

Data Science & ML Course

Lesson #18 Linear Regression

Ivanovitch Silva
November, 2018





- We are going to start by covering **linear regression**
 - One variable
 - Multiples variables
- We discuss the application of linear regression to **housing price prediction**
- Present the notion of a **cost function**
- Introduce the **gradient descent** method for learning.
- Refresher on **linear algebra concepts**.

Update from repository

```
git clone https://github.com/ivanovitchm/datascience2machinelearning.git
```

Or

```
git pull
```



Regression



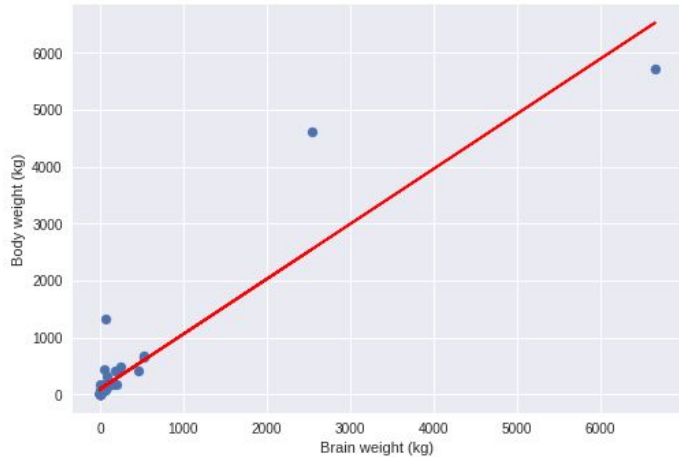
5 kilograms

200 kilograms

1.5 kilograms



????



1.5 kilograms



????

```

1 import pandas as pd
2 from sklearn import linear_model
3 import matplotlib.pyplot as plt
4
5 #read data
6 df = pd.read_fwf('brain_body.txt')
7 X = df[['Brain']]
8 y = df[['Body']]
9
10
11 #train model on data
12 model = linear_model.LinearRegression()
13 model.fit(X, y)
14
15 #visualize results
16 plt.scatter(X.values,y.values)
17 plt.plot(X.values,model.predict(X),color='red')
18 plt.xlabel('Brain weight (kg)')
19 plt.ylabel('Body weight (kg)')
20
21 plt.show()

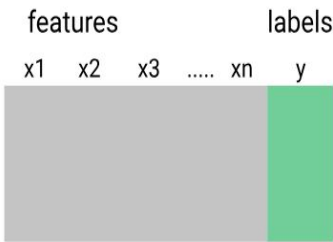
```


1_hello_world.ipynb



Training Process

Training Data



Select Features



Find the best model that best **approximates** training set.

Training is the most computationally intensive step.

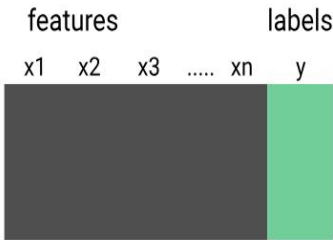
$O(n)^2$



$\hat{y} = 2.5x_1 + 11x_2 + 3.5x_3 + \dots + a_n x_n$

Testing Process

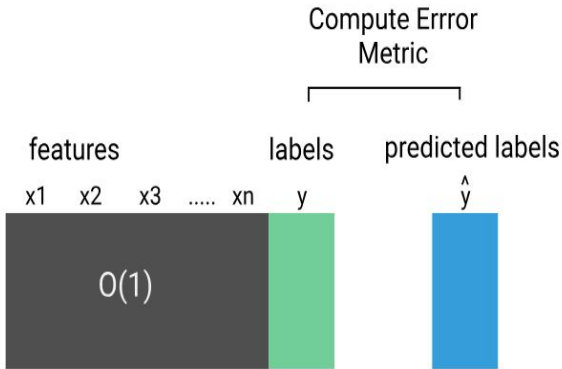
Test Data



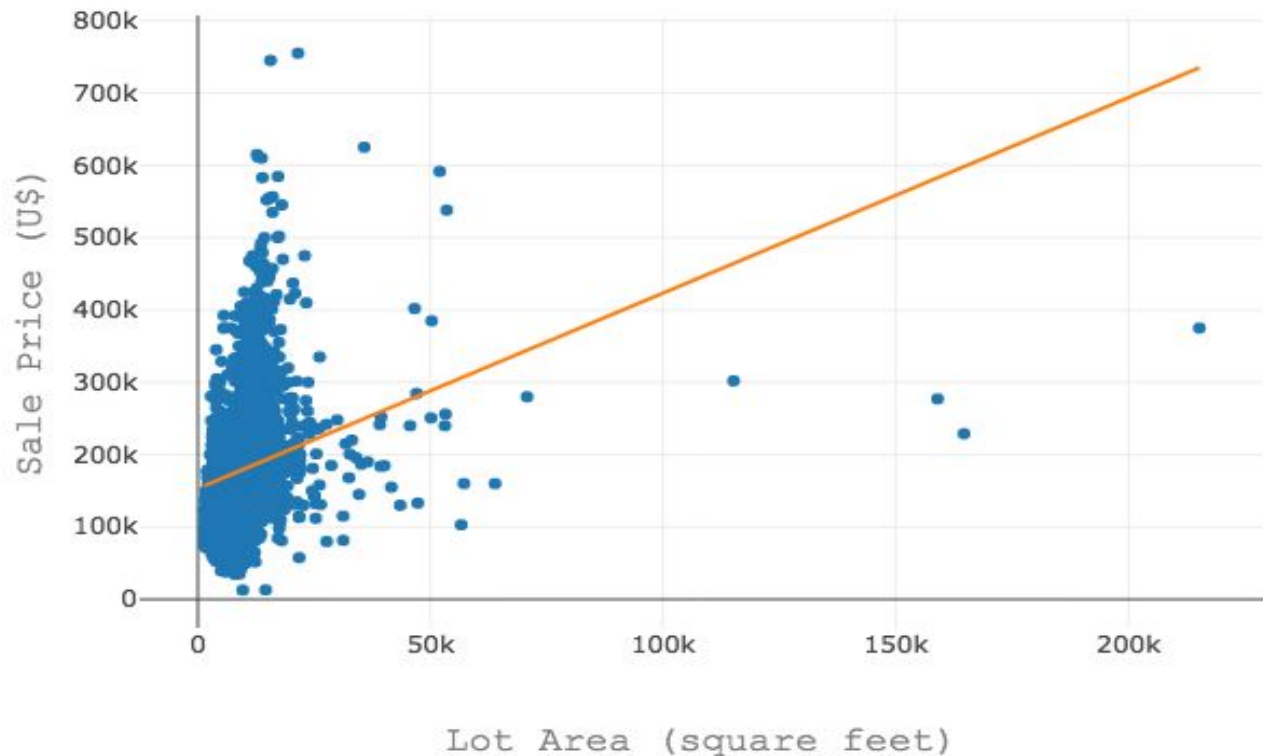
Predict label for an instance by plugging features into trained model.

