

# Course overview

---

1. Introduction to Linear Regression
2. Simple Linear Regression
3. Multiple Linear Regression
4. Evaluation of a Linear Regression Model
5. Practical Work



# What is Linear Regression ?

- Goal: predict real (continuous) valued outputs, by modeling how our observations that are associated with some features change as we change the values of these features.
- Example: training set of housing prices

$X$  = input, features, covariate, predictors

$y$  = output, target

$m$  = number of training examples

		Size (m2)	Built Year	Nb bathrooms	.....	Sale price (k\$)
$x^{(1)}$		200	2010	2	....	$y^{(1)} = 800$
$x^{(2)}$		300	1995	2	.....	$y^{(2)} = 750$
.....		.....	.....	.....	.....	.....

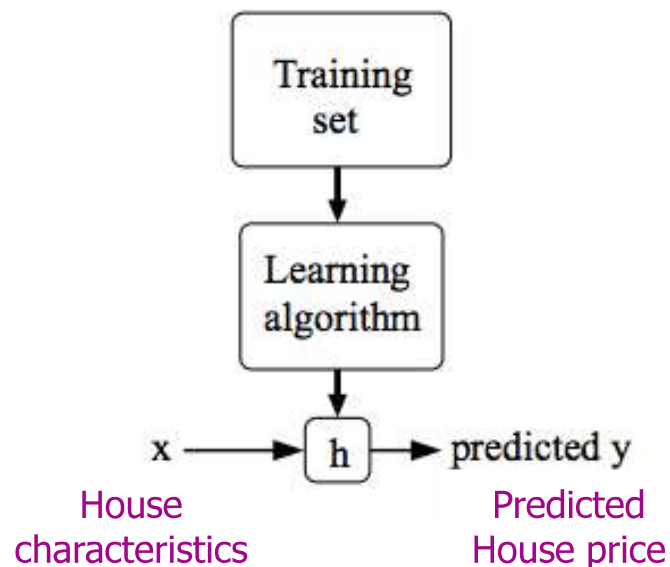
- Regression is about learning the relationship between  $X$  and  $y$ , and using it to predict the house price of new data

	250	2005	3	.....	???????
---	-----	------	---	-------	---------

# Model Representation

---

- $x^{(i)}$  denotes the “input” variables (house characteristics)
- $y^{(i)}$  denotes the “output” or target variable that we are trying to predict (price)
- A pair  $(x^{(i)}, y^{(i)})$  is called a training example
- A list of  $m$  training examples  $(x^{(i)}, y^{(i)}); i=1, \dots, m$ —is called a training set



$$\mathbf{h}: \mathbf{X} \rightarrow \mathbf{Y}$$

Hypothesis or function that takes as input the house's characteristics to estimate its price.

