



Universitat
de les Illes Balears

Machine Learning

Lesson 2: Supervised Learning

Linear Models (LMS, Logistic Regression, Perceptron)

Remember: a simple example...

Size (feet2)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

... for dataset notation

Dataset = $\{(\mathbf{x}^{(i)}, y^{(i)}); i = 1, \dots, m\}$

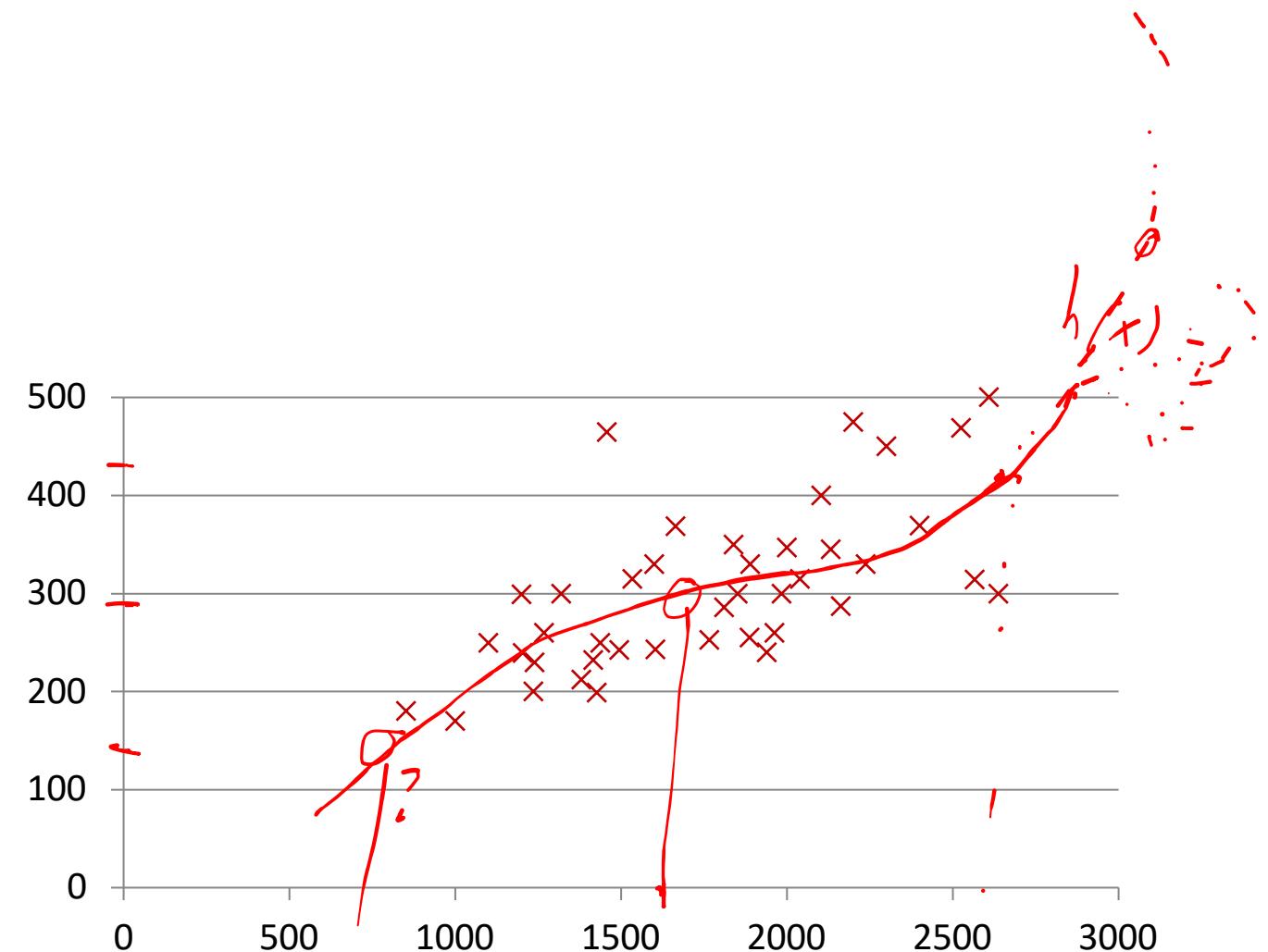
$(\mathbf{x}^{(i)}, y^{(i)})$ = *Training example*

$\mathbf{x}^{(i)}$ = “*input*” variable (*features*), $\mathbf{x}^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})$

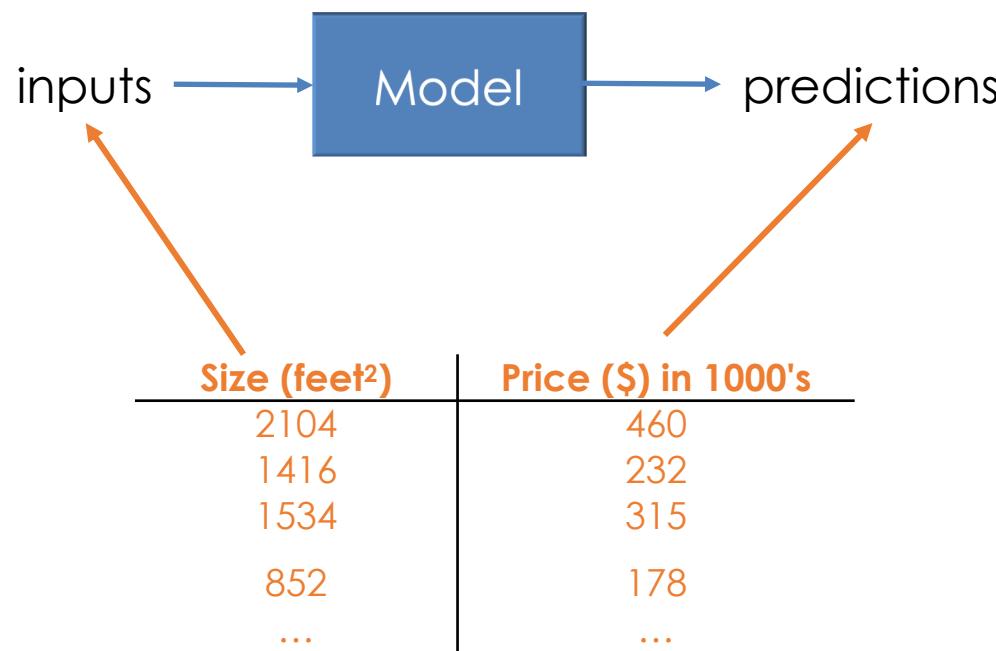
$y^{(i)}$ = “*output*” variable (*target*), $y \in \mathcal{Y}$

What price...?

Size (feet ²)	Price (\$) in 1000's
2104	460
1416	232
1534	315
852	178
...	...



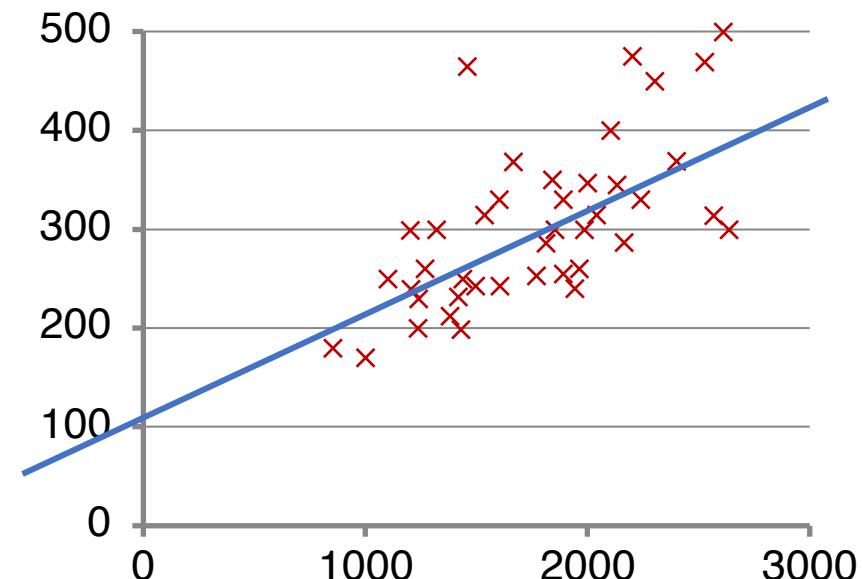
ML Model: regression



Linear regression

Hypothesis: the model is linear

$$h_{\theta}(x) = \theta_0 + \sum_{j=1}^n \theta_j x_j$$



Case n=1:

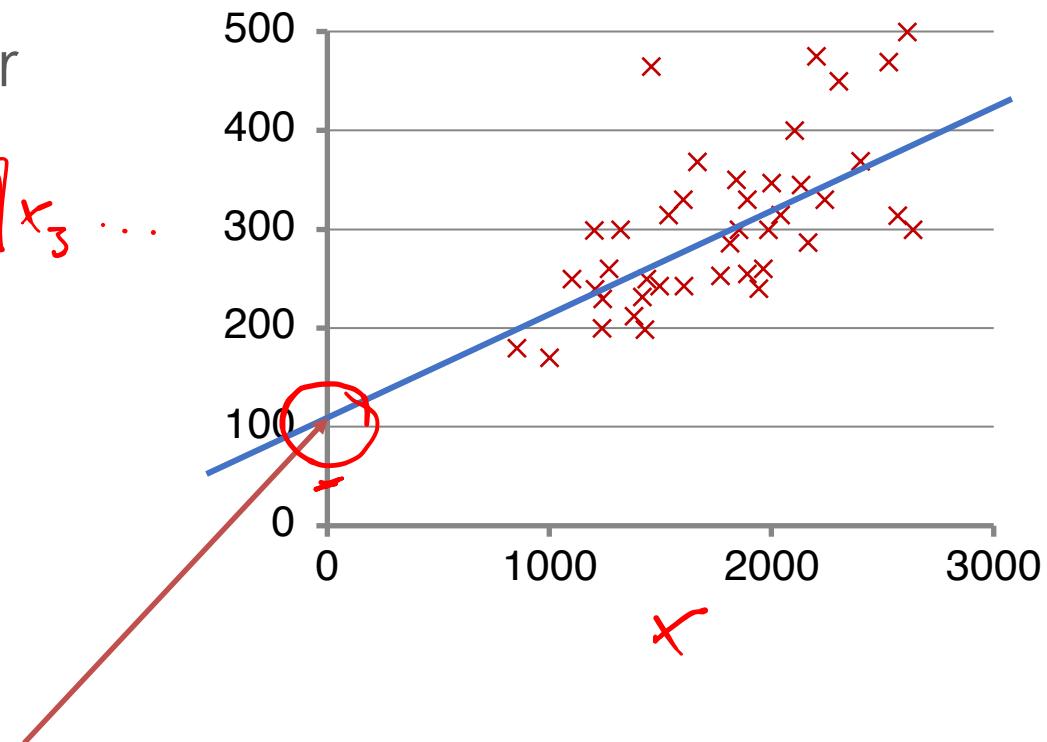
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Linear regression

Hypothesis: the model is linear

$$y = \underline{a} + b x_1 + k x_2 + \boxed{l} x_3 \dots$$
$$h_{\theta}(x) = \theta_0 + \sum_{j=1}^n \theta_j x_j$$

is called the **bias** because it represents
the intercept or prior (prejudice)



Linear regression

Hypothesis: the model is linear

$$h_{\theta}(x) = \theta_0 + \sum_{j=1}^n \theta_j x_j$$

$n = 3$

su | n. bengs | n. plos

m → n

$$\theta = [\theta_0 \dots \theta_n]$$
$$x = [1, \dots, x_n]$$
$$\theta^T x$$
$$\begin{bmatrix} \theta_0 \dots \theta_n \\ \vdots \\ x_n \end{bmatrix}$$

n = number of input dimensions

Generalization
(in matrix notation)

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

$x_0 = 1$ (Intercept term)

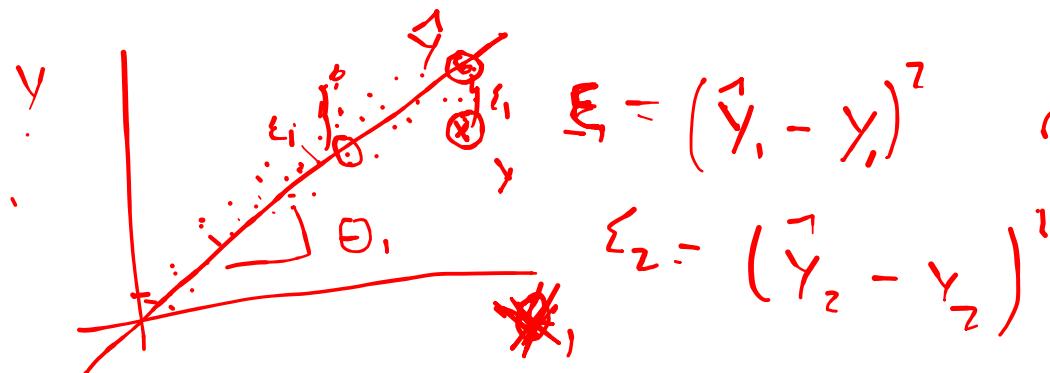
Loss function

Given a training set, how do we learn the parameters?