$\hat{y} = 2.5x_1 + 11x_2 + 3.5x_3 + \dots + a_n x_n$

Predicting is computationally cheap.



Linear Regression - Housing Price



● Lot Area vs Sale Price
— Linear regression

Lot Area	Overall Qual	Year Built	Yr Sold	SalePrice
31770	6	1960	2010	215000
11622	5	1961	2010	105000
14267	6	1958	2010	172000
11160	7	1968	2010	244000
13830	5	1997	2010	189900

One variable
Multiple variables



Reproduce this fig. using plotly and scikit-learn



Read sections 1 and 2
2. The linear Regression Model.ipynb



Linear Regression with One Variable

Notation:

- m - number of training examples
- X 's - input variable/features
- y 's - output variable/ target variable

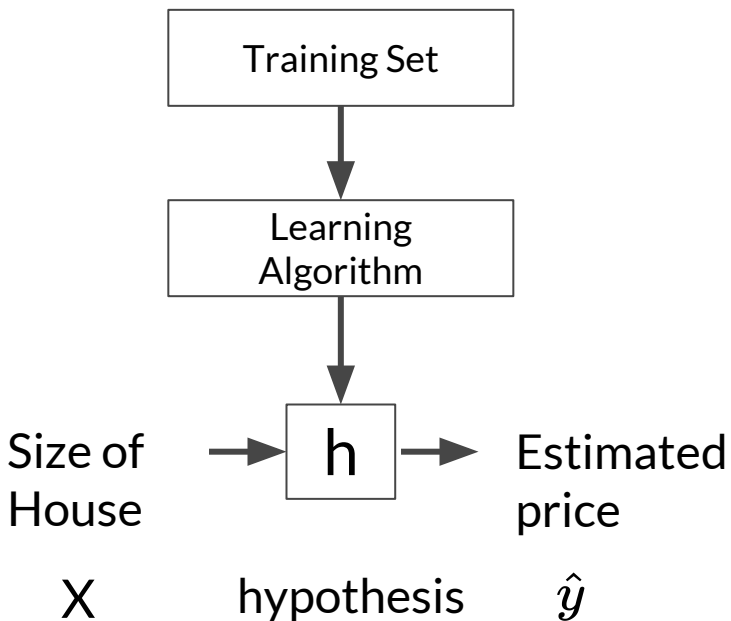
$$\begin{array}{ll} X^{(1)} = 31770 & y^{(1)} = 215000 \\ X^{(2)} = 11622 & y^{(2)} = 105000 \\ X^{(3)} = 14267 & y^{(3)} = 172000 \end{array}$$

$(X^{(i)}, y^{(i)}) = i^{\text{th}}$ training example

$m = 1465$

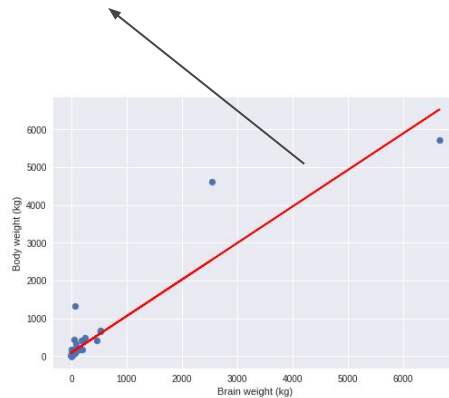
<div>X</div> <div>y</div>	
Lot Area SalePrice	
31770	215000
11622	105000
14267	172000
11160	244000
13830	189900

Model Representation (linear reg. one variable)



How do we represent h ?

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



Cost Function

$$f(x)$$

"minimize the error"

Cost Function (Linear Reg. One Var.)

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

θ_i = parameters

How to choose θ_i ?

$m = 1465$

Training Set

X

y

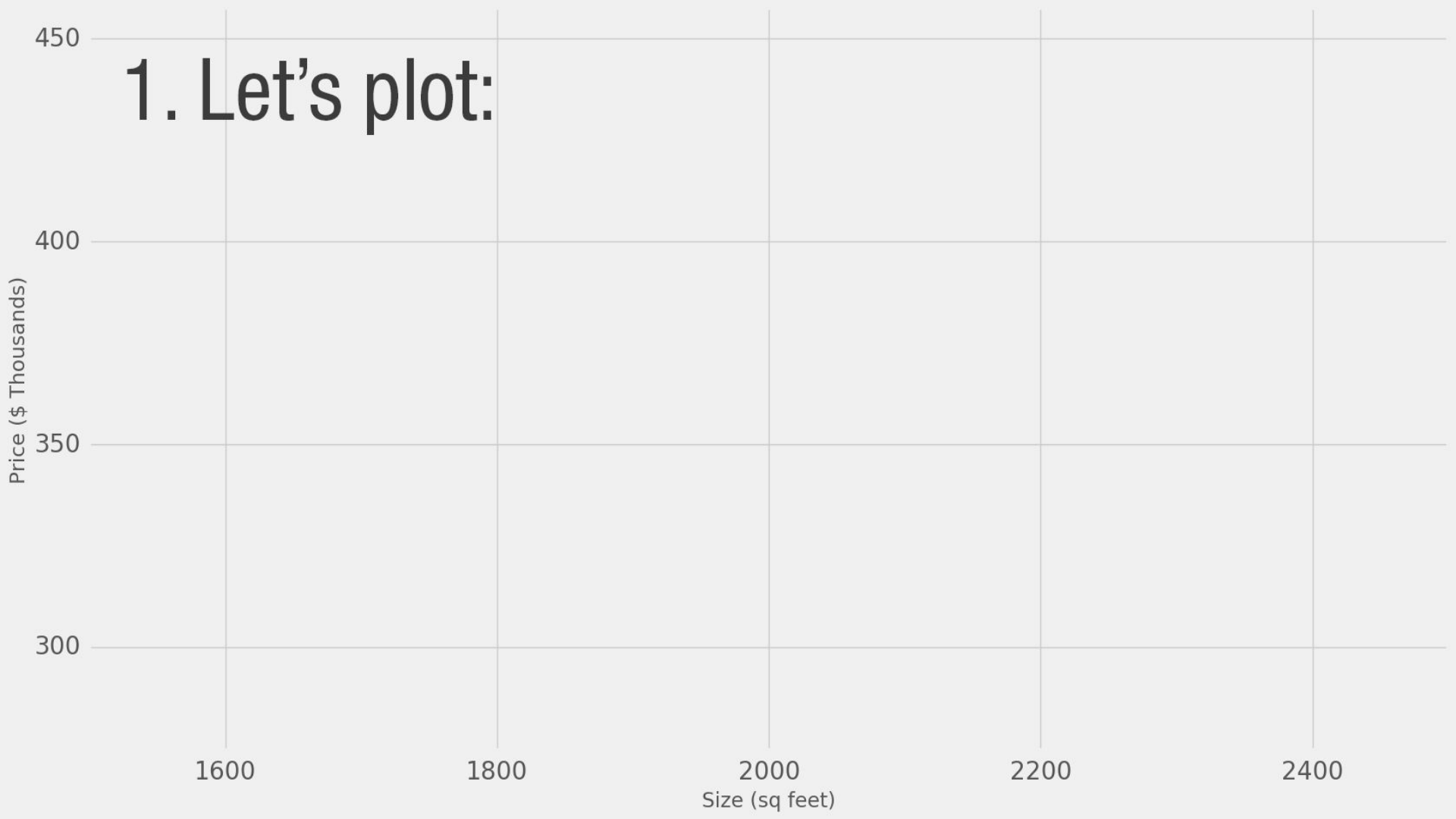
Lot Area	SalePrice
31770	215000
11622	105000
14267	172000
11160	244000
13830	189900

