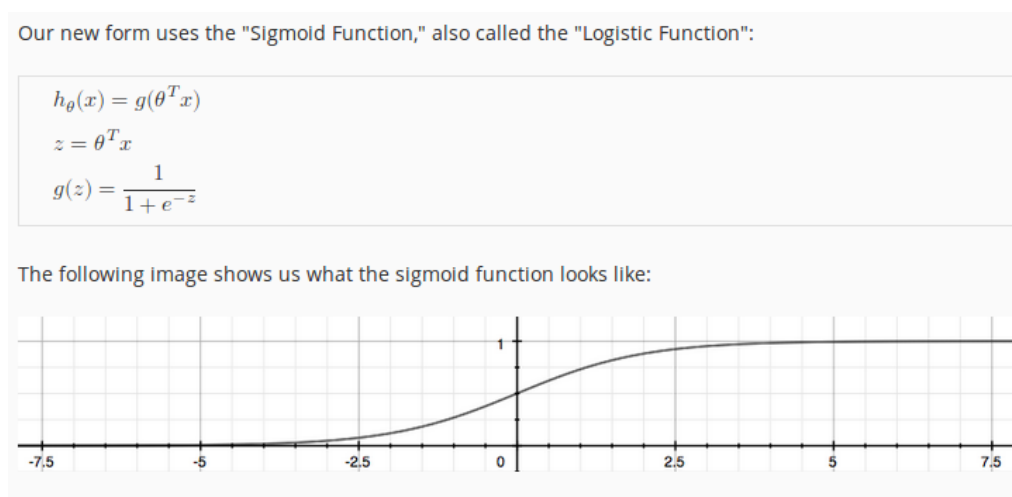


图 7: 逻辑函数表达式及图像

The function $g(z)$, shown here, maps any real number to the $(0, 1)$ interval, making it useful for transforming an arbitrary-valued function into a function better suited for classification.

逻辑回归模型的假设:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad (3)$$

$h_{\theta}(x)$ 的作用是,对于给定的输入变量,根据选择的参数计算输出变量=1 的可能性(estimated probability)即

$$h_{\theta}(x) = P(y = 1|x; \theta) \quad (4)$$

$h_{\theta}(x)$ will give us the probability that our output is 1.

$$h_{\theta}(x) = P(y = 1|x; \theta) = 1 - P(y = 0|x; \theta) \quad (5)$$

$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1 \quad (6)$$

例如，如果对于给定的 x ，通过已经确定的参数计算得出 $h_{\theta}(x) = 0.7$ ，则表示有 70% 的几率 y 为正向类，相应地 y 为负向类的几率为 $1-0.7=0.3$ 。

1.3 Decision Boundary

决策边界 (Decision Boundary) 可以更好的帮助我们理解假设函数在计算什么

首先回顾一下逻辑回归中假设函数的表达式：

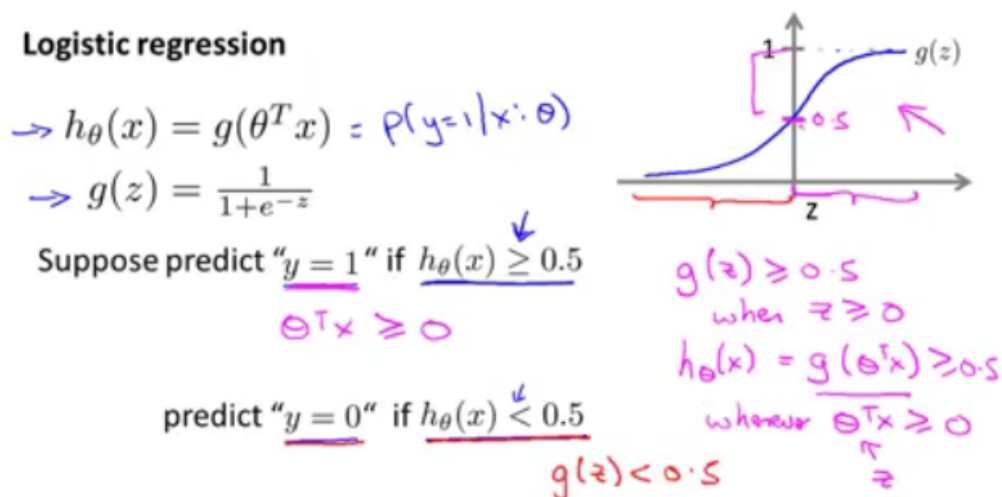


图 8: 逻辑函数回顾

根据上面绘制的S形函数的图像，我们可以知道：

当 $z=0$ 时， $g(z)=0.5$;