Basic idea: to make h_θ close to y

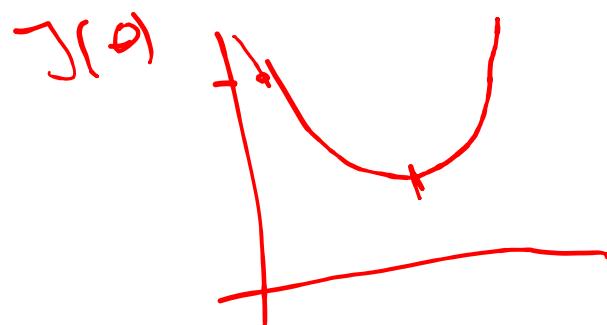
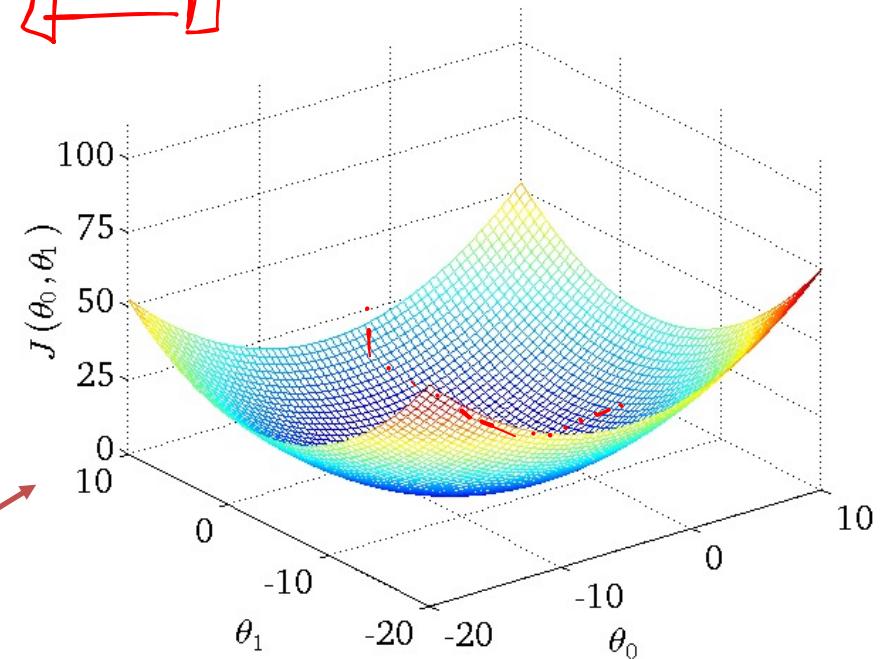
We've to define a function that measures, for each parameters' values, how close the $h_\theta(x^{(i)})$'s are to the corresponding $y^{(i)}$'s

Quadratic loss function



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

$JSE = \frac{1}{m} \sum_i (\hat{y}_i - y_i)^2$



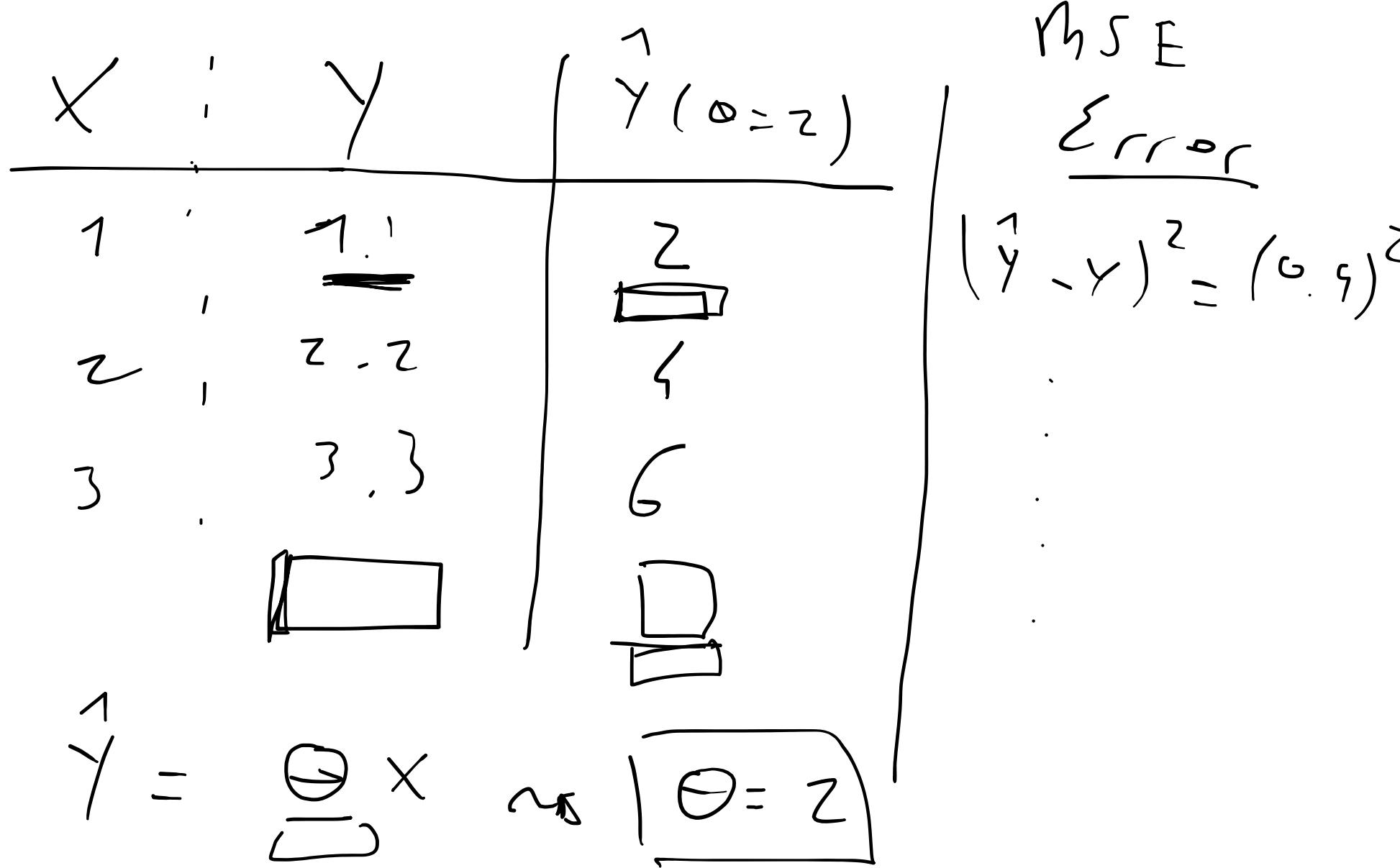
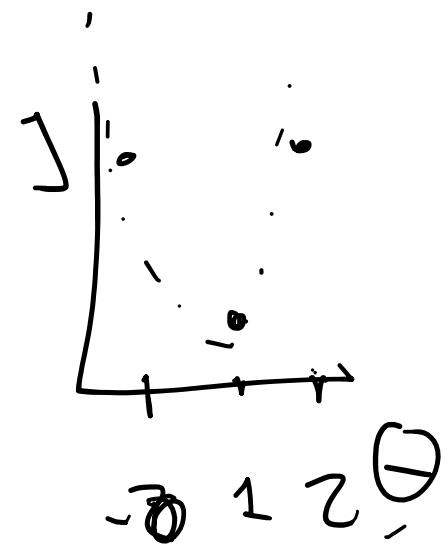
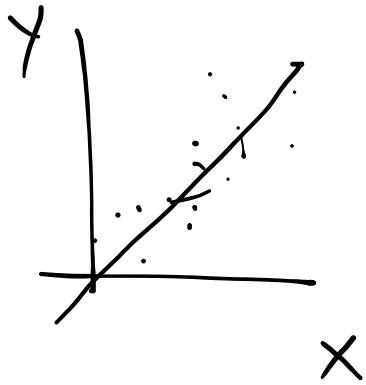
Case n=1:

$$h_\theta(x) = \theta_0 + \theta_1 x_1$$

$\theta_1 \cap \{y=1\}$

$x_1 = 1$

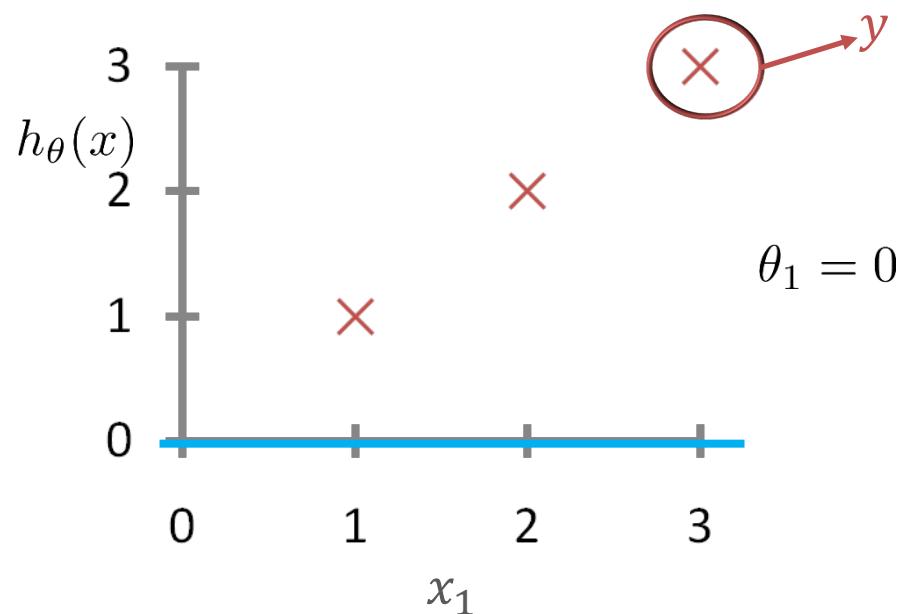




Quadratic loss function: example

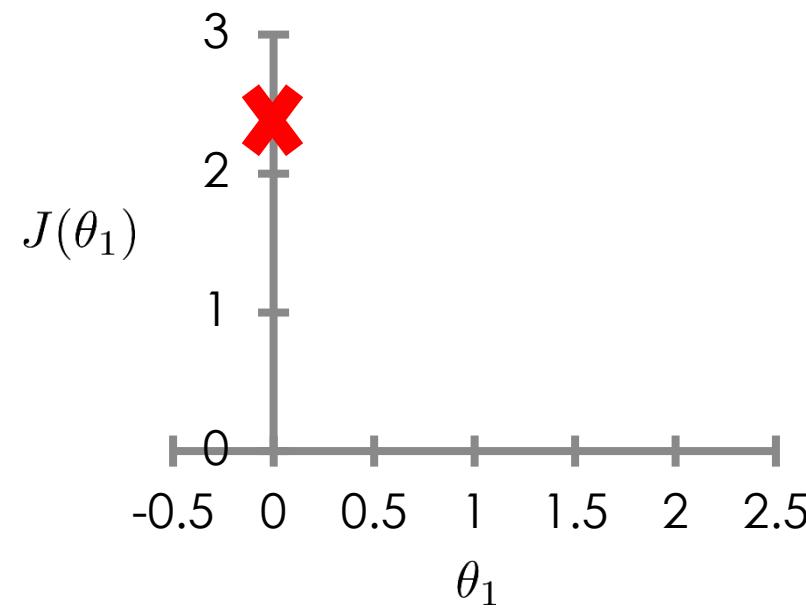
$$h_{\theta}(x) = \theta_1 x_1$$

(for a fixed θ_1 , this is a function of x)



