

Classification and Representation

- ✓ **Video:** Classification
8 min
- ✓ **Reading:** Classification
2 min
- ✓ **Video:** Hypothesis Representation
7 min
- ✓ **Reading:** Hypothesis Representation
3 min
- ✓ **Video:** Decision Boundary
14 min
- ✓ **Reading:** Decision Boundary
3 min

Logistic Regression Model

- ✓ **Video:** Cost Function
10 min
- ✓ **Reading:** Cost Function
3 min
- ✓ **Video:** Simplified Cost Function and Gradient Descent
10 min
- ✓ **Reading:** Simplified Cost Function and Gradient Descent
3 min
- ✓ **Video:** Advanced Optimization
14 min
- ✓ **Reading:** Advanced Optimization
3 min

Multiclass Classification

- ✓ **Video:** Multiclass Classification: One-vs-all
6 min
- ✓ **Reading:** Multiclass Classification: One-vs-all
3 min

Review

Solving the Problem of Overfitting

- ✓ **Video:** The Problem of Overfitting
9 min
- ✓ **Reading:** The Problem of



Cost Function

Note: [5:18 - There is a typo. It should be $\sum_{j=1}^n \theta_j^2$ instead of $\sum_{i=1}^n \theta_j^2$]

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:

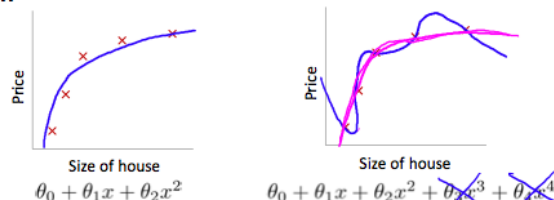
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

We'll want to eliminate the influence of $\theta_3 x^3$ and $\theta_4 x^4$. Without actually getting rid of these features or changing the form of our hypothesis, we can instead modify our **cost function**:

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

We've added two extra terms at the end to inflate the cost of θ_3 and θ_4 . Now, in order for the cost function to get close to zero, we will have to reduce the values of θ_3 and θ_4 to near zero. This will in turn greatly reduce the values of $\theta_3 x^3$ and $\theta_4 x^4$ in our hypothesis function. As a result, we see that the new hypothesis (depicted by the pink curve) looks like a quadratic function but fits the data better due to the extra small terms $\theta_3 x^3$ and $\theta_4 x^4$.

Intuition



Suppose we penalize and make θ_3, θ_4 really small.

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

$\theta_3 \approx 0$ $\theta_4 \approx 0$

We could also regularize all of our theta parameters in a single summation as:

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

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

Using the above cost function with the extra summation, we can smooth the output of our hypothesis function to reduce overfitting. If lambda is chosen to be too large, it may smooth out the function too much and cause underfitting. Hence, what would happen if $\lambda = 0$ or is too small?