当 $z>0$ 时， $g(z)>0.5$;

当 $z<0$ 时， $g(z)<0.5$;

又因为

$$z = \theta^T x \quad (7)$$

所以：

$\theta^T x \geq 0$ 时，预测 $y=1$;

$\theta^T x < 0$ 时，预测 $y=0$;

假设有一个模型：

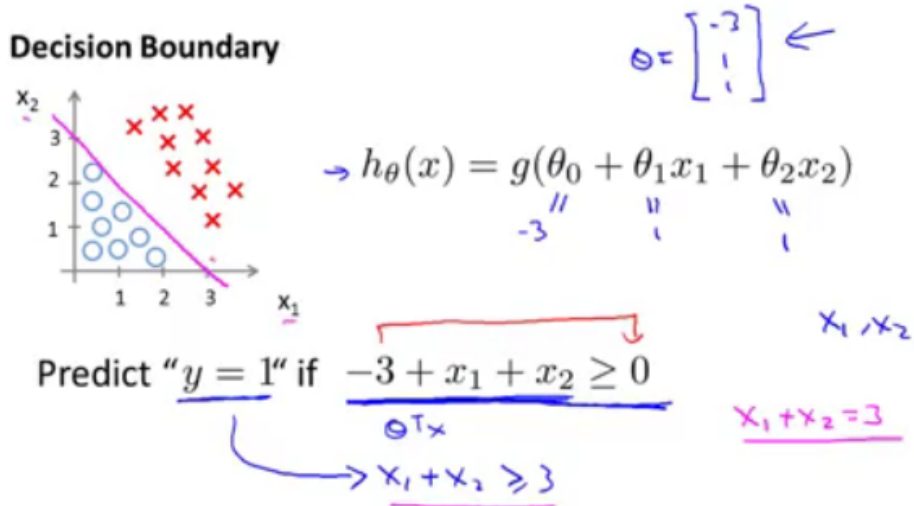


图 9: 线性决策边界

绘制直线 $x_1 + x_2 = 3$, 这条直线便是我们的模型的分界线, 将预测 $y=1$ 的区域和预测 $y=0$ 的区分隔开, 如图所示:

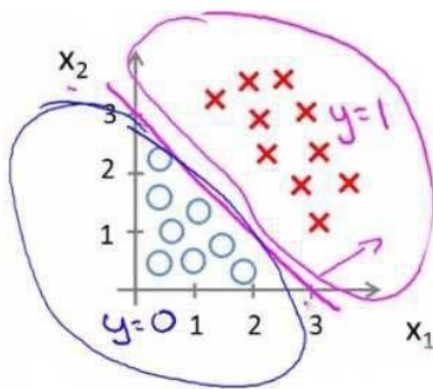


图 10: 线性决策边界图

另一种情况:

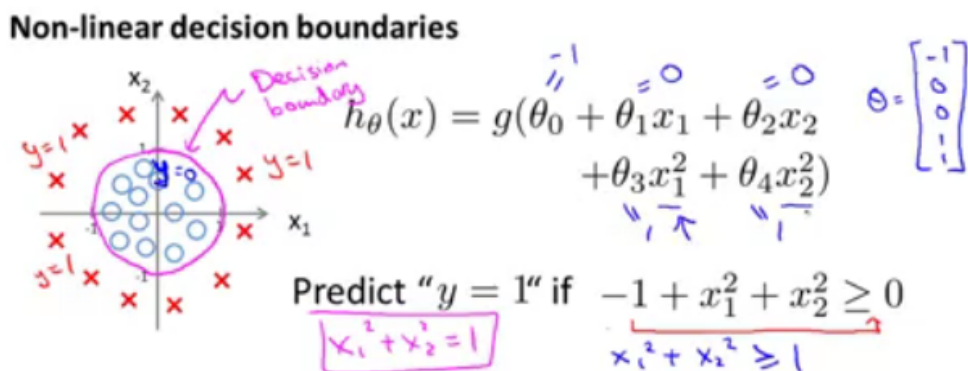


图 11: 非线性决策边界

需要二次方特征，得到的决策边界（判定边界）恰好是在原点且半径为1的圆形。

我们可以用非常复杂的模型来适应非常复杂的判定边界。例如：

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \dots) \quad (8)$$

我们可能会得到比椭圆等更复杂的判定边界：

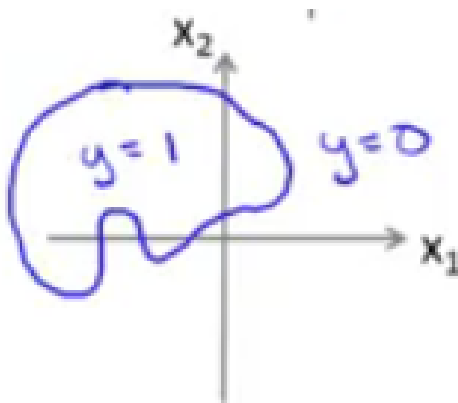


图 12: 更复杂的情况

2 Logistic Regression Model

2.1 Cost Function

介绍如何拟合逻辑回归模型的参数 θ ，具体来说，我们要定义来拟合参数的优化目标或者叫代价函数，这便是监督学习问题中的逻辑回归模型的拟合问题。

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m examples $x \in \begin{bmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{bmatrix} \quad x_0 = 1, y \in \{0, 1\}$

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

How to choose parameters θ ?

图 13: 如何选择参数 θ

Cost function

→ Linear regression: $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

logistic

→ $\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$

图 14: 线性代价函数讨论

Cost Function:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (9)$$

对于线性回归模型，我们定义的代价函数是所有模型误差的平方和。理论上来说，我们也可以对逻辑回归模型沿用这个定义，但是问题在于，当我们将 $h_{\theta} = \frac{1}{1+e^{-\theta^T x}}$ 带入到这样定义的代价函数中去，我们得到的代价函数将是一个非凸函数（non-convex function）。

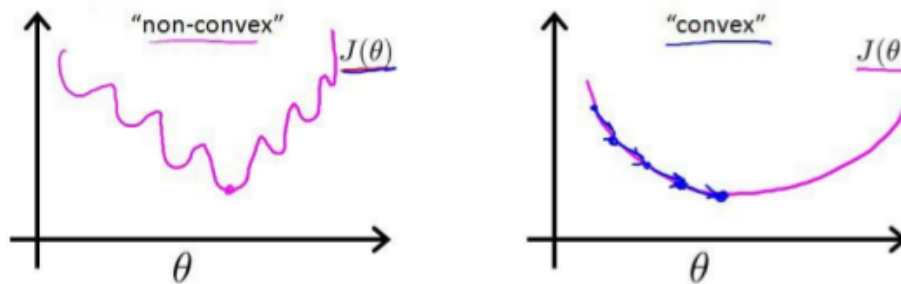


图 15: 非凸函数与凸函数

这就意味着我们的代价函数有许多局部最小值，这将影响梯度下降算法寻找全局最小值。

线性回归代价函数:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (10)$$

重新定义逻辑回归的代价函数：（这种方式可以保证 $J(\theta)$ 对于逻辑回归是凸函数）

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \quad (11)$$

其中:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \quad (12)$$

$h_{\theta}(x)$ 与 $Cost(h_{\theta}(x), y)$ 之间的关系如下图所示:

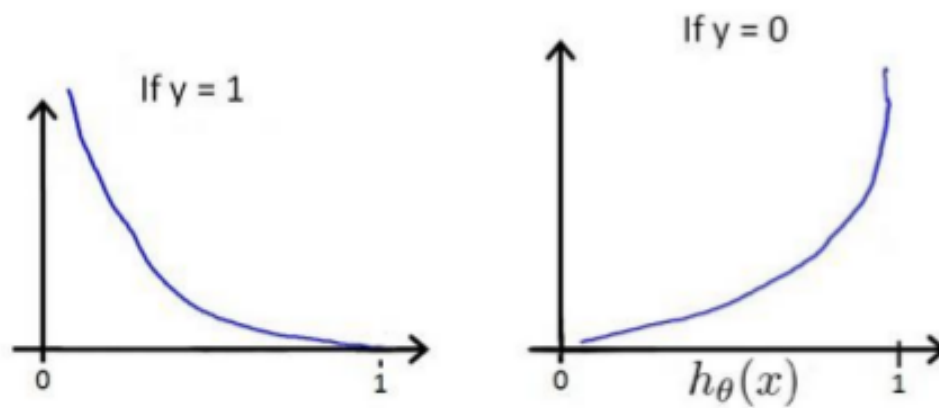


图 16: $h_{\theta}(x)$ 与 $Cost(h_{\theta}(x), y)$ 的关系

这样构建的函数的特点是:

- (1) 当实际的 $y=1$ 且 h_{θ} 也为1时的误差为0;
- (2) 当 $y=1$ 但不为1时误差随着的变小而变大;
- (3) 当实际的 $y=0$ 且 h_{θ} 也为0时代价为0;
- (4) 当 $y=0$ 但 h_{θ} 不为0时误差随着 h_{θ} 的变大而变大;

2.2 Simplified Cost Function and Gradient Descent

Logistic regression cost function

$$\begin{aligned} \rightarrow J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) \\ \rightarrow \text{Cost}(h_{\theta}(x), y) &= \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases} \\ \text{Note: } y &= 0 \text{ or } 1 \text{ always} \\ \rightarrow \text{Cost}(h_{\theta}(x), y) &= -y \log(h_{\theta}(x)) - (1-y) \log(1 - h_{\theta}(x)) \\ \text{If } y=1: \text{Cost}(h_{\theta}(x), y) &= -\log(h_{\theta}(x)) \leftarrow \\ \text{If } y=0: \text{Cost}(h_{\theta}(x), y) &= -\log(1 - h_{\theta}(x)) \end{aligned}$$

图 17: $h_{\theta}(x)$ 与 $\text{Cost}(h_{\theta}(x), y)$ 的关系

将构建的构造函数简化如下:

$$\text{Cost}(h_{\theta}(x)) = -y \log(h_{\theta}(x)) - (1-y) \log(1 - h_{\theta}(x)) \quad (13)$$

带入到代价函数中得到:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] \quad (14)$$

A vectorized implementation is:

$$h = g(X\theta) \quad (15)$$

$$J(\theta) = -\frac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h)) \quad (16)$$

Logistic regression cost function

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

To fit parameters θ :

$$\min_{\theta} J(\theta) \quad \text{Get } \theta$$

To make a prediction given new x :

$$\text{Output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}} \quad p(y=1 | x; \theta)$$

图 18: 梯度下降算法的导出

在得到这样一个代价函数以后，我们便可以用梯度下降算法来求得能使代价函数最小的参数 θ ，算法为：

repeat{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (17)$$

$$(simultaneously update all) \quad (18)$$

}

对 θ_j 求导后得出：

repeat{

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (19)$$

$$(simultaneously update all) \quad (20)$$

}

A vectorized implementation is:

$$\theta_j := \theta_j - \frac{\alpha}{m} X^T (g(X\theta) - \vec{y}) \quad (21)$$

在这个视频中，我们定义了单训练样本的代价函数，凸性分析的内容是超出我们这门课的范围的，但是可以证明我们所选的代价值函数会给我们一个凸优化问题。代价函数 $J(\theta)$ 会是一个凸函数，并没有局部最优值。

注意：

线性回归假设函数：

$$h_{\theta}(x) = \theta^T x \quad (22)$$

逻辑回归假设函数：

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad (23)$$

虽然得到的梯度下降算法表面看上去与线性回归的梯度下降算法一样，但是这里的逻辑回归假设函数与线性回归 $h_{\theta} = \theta^T x$ 的不同，所以实际上是不一样的。另外，在运行梯度下降算法之前，进行特征缩放依旧是非常有必要的，特征缩放也适用与逻辑回归。

逻辑回归是一种非常强大的，甚至是世界上使用最广泛的一种分类算法。

2.3 Advanced Optimization

利用高级优化算法，使通过梯度下降进行逻辑回归的速度大大提高，而且这也将使算法更加适合解决大型的机器学习问题，比如我们有数目庞大的特征量。

换个角度来看什么是梯度下降，我们有一个代价函数 $J(\theta)$ ，我们想要使其最小化，我们需要做的是编写代码，当输入参数 θ 时，它们会计算出两种东西： $J(\theta)$ 以及 J 关于 θ 的偏导数项。

Optimization algorithm

Cost function $J(\theta)$. Want $\min_{\theta} J(\theta)$.

Given θ , we have code that can compute

- $J(\theta)$
- $\frac{\partial}{\partial \theta_j} J(\theta)$ (for $j = 0, 1, \dots, n$)

Gradient descent:

Repeat {

- $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$

}

图 19: 梯度下降执行过程

假设我们已经完成了可以实现这两件事的代码，那么梯度下降所做的就是反复执行这些更新。

另一种考虑梯度下降的思路是：我们需要写出代码来计算 $J(\theta)$ 和这些偏导数，然后把它们插入到梯度下降中，然后它就可以为我们最小化这个函数。

对于梯度下降，不需要编写代码去计算代价函数 $J(\theta)$ ，只需要计算导数项，但是如果你希望代码还要能够监控这些 $J(\theta)$ 的收敛性，那就需要我们自己去编写代码来计算代价函数 $J(\theta)$ 和偏导数项 $\frac{\partial}{\partial \theta_j} J(\theta)$ ，所以在写完能够计算这两者的代码之后，我们就可以使用梯度下降。