$$J(\theta_1)$$

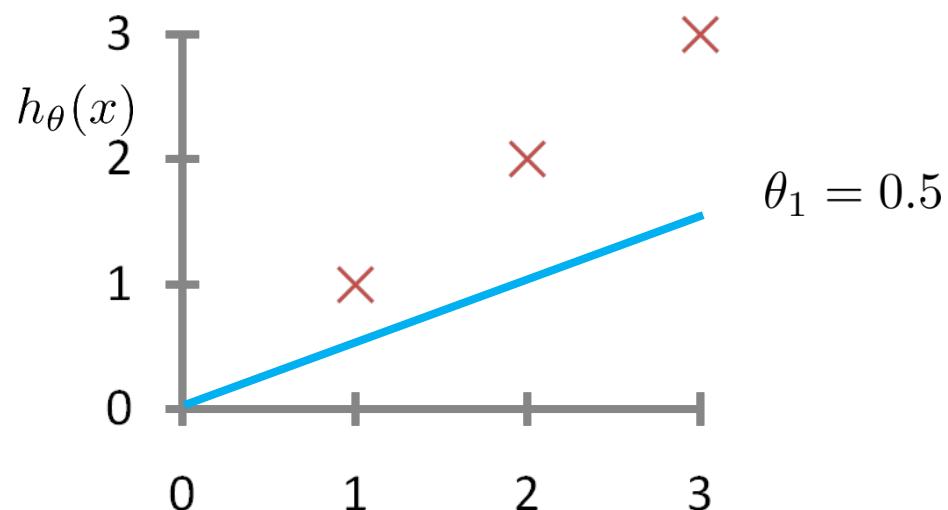
(function of the parameter θ_1)



Quadratic loss function: example

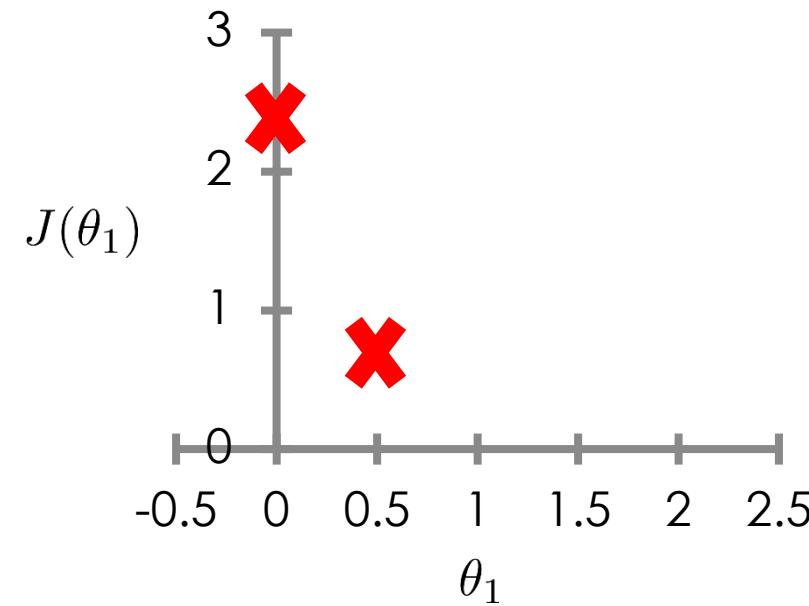
$$h_{\theta}(x) = \theta_1 x_1$$

(for a fixed θ_1 , this is a function of x)



$$J(\theta_1)$$

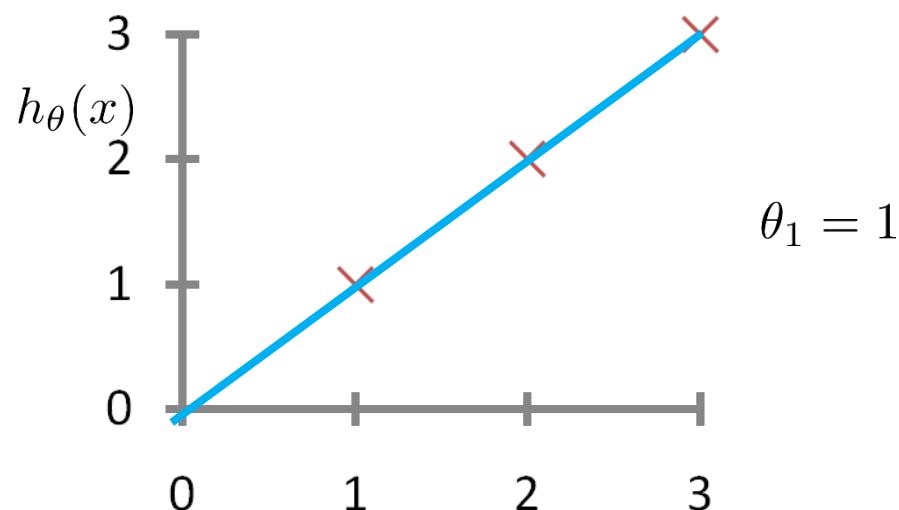
(function of the parameter θ_1)



Quadratic loss function: example

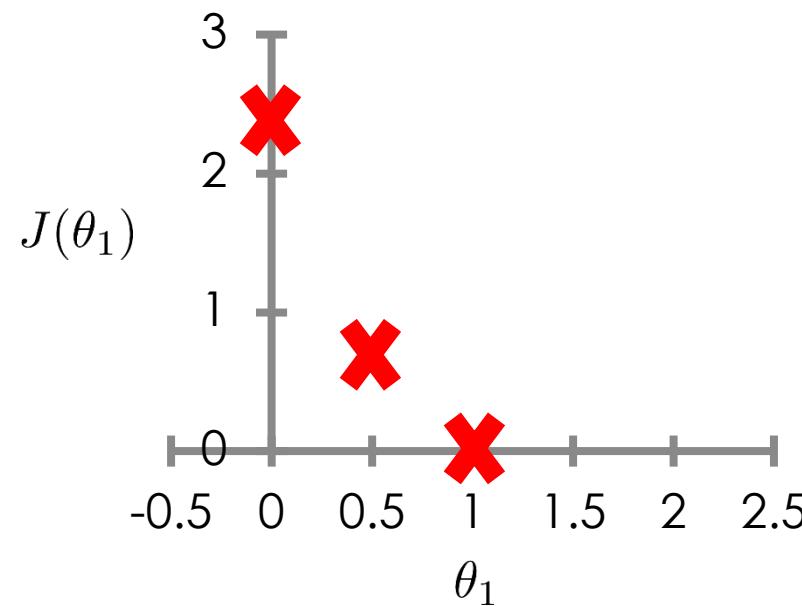
$$h_{\theta}(x) = \theta_1 x_1$$

(for a fixed θ_1 , this is a function of x)

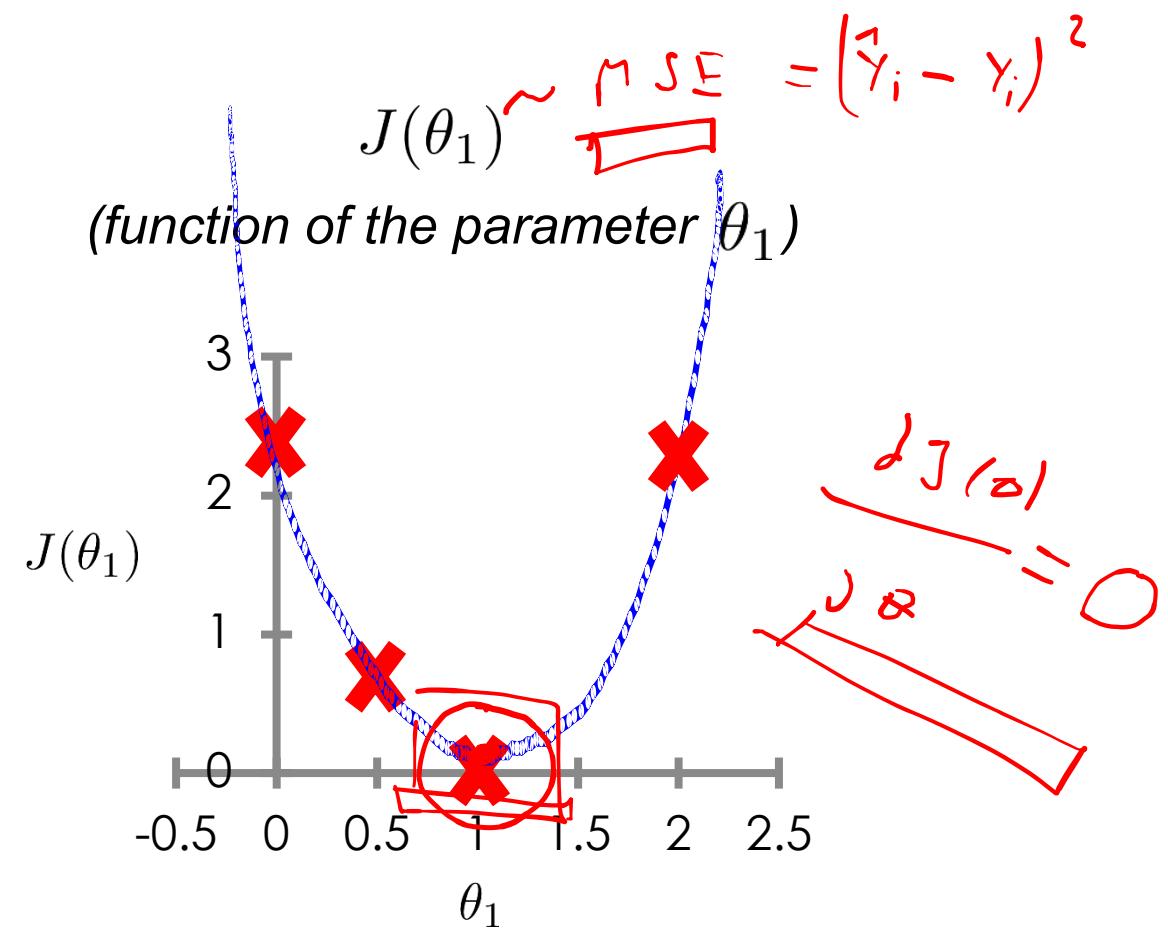
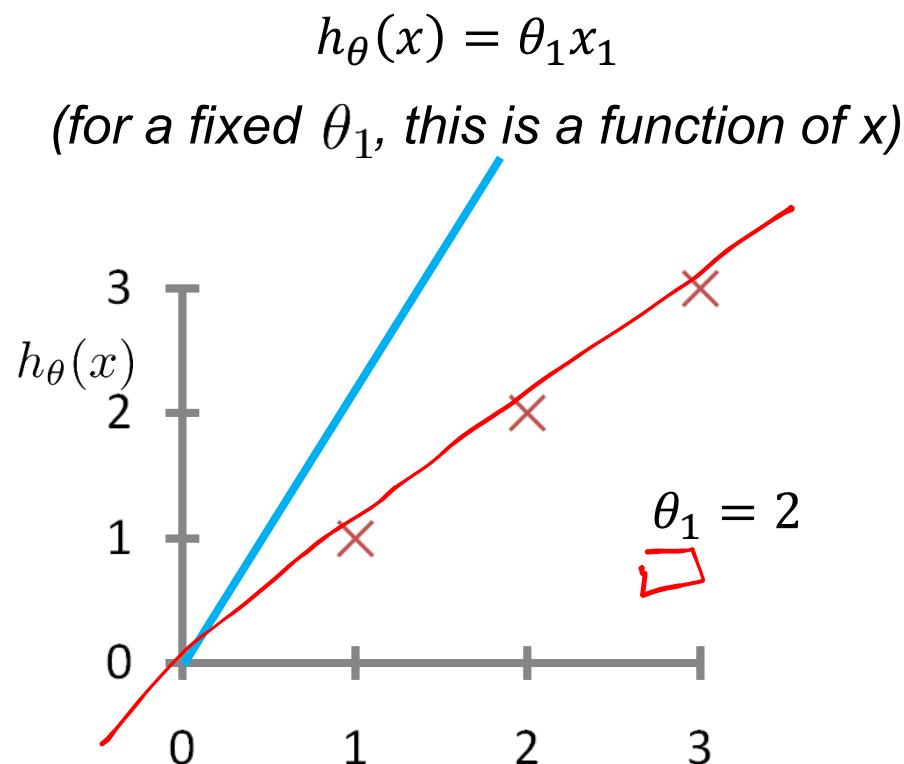


$$J(\theta_1)$$

(function of the parameter θ_1)