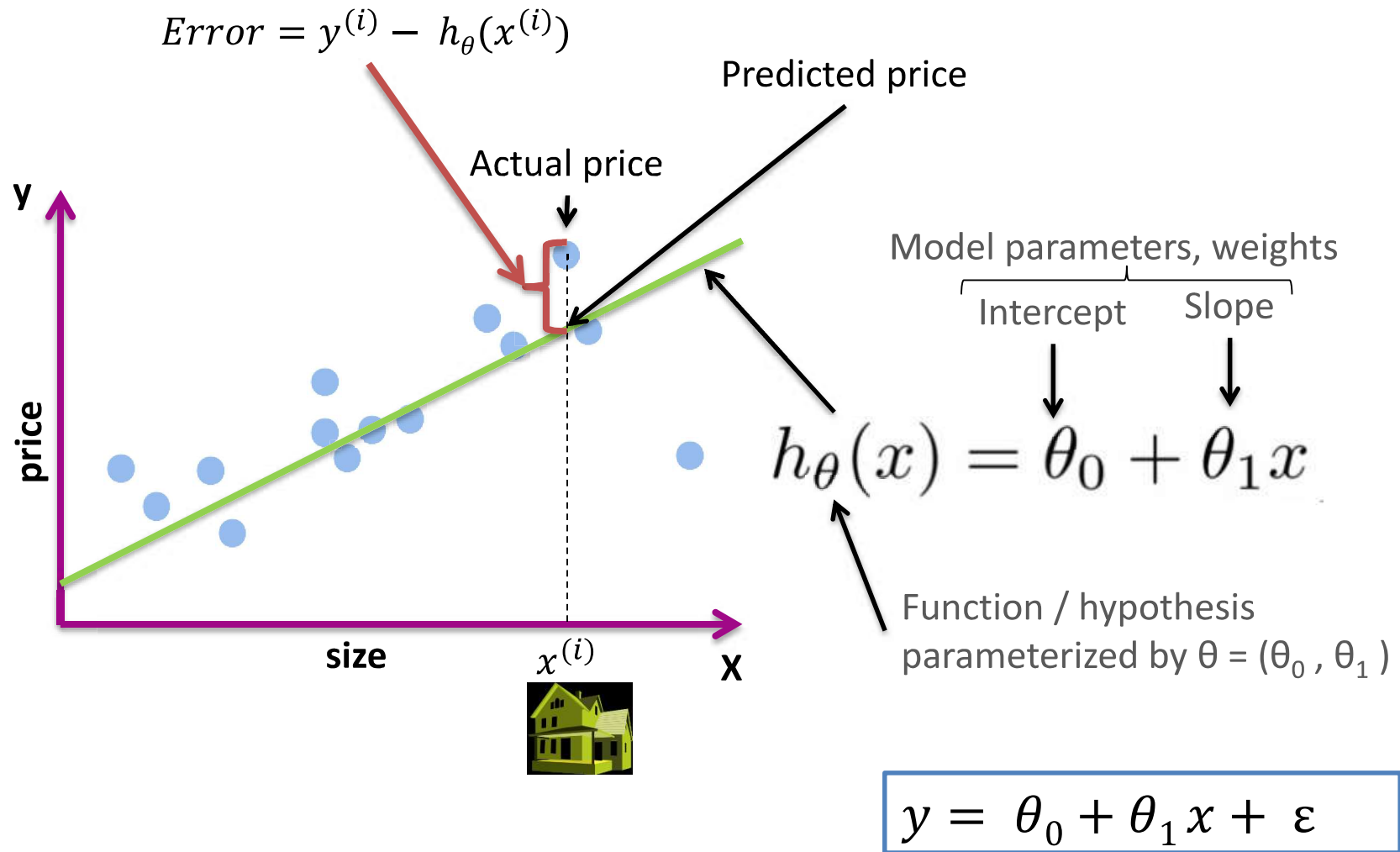
---

# **1.2**

## **Simple/Univariate Linear Regression**

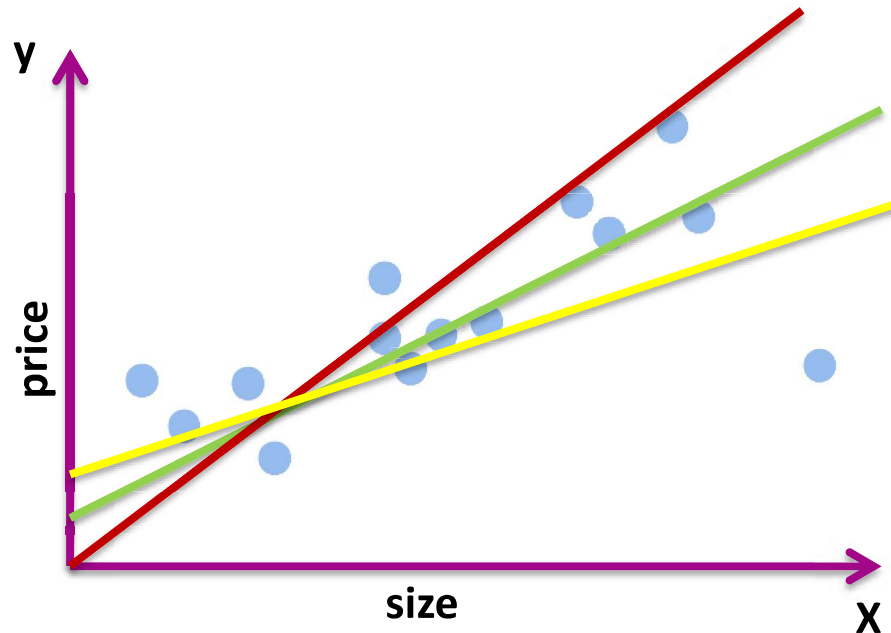
---

# Simple Linear Regression: Model



# Simple Linear Regression: Model Evaluation

---



Each line has different parameters  $\theta = (\theta_0, \theta_1)$   
→ different errors

- Which line is the best fit ?
- what should a good function  $h_{\theta}(x)$  minimize?
  - Sum error on all data points
  - Sum abs(error) on all data points
  - Sum error<sup>2</sup> on all data points

# Simple Linear Regression: Model Evaluation

- Sum error on all data points
- Sum  $\text{abs}(\text{error})$  on all data points
- Sum  $\text{error}^2$  on all data points

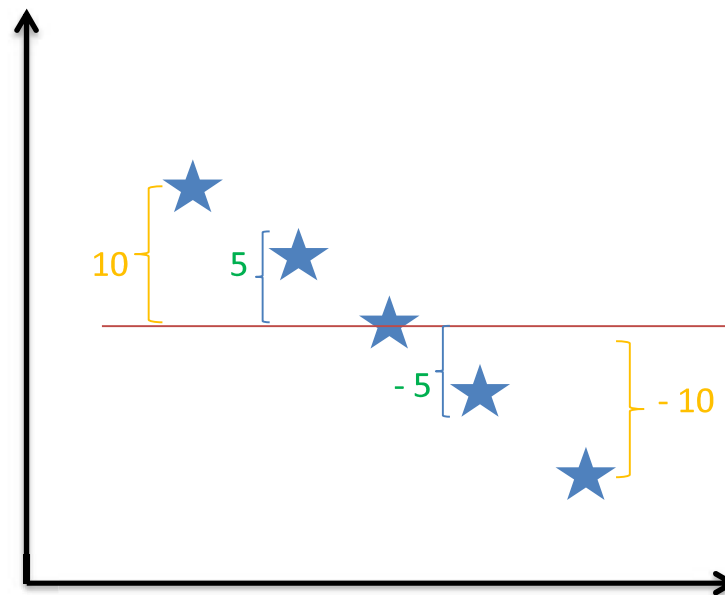
Sum error = 0



Sum  $\text{abs}(\text{error}) > 0$



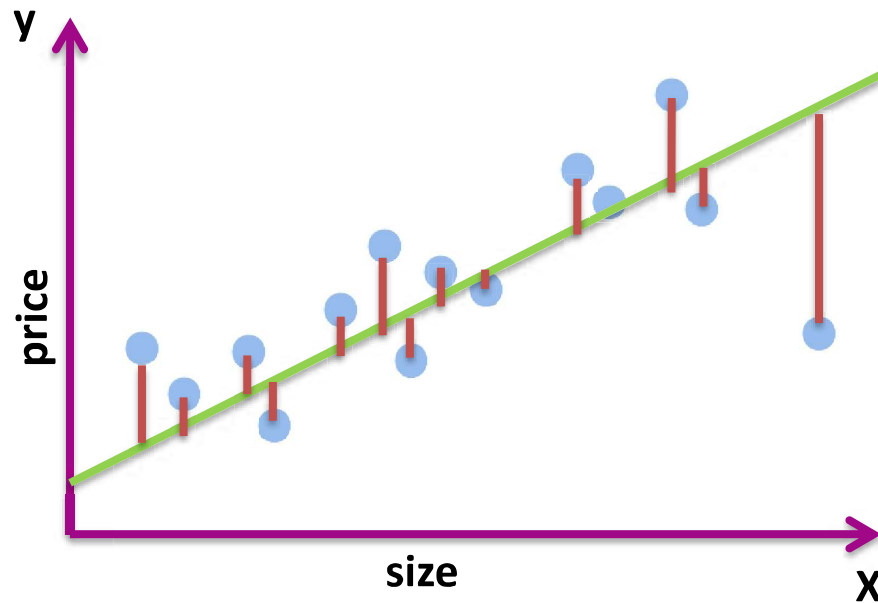
Sum  $\text{error}^2 > 0$



# Simple Linear Regression: Model Evaluation

---

- The Sum of Squared Errors (SSE) is also called Residual Sum of Squares (RSS)



