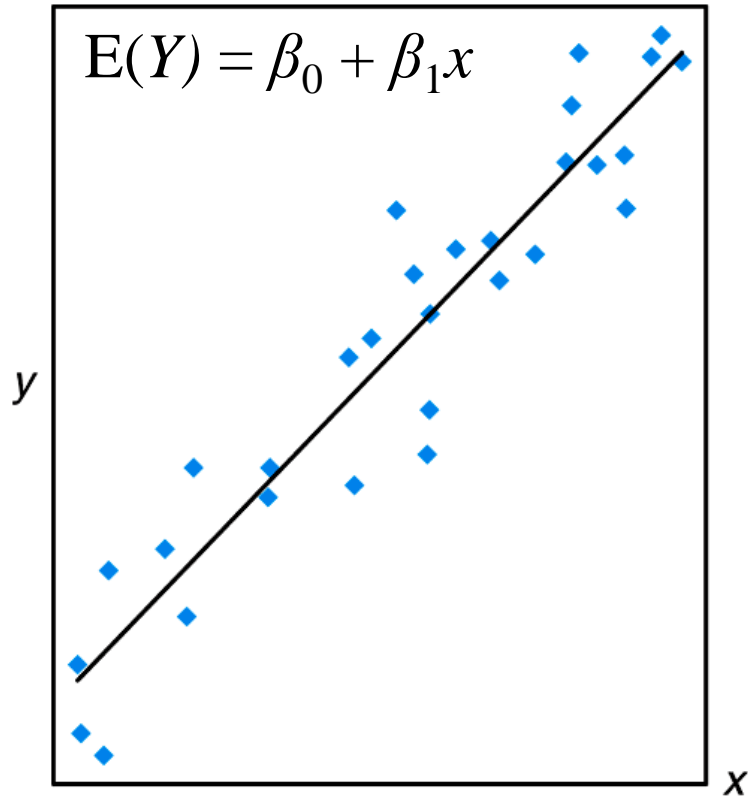


The Problem of Overfitting

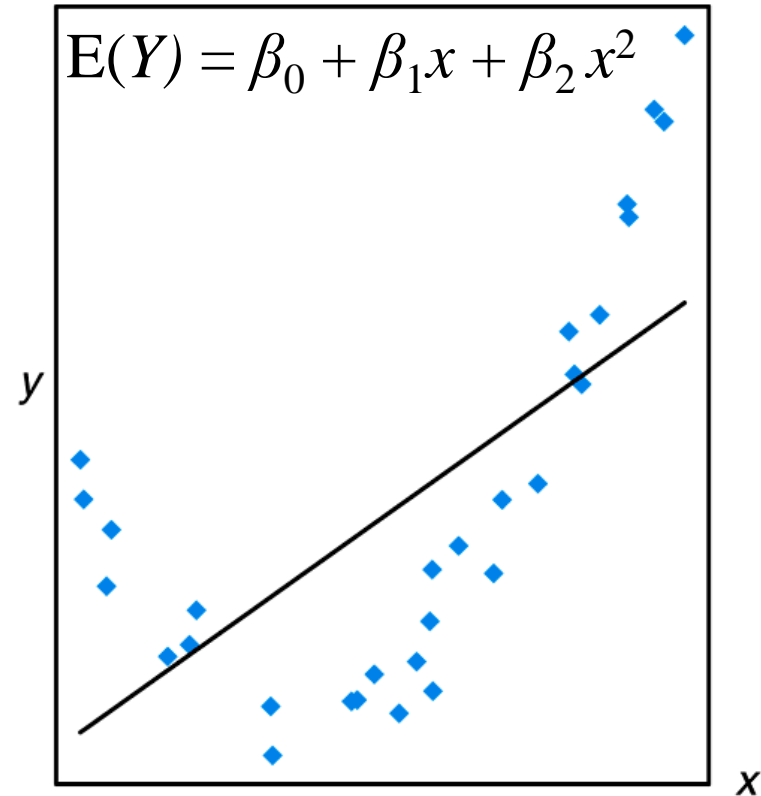
Solving the Problem of Overfitting

Regularization



(a)

(a) The relation between Y and x is linear.



(b)

(b) There is a second order relation between Y and x .

Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \textit{frontage} + \theta_2 \times \textit{depth}$$



Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

Area

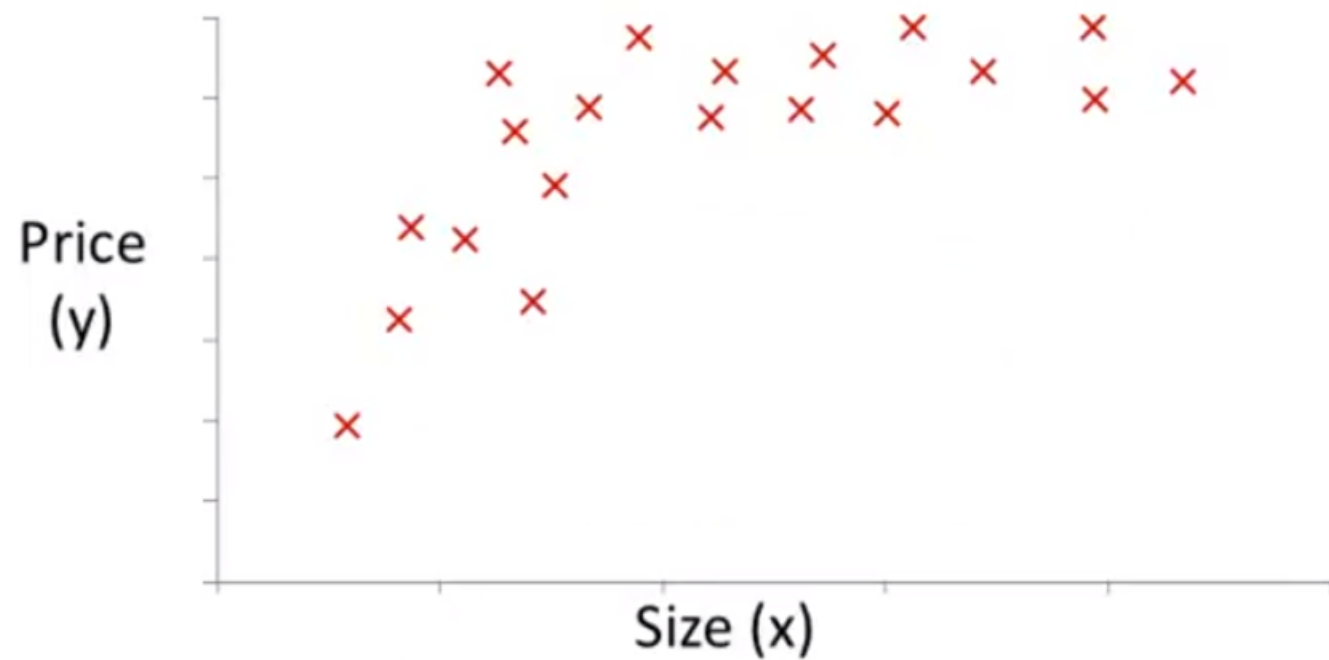
$$x = \underline{\text{frontage} \times \text{depth}}$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

