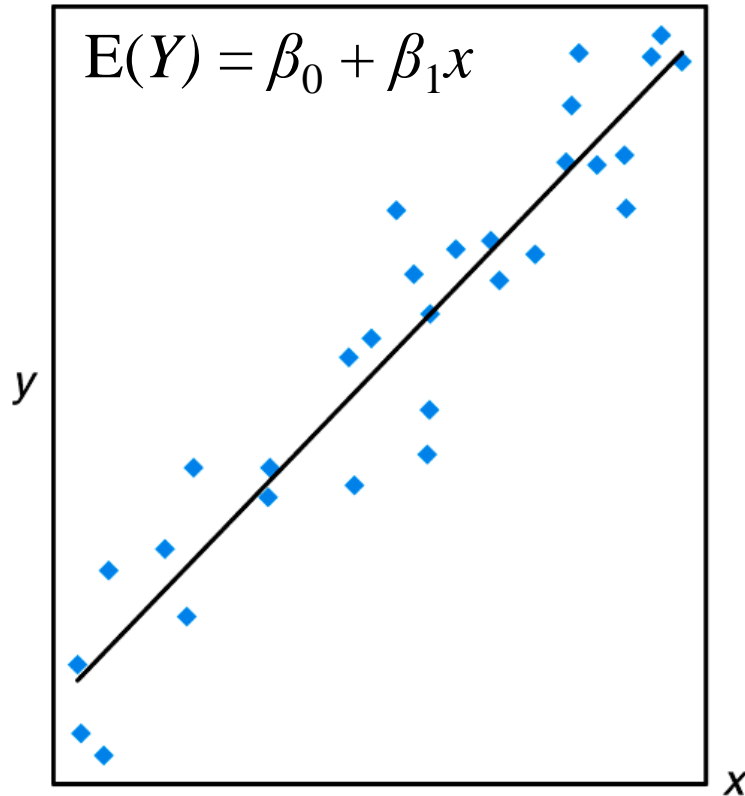


The Problem of Overfitting

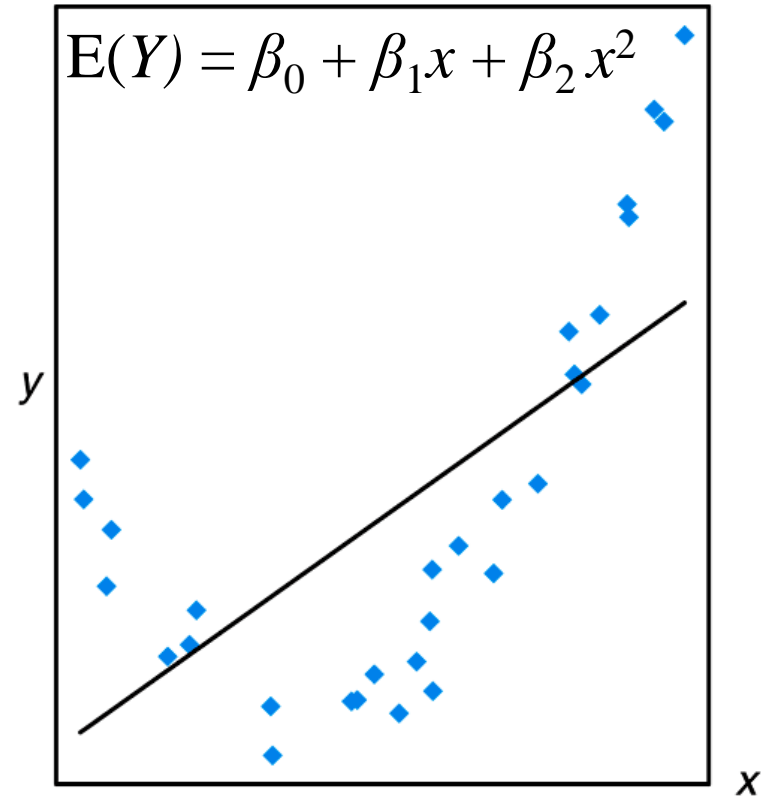
Solving the Problem of Overfitting

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

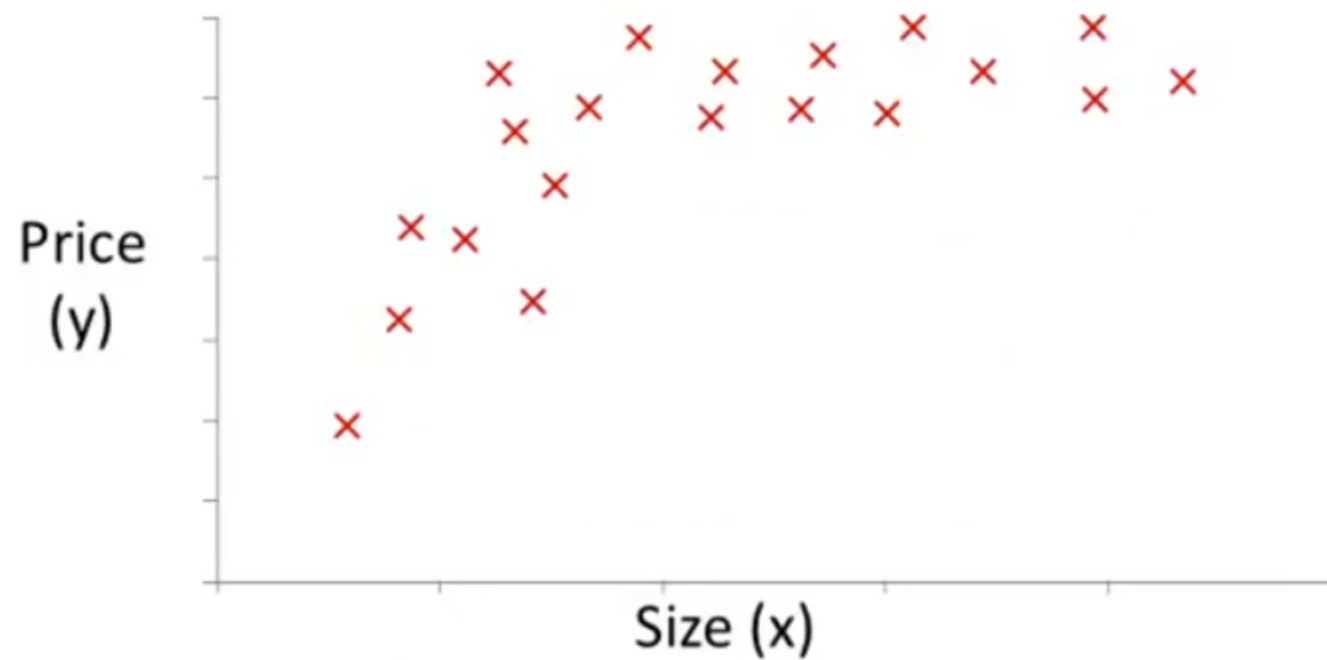
$$x = \underline{\text{frontage} \times \text{depth}}$$