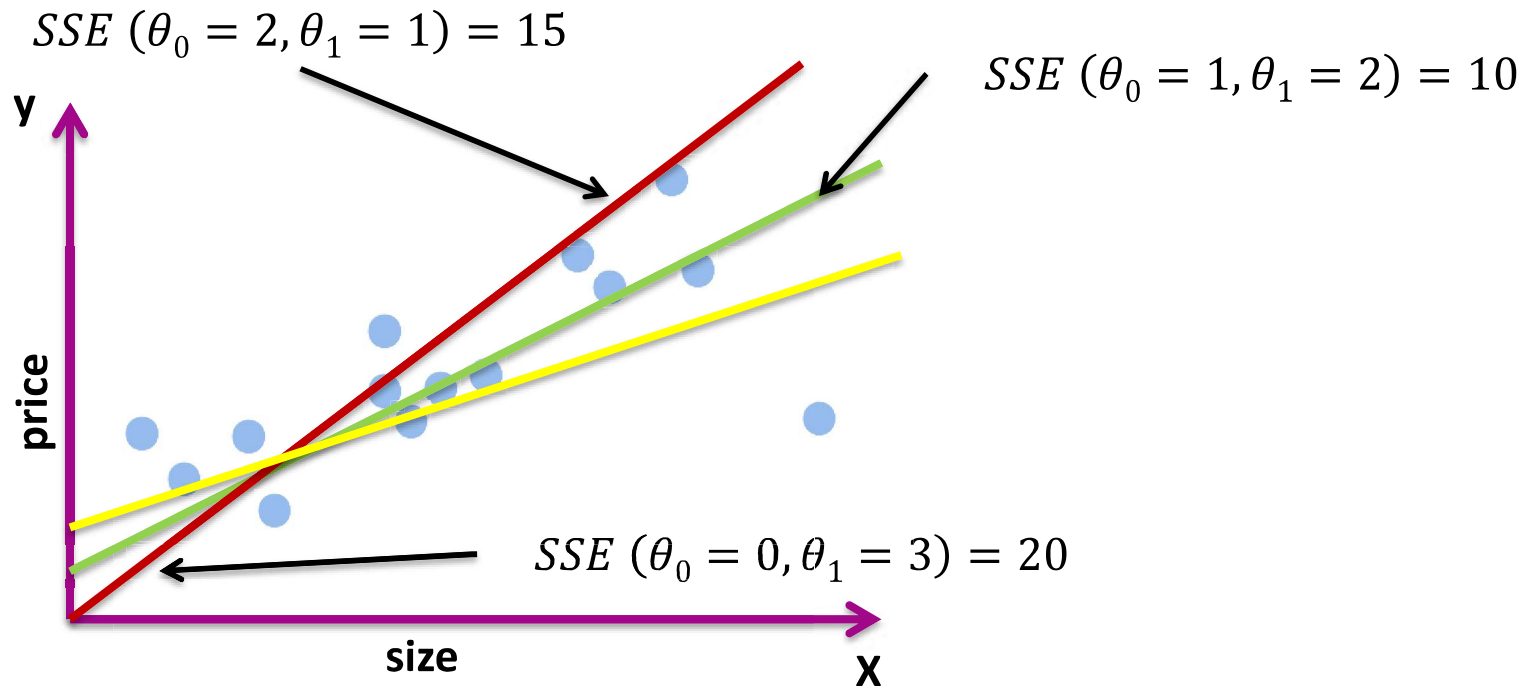
$$SSE(\theta_0, \theta_1) =$$

$$\begin{aligned} & (y^{(1)} - (\theta_0 + \theta_1 * x^{(1)}))^2 + \\ & (y^{(2)} - (\theta_0 + \theta_1 * x^{(2)}))^2 + \\ & \dots\dots\dots + \\ & (y^{(m)} - (\theta_0 + \theta_1 * x^{(m)}))^2 \end{aligned}$$



# Simple Linear Regression: Model Evaluation

---



- The green line is a better fit
- Let's see how to find the best line automatically

# Simple Linear Regression: Cost function

---

- The best hypothesis  $h_{\theta}(x)$  is the one that minimizes the cost function

The diagram shows the cost function formula  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ . Annotations include: 'Sum over all data points' with an arrow pointing to the summation symbol  $\sum$ ; 'Actual' with an arrow pointing to  $y^{(i)}$ ; 'Predicted' with an arrow pointing to  $h_{\theta}(x^{(i)})$ ; and an arrow pointing to the  $1/m$  term.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- $1/m$  - means we determine the average
- $1/2m$ , the 2 makes the math a bit easier, and doesn't change the weights  $\theta$  we determine at all (i.e. half the smallest value is still the smallest value!)
- The learning algorithm should find  $\theta^* = (\theta_0^*, \theta_1^*)$  that minimizes this cost
- Several algorithms:
  - Gradient descent
  - Ordinary least square (OLS): Used in the linear regression in python (sklearn)

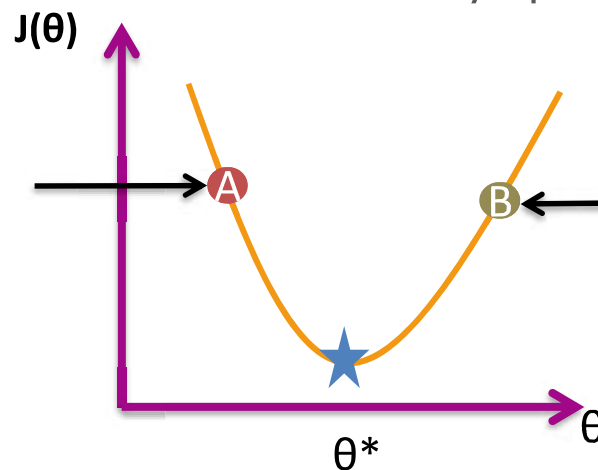
# Gradient Descent Algorithm: Intuition

- Let's assume that we have one parameter  $\theta$ , and that we want to minimize  $J(\theta)$

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

- GD starts by a random initial  $\theta$  and iteratively update it to get towards  $\theta^*$ .

In this case, the derivative  
(gradient)  $\partial J(\theta) / \partial \theta < 0$ .  
 $\theta^A - \alpha * \partial J(\theta) / \partial \theta > \theta^A$   
 $\theta^A$  is moving to the right  
 $\theta$  is increasing



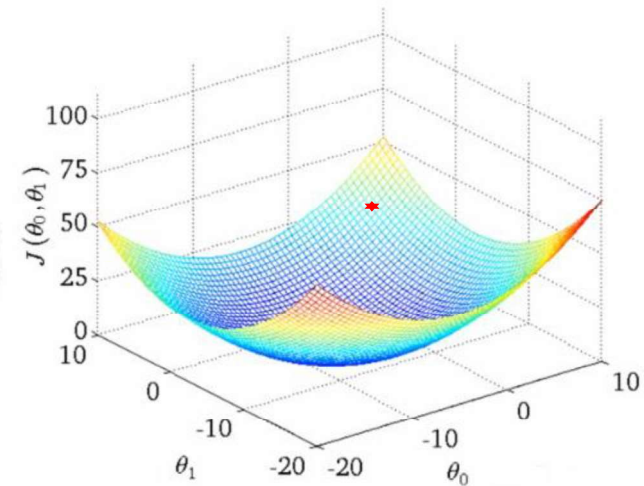
In this case, the derivative  
(gradient)  $\partial J(\theta) / \partial \theta > 0$ .  
 $\theta^B - \alpha * \partial J(\theta) / \partial \theta < \theta^B$   
 $\theta^B$  is moving to the left  
 $\theta$  is decreasing