Quadratic loss function: example



Analytic Linear regression

Dataset = X, \vec{y} , where X is a matrix $m \times (n + 1)$

$\eta s \in$

$$J(\theta) = \frac{1}{2m} (X\theta - \vec{y})^T (X\theta - \vec{y}) \rightarrow \sum (x\theta - y)^2$$

$$\vec{y} = X\theta$$

$$y = \boxed{w_0} x + \boxed{b}$$

$$X\theta = \vec{y}$$

$$X^T X\theta = X^T \vec{y}$$

$$\theta = (X^T X)^{-1} X^T \vec{y}$$

$$\frac{\partial J(\theta)}{\partial \theta} = 0 \quad \text{inv, que}$$

↳ 1) colinearity
2) computational costs

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

LMS learning algorithm

Loss function: Quadratic

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

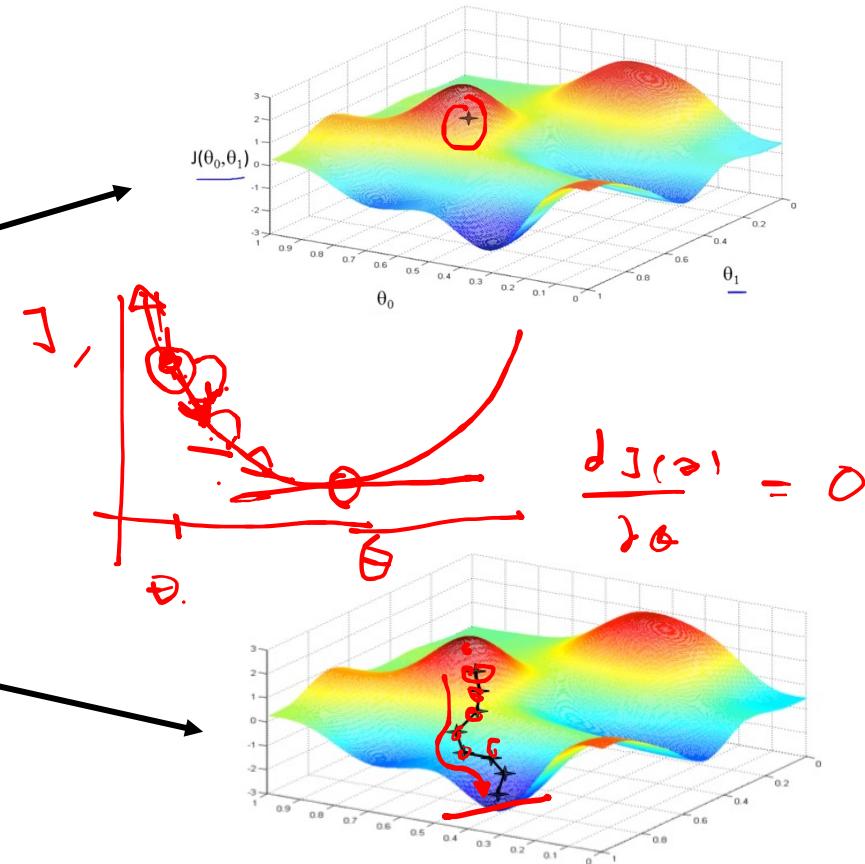
Goal: To minimize $J(\theta)$

$$\begin{matrix} m \\ | \\ j=0 \\ | \\ j=1 \\ | \\ \vdots \\ | \\ j=m \end{matrix}$$

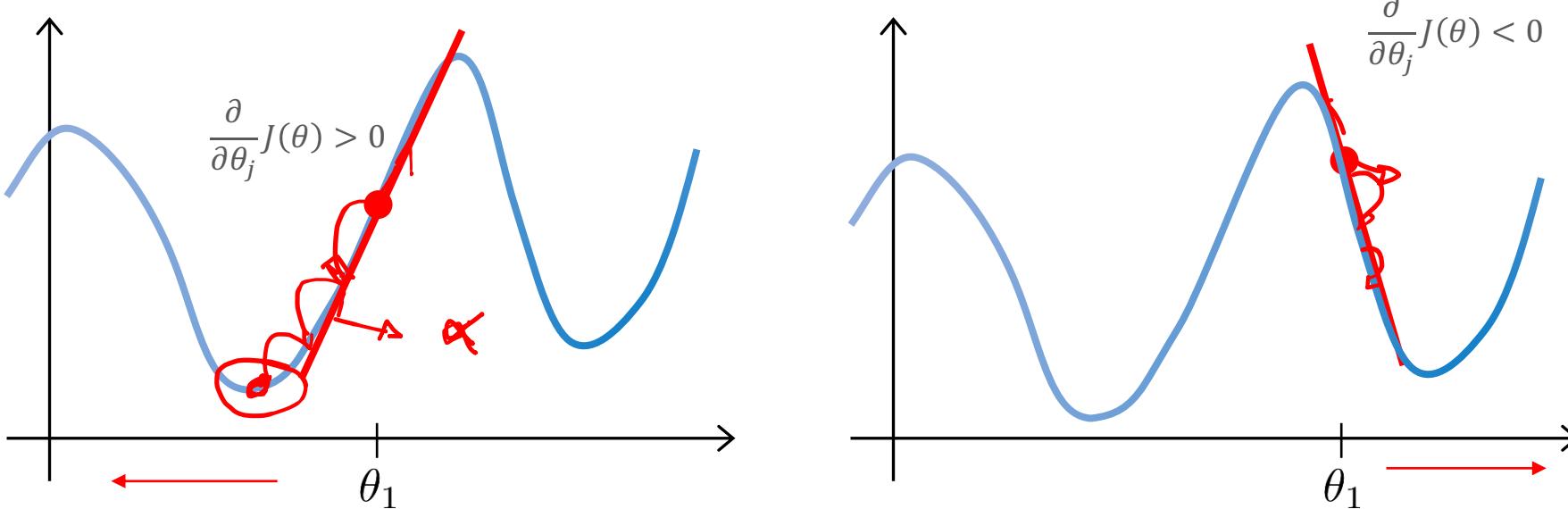
Gradient descent

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta)$$

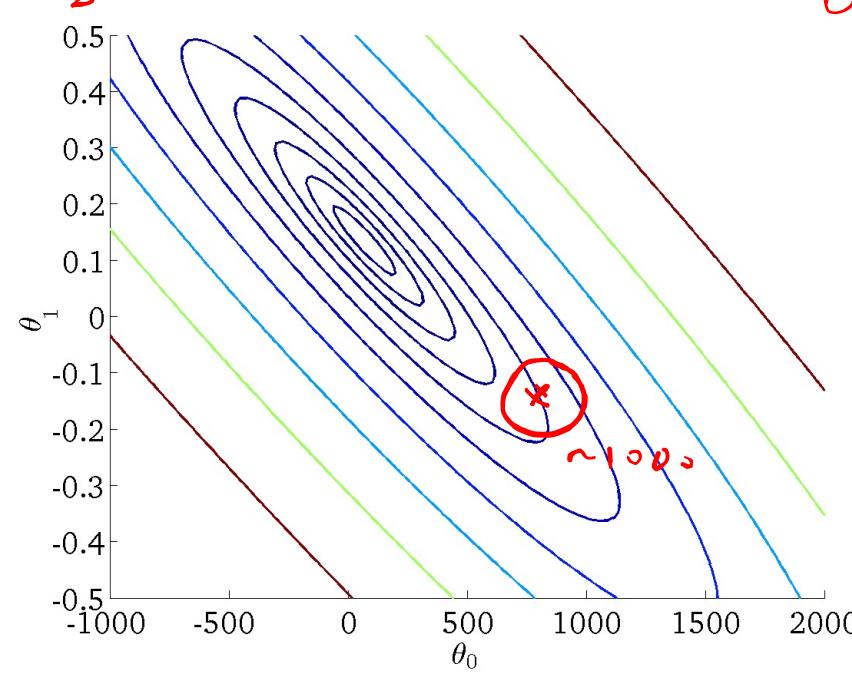
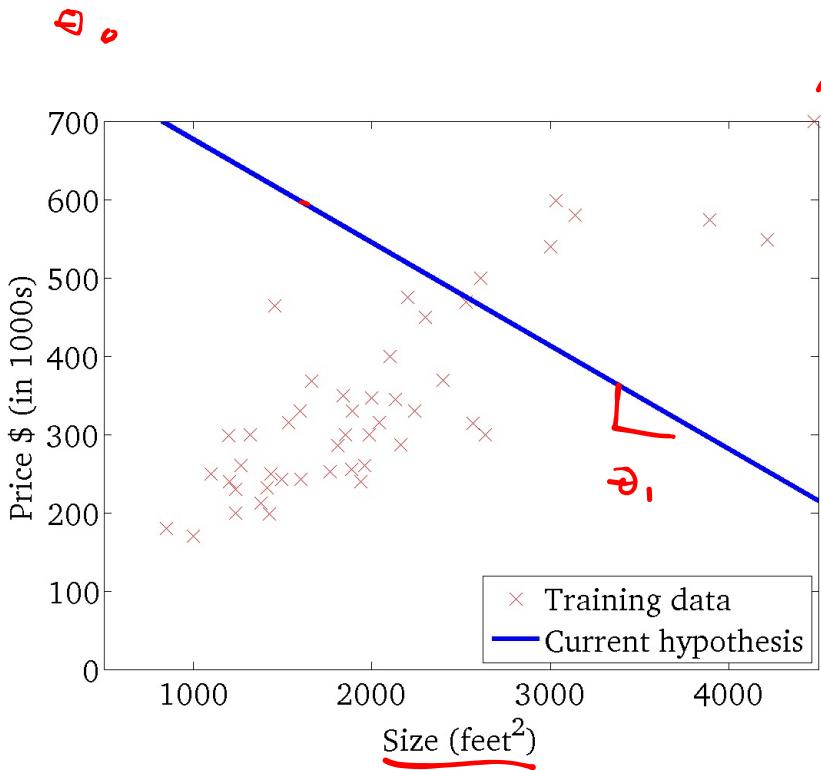
Learning rate