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

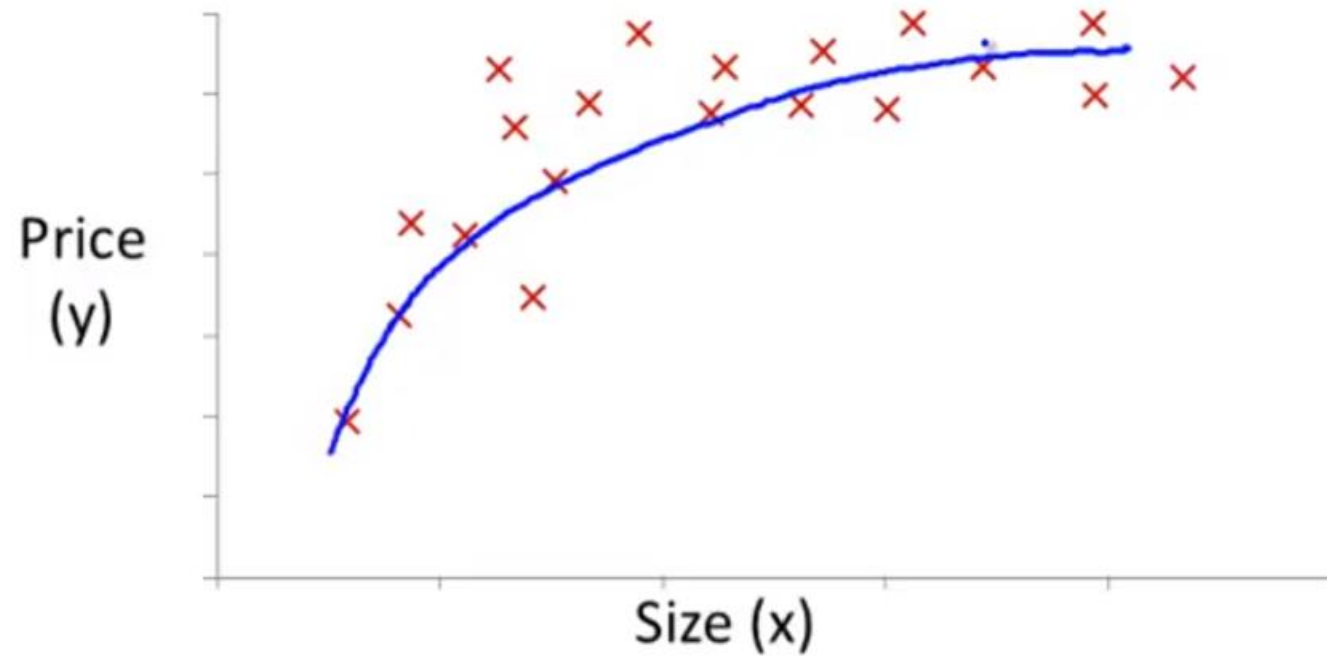
↖ land area



Polynomial regression

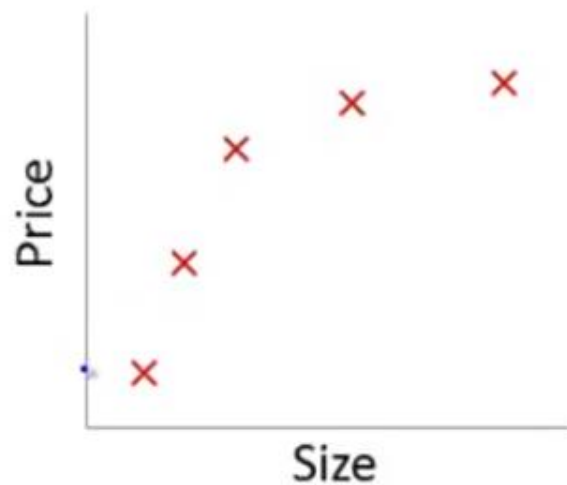


Polynomial regression

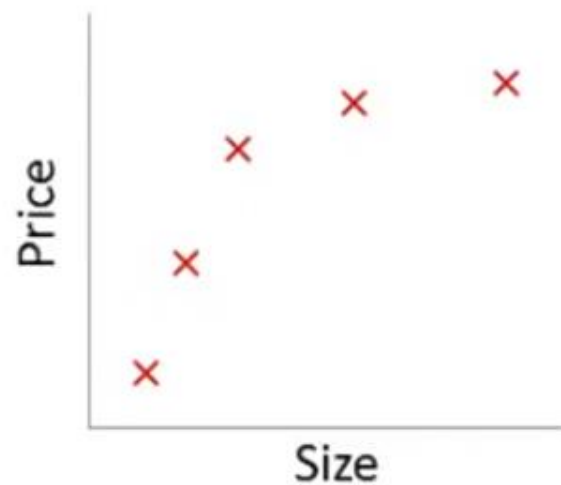


$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$

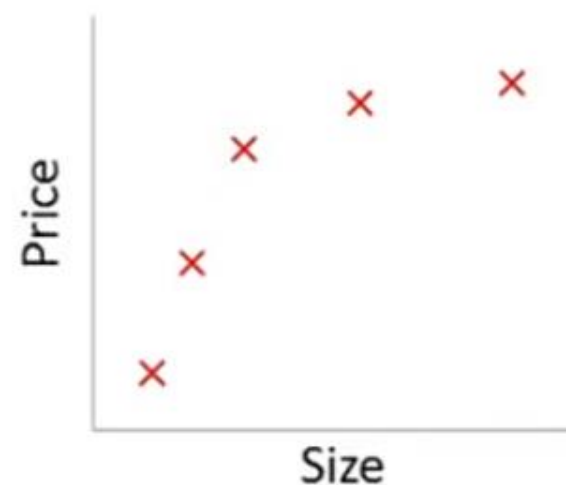
Example: Linear regression (housing prices)