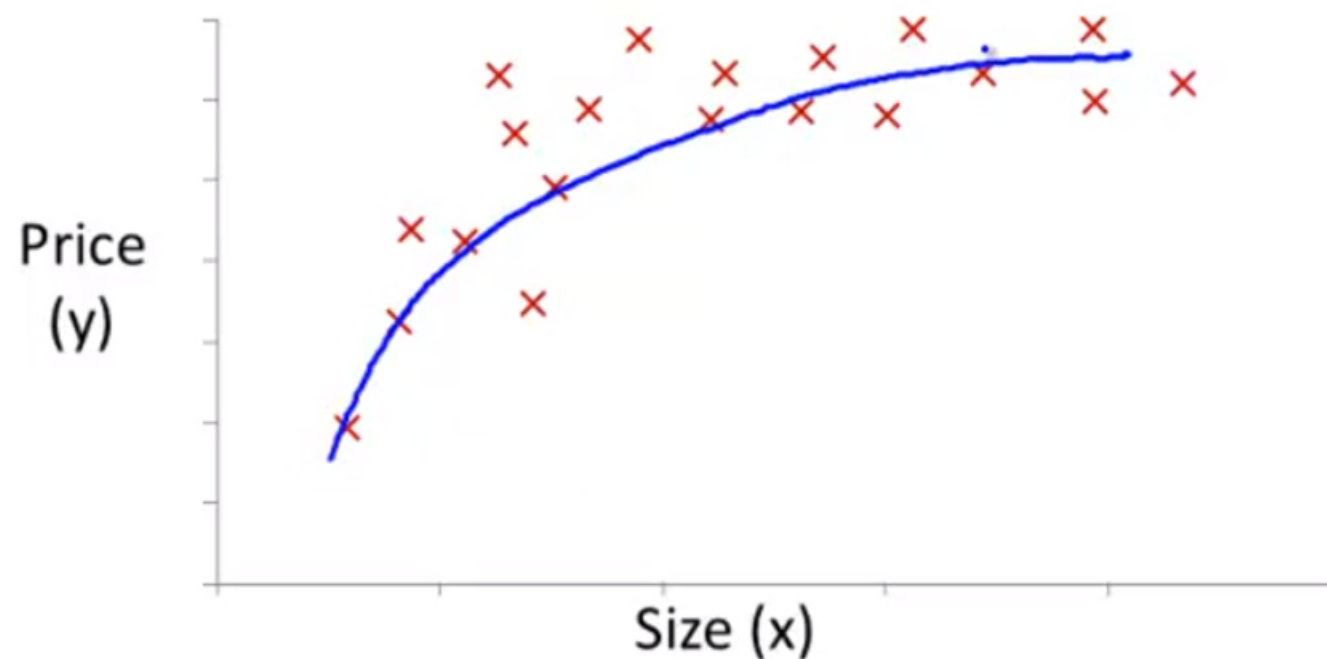
↖ land area



Polynomial regression

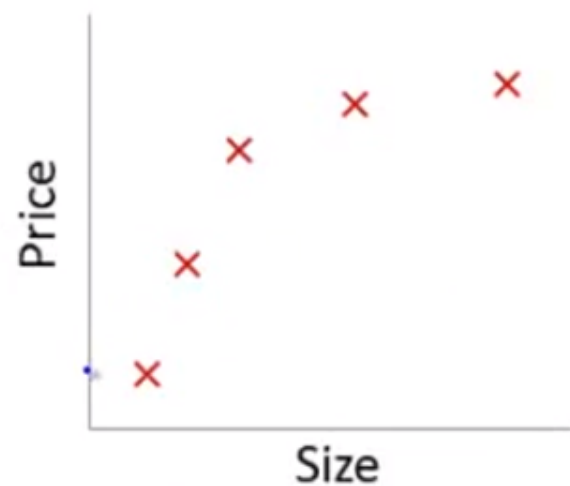


Polynomial regression

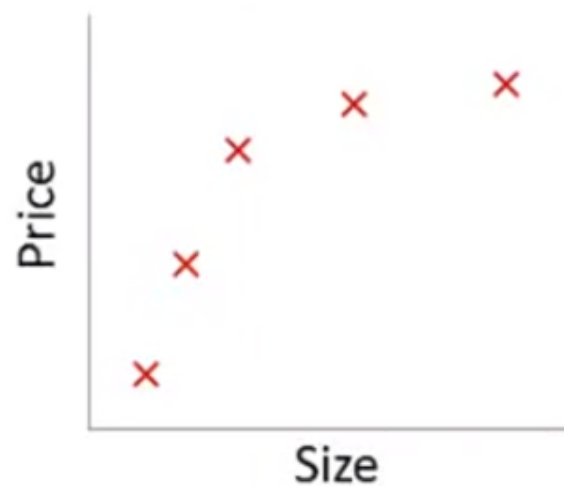


$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$

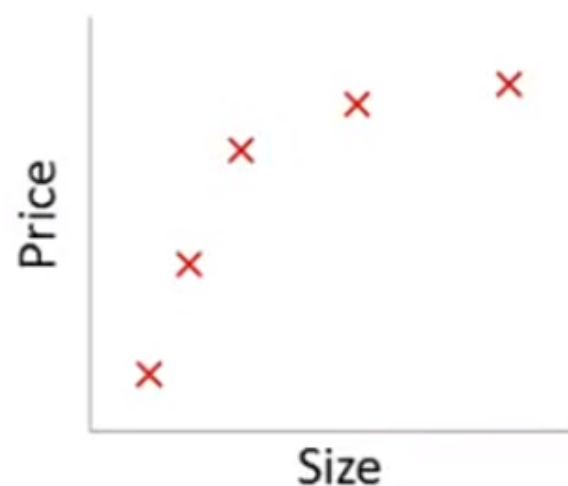
Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$

