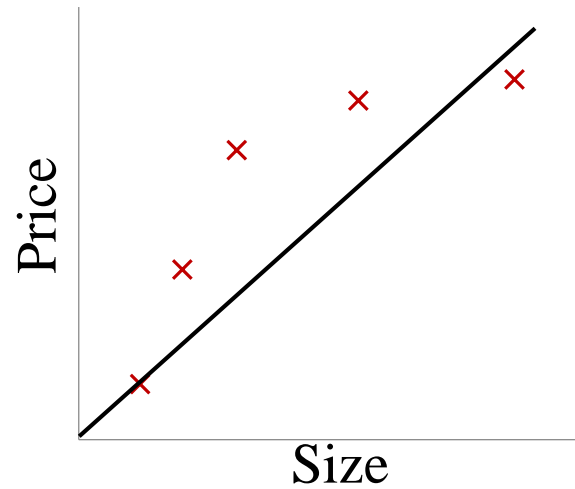
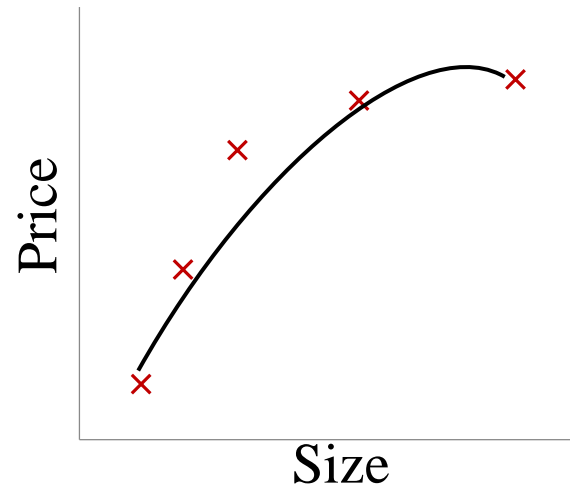


## Example: Linear regression (housing prices)



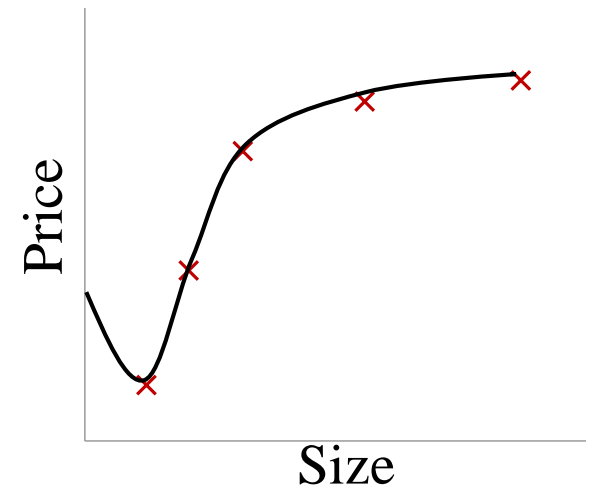
$$\theta_0 + \theta_1 x$$

Underfitting: “High Bias”



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

Correct Fit

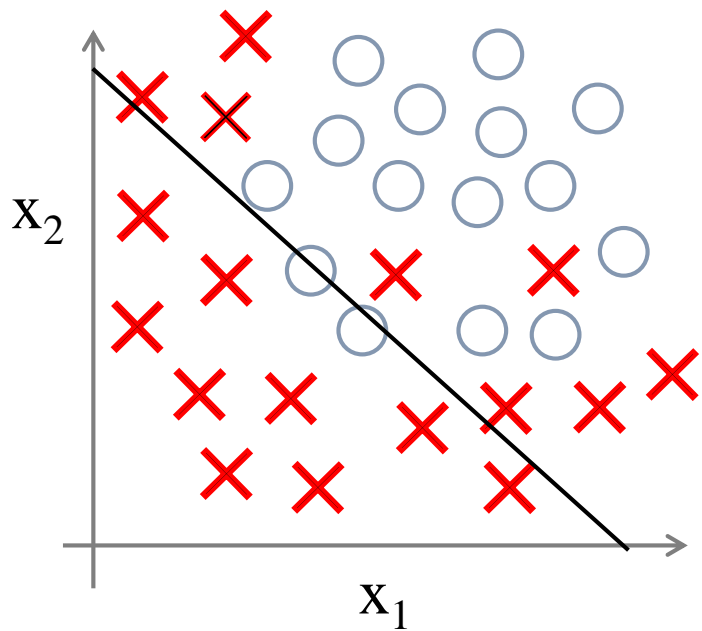


$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfitting: “High Variance”