- $\alpha$  is called the learning rate or the step size
  - Too small
    - Take baby steps  $\rightarrow$  Take too long to converge
  - Too large
    - Can overshoot the minimum  $\rightarrow$  fail to converge

# Simple Linear regression with GD

- Repeat until convergence :

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$



- If we calculate the derivatives, the expression becomes:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \quad \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

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

- Convergence means either:
  - The cost function is no longer changing by more than  $\epsilon$ .
  - The number of iterations is reached
- [Go further](#)

# Ordinary Least Square Algorithm

---

- The Ordinary least square (OLS) approach chooses  $\theta_0, \theta_1$  to minimize SSE

$$SSE = \sum_{i=1}^m (y^{(i)} - \theta_0 - \theta_1 * x^{(i)})^2$$

- In this method, we will minimize SSE by explicitly taking its derivatives with respect to the  $\theta_0, \theta_1$ , and setting them to zero.
- The minimizing values can be shown to be:

$$\hat{\theta}_1 = \frac{\sum_{i=1}^m (x^{(i)} - \bar{x})(y^{(i)} - \bar{y})}{\sum_{i=1}^m (x^{(i)} - \bar{x})^2}$$

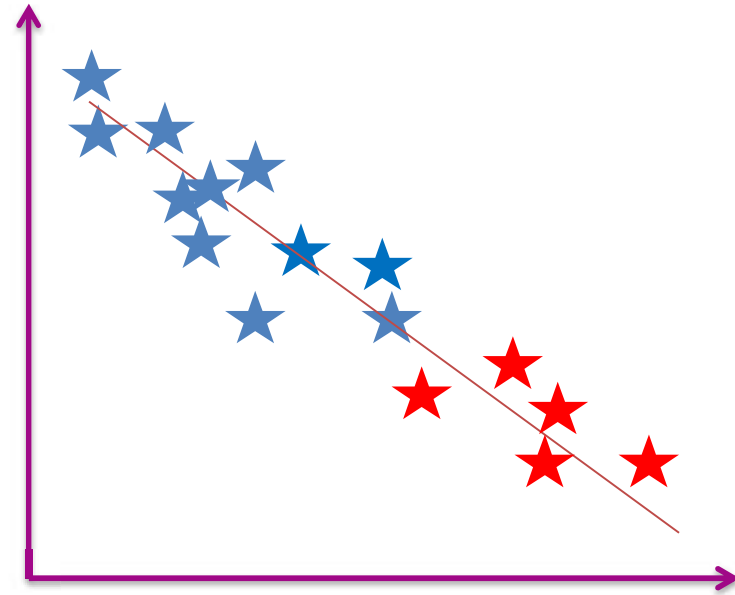
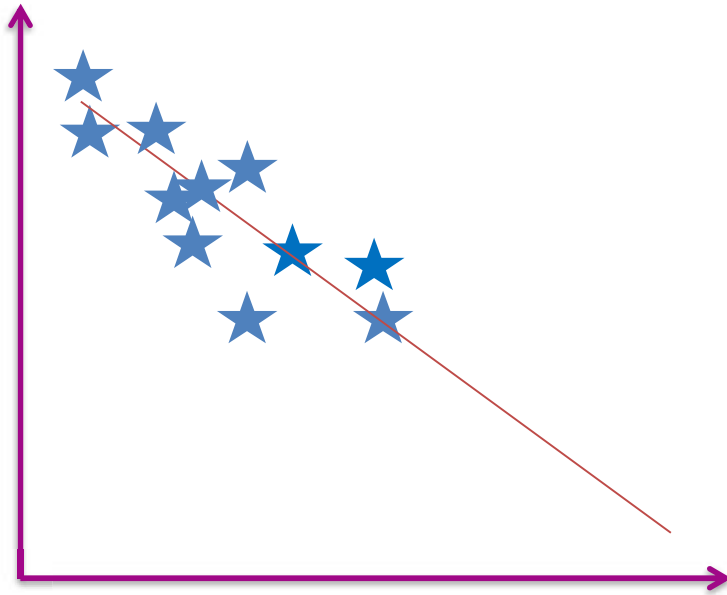
$$\hat{\theta}_0 = \bar{y} - \hat{\theta}_1 * \bar{x}$$

- Where  $\bar{y} = \frac{1}{m} \sum_{i=1}^m y^{(i)}$ ,  $\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$  are the sample means

# SSE is not perfect

---

- Which fit has larger SSE?



- Larger SSE doesn't necessarily mean worst fit → Need an other metric

# R<sup>2</sup> Statistic

---

- R<sup>2</sup> Answers the questions: **“how much of variability in the output (y) is explained by the change in the input (x)”**
- To calculate R<sup>2</sup>, we use the formula:

$$R^2 = \frac{TSS - SSE}{TSS} = 1 - \frac{SSE}{TSS}$$
$$TSS = \sum_{i=1}^m (y^{(i)} - \bar{y})^2$$

- An R<sup>2</sup> statistic that is close to 1 indicates that a large proportion of the variability in the response has been explained by the regression.
- A number near 0 indicates that the regression did not explain much of the variability in the response

---

# **1.2**

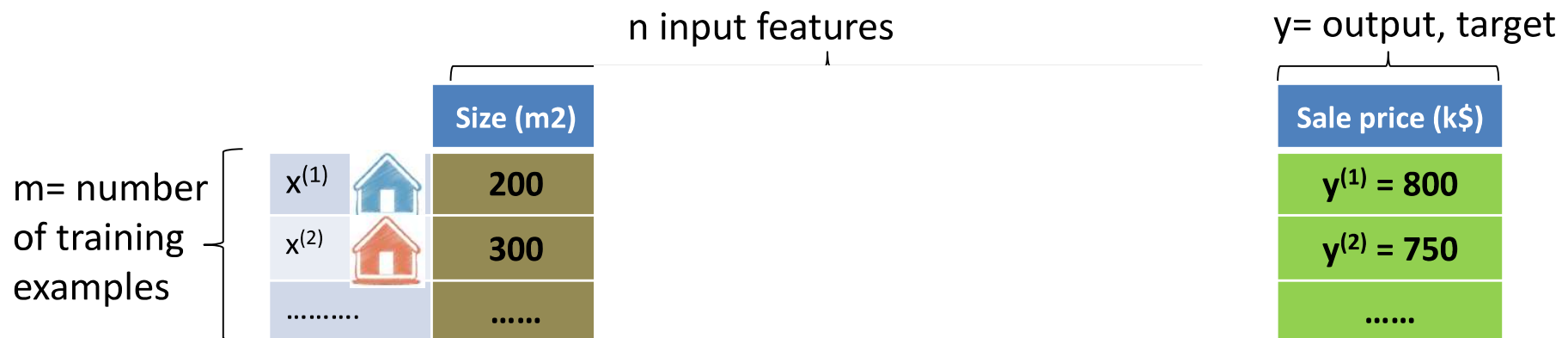
## **Multivariate Linear Regression**

---



# Multiple Linear Regression

- In simple linear regression, we use one feature  $x$  to predict  $(y)$
- In multiple linear regression, we have multiple features  $X=(x_1, x_2, \dots, x_n)$



