

Week 03

- Logistic Regression
 - Classification and Representation
 - Logistic Regression Model
 - Multiclass Classification
 - Regularization
 - Solving the Problem of Overfitting
-

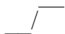
I. Logistic Regression

1. *Classification and Representation*

a. Classification

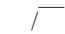
- Example: Email: Spam \rightarrow 0 or 1; Transaction: fraudulent;
- $y = 0$ as negative class;
- $y = 1$ as positive class;
 \Rightarrow binary classification.
- Use linear regression and map prediction based on threshold comparison.
 - Threshold classifier $h_{\theta}(x)$ can be > 0 or < 0
 - $x^{(i)}$ as a feature and y as binary value; \Rightarrow Classification problem **is like** regression problem on a discrete values.
 \Rightarrow This method doesn't work well because classification is not actually a linear function.

b. Hypothesis Representation

- In some cases, y as 0 or 1 is not good. Sometimes we need real values between 0 and 1;
- To fix this as $0 \leq h_{\theta}(x) \leq 1$, add it into a Logisitic function.
 - $h_{\theta}(x) = g(\theta^T x)$ and $z = \theta^T x$ so $g(z) = \frac{1}{1+e^{-z}}$
 - $g(z)$ converge to 0 when z is negative
 - $g(z)$ converge to 1 when z is positive;
 - $g(0) = 0.5$;
 - 

- z good to have values between 0 and 1;
- So $h_{\theta}(x)$ will give the **probability**.
- $h_{\theta}(x) = P(y = 1|x; \theta) = 1 - P(y = 0|x; \theta)$
- $\sum P = 1$

c. Decision

- To get discrete 0 or 1 classification, translate the output of the hypothesis function as follows:
 - $h_{\theta}(x) \geq 0.5 \Rightarrow y = 1$
 - $h_{\theta}(x) < 0.5 \Rightarrow y = 0$
- The logistic function g behaves as that when its input ≥ 0 , its output ≥ 0.5 :
 - $g(z) \geq 0.5$ when $z \geq 0$
- Remember:
 - $z = 0, e^0 = 1 \Rightarrow g(z) = 0.5$
 - $z = +\infty, e^{+\infty} \rightarrow \infty \Rightarrow g(z) = 1$
 - $z = -\infty, e^{-\infty} \rightarrow 0 \Rightarrow g(z) = 0$
 - 
- So if our input to g is θ^T , then that means:
 - $h_{\theta}(x) = g(\theta^T x) \geq 0.5$ when $\theta^T x \geq 0$
- So:
 - $\theta^T x \geq 0 \Rightarrow y = 1$
 - $\theta^T x < 0 \Rightarrow y = 0$
- Example:
 - choose parameters of θ
 - resolve θ in $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)$ so g is not necessarily linear when use x_1^2 or else *power*³;
 - $\theta = [5; 11; 0], y = 1$ if $5 - x_1 + 0x_2 \geq 0$ so it gives $x_1 \leq 5$
 \Rightarrow Decision boundary is a straight vertical line placed on the graph where $x_1 = 5$; at left $y = 1$; at right $y = 0$;
 - note: sigmoid function $g(z)$

2. Logistic regression

a. Cost function

Cannot use the same cost function like in linear regression; Logistic function will cause the output to be wavy, causing many local optima so it will not be a convex function.

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

$$\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x)); \text{ if } y = 0$$

When $y = 1$, we get the following plot for $J(\theta)$ vs $h_{\theta}(x)$:

- $0 \rightarrow +\infty$
- $1 \rightarrow \approx 0$

When $y = 0$, we get the following plot for $J(\theta)$ vs $h_\theta(x)$:

- $0 \rightarrow 0$
- $1 \rightarrow +\infty$

$$\begin{aligned} \text{Cost}(h_\theta(x), y) &= 0h_\theta(x) = y & \text{Cost}(h_\theta(x), y) &\rightarrow \infty \text{ if } y = 0 \text{ and } h_\theta(x) \rightarrow 1 \\ \text{Cost}(h_\theta(x), y) &\rightarrow \infty \text{ if } y = 1 \text{ and } h_\theta(x) \rightarrow 0 \end{aligned}$$

Note that writing the cost function in this way guarantees that $J(\theta)$ is convex for logistic regression.

b. Simplified Cost Function and Gradient Descent

i. Cost function

We can compress our cost function's two conditional cases into one case:

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

When $y = 0$ or $y = 1$ then one of the latter parts will be zero;

Full logistic regression cost function:

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

So minimize $J(\theta)$

A vectorized implementation is:

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

$$J(\theta) = 1/m \cdot (-y^T \log(h) - (1 - y)^T \log(1 - h))$$

ii. Gradient Descent

Remember that the general form of gradient descent is:

$$\text{Minimize } J(\theta) : \text{Repeat } \{ \theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta) \}$$

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

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

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

This algorithm is identical to the one we used in linear regression (where $h_\theta(x) = \theta^T x$). Simultaneously update all values in θ .