→ $\theta_0 + \theta_1 x$

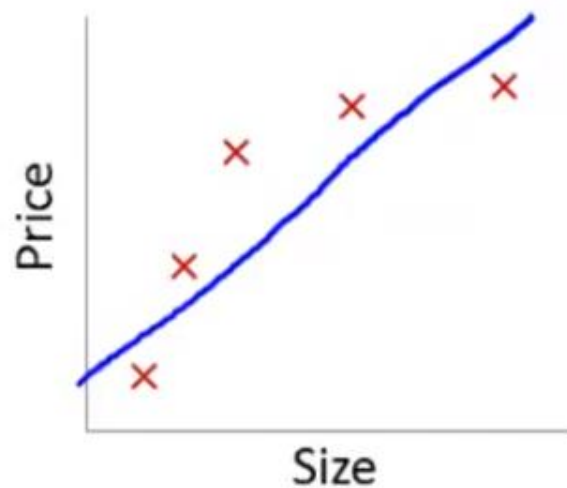


$\theta_0 + \theta_1 x + \theta_2 x^2$

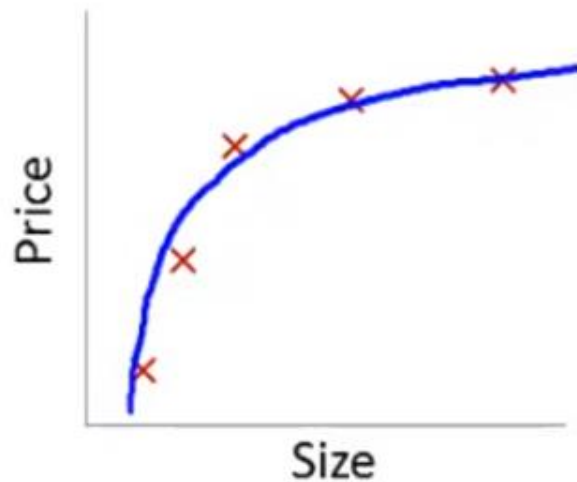


$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

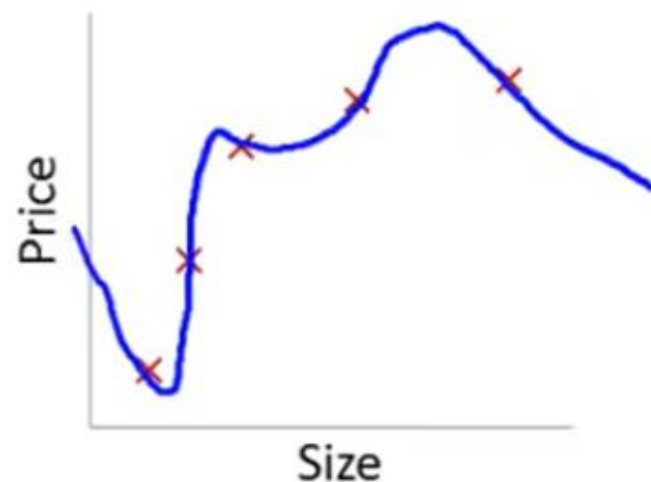
Example: Linear regression (housing prices)



→ $\theta_0 + \theta_1 x$

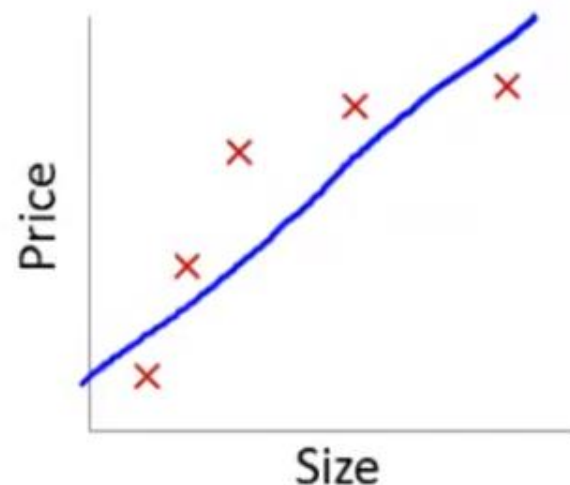


→ $\theta_0 + \theta_1 x + \theta_2 x^2$

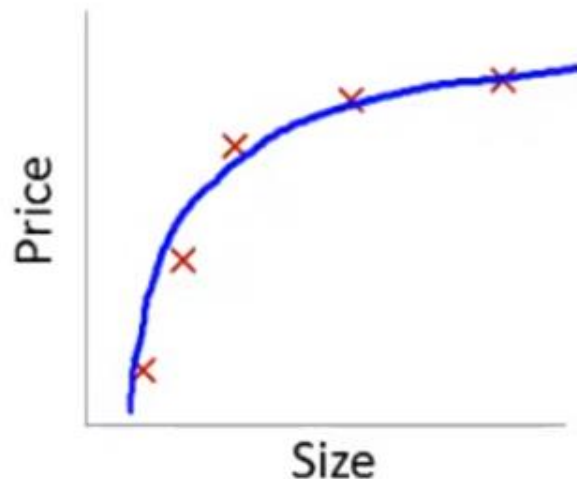


→ $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

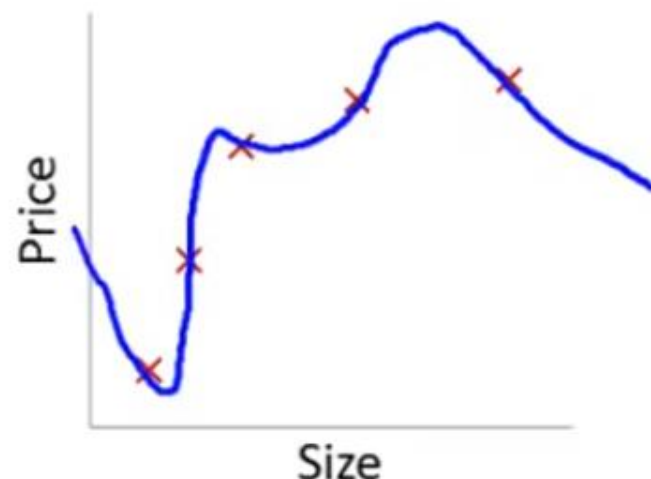
Example: Linear regression (housing prices)