- $X^{(i)}$  is an  $n$ -dimensional feature vector
- $X^{(1)} = (200, 2010, 2, \dots)^T$  the feature vector of the first training example.
- $x_j^{(i)}$  is the value of feature  $j$  in the  $i^{th}$  training example.  $x_2^{(1)} = 2010$

# Multiple Linear Regression

---

- In simple linear regression,  $h_{\theta}(x) = \theta_0 + \theta_1 x$
- In multiple linear regression,  $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$
- For convenience of notation, define  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- Rewrite in matrix notation

$\theta_0$		$x_0$	$h_{\theta}(x) = \sum_{i=1}^n \theta_i x_i = \theta^T x = \theta x^T$
$\theta_1$		$x_1$	
$\theta_2$		$x_2$	
....		....	
$\theta_n$		$x_n$	

$\theta, x \in R^{n+1}$

# Multiple Linear Regression

- Cost function:

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

## Simple Linear Regression, $n=1$

$$\begin{aligned} &\text{Repeat } \{ \\ &\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)} \\ &\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \\ &\quad \text{(simultaneously update } \theta_0, \theta_1 \text{)} \} \end{aligned}$$

## Multiple Linear Regression, $n>1$

$$\begin{aligned} &\text{Repeat } \{ \\ &\quad \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \\ &\quad \text{(simultaneously update } \theta_j \text{ for } j = 0, \dots, n) \\ &\} \\ &\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \\ &\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \\ &\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \\ &\dots \end{aligned}$$

# OLS for multiple features

---

- Cost function with matrix notations:

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

- Where

$$X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

- Derivatives with respect to  $\theta$   $\nabla_{\theta} J(\theta) = X^T X \theta - X^T y$

- By setting them to zero  $\theta = (X^T X)^{-1} X^T y$

# When to use OLS or GD ?

---

- The following is a comparison of GD and the OLS (normal equation):

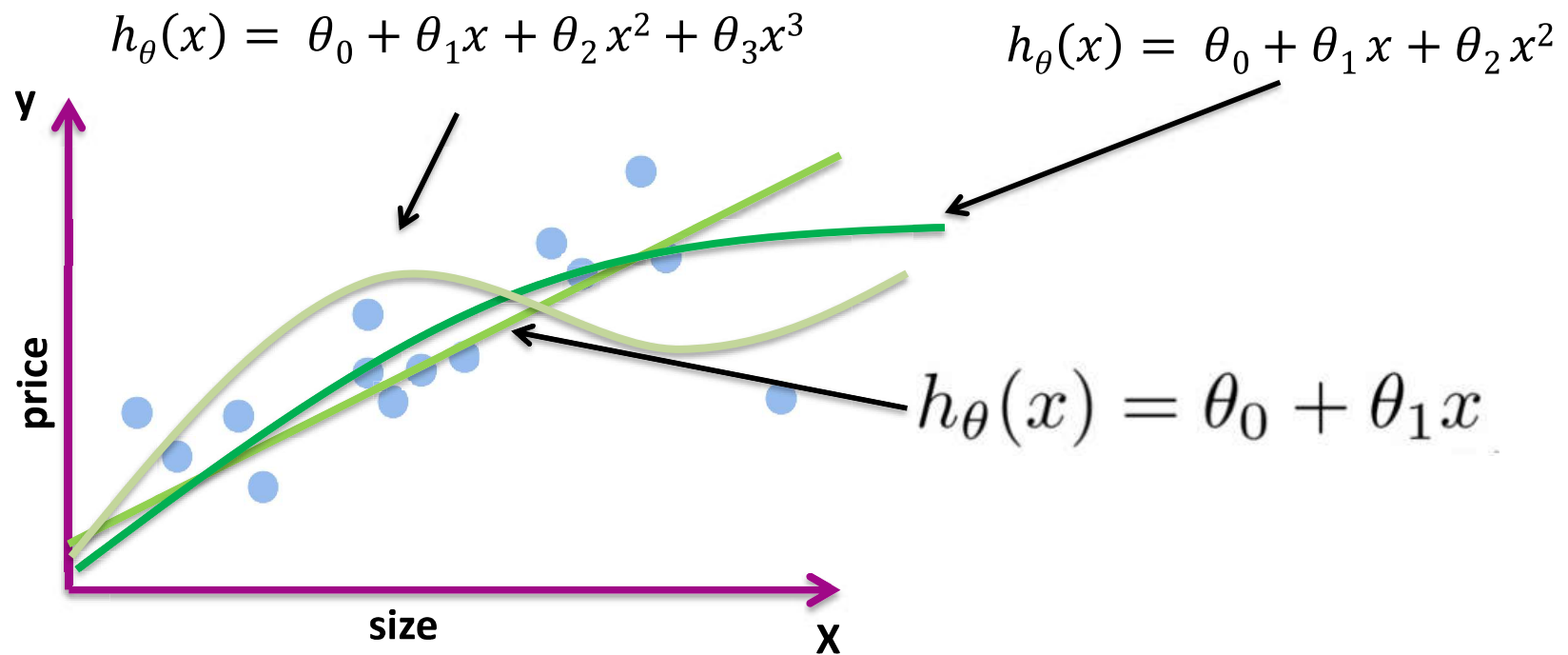
Gradient Descent	OLS (Normal Equation)
Need to choose alpha	No need to choose alpha
Needs feature scaling	No need for feature scaling
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$ , need to calculate inverse of $X^T X$
Works well when $n$ is large	Slow if $n$ is very large