Cost Function

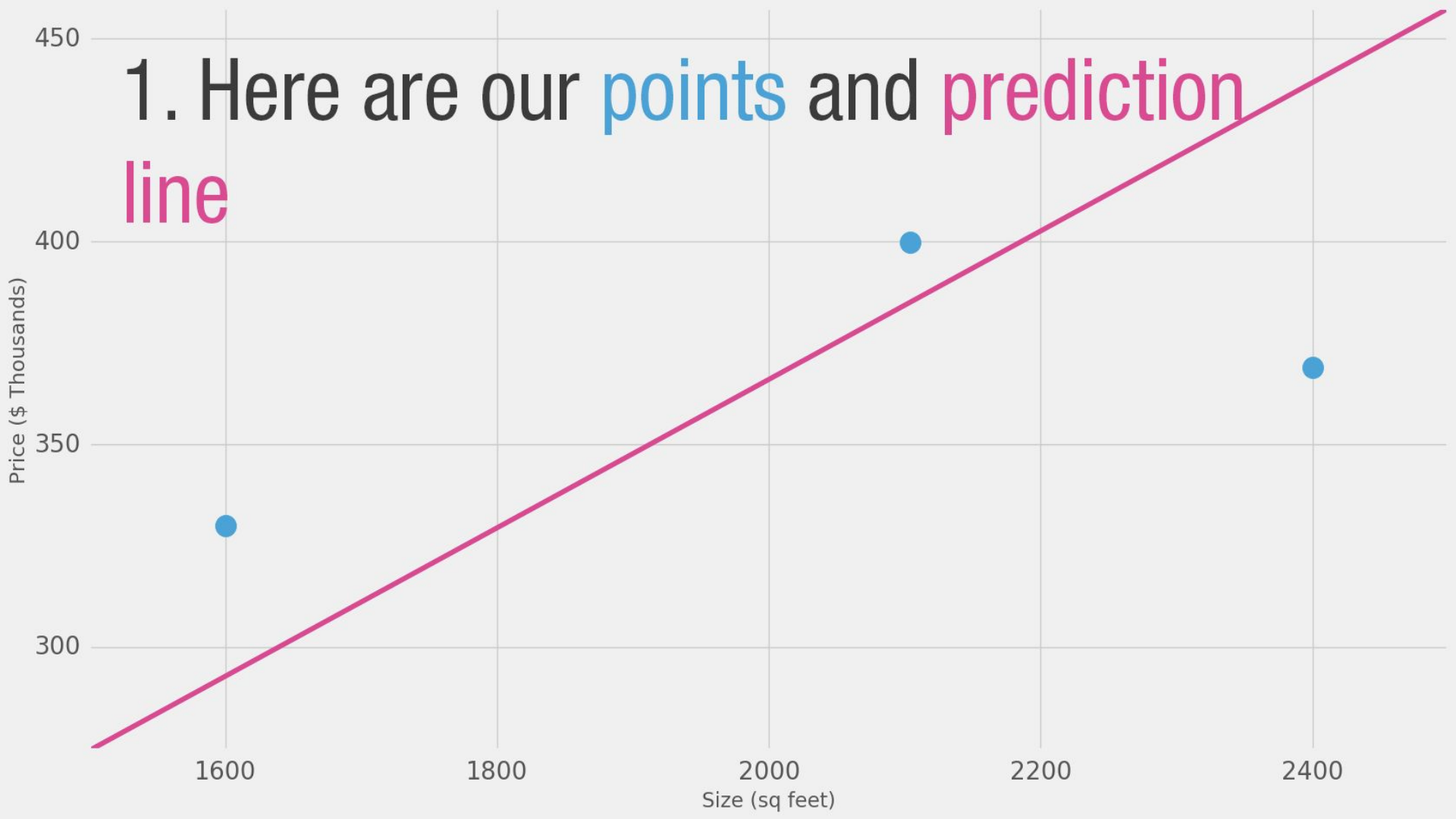
Intuition #01

(linear reg. One var)

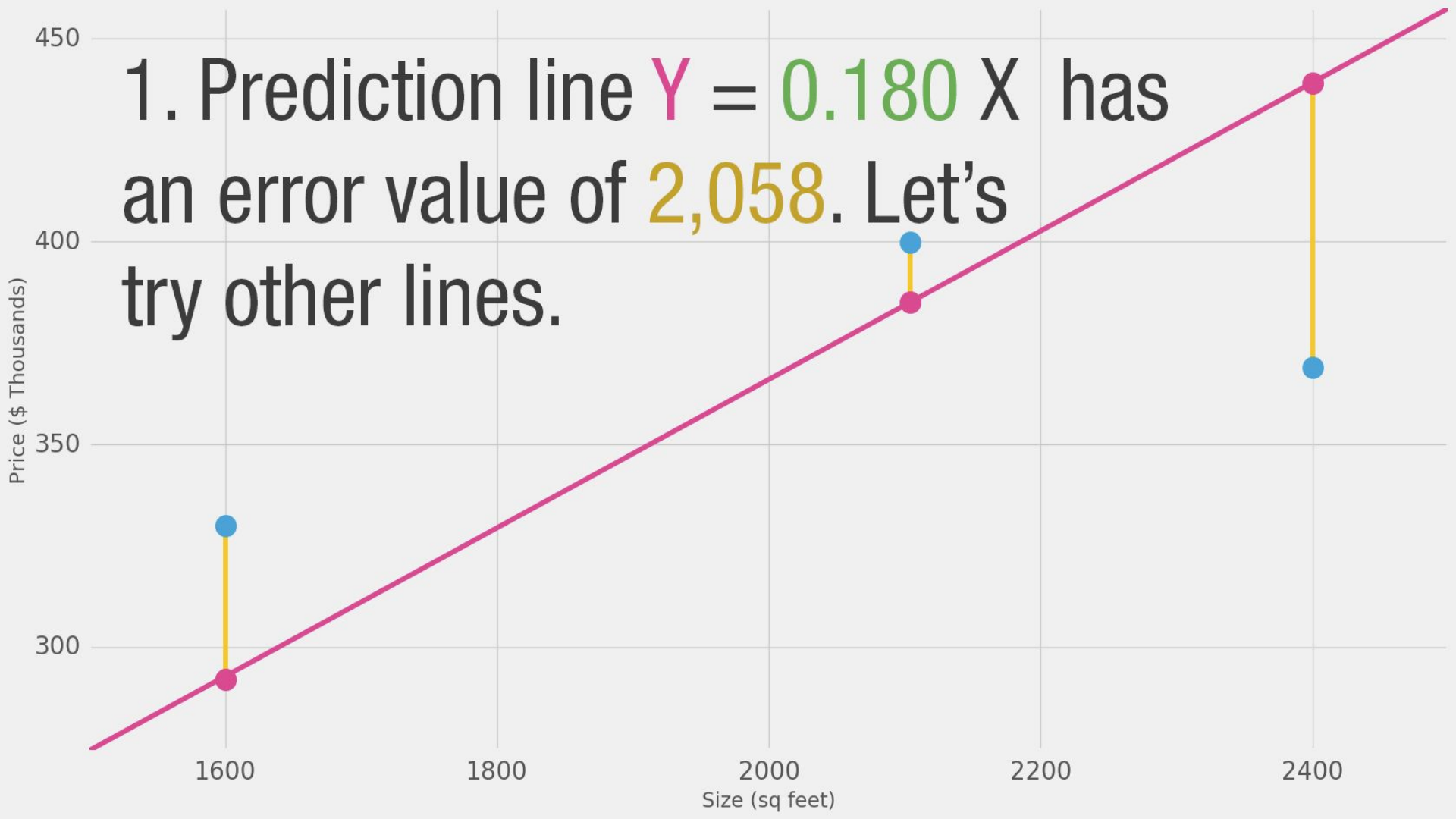
1. Let's plot:



1. Here are our **points** and **prediction**
line



1. Prediction line $Y = 0.180 X$ has an error value of 2,058. Let's try other lines.



Cost Function (square error function)

Training Set (m instances)

	Lot Area	SalePrice	
$x^{(1)}$	31770	215000	$y^{(1)}$
$x^{(2)}$	11622	105000	$y^{(2)}$
$x^{(3)}$	14267	172000	$y^{(3)}$
$x^{(4)}$	11160	244000	$y^{(4)}$
$x^{(5)}$	13830	189900	$y^{(5)}$

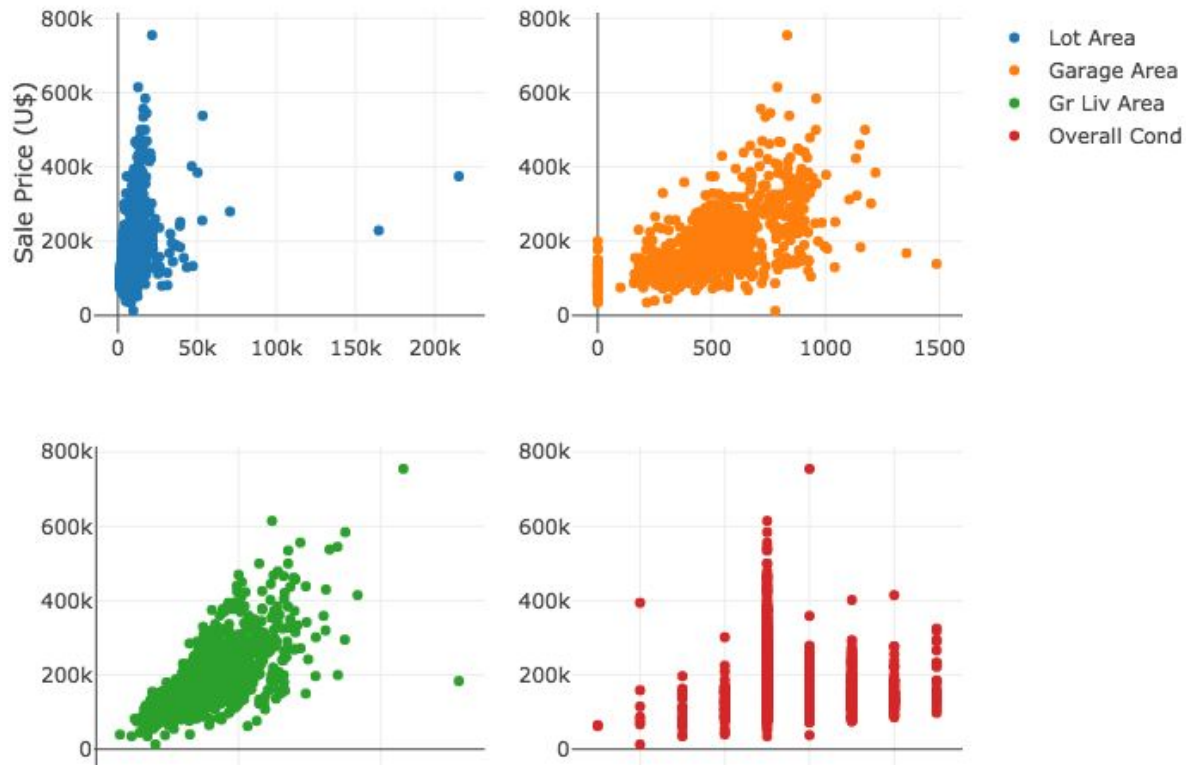
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]^2$$

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

Idea:

- choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples $(x^{(i)}, y^{(i)})$
- minimize (θ_0, θ_1)

Cost Function[step #01] - Select the feature x

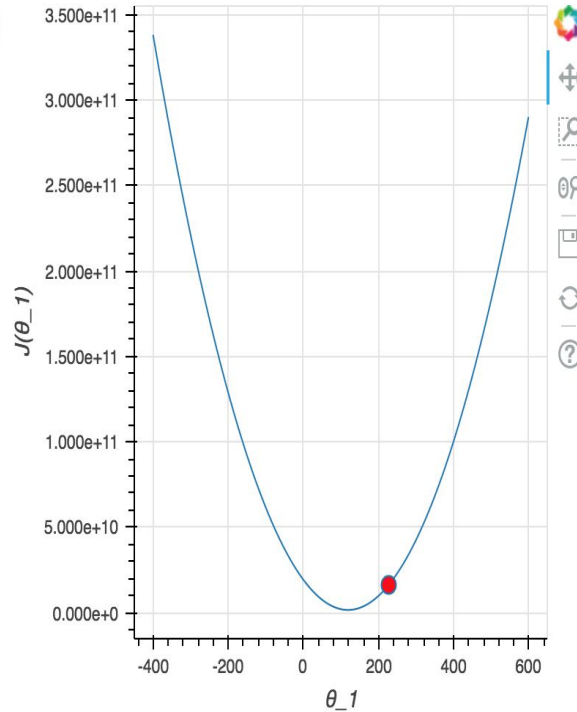
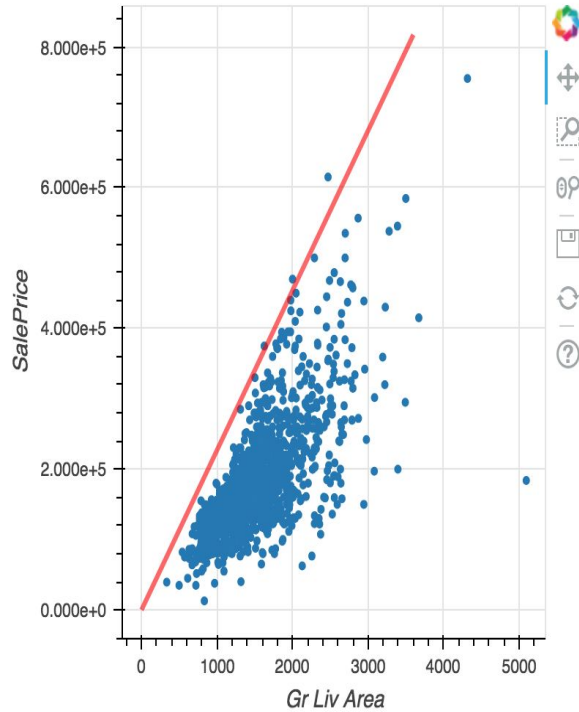