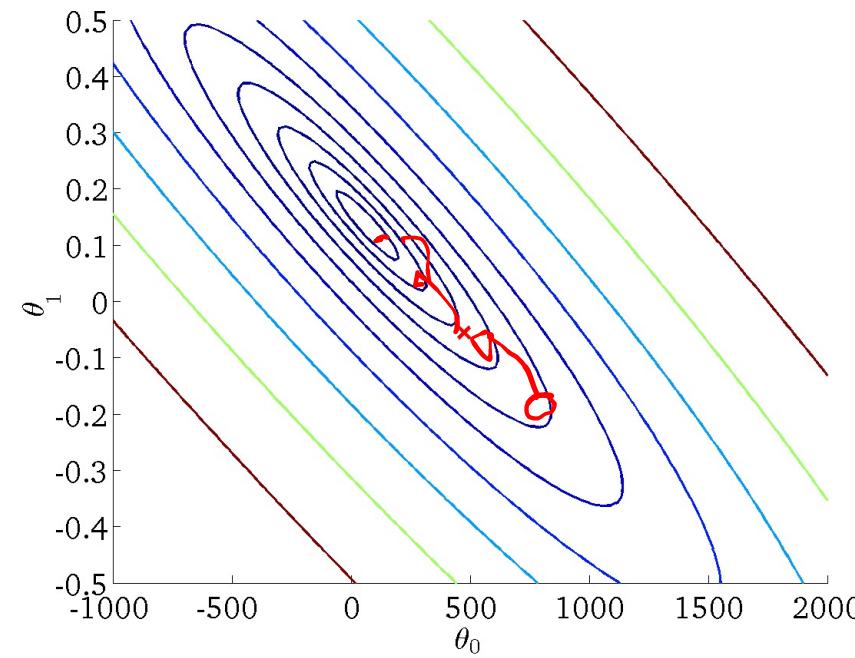
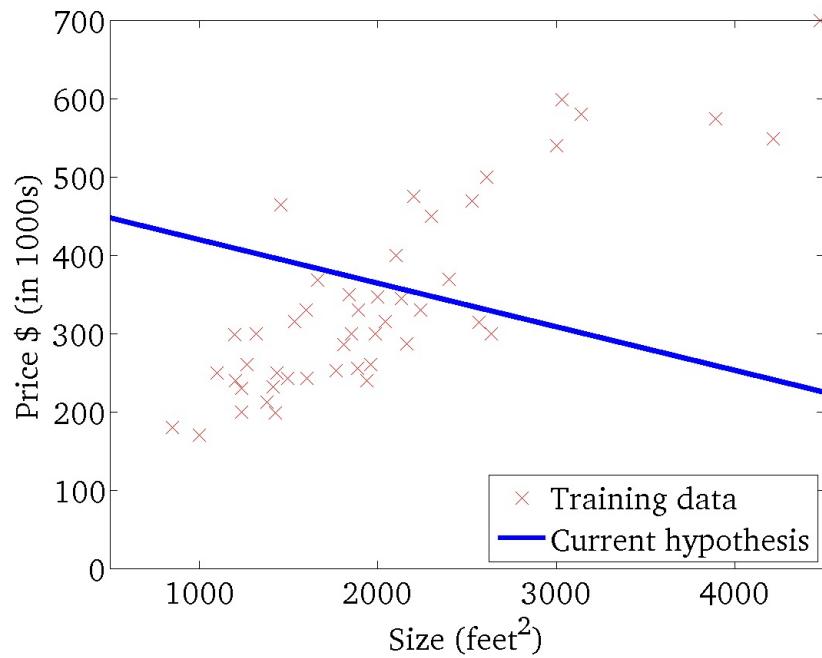
Gradient descent: basic idea



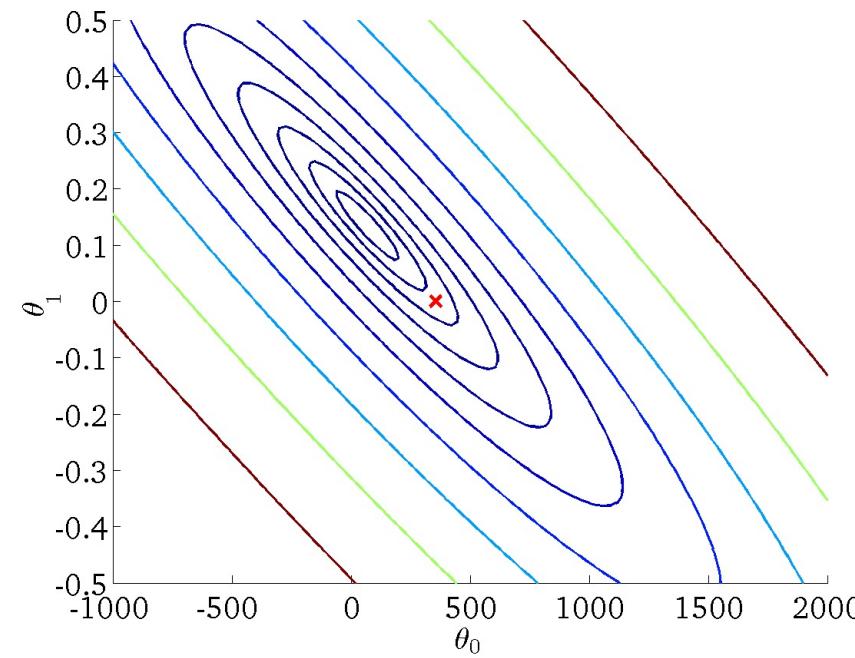
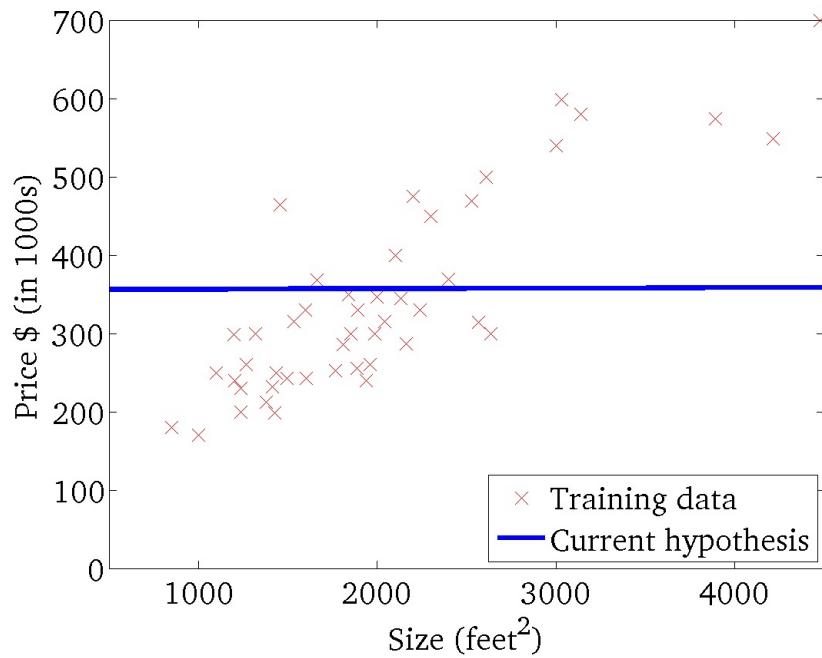
Gradient descent: Linear



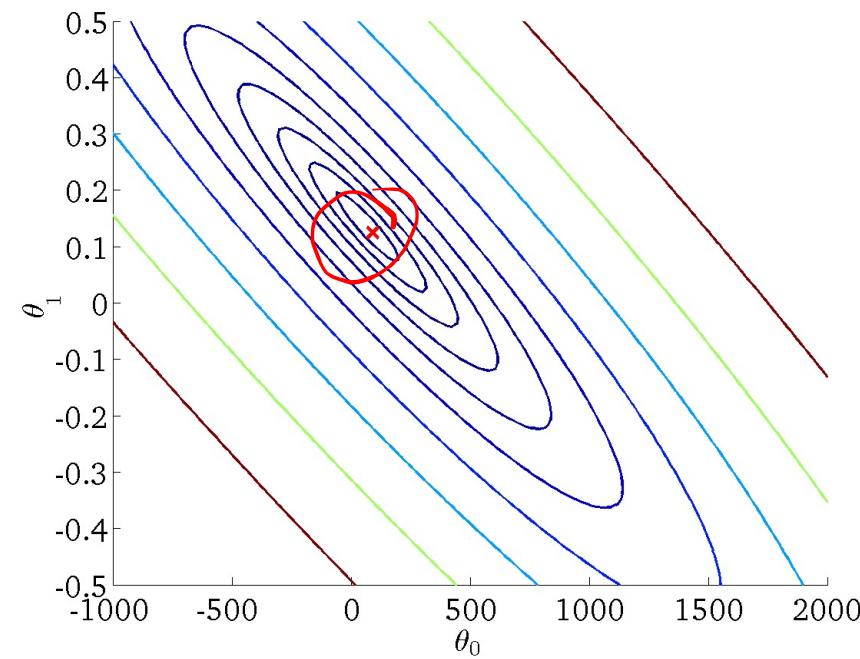
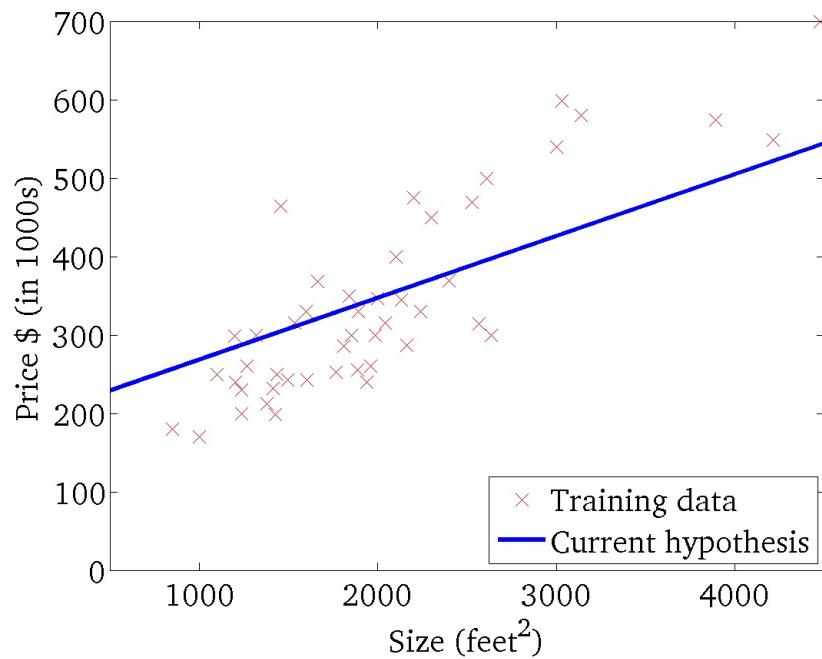
Gradient descent: Linear



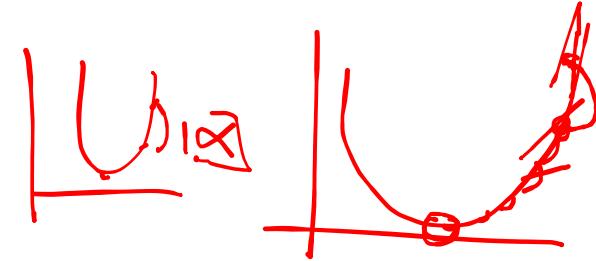
Gradient descent: Linear



Gradient descent: Linear



LMS algorithm: update rule



$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta) \longrightarrow$$

Exercise: Obtain the gradient for the quadratic cost function of:

$$h_\theta(X) = \theta_0 + \theta_1 X$$

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

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y) \cdot 1$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{1}{2m} \sum_{i=1}^m 2 \cdot (\theta_0 + \theta_1 x_i - y) \cdot x_i$$

$$\theta_j := \theta_j - \alpha \cdot (h_\theta(x) - y) x_j$$

This rule (also known as Widrow-Hoff learning rule) has several properties that seem natural and intuitive. For instance, the magnitude of the update is proportional to the error term ($y - h_\theta(x)$); thus, for instance, if we are encountering a training example on which our prediction nearly matches the actual value, then, we find that there is little need to change the parameters; in contrast, a larger change to the parameters will be made if our prediction has a large error.