**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).

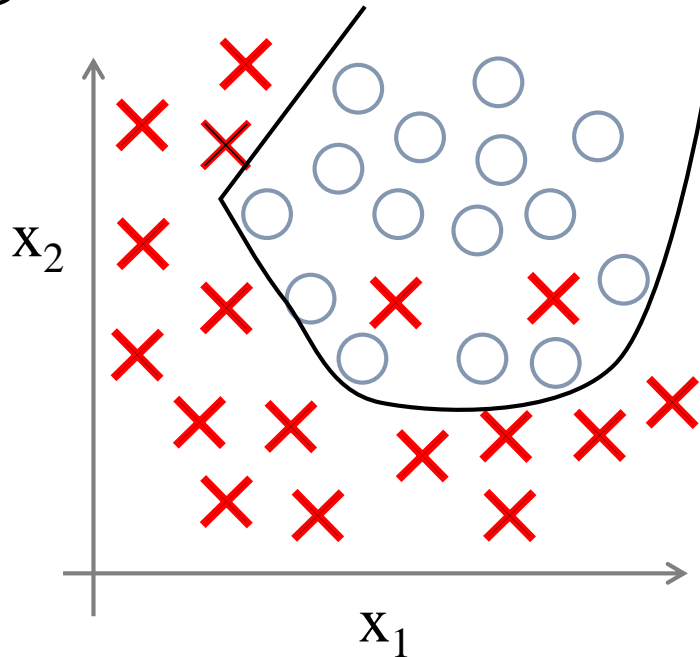
## Example: Logistic regression



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

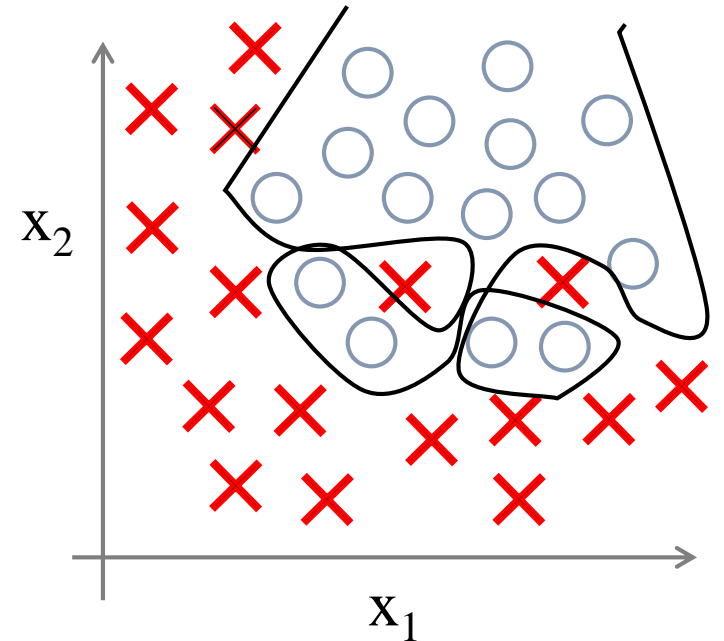
(  $g$  = sigmoid function)

Underfitting: “High Bias”



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

Correct 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:

## Options:

1. Reduce number of features.
  - Manually select which features to keep.
  - Model selection algorithm (later in course).
2. Regularization.
  - Keep all the features, but reduce magnitude/values of parameters  $\theta_j$
  - Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