Cost Function[step #01] - Select the feature x_1

```
1 train[['Garage Area', 'Gr Liv Area', 'Overall Cond', 'Lot Area', 'SalePrice']].corr()
```

	Garage Area	Gr Liv Area	Overall Cond	Lot Area	SalePrice
Garage Area	1.000000	0.473506	-0.145705	0.213122	0.625335
Gr Liv Area	0.473506	1.000000	-0.134157	0.248676	0.706364
Overall Cond	-0.145705	-0.134157	1.000000	-0.042415	-0.108979
Lot Area	0.213122	0.248676	-0.042415	1.000000	0.267714
SalePrice	0.625335	0.706364	-0.108979	0.267714	1.000000

$J(227) = 47472029566758$



Cost Function

Intuition #01

$(\theta_0 = 0)$

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

$$\hat{y} = h_{\theta}(x) = \theta_1 x$$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m [\theta_1 x^{(i)} - y^{(i)}]^2$$

A1: 227

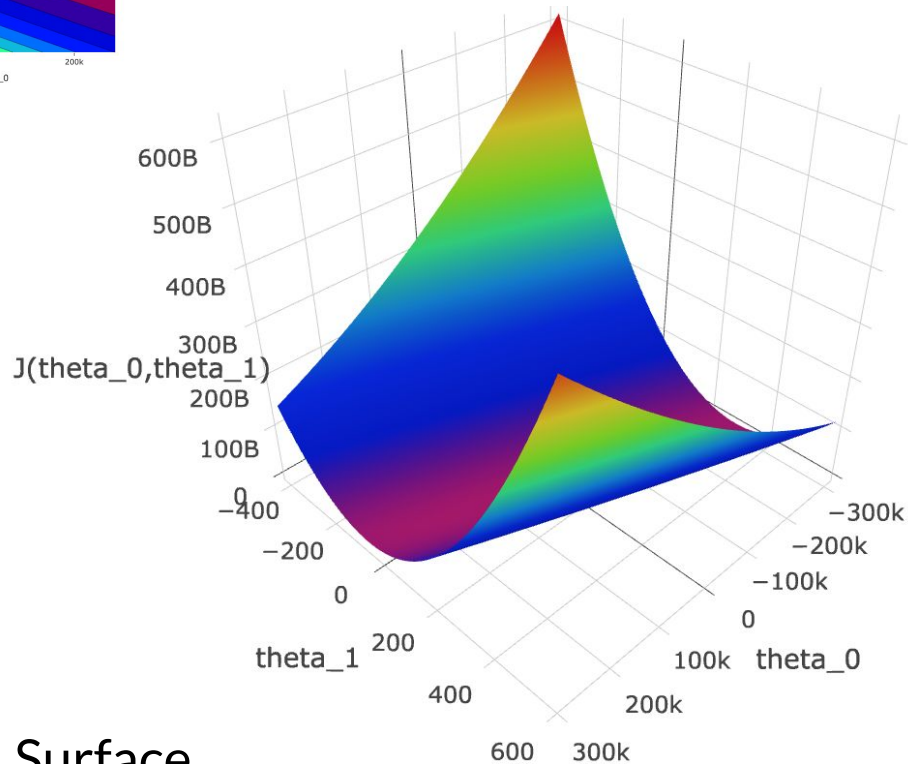
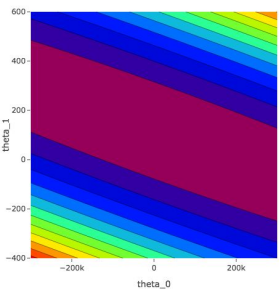


Cost Function

Intuition #02

(linear reg. One var)

Contour



Surface

Cost Function

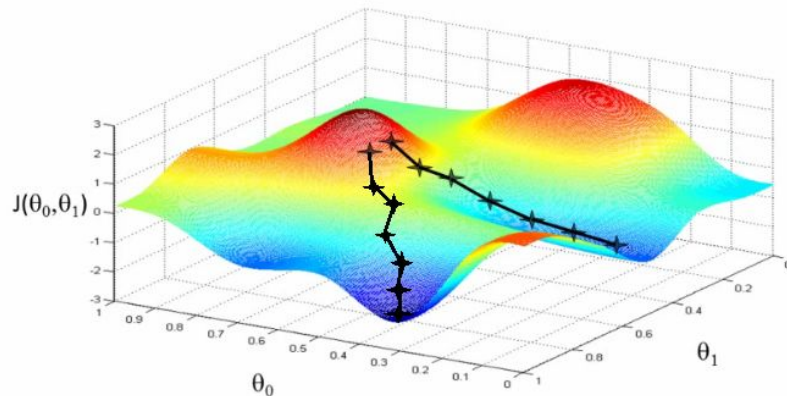
Intuition #02

$(\theta_0$ and θ_1 are defined)

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

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

Gradient Descent (linear reg. One var)