- $n=10^4$ - $10^5$  is usually the threshold of choosing GD over OLS
- Feature scaling helps converting the features to the same scale.
- For example, if  $x_i$  represents housing prices with a range of 100 to 2000 and a mean value of 1000, then,

$$x_i = \frac{price - 1000}{2000 - 100}$$

# Polynomial Regression

- Polynomial regression is a particular case of multiple regression where the features are powers of one single feature  $x$



- General model:

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots + \theta_p x^p + \varepsilon$$

---

# **1.3**

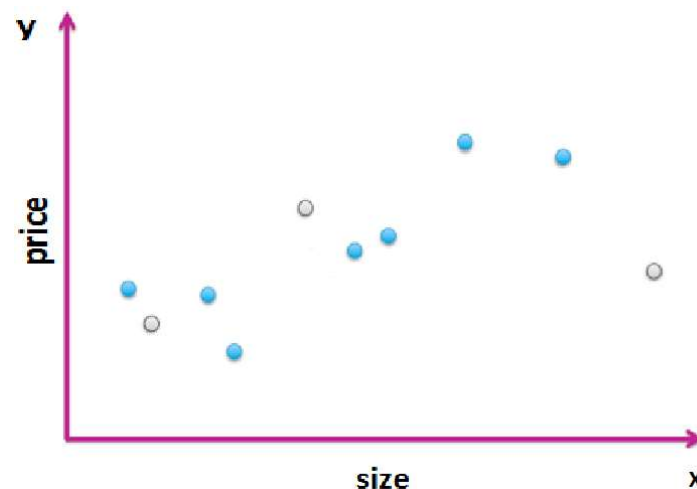
## **Assessing Performance**

---

# Assessing Performance

- This is about knowing how well the model will **generalize** to unseen data.
- One of the most used methods is splitting the data into train/test sets.

	Size	Price
Training set – 70%	2104	400
	1600	330
	2400	369
	1416	232
	3000	540
	1985	300
	1534	315
	1427	199
Test set 30%	1380	212
	1494	243



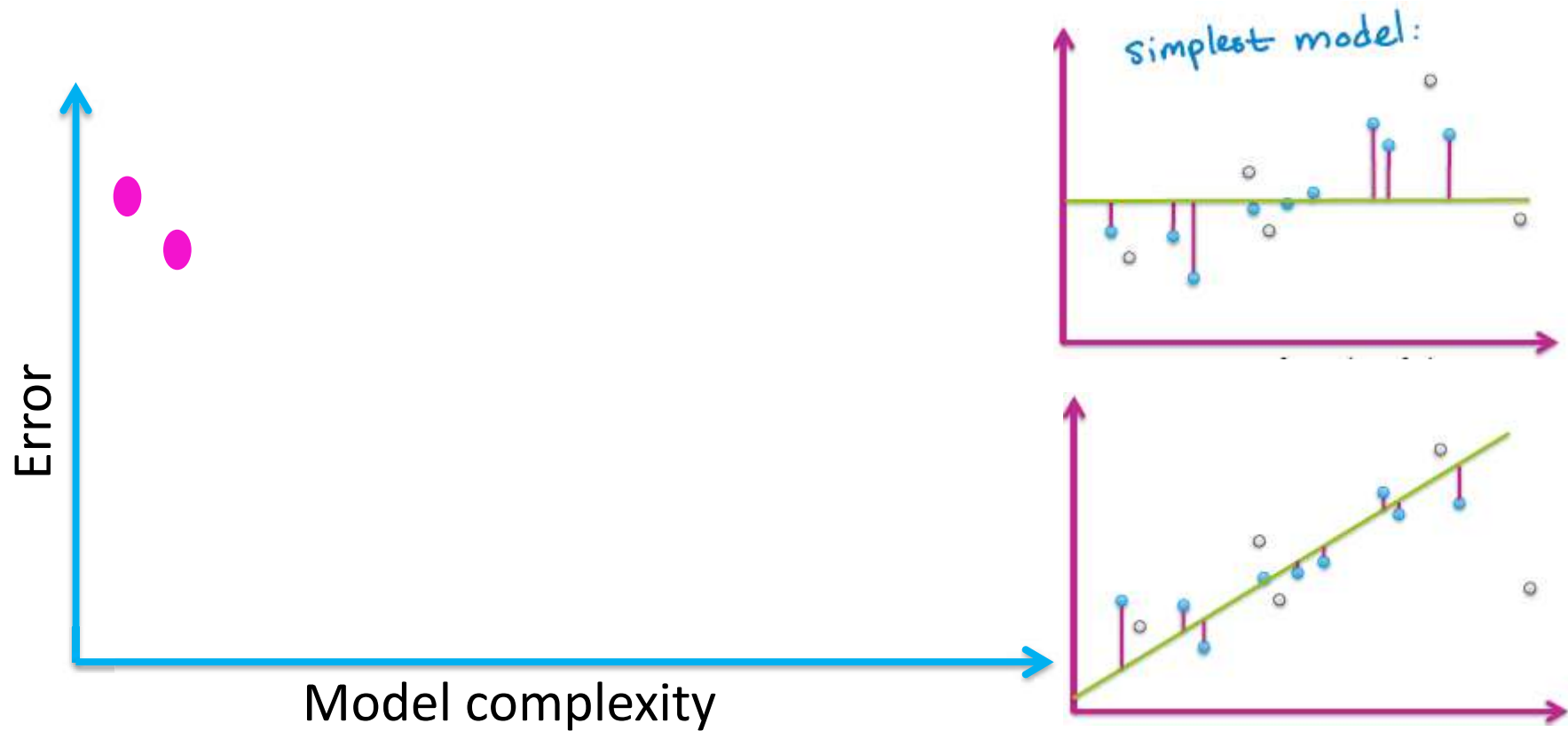
- Learn  $\theta$  from training data (minimizing training error  $J_{\text{train}}(\theta)$ )
- Compute test error  $J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_i^{m_{\text{test}}} (h_{\theta}(x^{(i)}) - y^{(i)})^2$