A vectorized implementation is: $\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - y)$

c. Advanced optimization

Cost function $J(\theta)$ and want to min $J(\theta)$ so do gradient descent as we have $\frac{d}{d\theta_i} J(\theta)$.

Optimization algorithms:

1. Gradient descent
2. Conjugate gradient
3. BFGS
4. L-BFGS

With the 3 latter ones, no to α , faster and more complex

Example:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}; J(\theta) = (\theta_1 - 5)^2 - (\theta_2 - 5)^2$$

$$\frac{d}{d\theta_1} J(\theta) = 2(\theta_1 - 5); \frac{d}{d\theta_2} J(\theta) = 2(\theta_2 - 5)$$

```
1 function [jVal, gradient] = costFn
2     jVal = (theta(1) - 5)^2 + (theta(2) - 5)^2
3     gradient = zeros(2, 1)
4     gradient(1) = 2*(theta(1) - 5)
5     gradient(2) = 2*(theta(2) - 5)
6
7 options = optimset('GradObj', 'on', 'MaxIter', 100);
8 initialTheta = zeros(2,1);
9 [optTheta, functionVal, exitFlag] = fminunc(@costFn, initialTheta, options);
```

3. Multiclass classification

Instead of $y = \{0,1\}$ we will expand our definition so that $y = \{0,1...n\}$.

Since $y = \{0,1...n\}$, we divide our problem into $n+1$ binary classification problems; in each one, we predict the probability that 'y' is a member of one of our classes.

$$y \in 0, 1...n$$

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

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

...

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

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

We are basically choosing one class and then lumping all the others into a single second class. Repeat this. Apply binary logistic regression to each case, and then use the hypothesis that returned the highest value as our prediction.

II. Regularization

1. The problem of overfitting

Example: Housing prices using linear regression

- $\theta_0 + \theta_1 x$ called underfit or high bias
- $\theta_0 + \theta_1 x + \theta_2 x^2$ called just right
- $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ called overfit or high variance

Plotted pass near the points

⇒ Overfit = many features + the learned hypothesis fit the training set ($J(\theta) \approx 0$) + fail to fit new examples.

$$j(\theta) = 1/2m \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Example: Logistic regression

- $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$ called underfit
- $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \dots)$
- $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \dots + \dots)$

Plotted create a boundary between points of 2 or more classes

Address overfitting:

1. Reduce number of features:

- Manually select which features to keep
- Model selection algorithm

2. Regularization:

- Keep all features but reduce values of θ_j
- Regularization works well when we have a lot of slightly useful features.

2. Cost function

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our function carry by increasing their cost.

Say we wanted to make the following function more quadratic: $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

We want to eliminate the influence of $\theta_3 x_3 + \theta_4 x_4$: ⇒ modify our **cost function**

$$\min_{\theta} 1/2m \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$$

The formula becomes: $J(\theta) = 1/2m [\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2]$

The lambda, is the **regularization parameter**. It determines how much the costs of our theta parameters are inflated.

3. Regularized Linear regression

Gradient Descent

We will modify our gradient descent function to separate out θ_0 from the rest of the parameters because we do not want to penalize θ_0 .

Repeat {

$$\theta_0 := \theta_0 - \alpha/m \sum_{i=1}^m h_{\theta}(x^{(i)} - y^{(i)})x_0^{(i)} \quad \theta_j := \theta_j - \alpha[(1/m \sum_{i=1}^m h_{\theta}(x^{(i)} - y^{(i)})x_j^{(i)}) + \lambda/m.\theta_j] \quad j \in 1..n$$

}

The term $\lambda/m.\theta_j$ performs the regulation. With some manipulation $\theta_j := \theta_j(1 - \alpha\lambda/m) - \alpha \dots$

The first part will be always less than 1. So reducing θ_j by some amount on every update gives us the same second part as before.

Normal Equation

Use the non-iterative normal equation. To add in regularization, the equation is the same as our original, except that we add another term inside the parentheses:

$$\theta = (X^T X + \lambda L)^{-1} X^T y \text{ where } L = \begin{bmatrix} 0 & & \\ & 1 & \\ & & \ddots \end{bmatrix}$$

L is a matrix with 0 at the top left and 1's down the diagonal, with 0's everywhere else.

4. Regularized logistic regression

We can regularize logistic regression in a similar way that we regularize linear regression.

As a result, we can avoid overfitting.

Recall that our cost function for logistic regression was:

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

We can regularize this equation by adding a term to the end:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

The sum in the added term means to explicitly exclude the bias term θ_0 . So the gradient will be like:

Repeat {

$$\theta_0 := \theta_0 - \alpha/m \sum_{i=1}^m h_{\theta}(x^{(i)} - y^{(i)})x_0^{(i)} \quad \theta_j := \theta_j - \alpha[(1/m \sum_{i=1}^m h_{\theta}(x^{(i)} - y^{(i)})x_j^{(i)}) + \lambda/m.\theta_j] \quad j \in 1..n$$