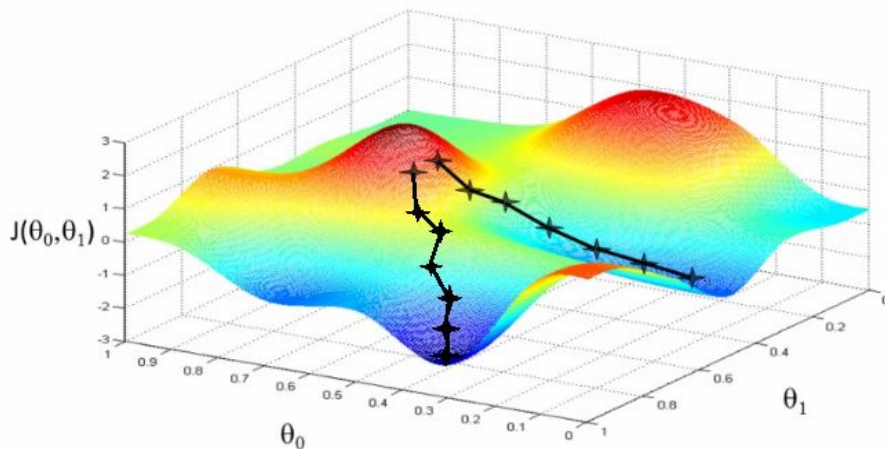


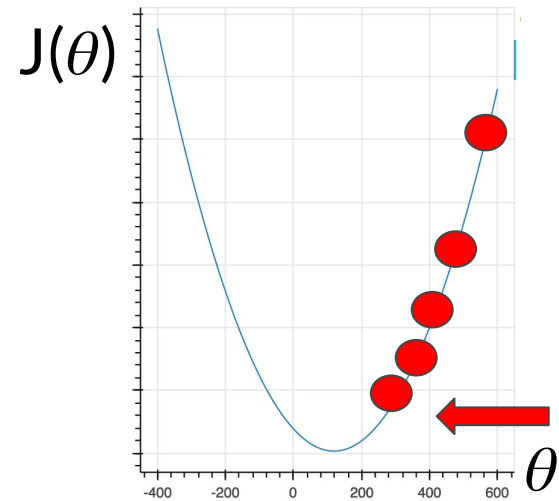
Algorithm - Idea

Have some function $J(\theta_0, \theta_1)$

Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Start with some θ_0, θ_1
Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum



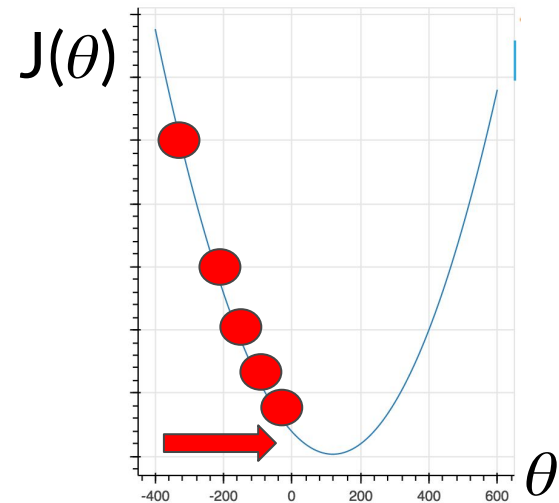


repeat until converge {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

α - learning rate



repeat until converge {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

Correct update

$$aux_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$aux_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 = aux_0$$

$$\theta_1 = aux_1$$

Incorrect update

$$aux_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 = aux_0$$

$$aux_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 = aux_1$$

repeat until converge {

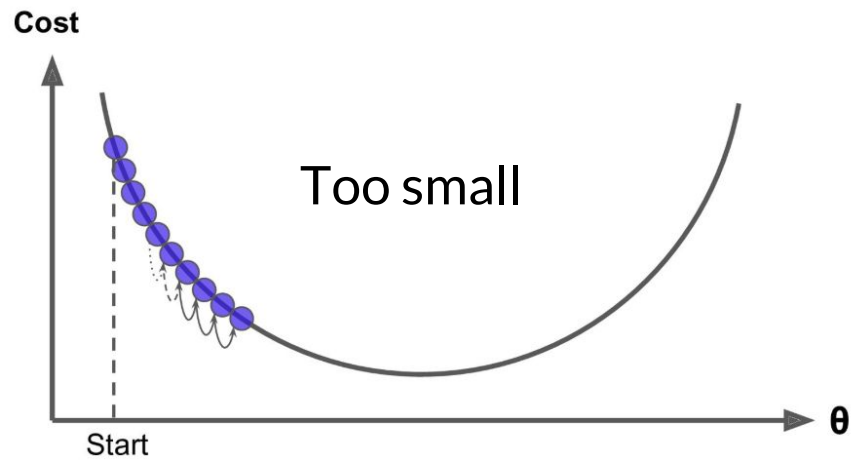
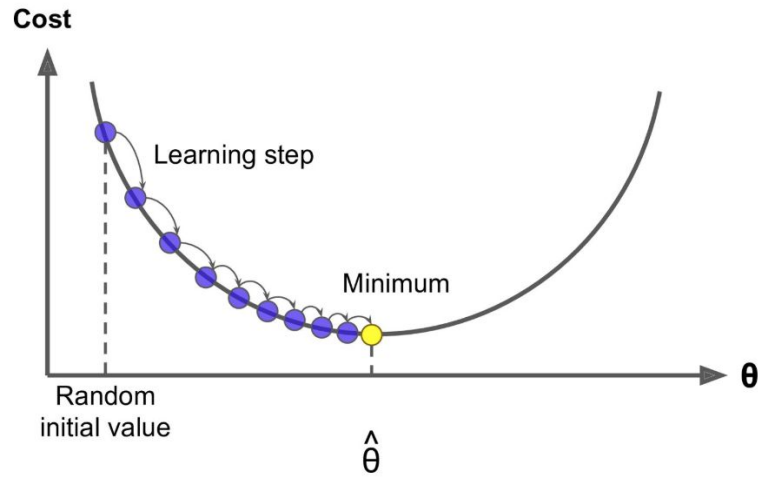
$$aux_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}]$$