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



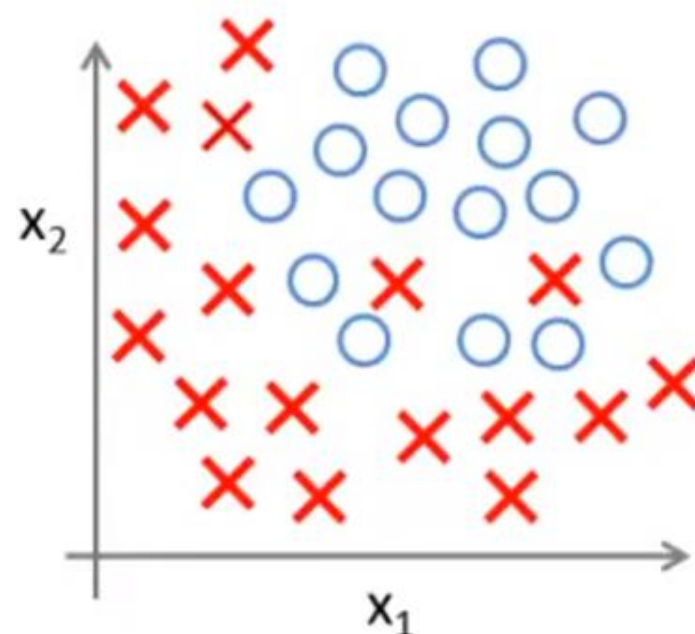
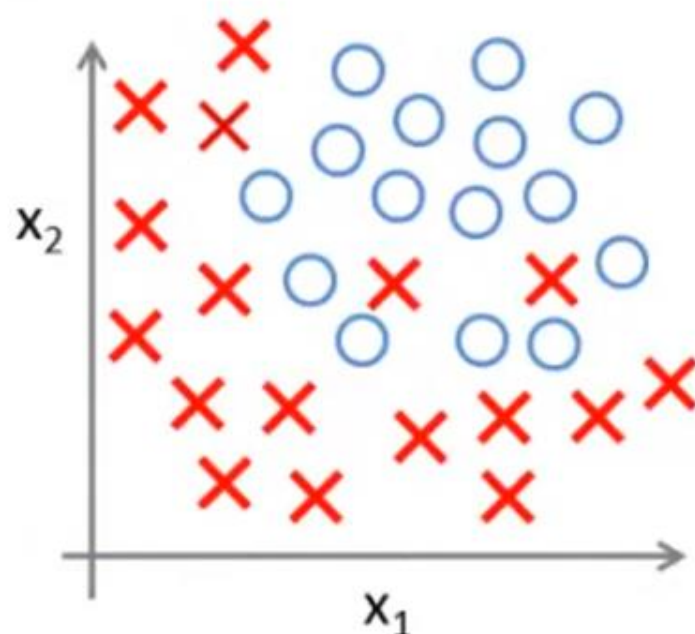
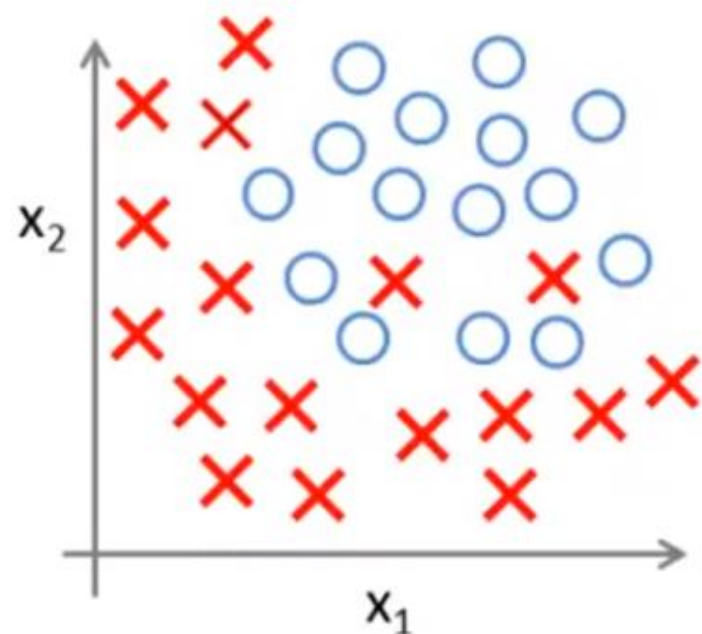
$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$
"Just right"



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

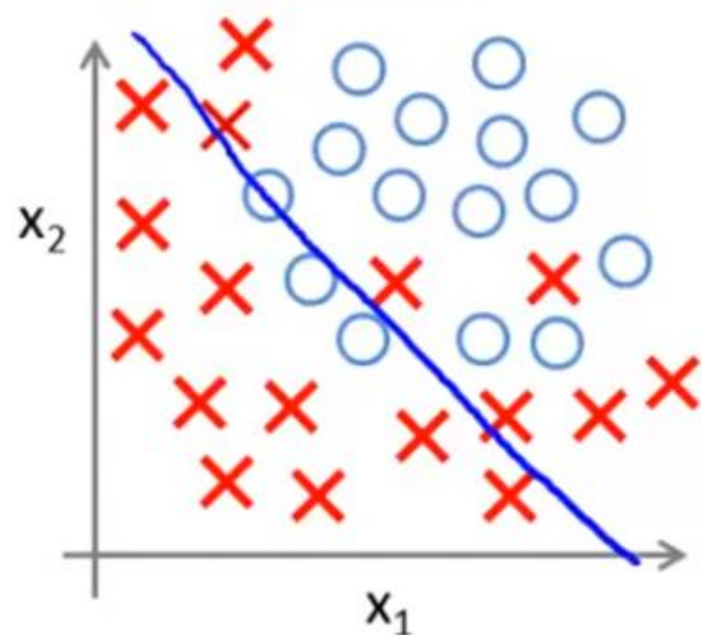


• $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$
 (g = sigmoid function)