$\lambda$  = regularization parameter

Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

assume  $\lambda = 1000$

Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

$\theta_3$  &  $\theta_4$  is going to decrease



Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

$\theta_3$  &  $\theta_4$  is going to decrease

Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

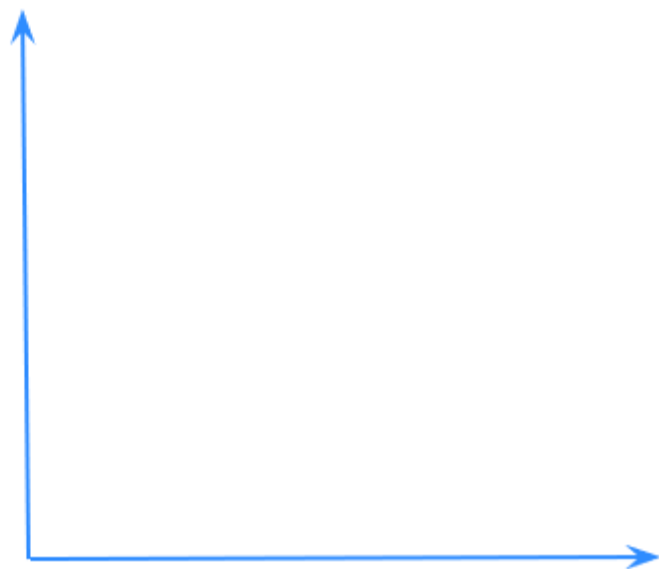
$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

$\theta_3$  &  $\theta_4$  is going to decrease



Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

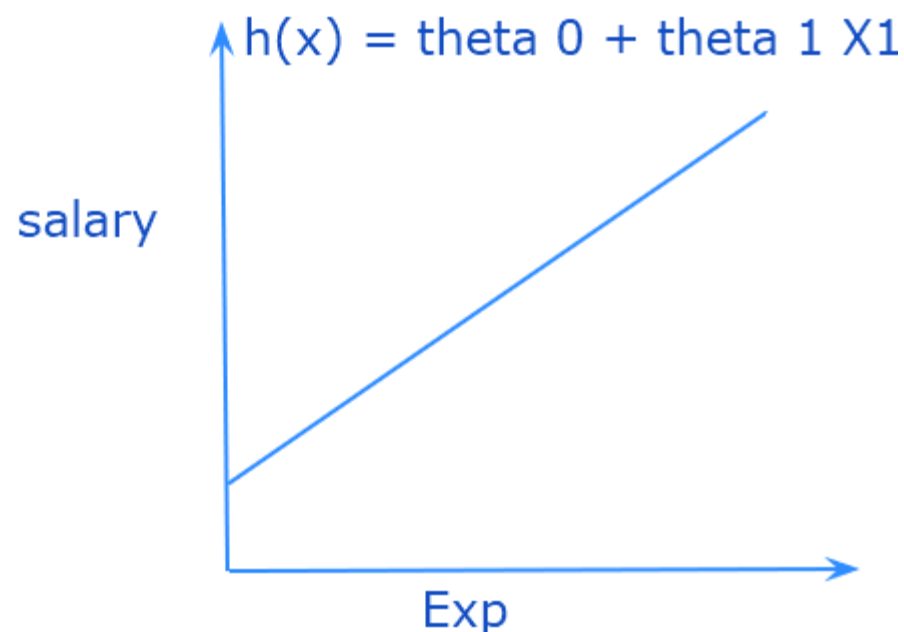
$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

$\theta_3$  &  $\theta_4$  is going to decrease



Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

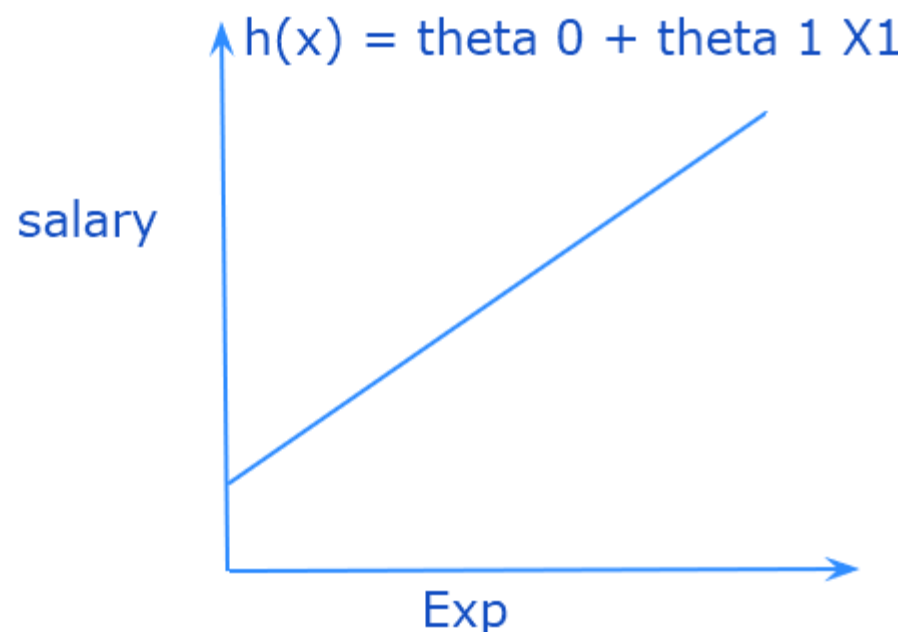
$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$   
 $\lambda(\text{slope})^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