$$aux_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] x^{(i)}$$

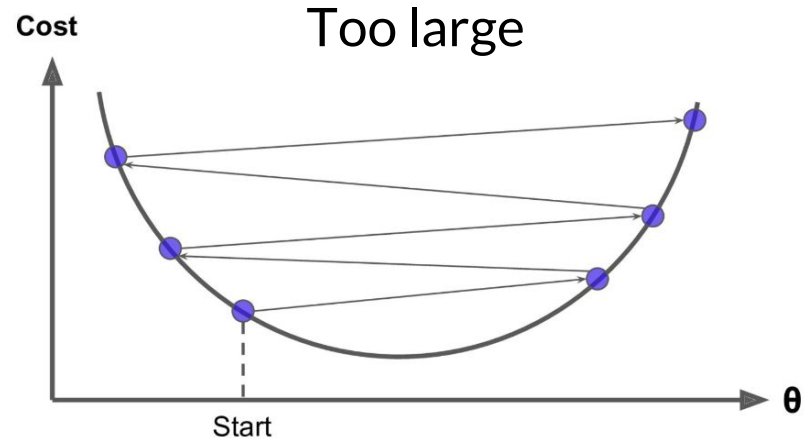
$$\theta_0 = aux_0$$

$$\theta_1 = aux_1$$

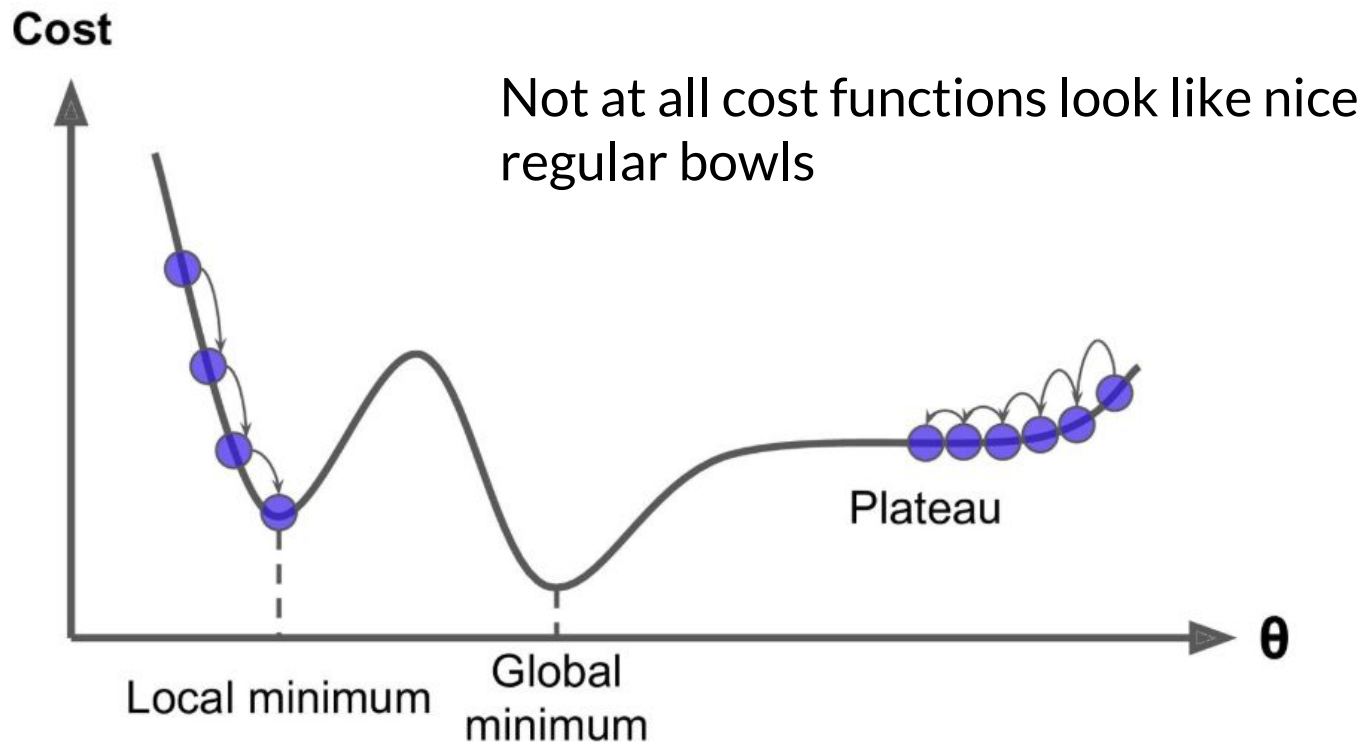
}



Learning rate
tradeoff



Gradient Descent Pitfalls



cost function & gradient descent from linear algebra perspective

Hypothesis

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

Gr Liv Area	SalePrice
2480	205000
1829	237000
2673	249000
1005	133500
1768	224900

[Export to plot.ly »](#)

$$\text{hypothesis} = \begin{bmatrix} 1 & 2480 \\ 1 & 1829 \\ 1 & 2679 \\ 1 & 1005 \\ 1 & 1768 \end{bmatrix} \times \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \begin{bmatrix} 2480 \theta_1 + \theta_0 \\ 1829 \theta_1 + \theta_0 \\ 2679 \theta_1 + \theta_0 \\ 1005 \theta_1 + \theta_0 \\ 1768 \theta_1 + \theta_0 \end{bmatrix}$$

```
def cost_function(X, y, theta):
    return np.sum(np.square(np.matmul(X, theta) - y)) / (2 * len(y))
```

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

Gr Liv Area	SalePrice
2480	205000
1829	237000
2673	249000
1005	133500
1768	224900

[Export to plot.ly »](#)

$$J(\theta_0, \theta_1) = \frac{1}{2 \times 5} \sum \left(\begin{bmatrix} 2480 \theta_1 + \theta_0 \\ 1829 \theta_1 + \theta_0 \\ 2679 \theta_1 + \theta_0 \\ 1005 \theta_1 + \theta_0 \\ 1768 \theta_1 + \theta_0 \end{bmatrix} - \begin{bmatrix} 205000 \\ 237000 \\ 249000 \\ 133500 \\ 224900 \end{bmatrix} \right)^2$$


```
def gradient_descent(X, y, alpha, iterations, theta):
    m = len(y)
    all_thetas = [theta]

    for i in range(iterations):
        t0 = theta[0] - (alpha / m) * np.sum(np.dot(X, theta) - y)
        t1 = theta[1] - (alpha / m) * np.sum((np.dot(X, theta) - y) * X[:,1])
        theta = np.array([t0, t1])
        all_thetas.append([t0,t1])

    return theta, np.array(all_thetas)
```

repeat until converge {

- Involves calculations over the full training set X at each gradient step
- It uses the whole batch of training data at every step.
- This is why the algorithm called **Batch Gradient Descent**