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

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

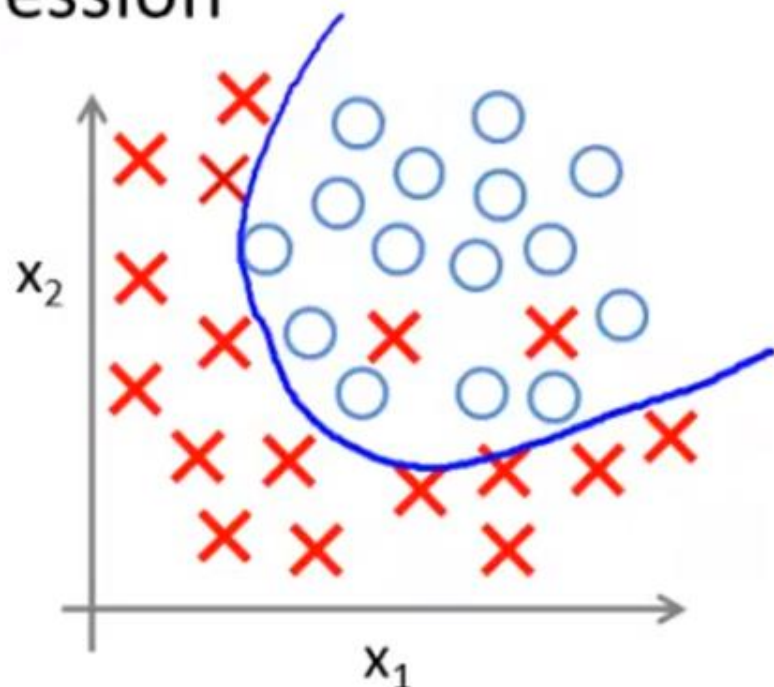
Example: Logistic regression



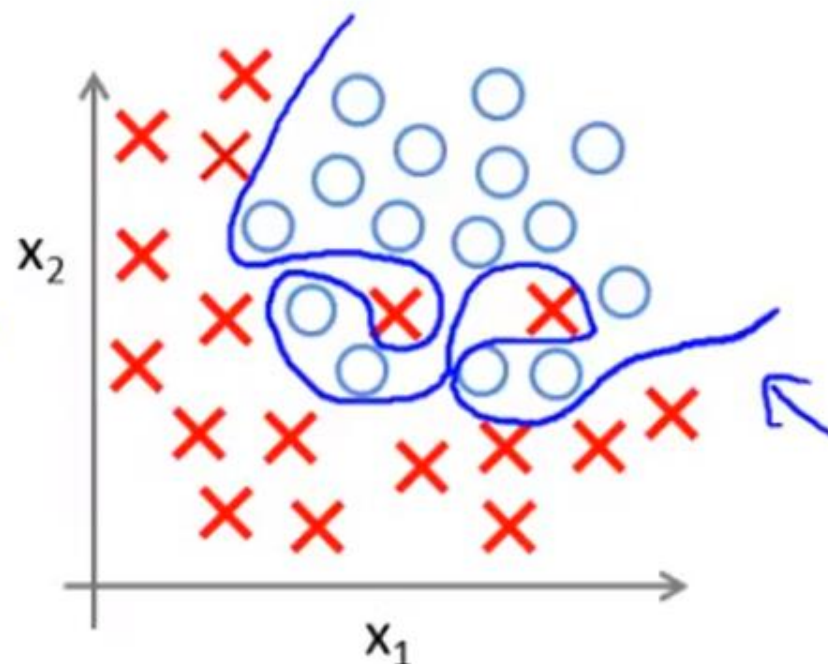
$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

(g = sigmoid function)

"Underfit"



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



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

"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

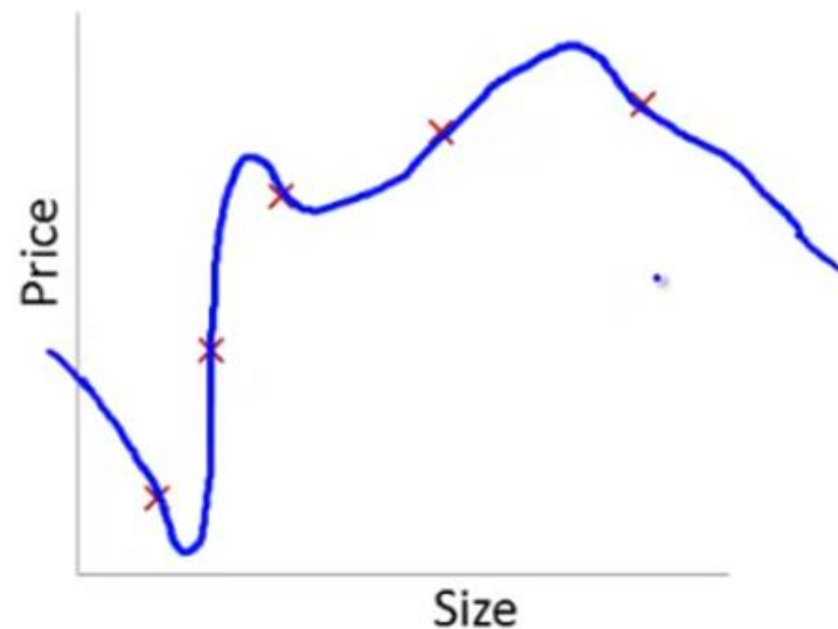
x_4 = age of house

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 .