$\theta_3$  &  $\theta_4$  is going to decrease



Regularization

Min =  $\frac{1}{2m} \sum_{i=1}^m (h(X)^i - y^i)^2$

parameters =  $\theta_j$

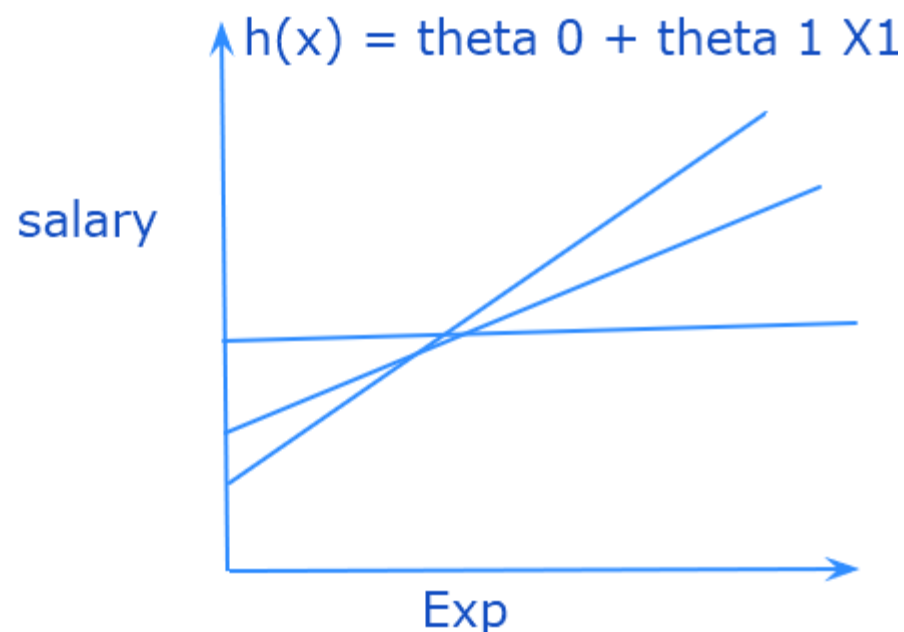
$\lambda$  = regularization parameter

Min = above terms are same +  $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$   
 $\lambda(\text{slope})^2$

assume  $\lambda = 1000$

Min = terms +  $1000 (\theta_3)^2 + 1000 (\theta_4)^2$

$\theta_3$  &  $\theta_4$  is going to decrease



## Regularized Linear Regression

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

## Regularized Linear Regression

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

## Regularized Logistic Regression

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

## Gradient Descent Algorithm

{

$\Theta_0 := \theta_0 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h(x)^i - y_i) X_0^i \right)$

$\Theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h(x)^i - y_i) X_j^i - \lambda \theta_j / m \right]$

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

}



## Gradient Descent Algorithm

{

$\Theta_0 := \theta_0 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (h(x)^i - y_i) X_0^i \right)$

$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h(x)^i - y_i) X_j^i - \frac{\lambda}{m} \theta_j \right]$

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

}

penalize the features that has higher value of slope  
 $\lambda = 0$  to +ve values

## Gradient Descent Algorithm

{

$\Theta_0 := \theta_0 - \alpha (1/m) \sum_{i=1}^m (h(x)^i - y_i) X_0^i$

$\Theta_j := \theta_j - \alpha [1/m \sum_{i=1}^m (h(x)^i - y_i) X_j^i - \lambda/m \theta_j]$

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

}

penalize the features that has higher value of slope  
 $\lambda = 0$  to +ve values

