

CS1138

Machine Learning

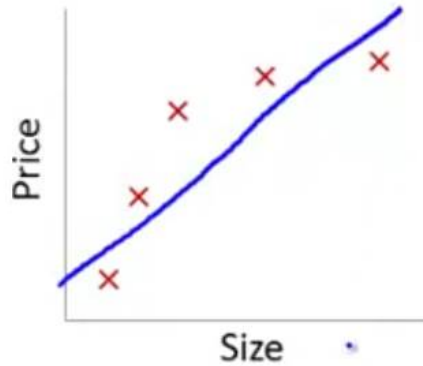
Lecture : Regularization: Problem of Overfitting

(Slide Credits: Andrew Ng)

Arpan Gupta

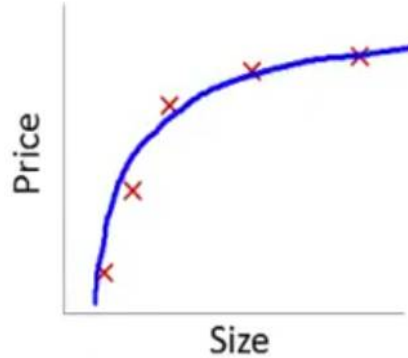
Problem of Overfitting

Example: Linear regression (housing prices)



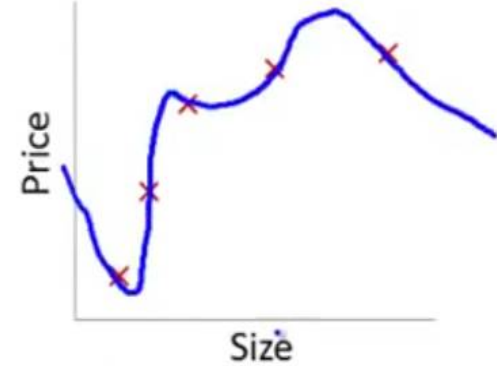
→ $\theta_0 + \theta_1 x$

Underfitting.
High Bias.



→ $\theta_0 + \theta_1 x + \theta_2 x^2$

“Just right”



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

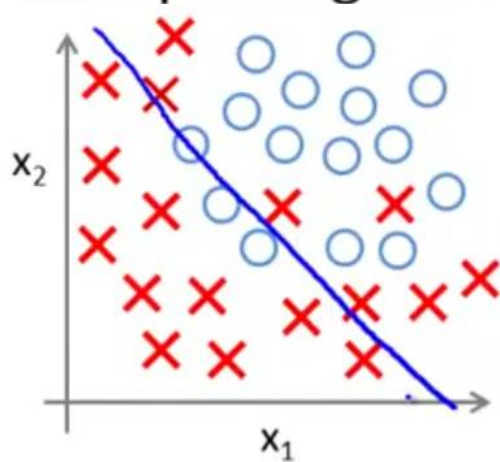
Overfitting.
High Variance

Problem of Overfitting

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

Problem of Overfitting

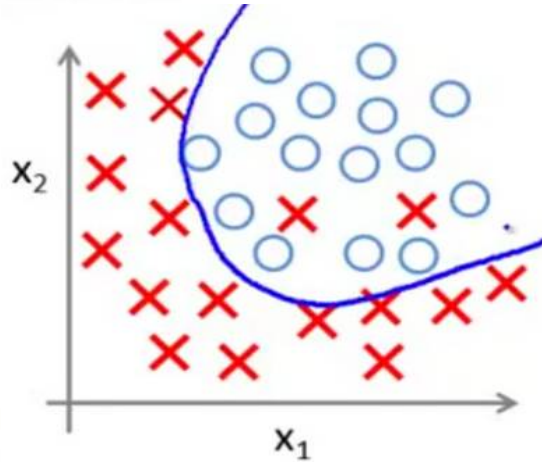
Example: Logistic regression



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

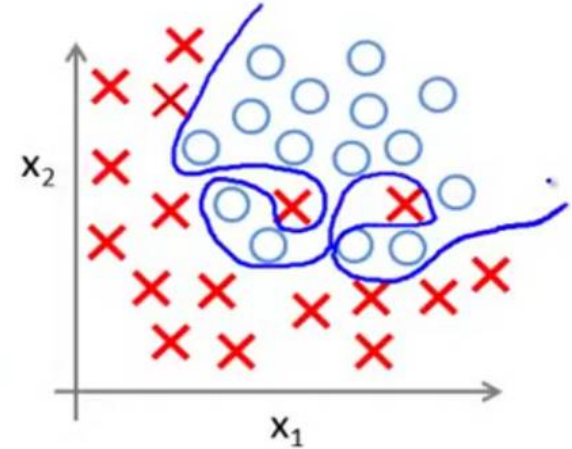
(g = sigmoid function)

Underfitting.
High Bias.