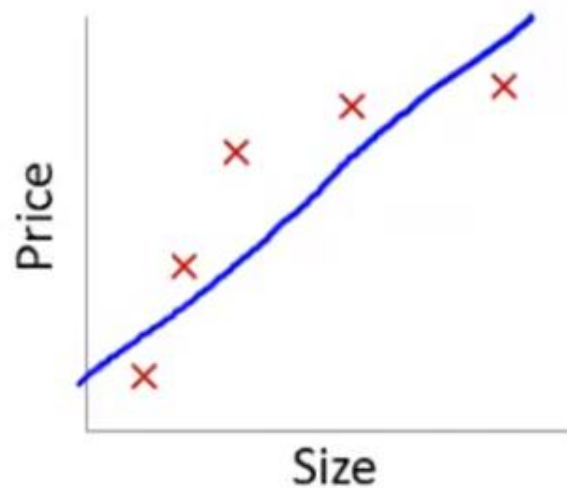


$\theta_0 + \theta_1 x + \theta_2 x^2$

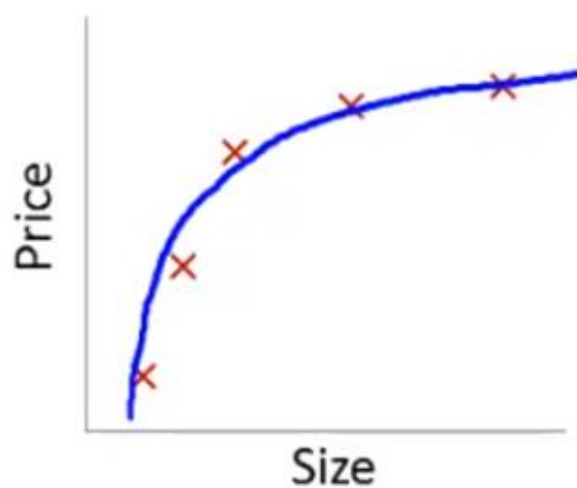


$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

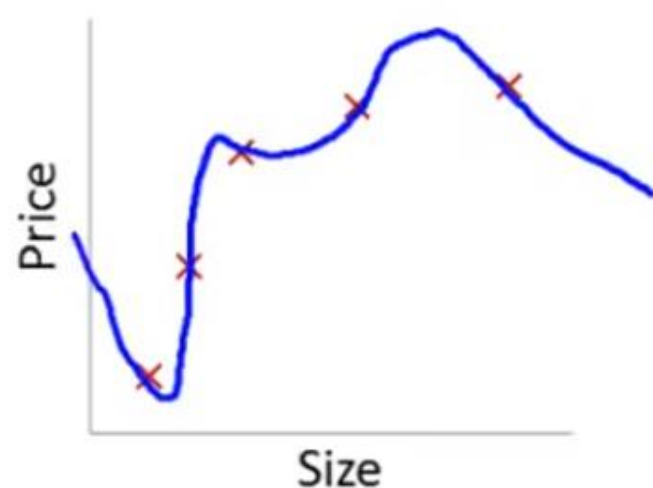
Example: Linear regression (housing prices)



→ $\theta_0 + \theta_1 x$

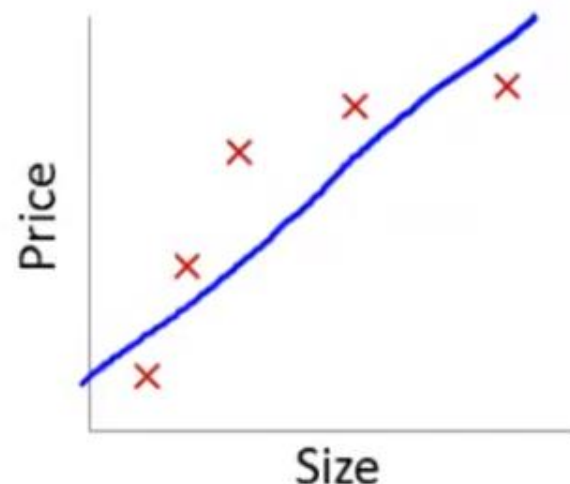


→ $\theta_0 + \theta_1 x + \theta_2 x^2$

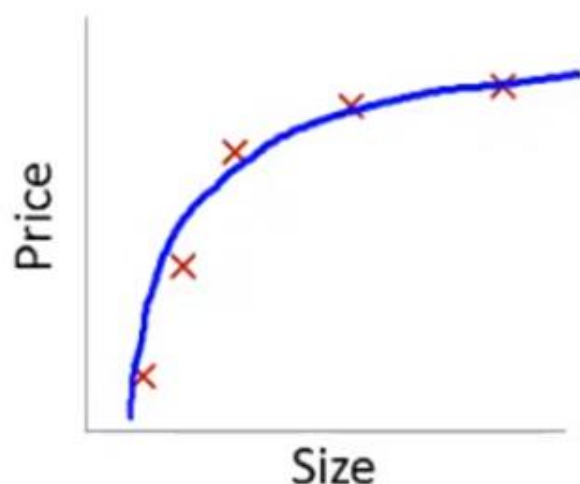


→ $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

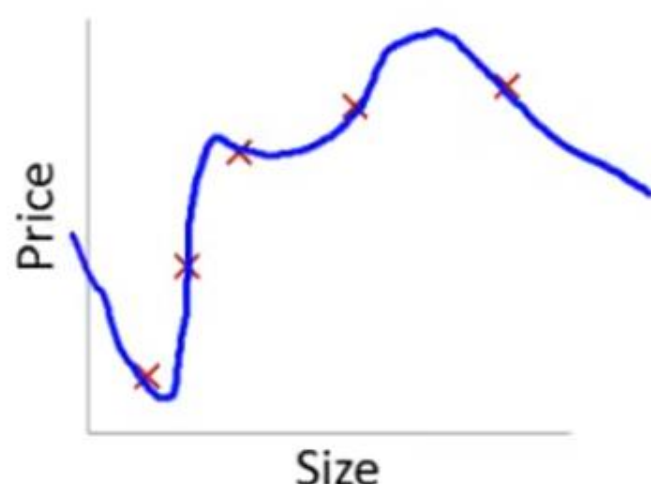
Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$
"Underfit" "High bias"



