

# Regularization

Zongheng Yang

May 19, 2012

## 1 Overfitting and Solutions

Overfitting: fits training data well, yet fails to generalize and has poor performance on test data.

Solutions:

1. Reduce # of features (manually select features; use model selection algorithms).
2. Regularization: reduce magnitudes of all  $\theta$ , works well when we have a lot of features. If  $\theta$ 's become small, the cost function becomes "smoother".

## 2 Regularized Linear Regression

In order to minimize the magnitudes of all  $\theta$ , the cost function is added with a new term:

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

The gradient descent:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}; \quad (2)$$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} - \frac{\lambda}{m} \theta_j \right] \quad (3)$$

$$= \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}. \quad (4)$$

$(1 - \alpha \frac{\lambda}{m})$  is usually slightly smaller than 1, so essentially in regularization  $\theta_j$  is adjusted to a smaller value.)

The normal equation: if  $\lambda > 0$ ,

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y. \quad (5)$$

Notice that if  $m(\#examples) \leq n(\#features)$ , the original  $X^T X$  is non-invertible. However, regularization here guarantees invertibility.

### 3 Regularized Logistic Regression

Cost function:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2. \quad (6)$$

The gradient descent has the same form as that of regularized linear regression.

### 4 $\theta_0$ and $\lambda$

$\theta_0$  is usually not regularized.

If  $\lambda$  is too large  $\rightarrow$  all  $\theta$ 's are approximately zeros except  $\theta_0 \rightarrow h_{\theta}(x) = \theta_0$ ; this is called **underfitting**.