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

Just right fit

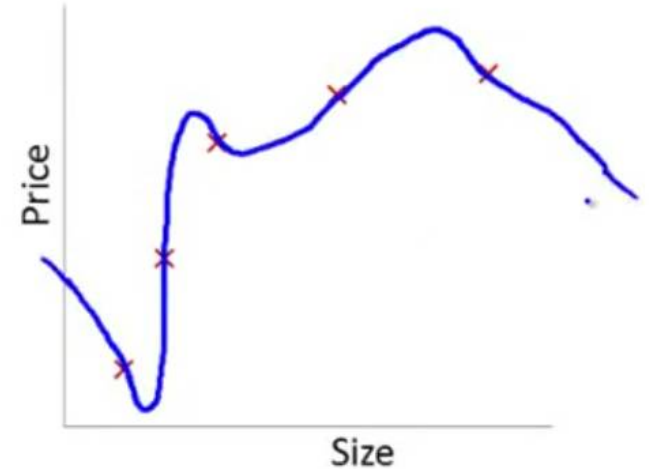


$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots)$$

Overfitting.
High Variance

Addressing Overfitting

- If overfitting is occurring, what to do to address it.
 - For 1D or 2D data, we can just plot and see which degree of polynomial best fits the data.
 - But it doesn't always work.



Addressing Overfitting (Contd.)

- We may have a learning problem when there are a lot of features.
 - Becomes much harder to visualize for more features.
 - **If we have lot of features and very little data, then overfitting can occur.**

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

x_4 = age of house

x_5 = average income in neighborhood

x_6 = kitchen size

\vdots

x_{100}

Addressing Overfitting (contd.)

- Options:

- 1) Reduce number of features

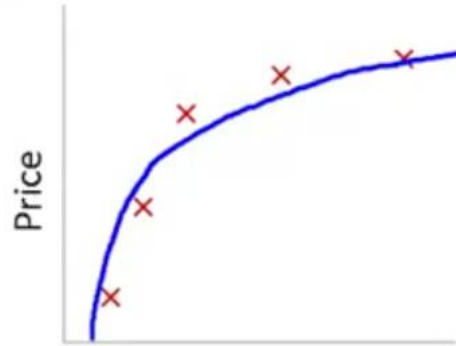
- a) Manually select which features to keep
- b) Use a model selection algorithm

- 2) Regularization

- a) Keep all features, but reduce magnitude/values of parameters Θ_j .
- b) Works well when we have a lot of features, each of which contributes a bit to predicting \mathbf{y} .

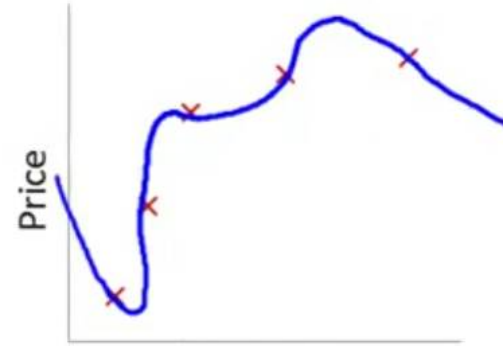
Regularization: Cost Function

Intuition



Size of house

$$\theta_0 + \theta_1 x + \theta_2 x^2$$



Size of house

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

Suppose we penalize and make θ_3, θ_4 really small.

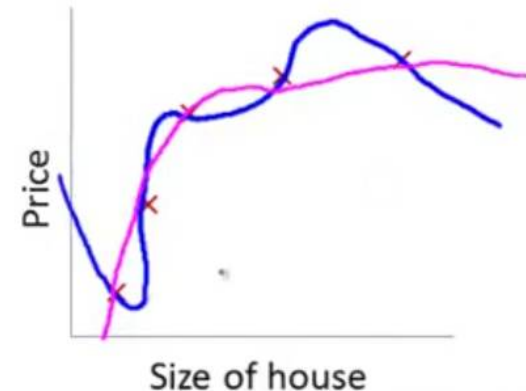
$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

Regularization: Cost Function

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

When we minimize this new function, then we have to make θ_3 and θ_4 very small.