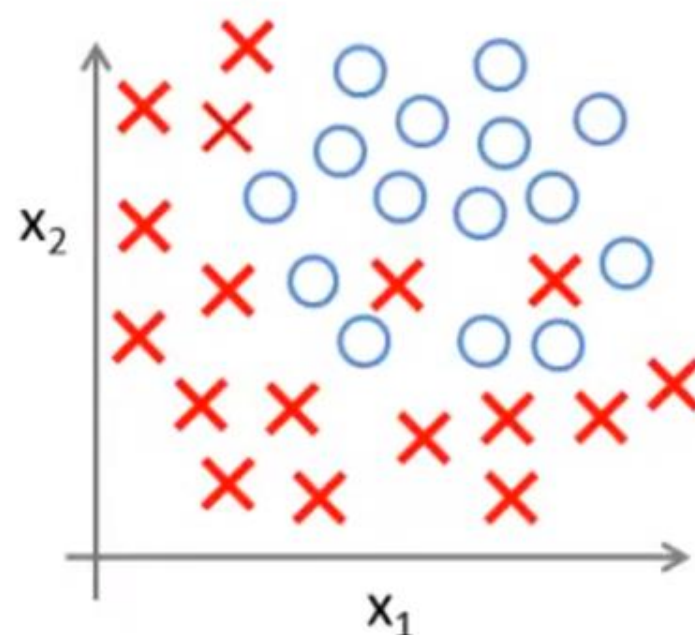
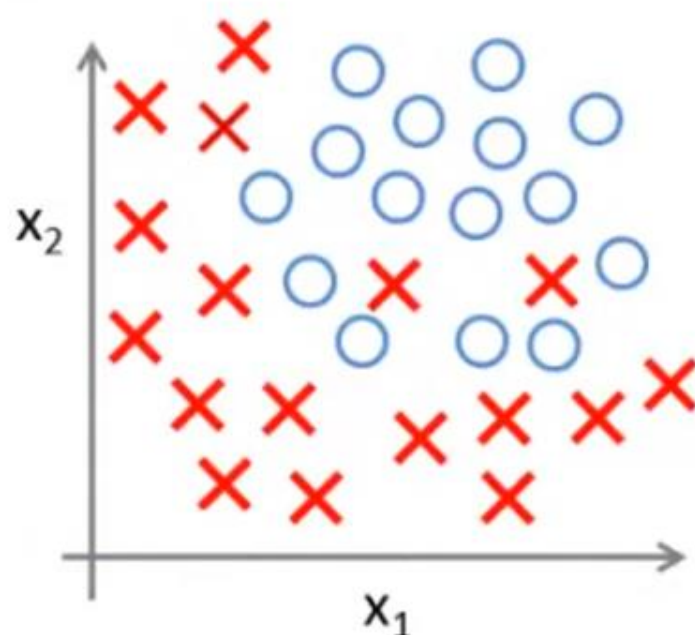
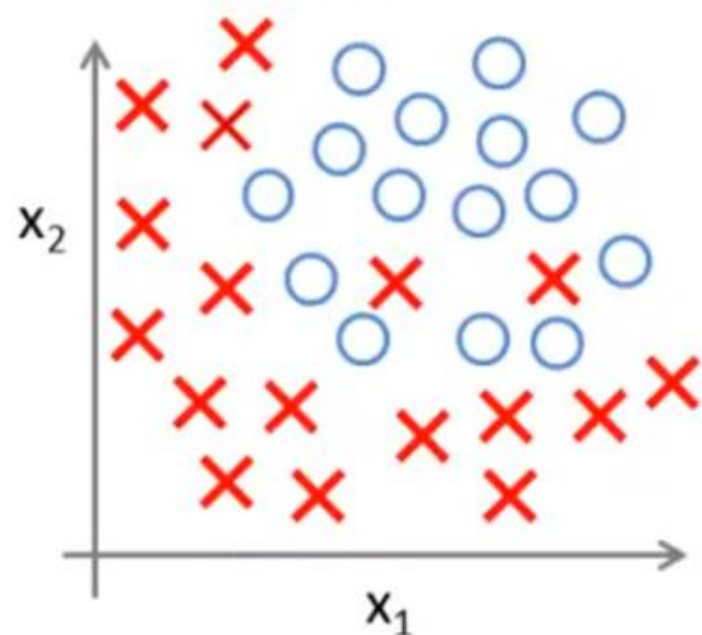
$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$
"Just right"



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

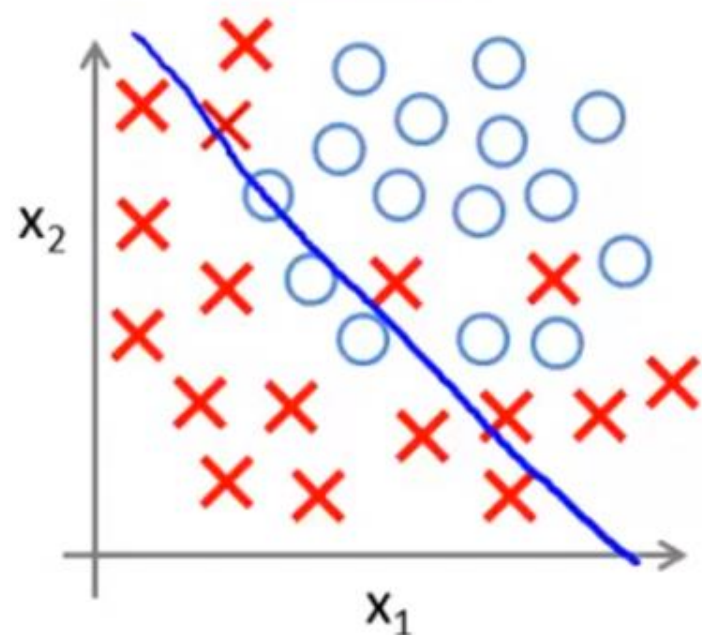


• $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$
 (g = sigmoid function)