# Training error vs. model complexity

---

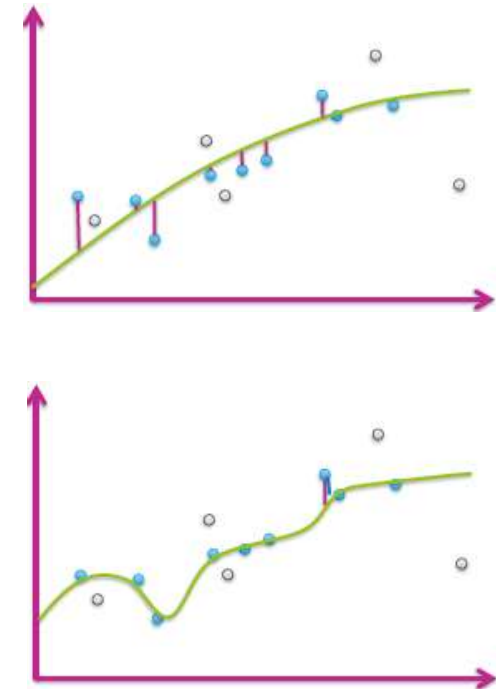
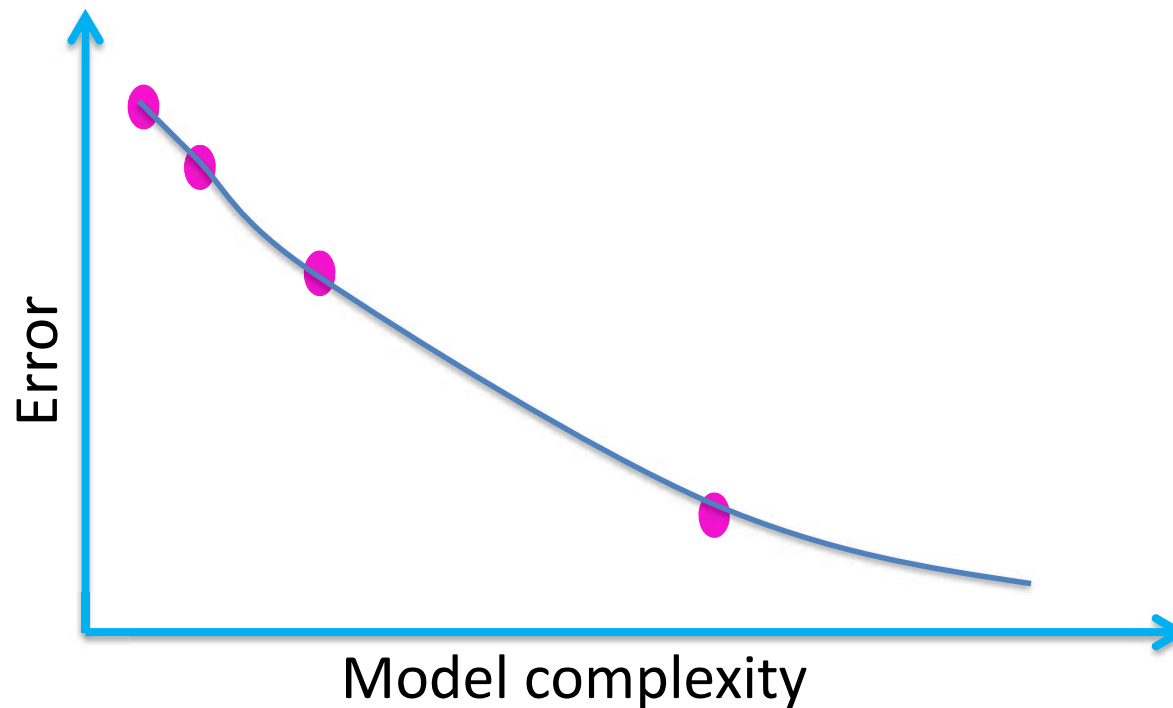
- Let's take polynomial regression as an example



# Training error vs. model complexity

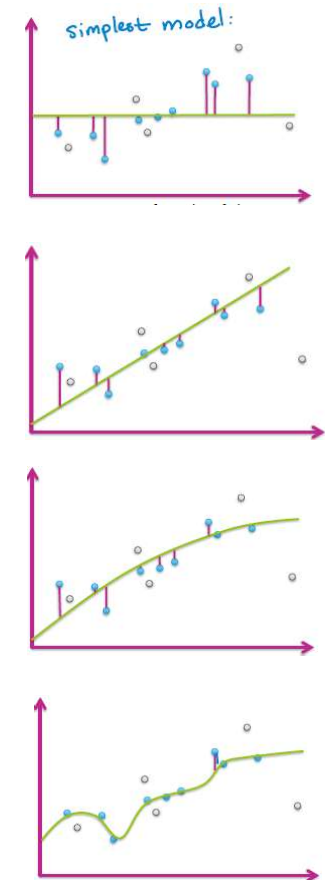
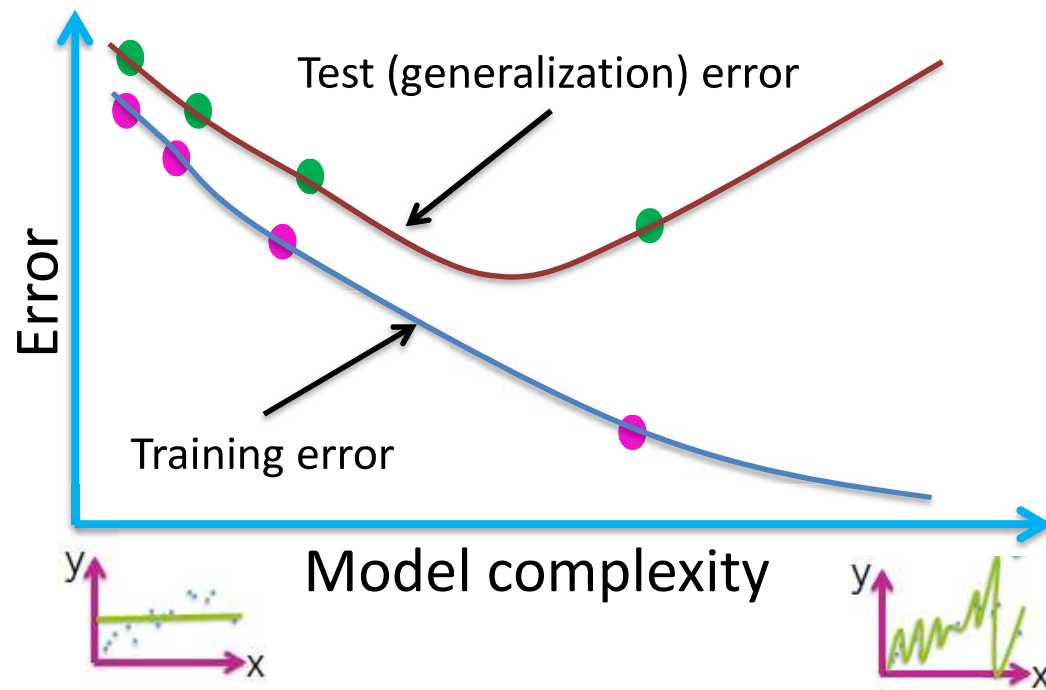
---