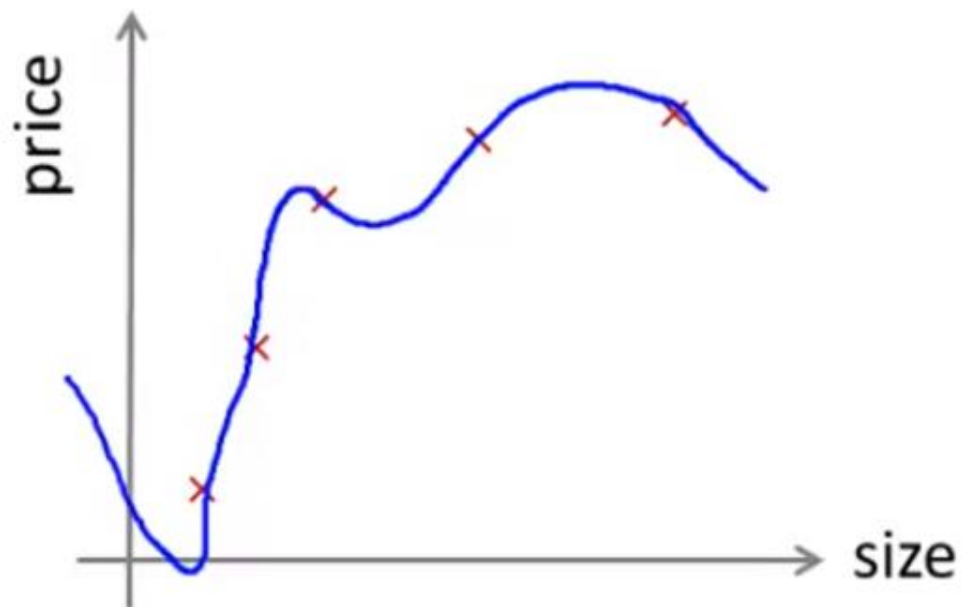
Later...

Evaluating a Hypothesis

Train Set, Test Set

Evaluating a Learning Algorithm

Advice for Applying Machine Learning

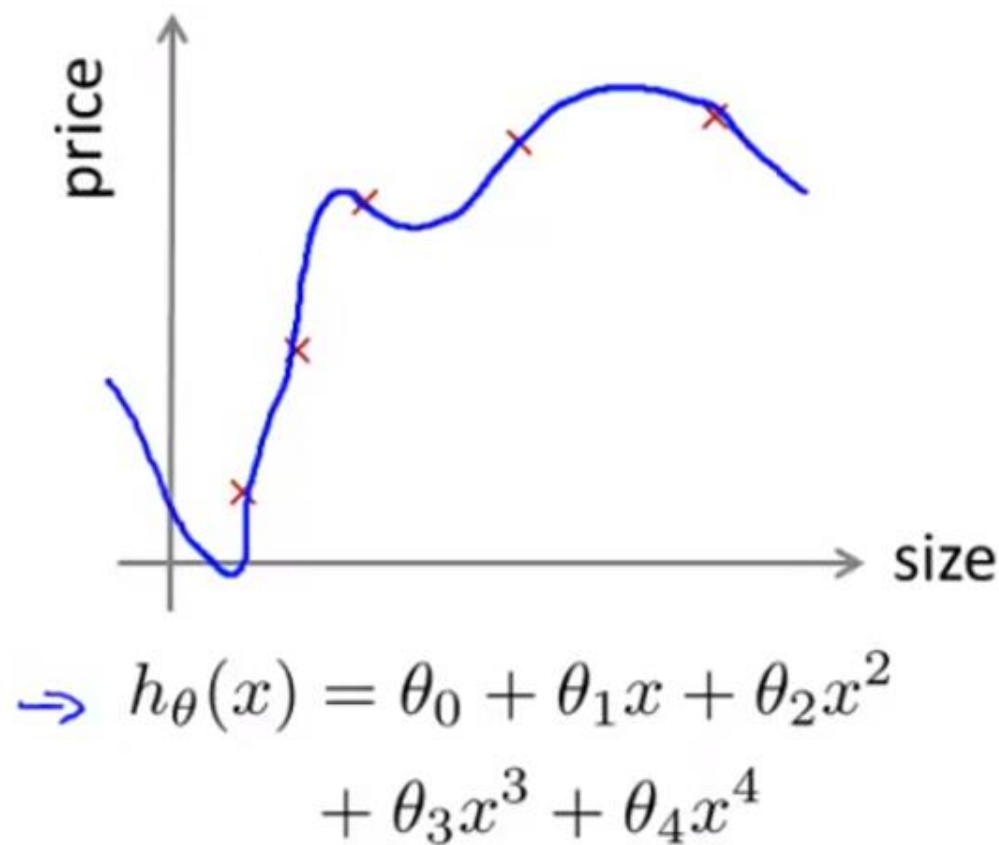


→
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

How to plot many features?

Evaluating your hypothesis

Fails to generalize to new examples not in training set.



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}

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	} Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	
<hr/>		
1427	199	} Test set
1380	212	
1494	243	

20%

30%

$(x^{(1)}, y^{(1)})$
 $(x^{(2)}, y^{(2)})$
 \vdots
 $(x^{(m)}, y^{(m)})$

$(x_{test}^{(1)}, y_{test}^{(1)})$
 $(x_{test}^{(2)}, y_{test}^{(2)})$
 \vdots
 $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

Exercise

- Suppose an implementation of linear regression (without regularization) is badly overfitting the training set.
- In this case, we would expect:
 - The training error $J(\theta)$ to be **low** and the test error $J_{\text{test}}(\theta)$ to be **high**
 - The training error $J(\theta)$ to be **low** and the test error $J_{\text{test}}(\theta)$ to be **low**
 - The training error $J(\theta)$ to be **high** and the test error $J_{\text{test}}(\theta)$ to be **low**
 - The training error $J(\theta)$ to be **high** and the test error $J_{\text{test}}(\theta)$ to be **high**

Procedure

- Training/testing procedure for linear regression
 - Learn parameter θ from training data by minimizing the training error $J(\theta)$
 - Compute the test error for linear regression:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left(\underbrace{h_{\theta}(x_{\text{test}}^{(i)})}_{\uparrow} - y_{\text{test}}^{(i)} \right)^2$$

- Compute the test error for **logistic regression**

$$\rightarrow \underline{J_{\text{test}}(\theta)} = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log h_{\theta}(x_{\text{test}}^{(i)})$$