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

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

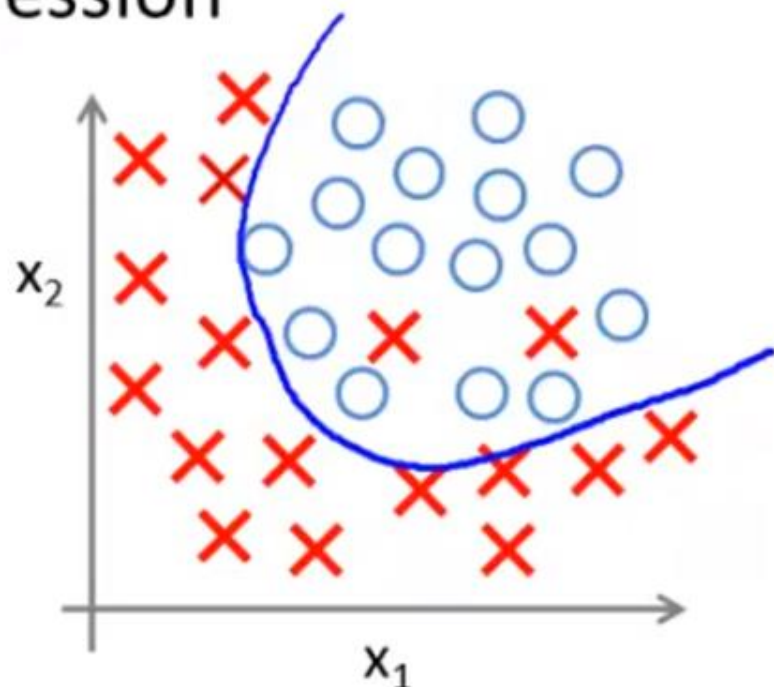
Example: Logistic regression



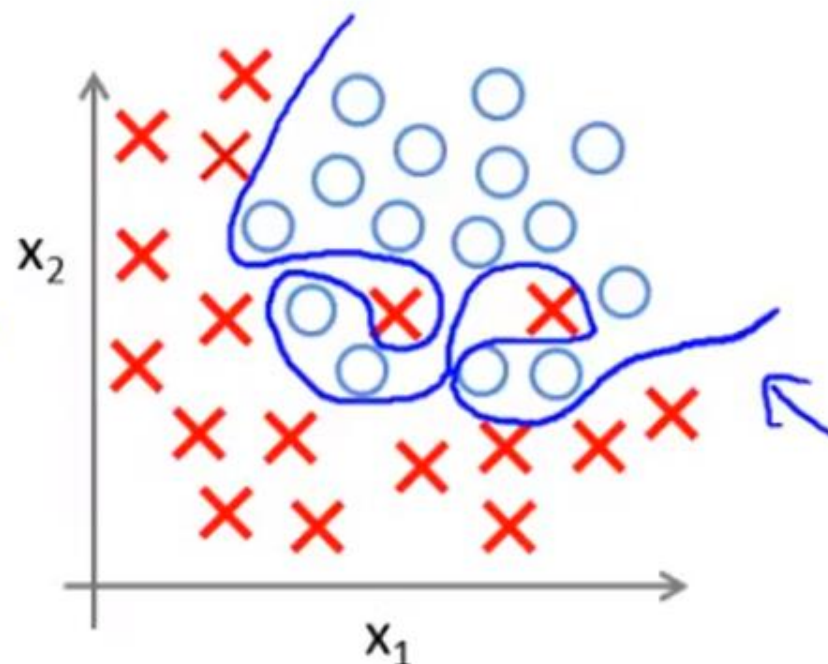
$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(g = sigmoid function)

"Underfit"



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



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

"Overfit"

Exercise

- Consider the medical diagnosis problem of classifying tumors as malignant or benign. If a hypothesis $h(x)$ has overfit the training set, it means that:
 - It makes accurate predictions for examples in the training set and generalizes well to make accurate predictions on new, previously unseen examples.
 - It does not make accurate predictions for examples in the training set, but it does generalize well to make accurate predictions on new, previously unseen examples.
 - It makes accurate predictions for examples in the training set, but it does not generalize well to make accurate predictions on new, previously unseen examples.
 - It does not make accurate predictions for examples in the training set and does not generalize well to make accurate predictions on new, previously unseen examples.

Addressing overfitting:

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

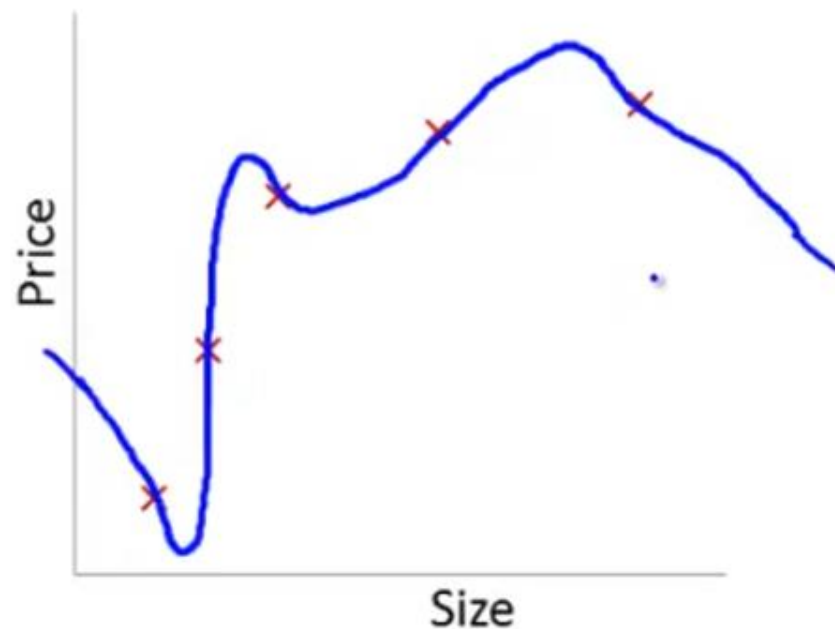
x_4 = age of house

x_5 = average income in neighborhood

x_6 = kitchen size

⋮

x_{100}



Addressing overfitting:

Options:

1. Reduce number of features.
 - Manually select which features to keep.
 - Model selection algorithm (later in course).

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 θ_j .

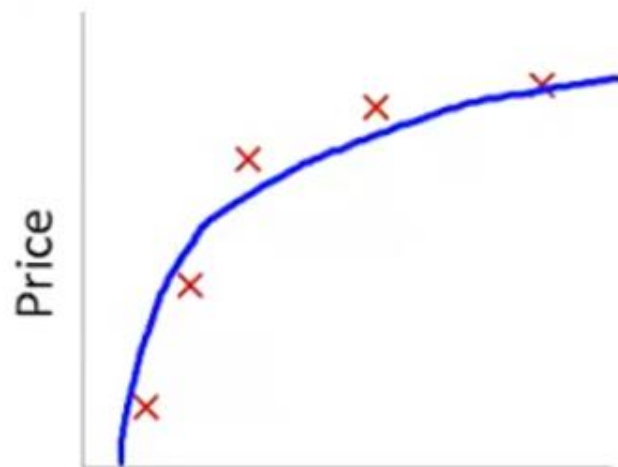
— Works well when we have a lot of features, each of which contributes a bit to predicting y .

