

# Regularization

• Small value of parameter  $\Rightarrow$  simpler hypothesis  
 ↓  
 less overfitting

• ex)  $J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i]^2 + \lambda \sum_{j=1}^n \theta_j^2$

regularization term

if  $\lambda$  too large then underfit  
 high penalty  $\rightarrow$  as if model don't contain  $\theta_1, \theta_2, \theta_3, \dots$

$\therefore$  Now the problem is finding the optimal  $\lambda$  that minimizes  $J(\theta)$   
 penalize  $\theta_1, \dots, \theta_n$  but not  $\theta_0$

cost fun  
 ↓  
 $J(\theta)$

(ex) Regularized linear regression:

Gradient descent  $\theta_j = \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x_i) - y_i] x_{ij}$

$1 - \alpha \frac{\lambda}{m} < 1$  ex) 0.99  
 $\theta_j$  gets replaced by  $0.99 \theta_j \Rightarrow$  make  $\theta_j$  smaller  
 $\rightarrow$  regular gradient descent  
 $\rightarrow$  So basically you are just shrinking the parameter by a bit when you do regularization

Normal Equation

$$\theta = \left( X^T X + \lambda \begin{bmatrix} 0 & 1 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & 1 & \ddots \end{bmatrix} \right)^{-1} X^T y \Rightarrow \text{yield global optimum}$$

guaranteed to be invertible matrix

(ex) Regularized Logistic Regression

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots)$$

with regularization.  
 even lots of features  $\theta_1, \dots, \theta_n$  is OK