Training/testing procedure for **logistic regression**

→ - Learn parameter θ from training data

m_{test}

- Compute test set error:

→
$$\underline{J_{test}(\theta)} = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y_{test}^{(i)} \log h_{\theta}(x_{test}^{(i)}) + (1 - y_{test}^{(i)}) \log h_{\theta}(x_{test}^{(i)})$$

- Misclassification error (0/1 misclassification error):

$$err(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq \underline{0.5}, y = \underline{0} \\ & \text{or if } h_{\theta}(x) < \underline{0.5}, y = \underline{1} \end{cases} \text{ error}$$

0 otherwise

$$\text{Test error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_{\theta}(x_{test}^{(i)}), y^{(i)}).$$

Model Selection and Train/Validation/Test Sets

Evaluating a Learning Algorithm

Advice for Applying Machine Learning

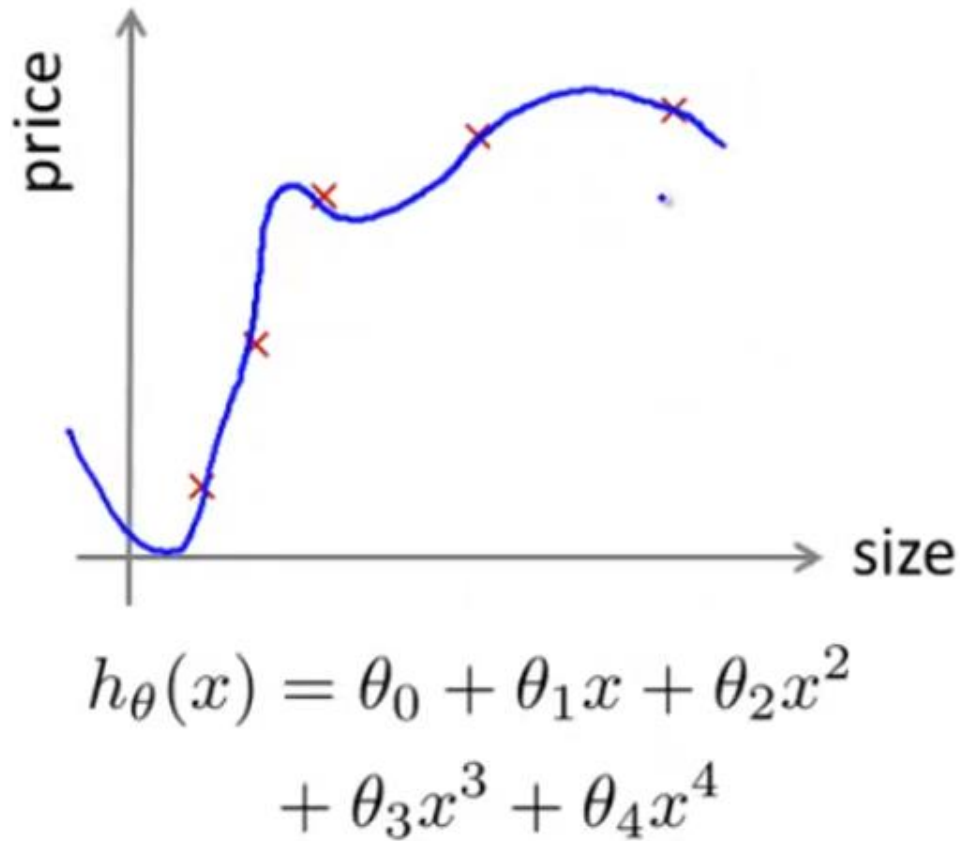
Introduction

- Suppose you are left to decide what degree of polynomial to fit to a data set.
- So that what features to include that gives you a learning algorithm.
- Or suppose you'd like to choose the regularization parameter λ for learning algorithm
- These are called **model selection** problems.

Introduction

- We've already seen a lot of times the problem of overfitting, in which just because a learning algorithm fits a training set well, that doesn't mean it's a good hypothesis.
- More generally, this is why the training set's error is not a good predictor for how well the hypothesis will do on new example.

Overfitting example



Once parameters $\theta_0, \theta_1, \dots, \theta_4$ were fit to some set of data (training set), the error of the parameters as measured on that data (the training error $J(\theta)$) is likely to be lower than the actual generalization error.

Therefore it's an optimistic estimate for the real life error.

Model selection

1. $h_{\theta}(x) = \theta_0 + \theta_1 x$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$

$d = \text{degree of polynomial}$ ↓

Model selection

$$d=1 \quad 1. \quad \underline{h_{\theta}(x) = \theta_0 + \theta_1 x} \quad \longrightarrow \quad \Theta^{(1)} \quad \longrightarrow \quad J_{\text{test}}(\Theta^{(1)})$$

$$d=2 \quad 2. \quad \underline{h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2} \quad \longrightarrow \quad \Theta^{(2)} \quad \longrightarrow \quad J_{\text{test}}(\Theta^{(2)})$$

$$d=3 \quad 3. \quad \underline{h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3} \quad \longrightarrow \quad \Theta^{(3)} \quad \longrightarrow \quad J_{\text{test}}(\Theta^{(3)})$$

⋮

⋮

⋮

$$d=10 \quad 10. \quad \underline{h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}} \quad \longrightarrow \quad \Theta^{(10)} \quad \longrightarrow \quad J_{\text{test}}(\Theta^{(10)})$$

Choose $\underline{\theta_0 + \dots + \theta_5 x^5}$ ←

How well does the model generalize? Report test set error $J_{\text{test}}(\theta^{(5)})$.

→ $d = \text{degree of polynomial}$

Model selection

- $d=1$ 1. $\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \Theta^{(1)} \rightarrow J_{\text{test}}(\Theta^{(1)})$
- $d=2$ 2. $\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \Theta^{(2)} \rightarrow J_{\text{test}}(\Theta^{(2)})$
- $d=3$ 3. $\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \Theta^{(3)} \rightarrow J_{\text{test}}(\Theta^{(3)})$
- \vdots
- $d=10$ 10. $\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \Theta^{(10)} \rightarrow J_{\text{test}}(\Theta^{(10)})$

Choose $\theta_0 + \dots + \theta_5 x^5 \leftarrow$

How well does the model generalize? Report test set error $J_{\text{test}}(\theta^{(5)})$.