Cost Function

Solving the Problem of Overfitting

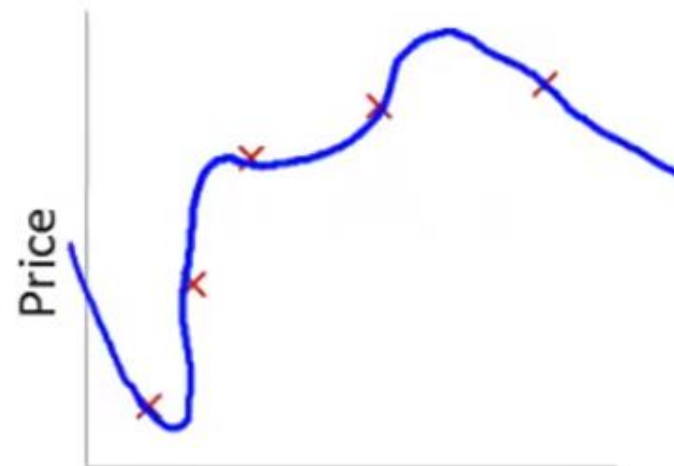
Regularization

Intuition



Size of house

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

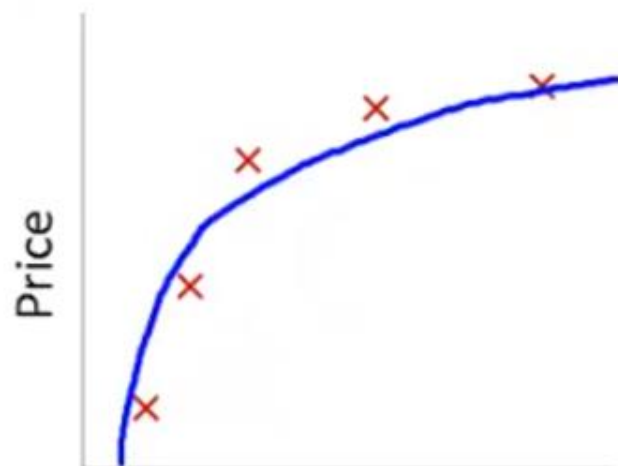


Size of house

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

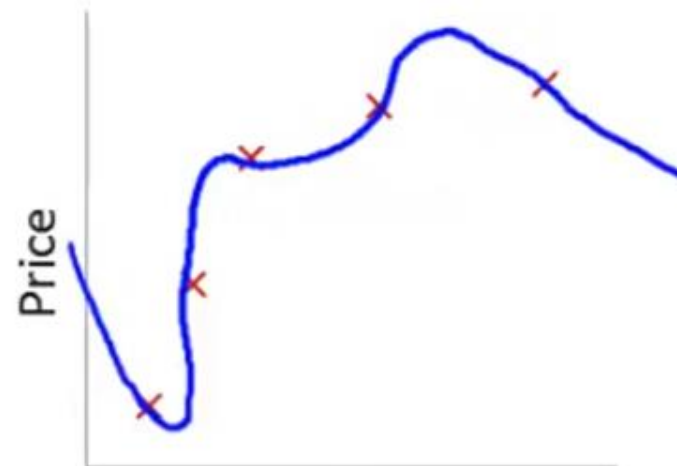
Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

Intuition



Size of house

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



Size of house

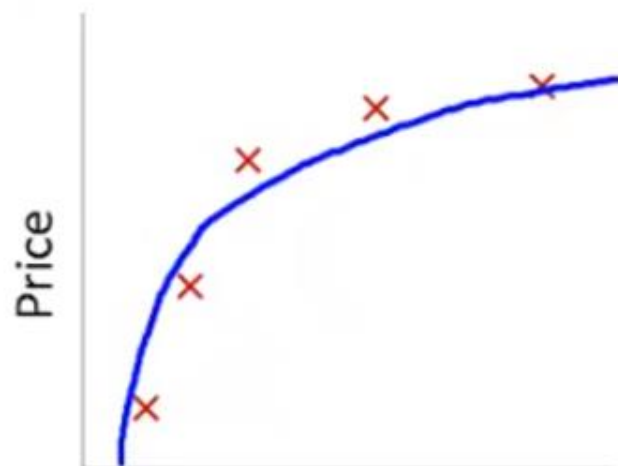
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make θ_3, θ_4 really small.

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

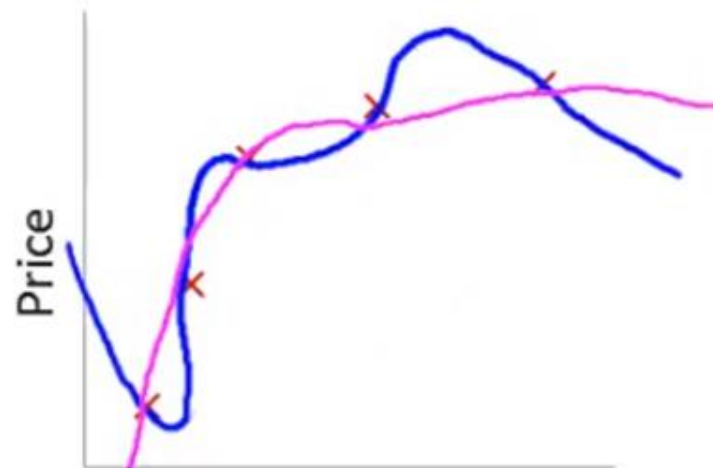
Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

Intuition



Size of house

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



Size of house

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

↑ ↑

Suppose we penalize and make θ_3, θ_4 really small.

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

$\underline{\theta_3 \approx 0} \qquad \underline{\theta_4 \approx 0}$

Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

Regularization.

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

- "Simpler" hypothesis
- Less prone to overfitting

$\rightarrow \boxed{\theta_3, \theta_4}$
 ≈ 0

Housing:

- Features: $\underline{x}_1, \underline{x}_2, \dots, x_{100}$
- Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

- Which one to minimize???

Regularization.

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

- “Simpler” hypothesis
- Less prone to overfitting

$\rightarrow \theta_3, \theta_4$
 $\nearrow \approx 0$

Housing:

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

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

~~$\theta_1, \theta_2, \theta_3, \dots, \theta_{100}$~~

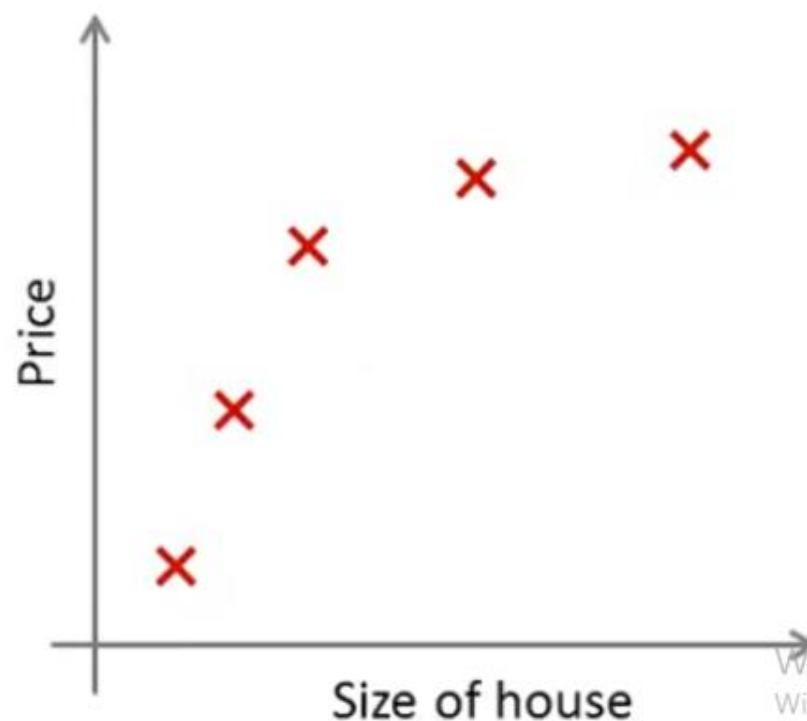
Regularization.