We will end up with θ_3 and θ_4 close to 0. Similar to getting rid of the higher order terms in the hypothesis function.



$$\underline{\theta_0 + \theta_1 x + \theta_2 x^2} + \cancel{\theta_3 x^3} + \cancel{\theta_4 x^4}$$

Regularization

Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$

- “Simpler” hypothesis
- Less prone to overfitting

Housing:

- Features: x_1, x_2, \dots, x_{100}
- Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

We may not know which parameters to shrink, so we select all of them.

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

Regularization

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

- The convention is to penalize Θ_1 to Θ_n and not penalize Θ_0 , i.e., summation j from 1 to n . But in practice, it does not make much difference.
- If we penalize Θ_0 (intercept $\rightarrow 0$) means that the hyperplane passes through the origin.
- The second term is the **regularization term** and **λ is the regularization parameter**.

Regularization

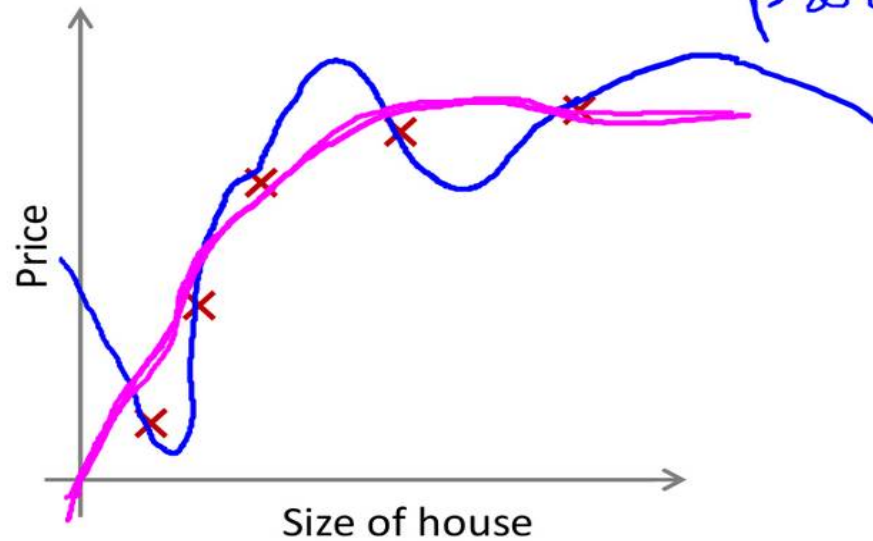
$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

- This objective function ensures two goals, corresponding to the two summation terms.
 - **Fitting training set well:** The first term (squared error $(\mathbf{h}(\mathbf{x}) - \mathbf{y})^2$) tells that we want to fit the training set well.
 - **Keep the parameters small:** The second term (regularization term) tells that the chosen Θ values should be small.
- λ controls the tradeoff between the above two goals (fitting the training set well and keeping the parameters small), thereby, keeping the hypothesis relatively “simple” and avoid overfitting.

Regularization.

$$\rightarrow J(\theta) = \frac{1}{2m} \left[\underbrace{\sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2}_{\text{data fit}} + \underbrace{\lambda \sum_{j=1}^n \theta_j^2}_{\text{regularization parameter}} \right]$$