$$\frac{\partial J(g(x))}{\partial x} = g'(g(x)) \cdot g'(x)$$

Exercise 2. Using the Gradient derived above, calculate the first iteration of the gradient descent

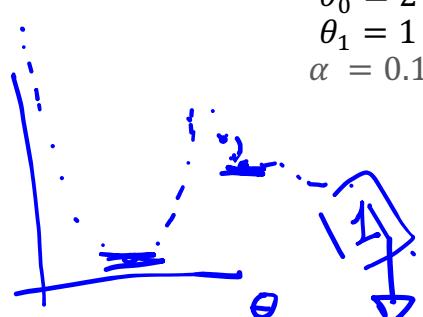
Initial Model

$$\hat{y} = \theta_0 + \theta_1 x$$

$$\theta_0 = 2$$

$$\theta_1 = 1$$

$$\alpha = 0.1$$

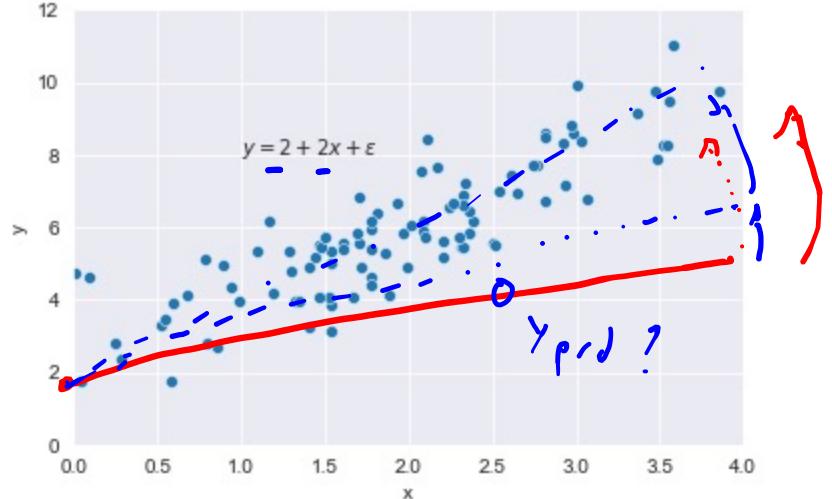


Grad Desc Algorithm

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\theta_0^j := \theta_0^j - \alpha \left((\theta_0 + \theta_1 x_i - y) x_i \right)$$

$$\theta_1^j := \theta_1^j - \alpha \left((\theta_0 + \theta_1 x_i - y) \right)$$



Sample	x_0	x_1	y_{real}	θ_0^j	θ_1^j	y_{pred}	error_{te}	Grad_{θ_0}	Grad_{θ_1}	Learn_rate	θ_0^{j+1}	θ_1^{j+1}
0	1.0	2.50	5.58	2	1	4.5	-1.08	-1.08	-2.7	0.1	2.108	1.27
1	1.0	1.86	5.30	2.108	1.27							
2	1.0	2.65	6.96									
3	1.0	3.52	8.24									
4	1.0	1.77	5.38									

Learning rate

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta)$$

Learning rate

$\hookrightarrow \alpha \sim \text{precise/error}$

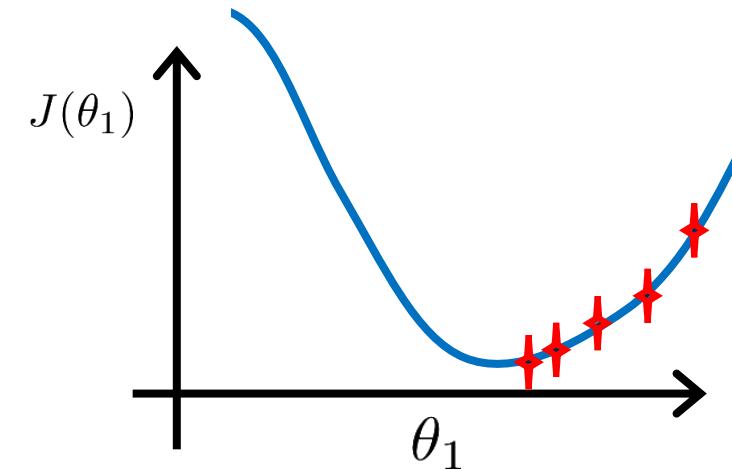
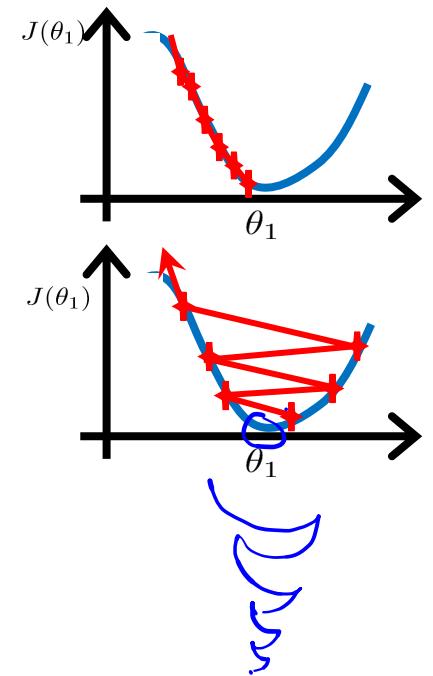
Gradient descent can be susceptible to **local minima** in general, the optimization problem we have posed here for linear regression has only one global, and no other local, optima; thus gradient descent always converges (assuming the learning rate α is not too large) to the global minimum.

As we approach a local minimum, gradient descent will automatically take smaller steps, being $[0,1]$ smooths the derivative by reducing the jump stride.

In Practice:

If α is too small, gradient descent can be slow.

If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge



Learning rate

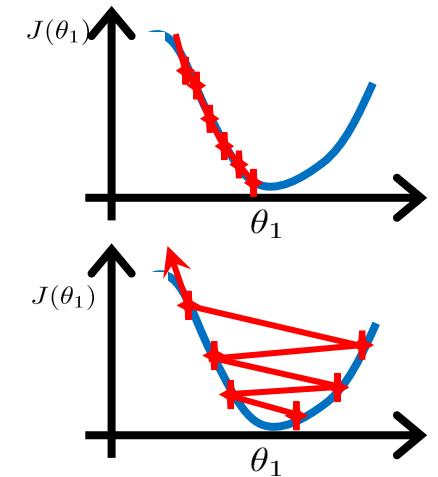
$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta)$$

Learning rate

In Practice:

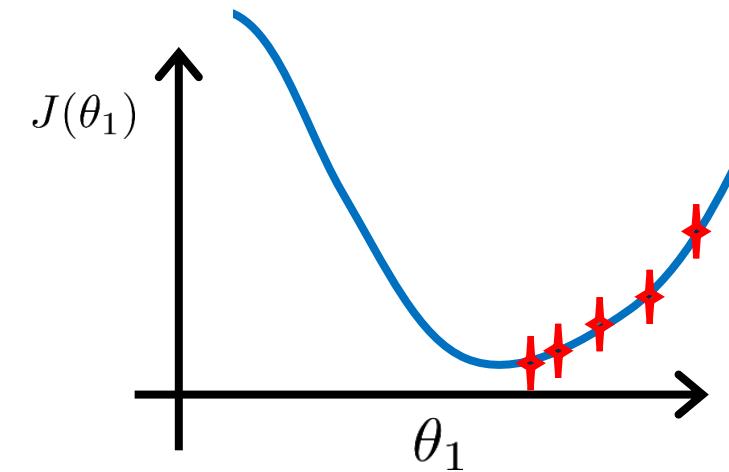
If α is too small, gradient descent can be slow.

If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge



Gradient descent can be susceptible to **local minima** in general, the optimization problem we have posed here for linear regression has only one global, and no other local, optima; thus gradient descent always converges (assuming the learning rate α is not too large) to the global minimum.

As we approach a local minimum, gradient descent will automatically take smaller steps, being $[0,1]$ smooths the derivative by reducing the jump stride.



LMS algorithm: incremental rule

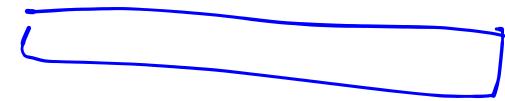
```
Loop{  
    for i=1 to m, {  
        }  $\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$   
    }  
}
```

When the training set is large, **stochastic gradient descent** is often preferred over batch gradient descent

Whereas batch gradient descent has to scan through the entire training set before taking a single step—a costly operation if m is large—stochastic gradient descent can start making progress right away, and continues to make progress with each example it looks at. Often, stochastic gradient descent gets θ “close” to the minimum much faster than batch gradient descent. (Note however that it may never “converge” to the minimum, and the parameters will keep oscillating around the global minimum; but in practice, most of the values near the minimum will be reasonably good approximations to the true minimum.