Ridge regression & lasso regression

ridge regression

$\lambda(\text{slope})^2$

Lasso regression

$\lambda(\text{absolute value of the slope})$

ridge regression

$\lambda(\text{slope})^2$

Lasso regression

$\lambda(\text{absolute value of the slope})$

ridge regression

$\lambda(\text{slope})^2$       closely equal to zero

Lasso regression

$\lambda(\text{absolute value of the slope}) = \text{zero}$

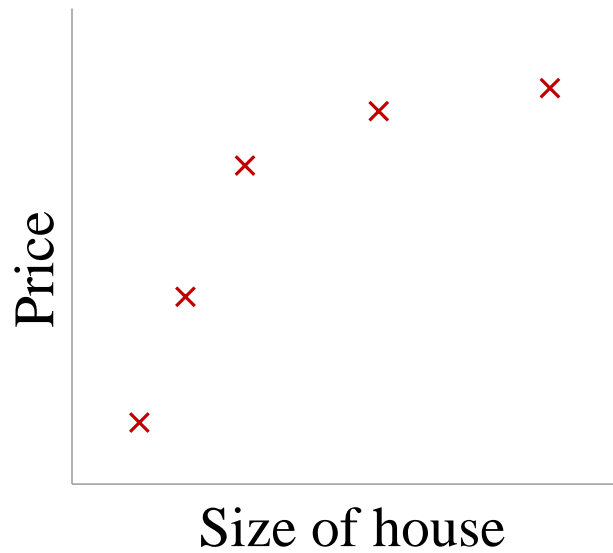
ridge regression    when most predictors impact the response

$\lambda(\text{slope})^2$     closely equal to zero

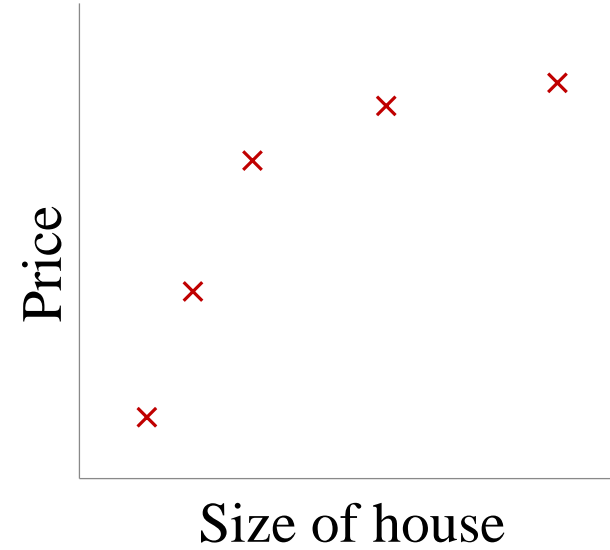
Lasso regression    when only a few predictors actually influence the response

$\lambda(\text{absolute value of the slope}) = \text{zero}$

# Intuition



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



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

Suppose we penalize and make  $\theta_3, \theta_4$  really small.

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

## 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:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

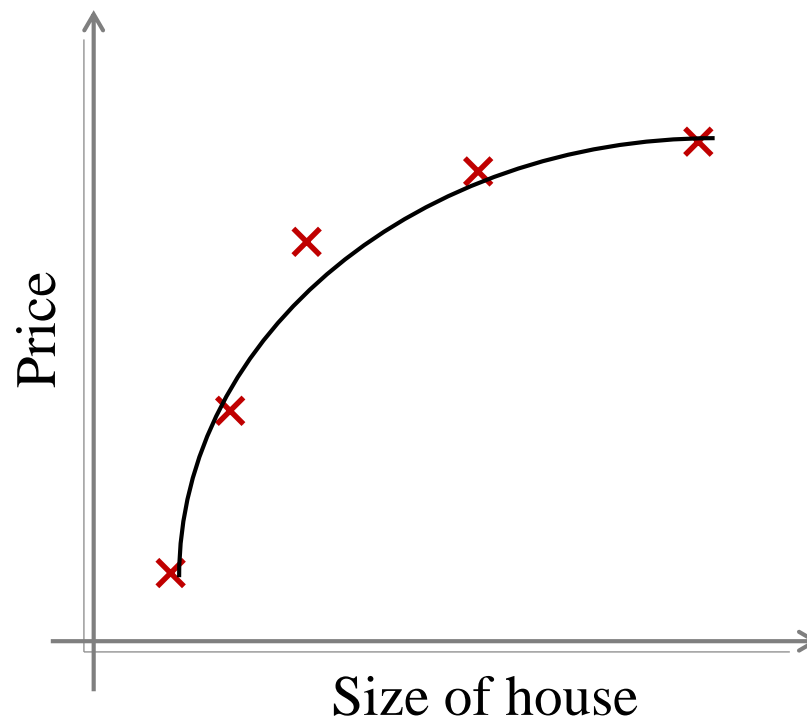
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$