共轭梯度法，BFGS（变尺度法）和L-BFGS（限制变尺度法）就是其中一些更高级的优化算法，它们需要有一种方法来最小化代价函数 $J(\theta)$ 。

三种方法的优缺点：

Optimization algorithm

Given θ , we have code that can compute

$$\begin{aligned} & - J(\theta) \leftarrow \\ & - \frac{\partial}{\partial \theta_j} J(\theta) \leftarrow \quad (\text{for } j = 0, 1, \dots, n) \end{aligned}$$

Optimization algorithms:

- - Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:

- No need to manually pick α
- Often faster than gradient descent.

Disadvantages:

- More complex

图 20: 高级优化算法优缺点

实际上,我不会建议你们编写自己的代码来计算数据的平方根,或者计算逆矩阵,因为对于这些算法,我还是会建议你直接使用一个软件库,比如说,要求一个平方根,我们所能做的就是调用一些别人已经写好用来计算数字平方根的函数。幸运的是现在我们有 Octave 和与它密切相关的 MATLAB 语言可以使用。

Octave 有一个非常理想的库用于实现这些先进的优化算法,所以,如果你直接调用它自带的库,你就能得到不错的结果。我必须指出这些算法实现得好或不好是有区别的,因此,如果你正在你的机器学习程序中使用一种不同的语言,比如如果你正在使用 C、C++、Java 等等,你可能会想尝试一些不同的库,以确保你找到一个能很好实现这些算法的库。因为在 L-BFGS 或者等高线梯度的实现上,表现得好与不太好是有差别的,因此现在让我们来说明:

如何使用这些算法:

Example: $\min_{\theta} J(\theta)$

$$\begin{aligned} \rightarrow \theta &= \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_1=5, \theta_2=5. \\ \rightarrow J(\theta) &= (\theta_1 - 5)^2 + (\theta_2 - 5)^2 \\ \rightarrow \frac{\partial}{\partial \theta_1} J(\theta) &= 2(\theta_1 - 5) \\ \rightarrow \frac{\partial}{\partial \theta_2} J(\theta) &= 2(\theta_2 - 5) \end{aligned}$$

图 21: 举例

运行一个像这样的Octave函数:

```
function [jVal, gradient]
    = costFunction(theta)
jVal = (theta(1)-5)^2 + ...
      (theta(2)-5)^2;
gradient = zeros(2,1);
gradient(1) = 2*(theta(1)-5);
gradient(2) = 2*(theta(2)-5);
```

图 22: Octave函数

这样就计算出这个代价函数,函数返回的第二个值是梯度值,梯度值应该是一个 2×1 的向量,梯度向量的两个元素对应这里的两个偏导数项,运行这个 `costFunction` 函数后,你就可以调用高级的优化函数,这个函数叫 `fminunc`,它表示 Octave 里无约束最小化函数。调用它的方式如下:

```
> options = optimset('GradObj', 'on', 'MaxIter', '100');
> initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] ...
    = fminunc(@costFunction, initialTheta, options);
```

图 23: 调用高级优化函数

你要设置几个 options,这个 options 变量作为一个数据结构可以存储你想要的 options,所以 `GradObj` 和 `On`,这里设置梯度目标参数为打开(on),这意味着你现在确实要给这个算法提供一个梯度,然后设置最大迭代次数,比方说 100,我们给出一个 θ 的猜测初始值,它是一个 2×1 的向量,那么这个命令就调用 `fminunc`,这个 @ 符号表示指向我们刚刚定义的 `costFunction` 函数的指针。如果你调用它,它就会使用众多高级优化算法中的一个,当然你也可以把它当成梯度下降,只不过它能自动选择学习速率 α ,你不需要自己来做。然后它会尝试使用这些高级的优化算法,就像加强版的梯度下降法,为你找到最佳的 θ 值。

总结: 需要学到的主要内容是: 写一个函数,它能返回代价函数值,梯度值,因此要把这个应用到逻辑回归中,甚至线性回归中,你可以把这些优化算法用于线性回归,你需要做的就是输入合适的代码来计算这里的这些东西。

3 Multiclass Classification

3.1 Multiclass Classification:One-vs-all

4 Solving the Problem of Overfitting

4.1 The Problem of Overfitting

4.2 Cost Function

4.3 Regularized Linear Regression

4.4 Regularized Logistic Regression