Problem: $J_{\text{test}}(\theta^{(5)})$ is likely to be an optimistic estimate of generalization error. I.e. our extra parameter ($d = \text{degree of polynomial}$) is fit to test set.

Evaluating your hypothesis

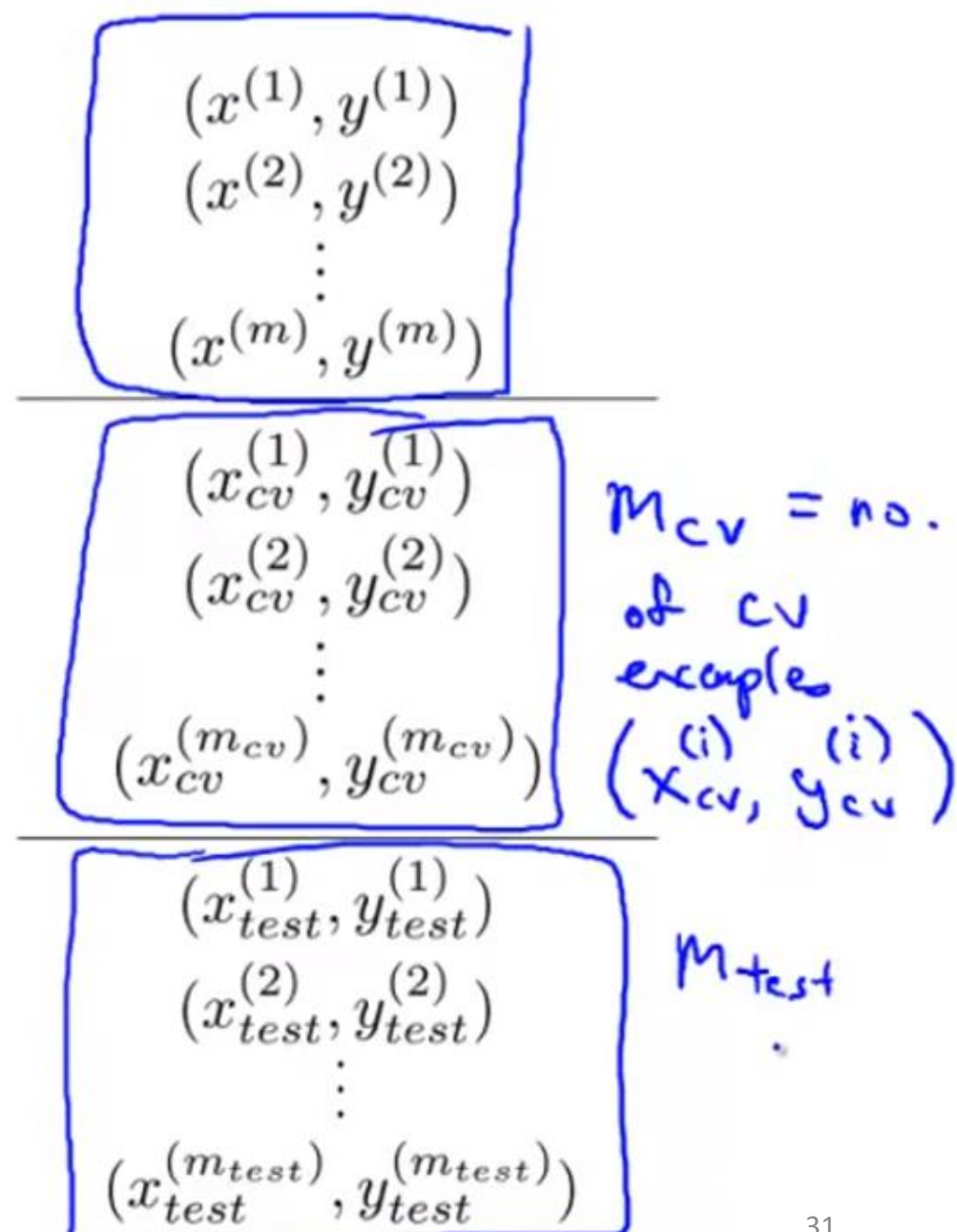
Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243

Evaluating your hypothesis

Dataset:

	Size	Price	
	2104	400	60% Training set
	1600	330	
	2400	369	
	1416	232	
	3000	540	
	1985	300	
	1534	315	20% Cross validation set (CV)
	1427	199	
	1380	212	20% Test set
	1494	243	



Train/validation/test error

Training error:

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

Cross Validation error:

$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

$$\rightarrow J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Model selection

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

2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_3 x^3$

\vdots

10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \cdots + \theta_{10} x^{10}$

Model selection CORRECTED

1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)}$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)}$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)}$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)}$

Find theta's using the **TRAINING** set,
i.e., find thetas of all models
that minimizes
the error of the **TRAINING** set.

Model selection

1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(4)})$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{cv}(\theta^{(10)})$

Model selection

$$\begin{array}{ll}
 d=1 & 1. \quad h_{\theta}(x) = \theta_0 + \theta_1 x \quad \longrightarrow \quad \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)}) \\
 d=2 & 2. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \quad \longrightarrow \quad \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)}) \\
 d=3 & 3. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \quad \longrightarrow \quad \theta^{(3)} \rightarrow J_{cv}(\theta^{(4)}) \\
 & \vdots \\
 d=10 & 10. \quad h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \quad \longrightarrow \quad \theta^{(10)} \rightarrow J_{cv}(\theta^{(4)})
 \end{array}$$

$\underline{d=4} \quad \nearrow$

Pick $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4 \leftarrow$

Estimate generalization error for test set $J_{test}(\theta^{(4)})$ \leftarrow

Exercise

- Consider the model selection procedure where we choose the degree of polynomial using a cross validation set. For the final model (with parameters θ), we might generally expect $J_{CV}(\theta)$ to be lower than $J_{test}(\theta)$
 - An extra parameter (d , the degree of the polynomial) has been fit to the cross validation set.
 - An extra parameter (d , the degree of the polynomial) has been fit to the test set.
 - The cross validation set is usually smaller than the test set.
 - The cross validation set is usually larger than the test set.