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

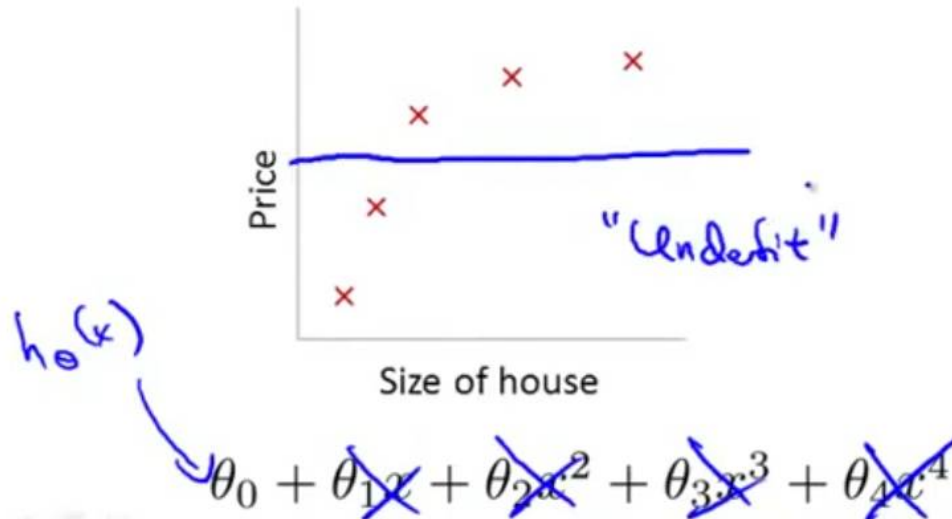


Effect of λ

In regularized linear regression, we choose θ to minimize

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

What if λ is set to an extremely large value (perhaps too large for our problem, say $\lambda = 10^{10}$)?



$$\begin{aligned} &\theta_1, \theta_2, \theta_3, \theta_4 \\ &\theta_1 \approx 0, \theta_2 \approx 0 \\ &\theta_3 \approx 0, \theta_4 \approx 0 \\ &\underline{h_{\theta}(x) = \theta_0} \end{aligned}$$

Regularized Linear Regression

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

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

Gradient Descent for regularized linear regression

Repeat {

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

$(j = 0, 1, 2, 3, \dots, n)$

}

Gradient Descent for regularized linear regr

Repeat {

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

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

$(j = \text{X}, 1, 2, 3, \dots, n)$

}

Gradient Descent for regularized linear regr

Repeat {

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

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

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

$(j = \text{~~0~~, 1, 2, 3, \dots, n)$

$$\frac{\partial}{\partial \theta_j} J(\theta) \quad \text{regularized}$$

}

Gradient Descent for regularized linear regr

Repeat {

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

$$\frac{2}{2\theta_0} J(\theta)$$

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

$(j = \text{~~0~~, 1, 2, 3, \dots, n)$

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

Gradient Descent for regularized linear regr

Repeat {

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

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

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

($j = \cancel{0}, 1, 2, 3, \dots, n$)

$$- \frac{\lambda}{m} \theta_j$$

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

$$\rightarrow J(\theta)$$

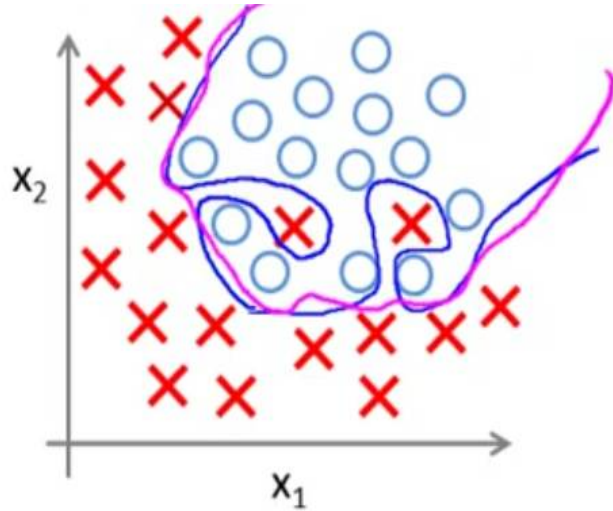
$$\theta_j^2$$

$$1 - \alpha \frac{\lambda}{m} < 1$$

$$0.99$$

$$\theta_j \times 0.99$$

Regularized Logistic Regression



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

Cost function:

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad | \quad \theta_1, \theta_2, \dots, \theta_n$$

Gradient Descent for regularized logistic regression

Repeat {

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

$(j = 0, 1, 2, 3, \dots, n)$

}

Gradient Descent for regularized logistic regression

Repeat {

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

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

$(j = \text{~~x~~, } \underline{1, 2, 3, \dots, n})$
 $\theta_1 \dots \theta_n$

}

Similar to linear regression but hypothesis is different.

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

Gradient Descent for regularized logistic regression

Repeat {

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

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

$(j = \text{red X}, 1, 2, 3, \dots, n)$
 $\theta_1, \dots, \theta_n$

}

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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

Ridge and Lasso

- Cost function where we take square of all parameters, is also known as **Ridge** regression.

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

- If we take the absolute values of parameters Θ , then it is called **Lasso** regression.

End of Lecture