LMS algorithm: batch rule

epoch

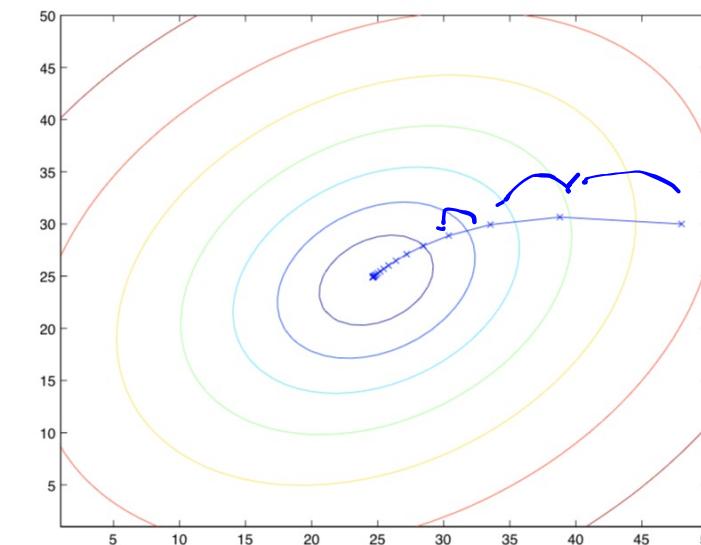


Repeat until convergence {

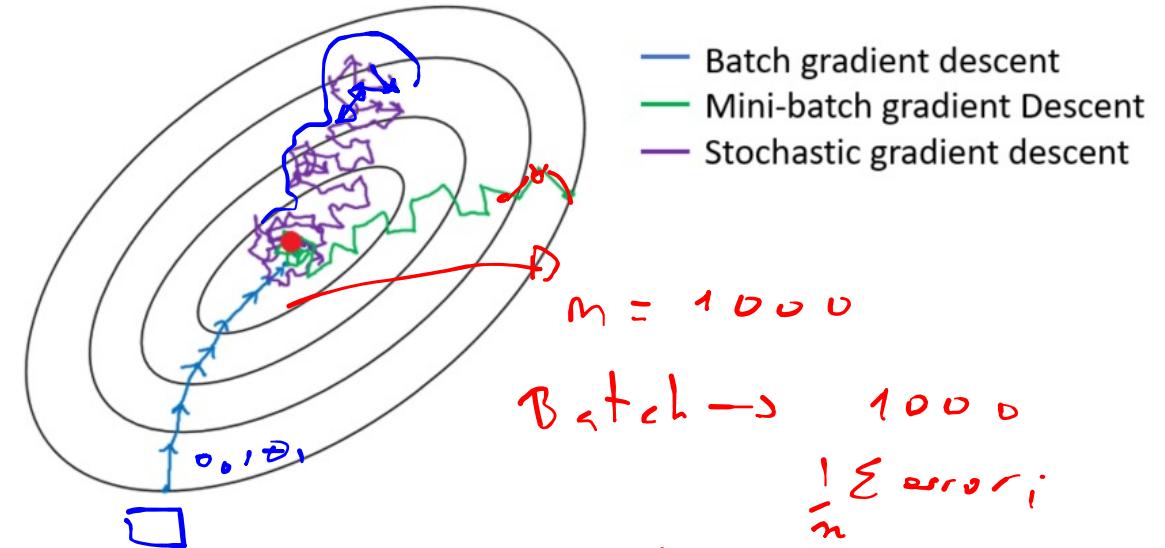
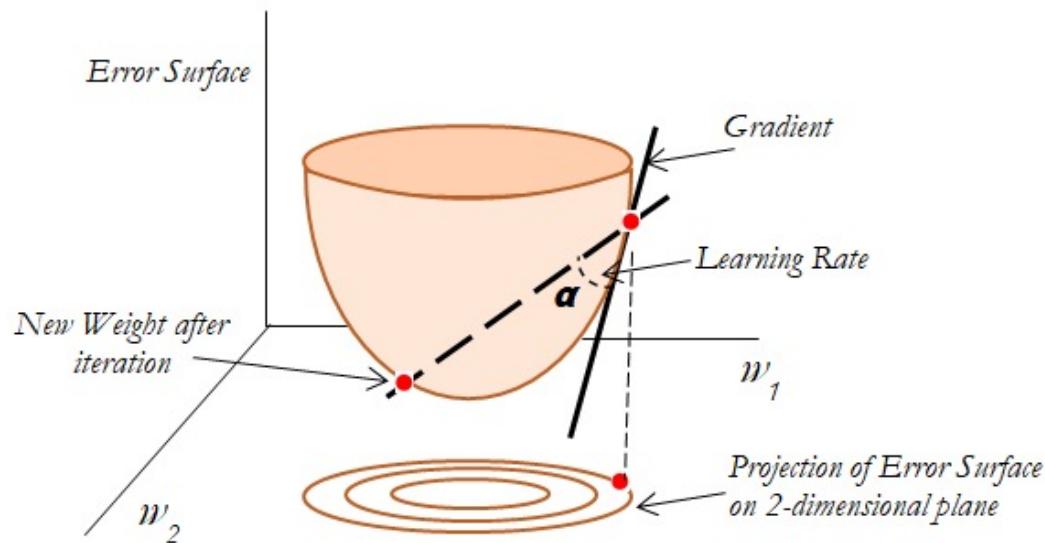
$$\theta_j := \theta_j + \alpha \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

}

(for every j)



Stochastic Gradient Descent with Batch size “1”



Gradient descent

minimize an objective function $J(\theta)$ by updating the parameters in the opposite direction of the gradient of the objective function.

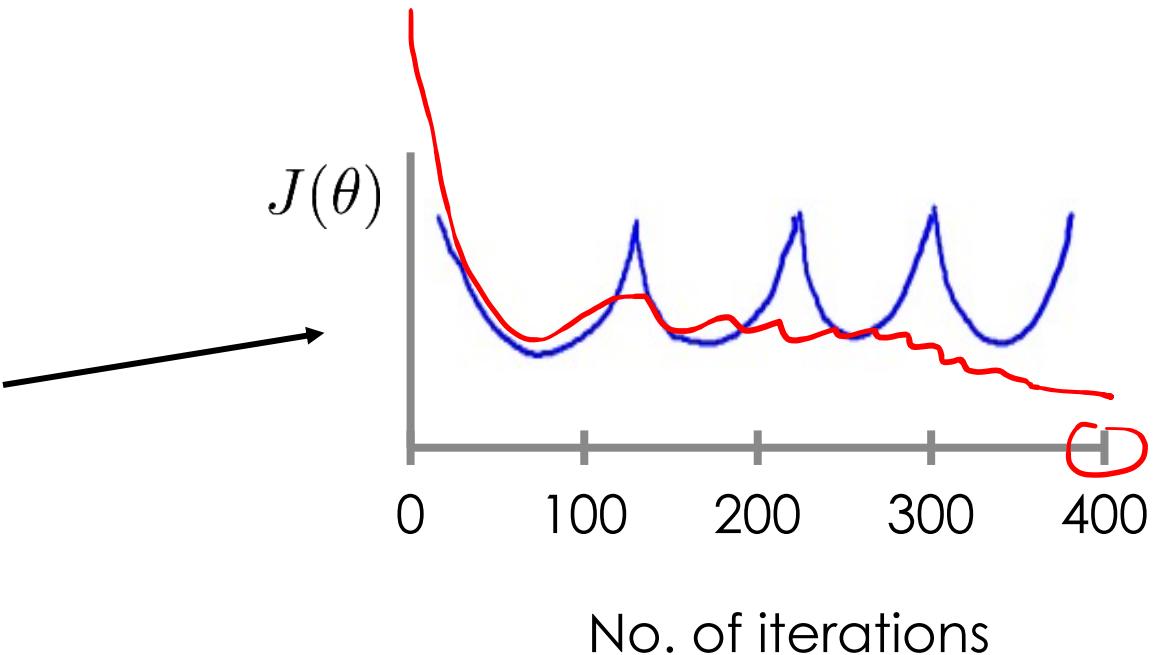
The learning rate determines the size of the steps taken to reach the minimum

$\rightarrow i.i.d$

- Batch gradient descent: all training observations utilized in each iteration
- SGD: one observation per iteration
- Mini batch gradient descent: size of about 50 training observations for each iteration

LMS in practice

It is necessary to display
the cost function



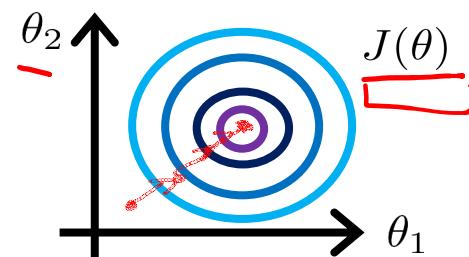
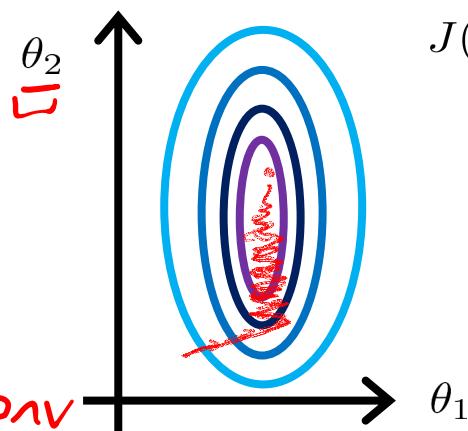
Simultaneous update

```
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ 
temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ 
 $\theta_0 := \text{temp0}$ 
 $\theta_1 := \text{temp1}$ 
```

Data scaling / normalization

1) $\|x\|_{\text{L}_2}$, ω_{inv}

2) Expl.



Summary of Linear models

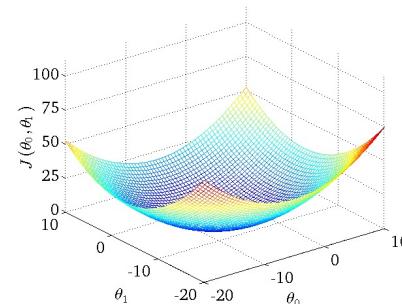
Linear regression

Hypothesis:
the model is linear

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

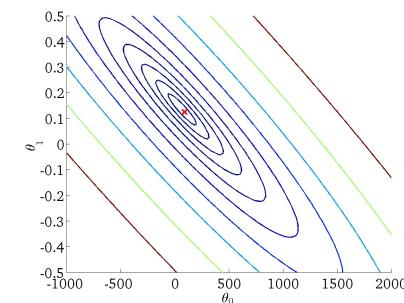
Gradient descent

Idea: to make h_{θ} close to y by means *minimizing a cost function*



LMS

Using the gradient to find the minimum (batch & incremental)



Linear regression revisited

Hypothesis: the model is linear

n = number of input dimensions

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

Generalization

(in matrix notation)

Case n=1:

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

$x_0=1$ (Intercept term)

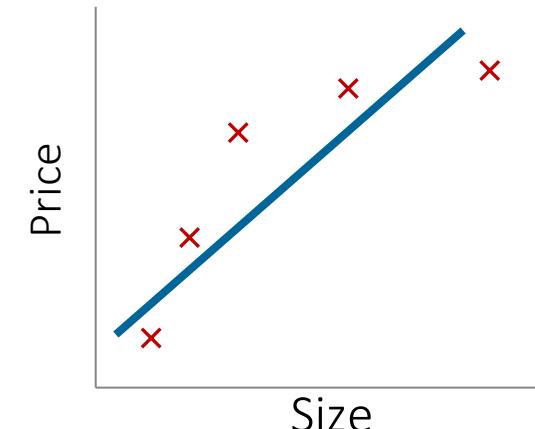
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

Bias-Variance Tradeoff

A learning algorithm is **biased** for a particular input (x) if, when trained on different data sets, it is systematically incorrect when predicting the correct output for x .

Therefore, the **bias** is an error from erroneous assumptions in the learning algorithm.

High bias can cause an algorithm to miss the relevant relations between features and target outputs (**underfitting**).



$$\theta_0 + \theta_1 x$$

Bias

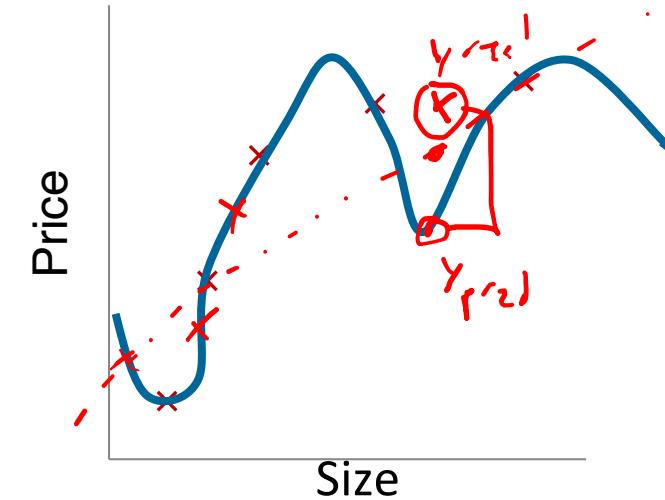
(Underfitting)

Bias-Variance Tradeoff

A learning algorithm has high **variance** for a particular input x if it predicts different output values when trained on different training sets.

The **variance** is an error from sensitivity to small fluctuations in the training set.

High variance can cause an algorithm to model the random noise in the training data, rather than the intended outputs (**overfitting**).



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

Variance
(Overfitting)

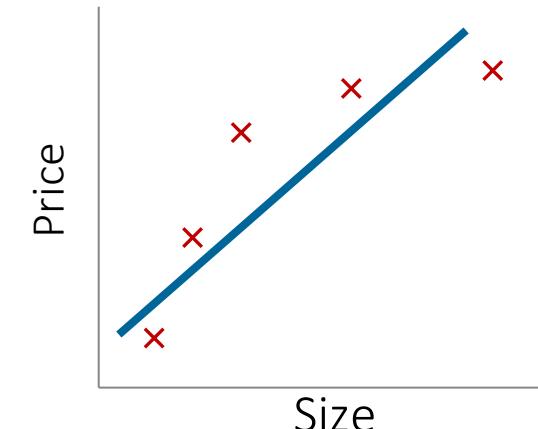


Bias-Variance Tradeoff

A learning algorithm is **biased** for a particular input (x) if, when trained on different data sets, it is systematically incorrect when predicting the correct output for x .

Therefore, the **bias** is an error from erroneous assumptions in the learning algorithm.

High bias can cause an algorithm to miss the relevant relations between features and target outputs (**underfitting**).