# Regularization.

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

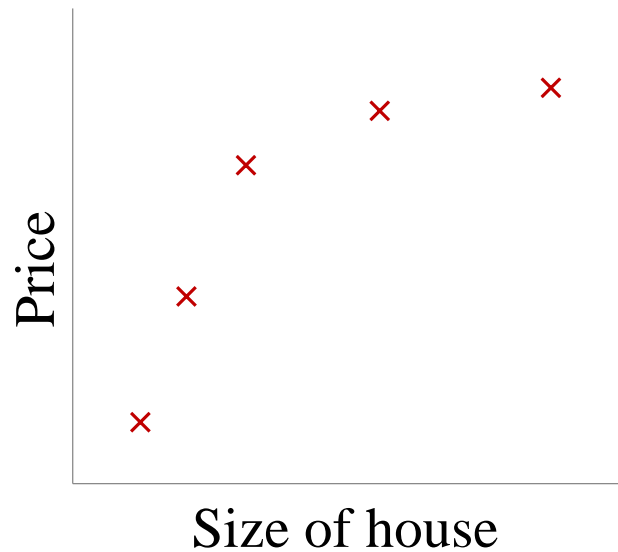
$$\min_{\theta} J(\theta)$$



In regularized linear regression, we choose  $\theta$  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  $\lambda$  is set to an extremely large value (perhaps far too large for our problem, say  $\lambda = 10^{10}$ )?



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

# 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

Repeat {

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

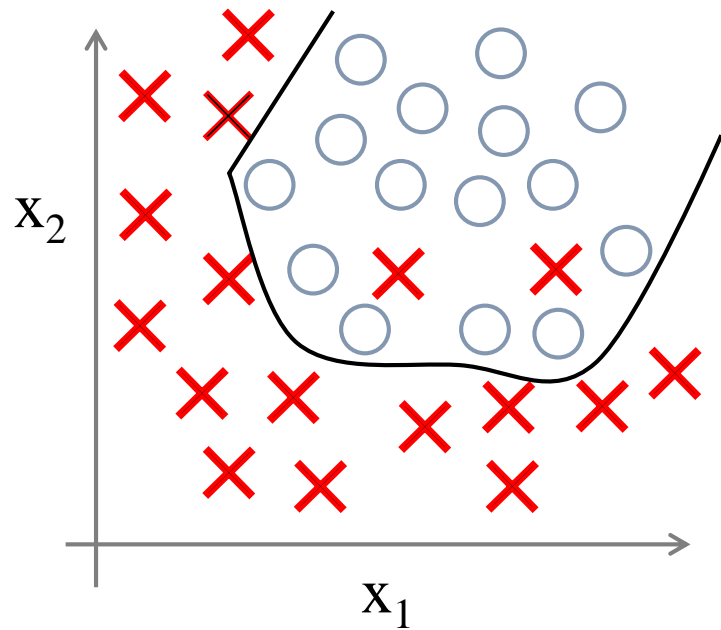
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$(j = \text{~~8~~, 1, 2, 3, \dots, n})$

}

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

# 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:

$$J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\text{lambda}}{2m} \text{summation } i = 1 \text{ to } n \text{ theta}_{j2}$$

# Gradient descent

Repeat {

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

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

$(j = \text{~~8~~, 1, 2, 3, \dots, n)$

}

# RIDGE & LASSO REGRESSION

- Both the techniques are adopted to **overcome the overfitting issue**.
- Lasso regression stands for Least Absolute Shrinkage and Selection Operator & it tends to make coefficients to absolute zero as compared to Ridge which never sets the value of coefficient to absolute zero.
- **Limitation of Ridge Regression:** Ridge regression decreases the complexity of a model but does not reduce the number of variables since it never leads to a coefficient been zero rather only minimizes it. Hence, this model is not good for feature reduction.
- **Limitation of Lasso Regression:** Lasso sometimes struggles with some types of data. Lasso will pick at most  $n$  predictors as non-zero, even if all predictors are relevant (or may be used in the test set).
- If there are two or more **highly collinear variables then LASSO regression select one of them randomly** which is not good for the interpretation of data.