- Training error decreases with increasing model complexity



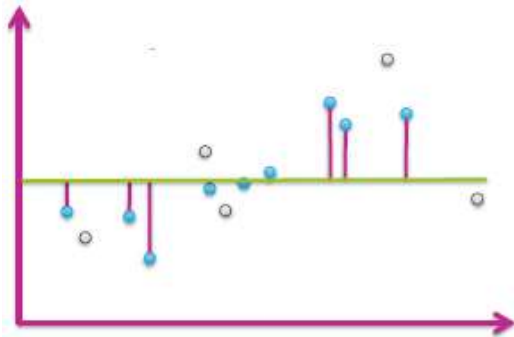
# Testing error vs. model complexity

- Test error can be used as an approximation of the generalization error



# Bias – Variance tradeoff

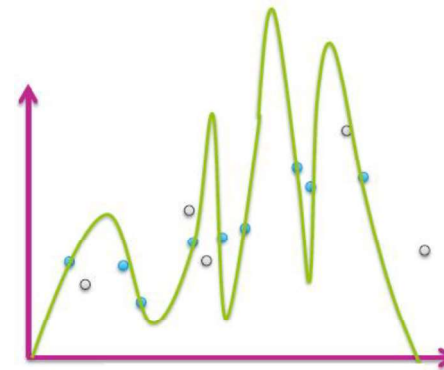
---



- High bias model

The model does not capture enough  
The structure of the training set  
**Parameters tend to be small**

**UNDERFIT**

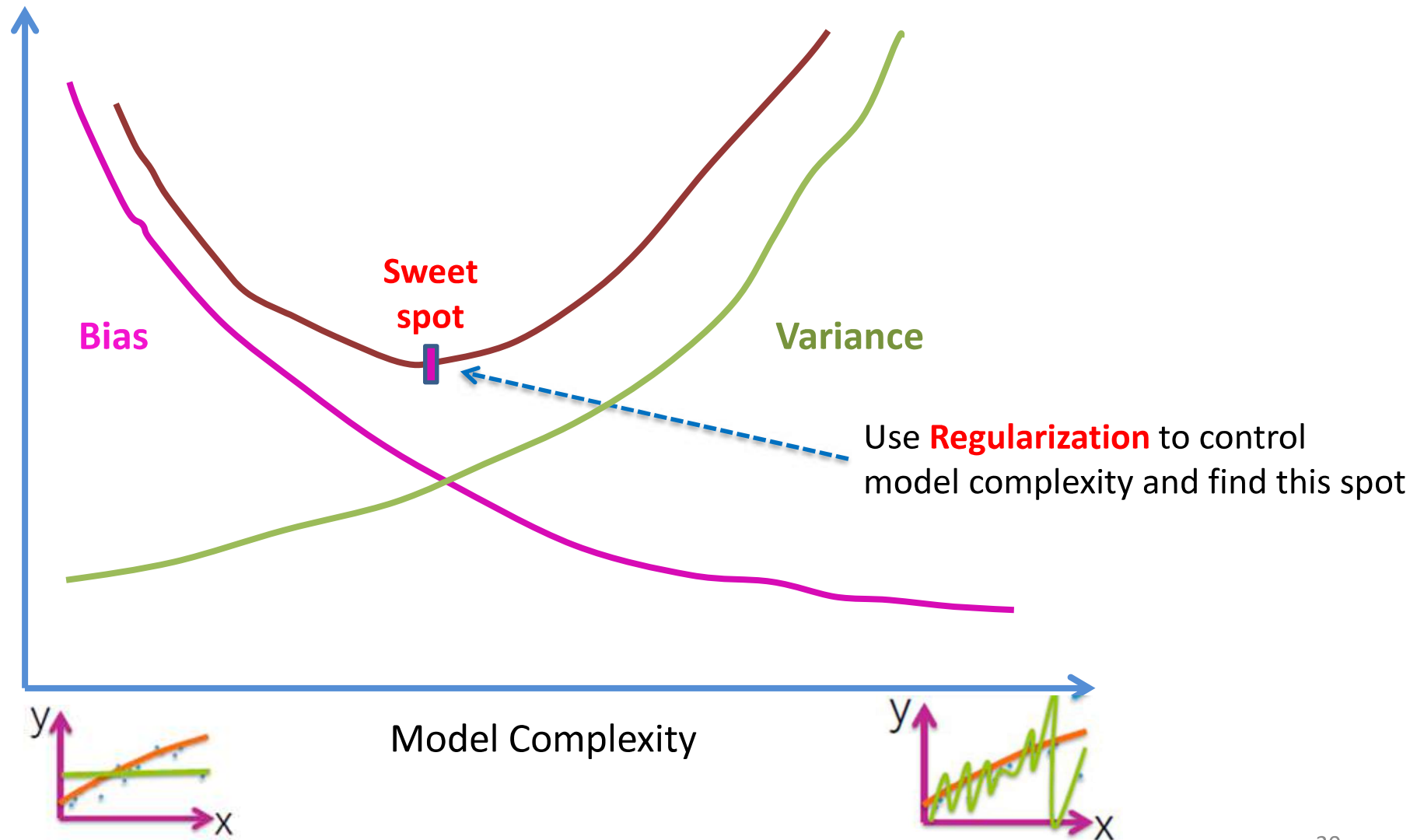


- High variance model

The model is too specific to the structure  
Of the training set  
**Parameters tend to be very large**

**OVERFIT**

# Bias – Variance tradeoff



---

# **1.4**

## **Regularization**

---

# Regularization

---

- It's about finding balance between:
  - How well the model fits the data
  - The magnitude of coefficients
- This is achieved by incorporating a penalty on weights  $\theta$  in the cost function

- Ridge Regression ( $L_2$  regularization)

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n (\theta_j)^2 \right]$$

- Lasso Regression ( $L_1$  regularization)

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=0}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n |\theta_j| \right]$$

- $\lambda$  is the regularization parameter:
  - Ridge: Encourages small weights  $\theta$  but not exactly 0
  - Lasso: "Shrink" some weights  $\theta$  exactly to 0

# GD with Regularization

---

- Reminder of the L2 regularization cost function:

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

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

- With regularization: 
$$\theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$
- $\alpha, \lambda$  are learning parameters to choose manually
- In practice:  $(1 - \alpha\lambda/m)$  is between 0.99 and 0.95



# Other Linear Regression Models

---

- Number of visiting customers to a website
- Product demand, inventory, failure, ...
- Stock pricing
- Insurance claims severity

“Remember that all models are wrong;  
the practical question is how wrong do they  
have to be to not be useful.”

George Box, 1987