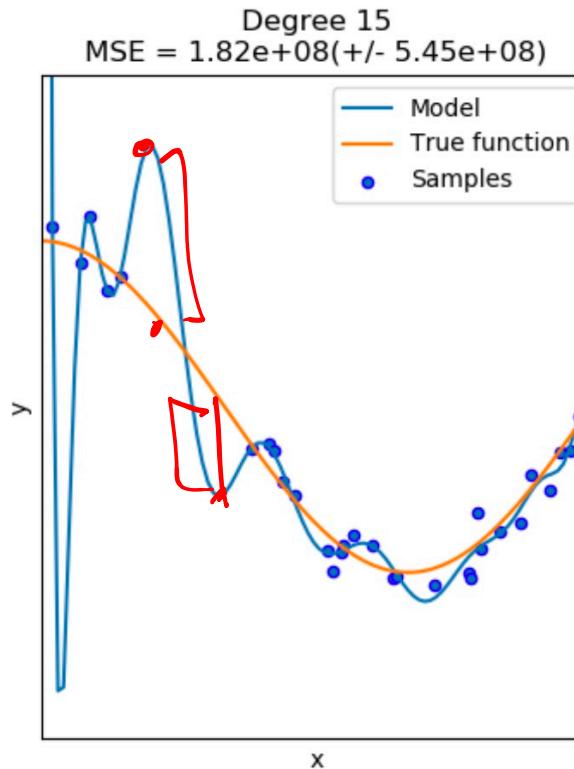
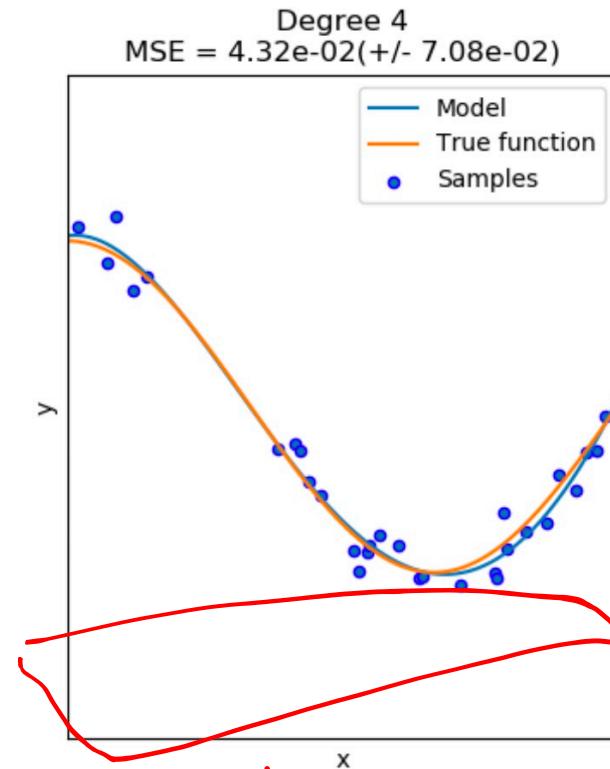
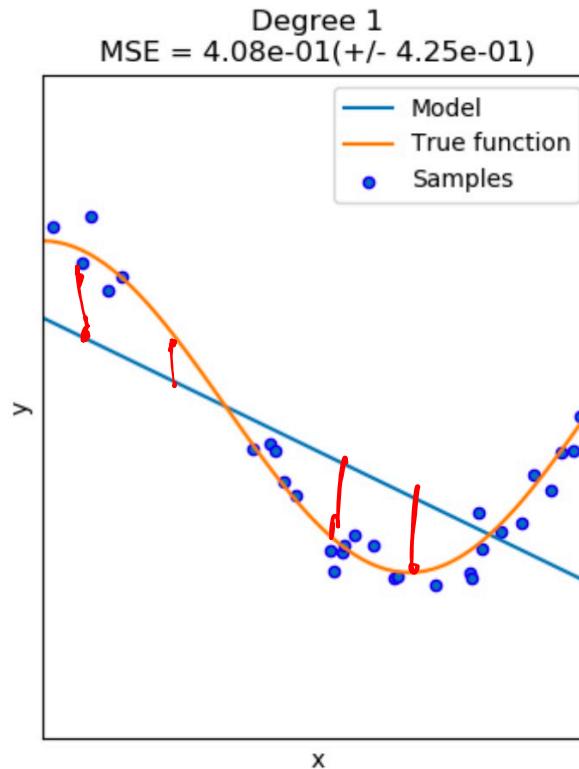


$$\theta_0 + \theta_1 x$$

Bias

(Underfitting)

Bias-Variance Tradeoff



Bias-Variance Tradeoff

$$Err(x) = \mathbb{E} \left[(y_{test} - \hat{y})^2 \right] = \underbrace{\left(E[\hat{f}(x)] - f(x) \right)^2}_{\text{Bias}} + \underbrace{E \left[(\hat{f}(x) - E[\hat{f}(x)])^2 \right]}_{V_{ff}}$$

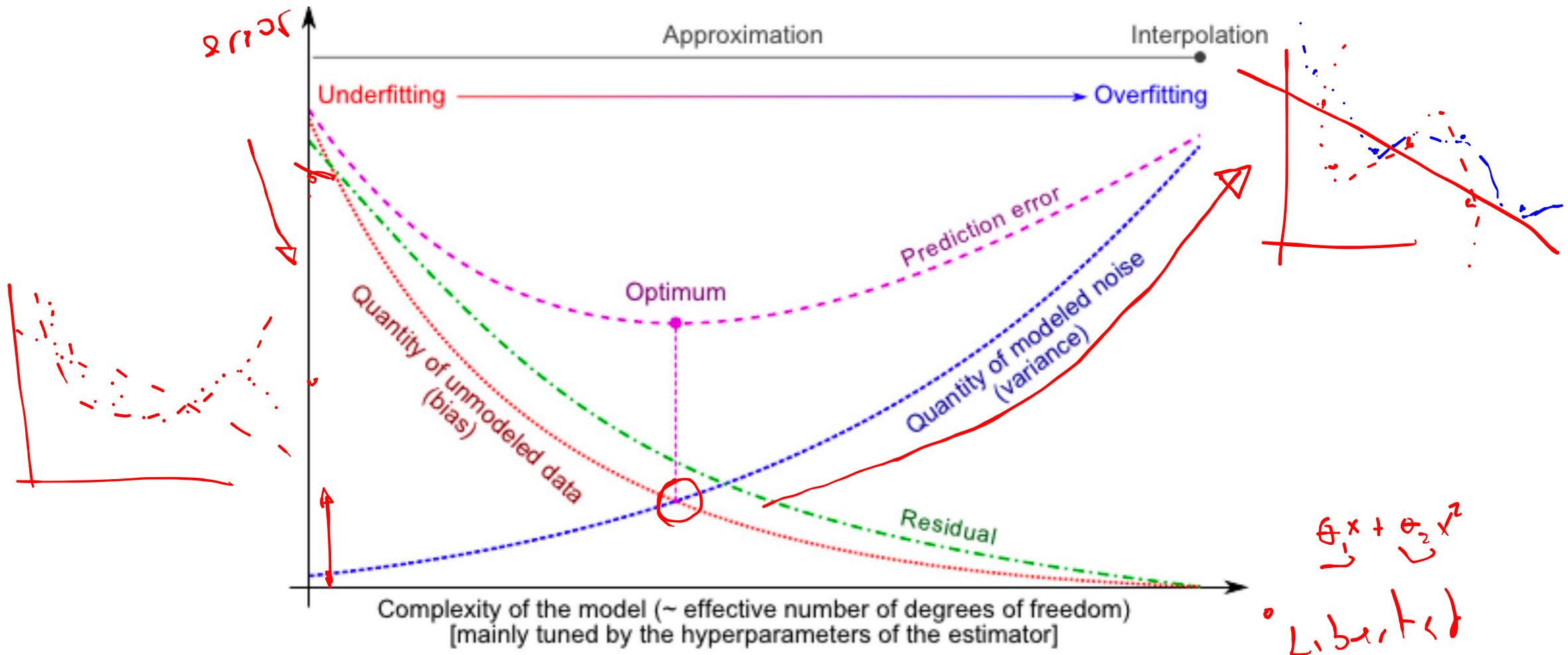
σ_e^2 is circled in red.

Bies: La diferencia entre el valor predicho por el modelo y el valor real.

Varianza: El error producido debido a la sensibilidad del modelo con respecto a los datos de entreno.

Error irreducible: ruido que no conseguiremos modelizar

Bias-Variance Tradeoff



Bias-Variance Tradeoff

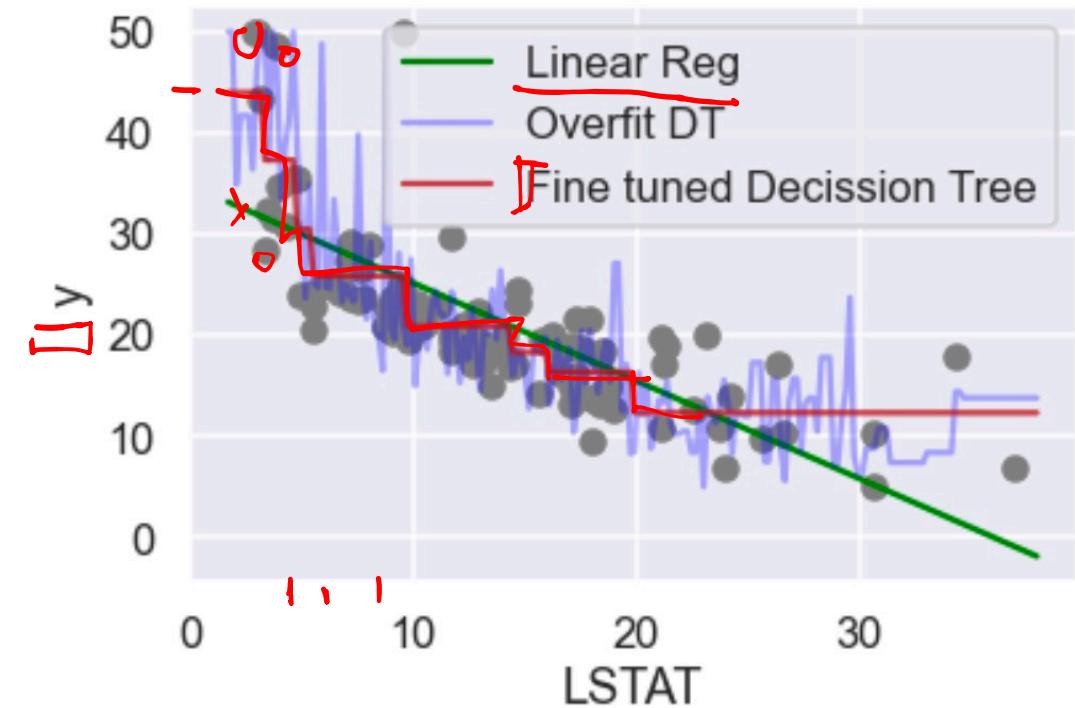
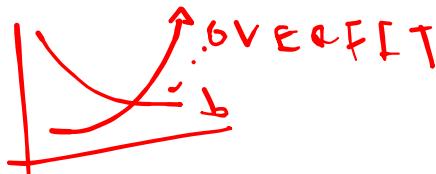
paramétricos $\rightarrow y = mx + b$

Alto bies: Simplificaciones sobre la forma de los datos, menos flexibles: Linear models

Baja Varianza: Pequeños cambios en las predicciones: Linear models

Bajo bies: Pocos supuestos sobre la forma de los datos, más flexibles: Decision trees, knn, ...

Alta Varianza: Grandes cambios en las predicciones: Decision trees, knn, ...



Bias-Variance Tradeoff

La única forma de comprobar que nuestro modelo generaliza bien es guardando una parte de los datos para medir el error *out-of-sample*, es decir, el error sobre los datos que no ha visto el modelo durante la fase de entrenamiento.

Con tal de poder elegir el mejor modelo y sus parámetros se dividirá inicialmente el dataset en dos subsets:

- subset de entrenamiento: normalmente se utiliza el 80% del tamaño de la muestra
- subset de testeo: el 20% restante.