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

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

This is also called the Ridge regression.

regularization
parameter



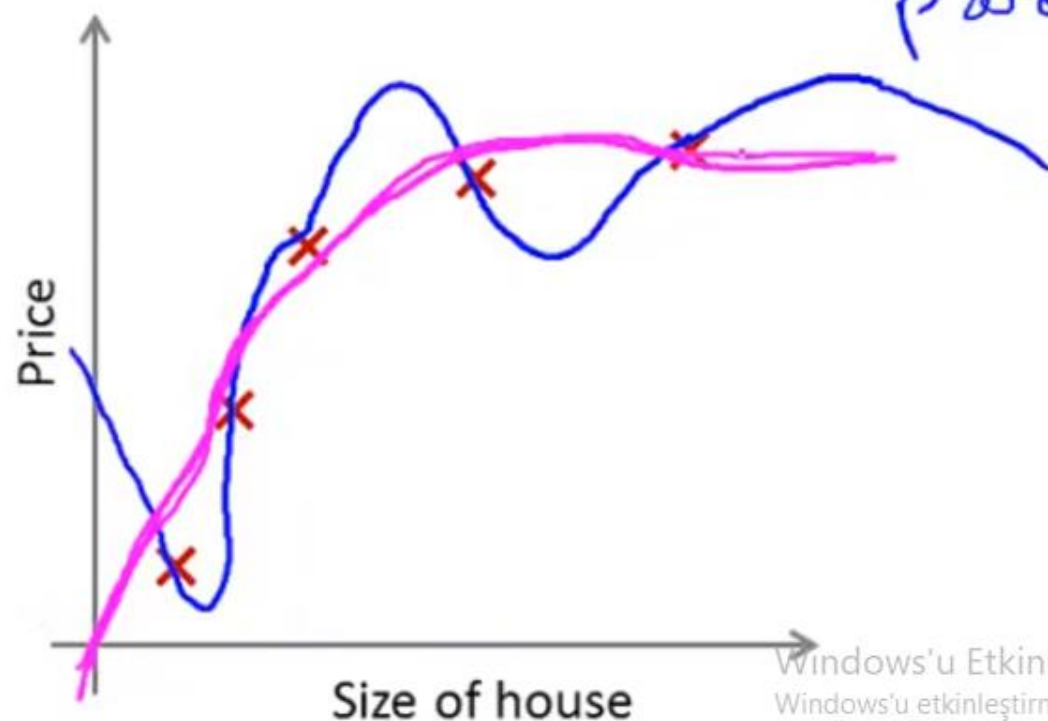
Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

Regularization.

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

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

regularization parameter

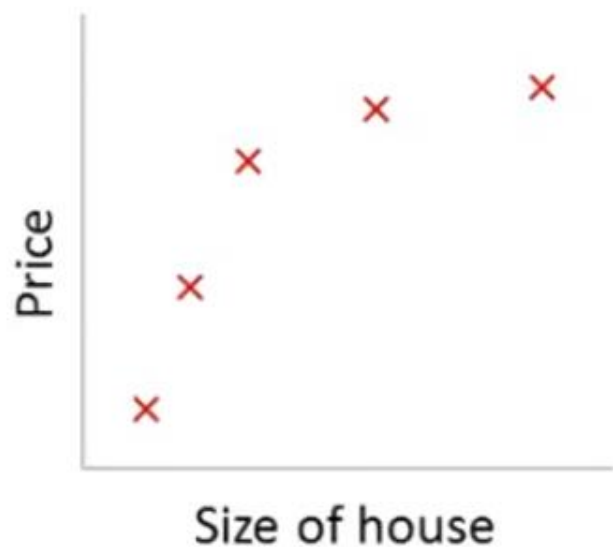


Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

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 far too large for our problem, say $\lambda = 10^{10}$)?



$h_{\theta}(x)$

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

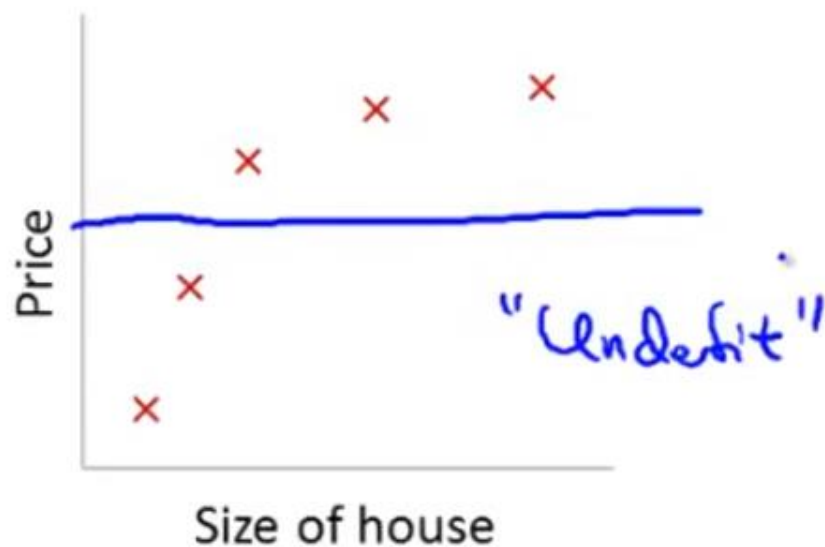
$\theta_1, \theta_2, \theta_3, \theta_4$
 $\theta_1 \approx 0, \theta_2 \approx 0$
 $\theta_3 \approx 0, \theta_4 \approx 0$

Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

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 far too large for our problem, say $\lambda = 10^{10}$)?



$$\theta_1, \theta_2, \theta_3, \theta_4$$

$$\theta_1 \approx 0, \theta_2 \approx 0$$

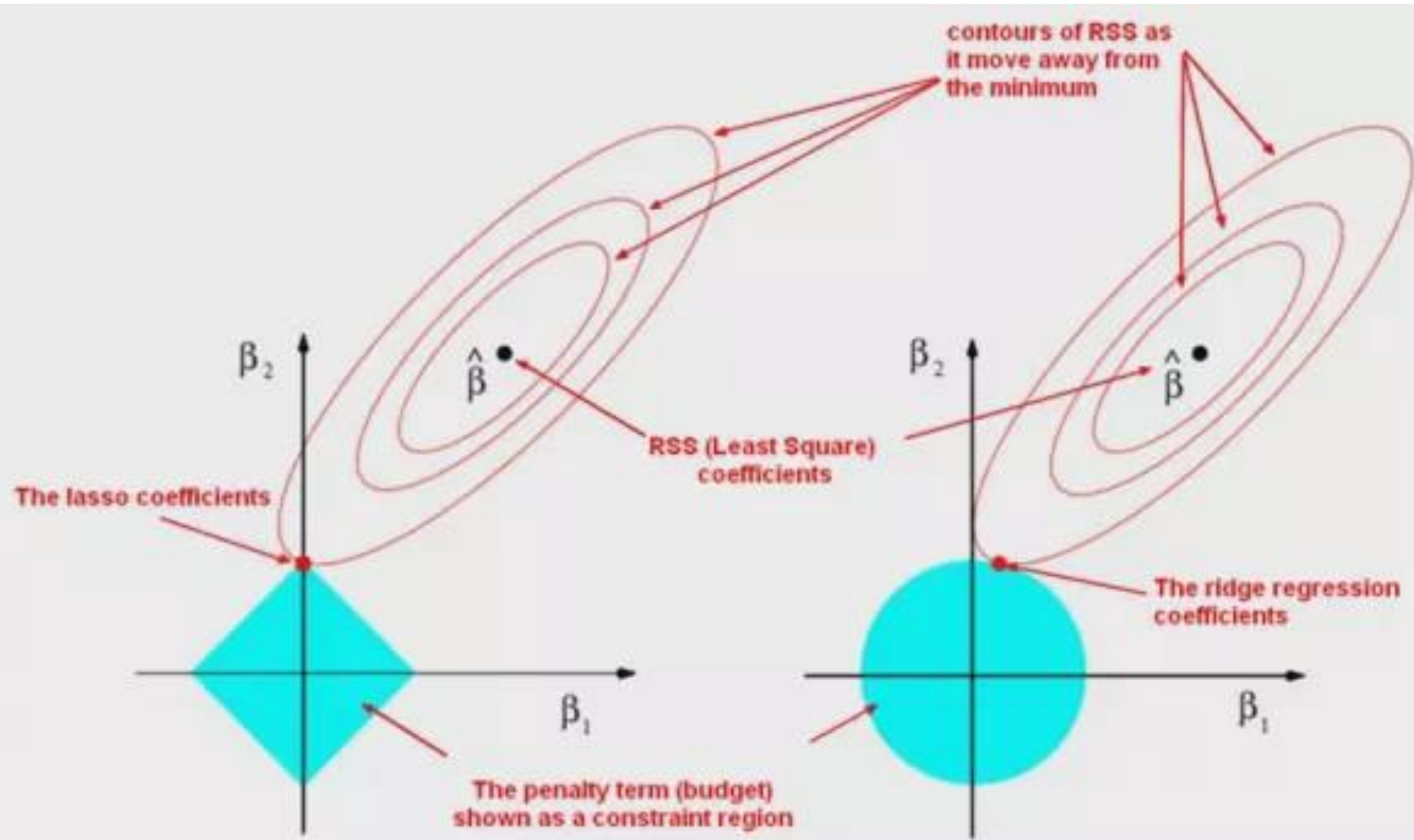
$$\theta_3 \approx 0, \theta_4 \approx 0$$

$$h_{\theta}(x) = \theta_0$$

Windows'u Etkinleştirin
Windows'u etkinleştirmek için Ayarlar'a gidin.

$h_{\theta}(x)$

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



LASSO

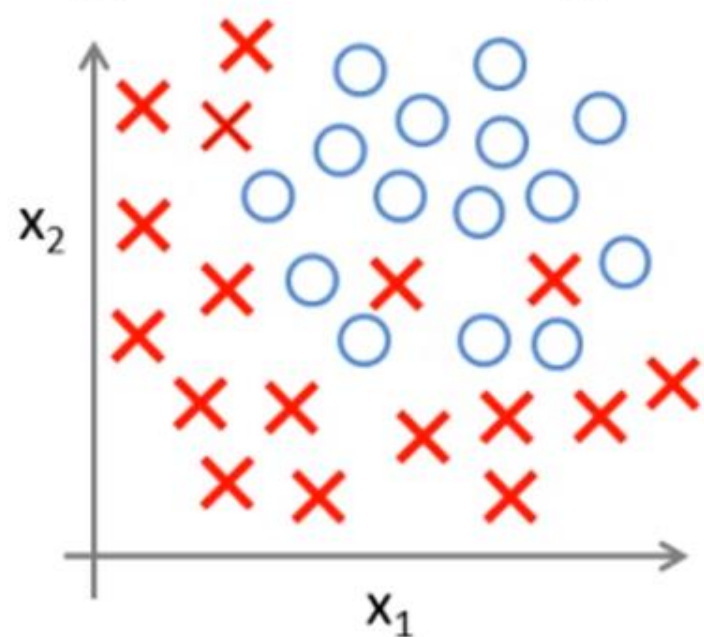
RIDGE REGRESSION

Regularized Logistic Regression

Solving the Problem of Overfitting

Regularization

Regularized logistic regression.



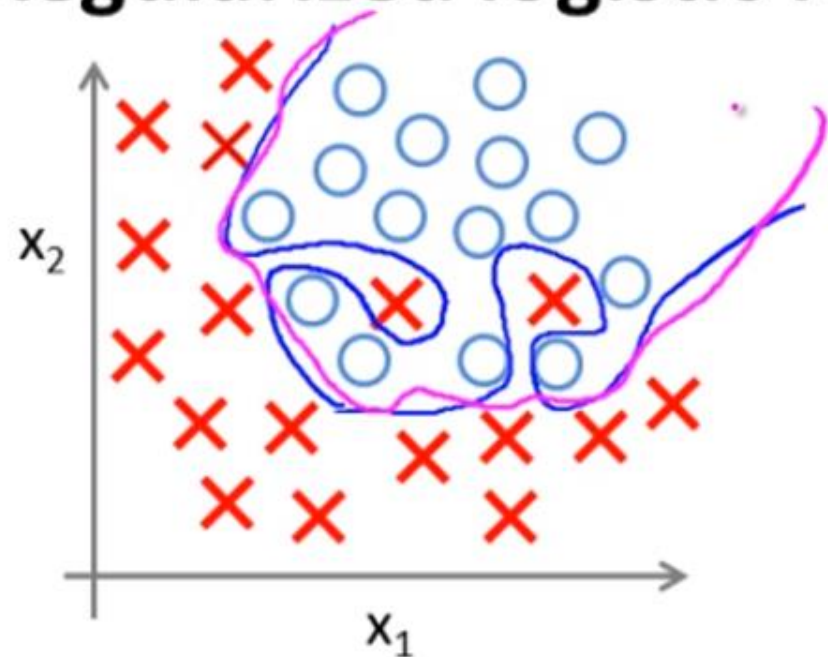
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

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

Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.

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:

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

Windows'u Etkinleştir
Windows'u etkinleştirmek için Ayarlar'a gidin.