$$aux_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)})]$$

$$aux_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)})] x^{(i)}$$

$$\theta_0 = aux_0$$

$$\theta_1 = aux_1$$

$$\}$$

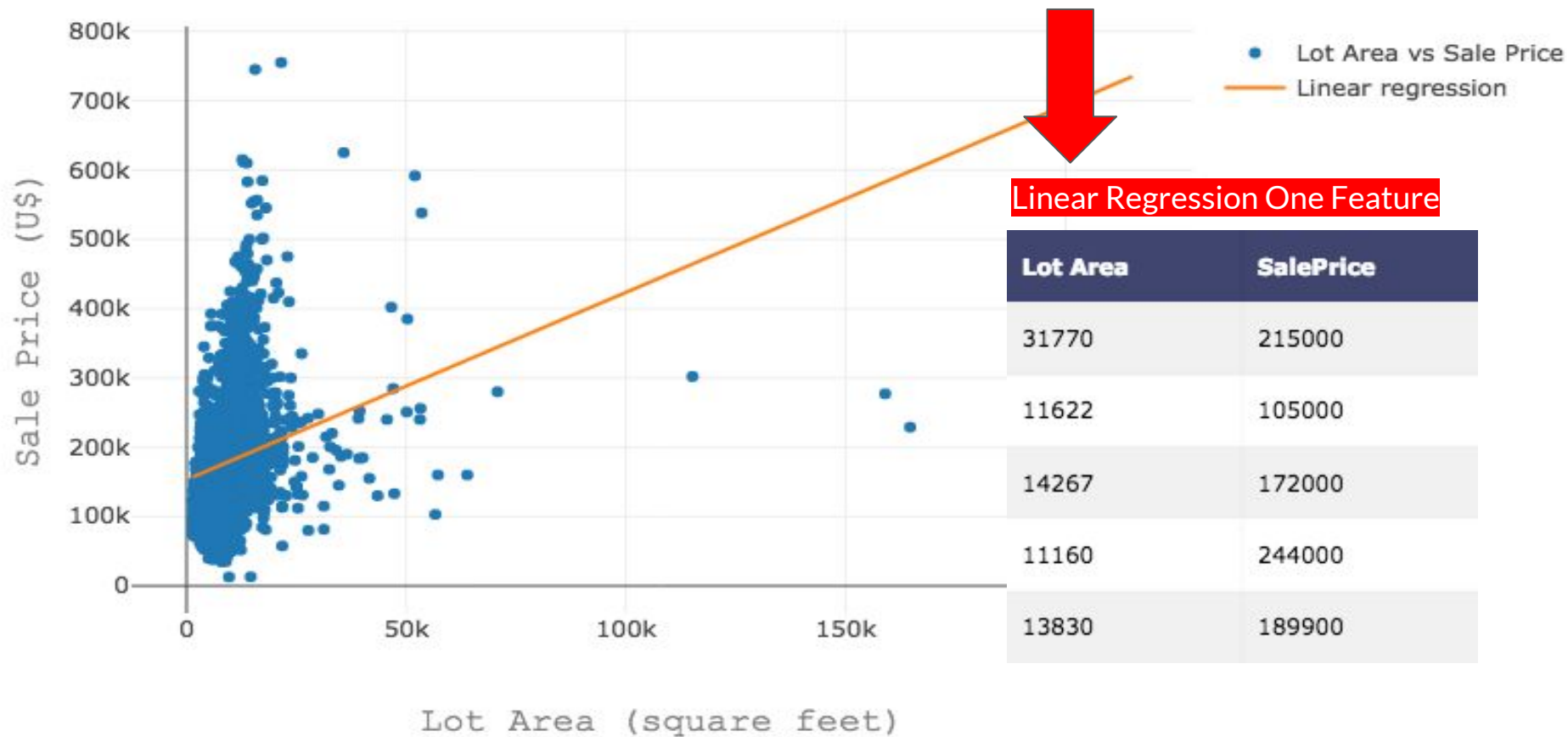
2_The linear Regression Model.ipynb



PREVIOUSLY ON...

lesson #18

Linear Regression - Housing Price



Training Set (m instances)

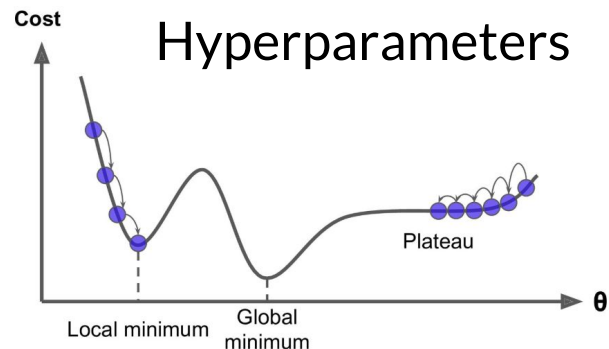
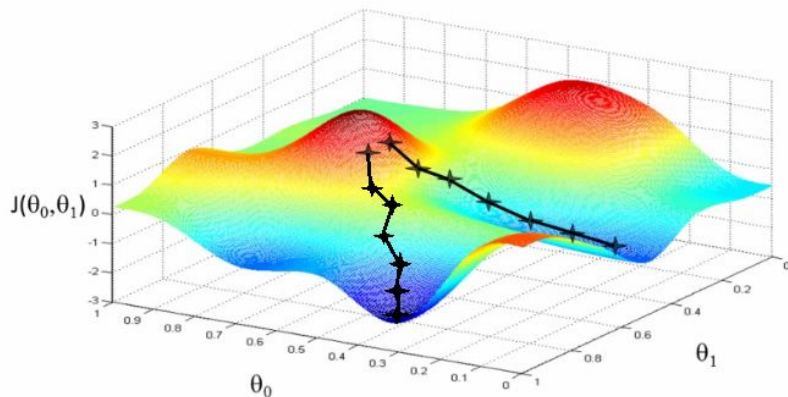
	Lot Area	SalePrice	
$x^{(1)}$	31770	215000	$y^{(1)}$
$x^{(2)}$	11622	105000	$y^{(2)}$
$x^{(3)}$	14267	172000	$y^{(3)}$
$x^{(4)}$	11160	244000	$y^{(4)}$
$x^{(5)}$	13830	189900	$y^{(5)}$

Cost Function

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

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

Gradient Descent





- We are going to start by covering linear regression
 - Multiple variables
- We discuss the application of linear regression to **housing price prediction**

Linear Regression with Multiple Variables

Notation:

- m - number of training examples
- n - number of features
- $x^{(i)}$ - input features of i^{th} training example
- $x_j^{(i)}$ - value of feature j in i^{th} training example
- $y^{(i)}$ - target value of i^{th} training examples

$$x^{(2)} = \begin{bmatrix} 11622 \\ 5 \\ 1961 \\ 2010 \\ 105000 \end{bmatrix}$$

$$x_3^{(2)} = 1961$$

$n = 4$

x_1	x_2	x_3	x_4	y
Lot Area	Overall Qual	Year Built	Yr Sold	SalePrice
31770	6	1960	2010	215000
11622	5	1961	2010	105000
14267	6	1958	2010	172000
11160	7	1968	2010	244000
13830	5	1997	2010	189900

$m = 5$

Hypothesis (previously)

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

Multivariable case

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

For convenience of notation, define $x_0=1$. In other words:

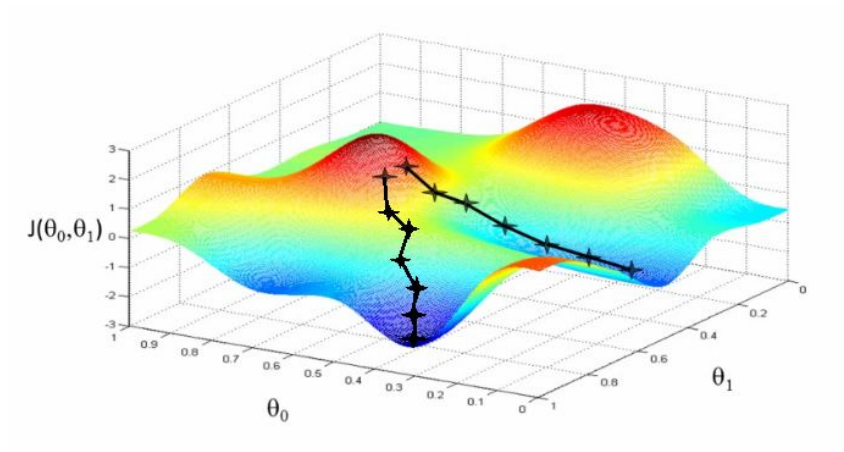
$$x_0^{(i)}=1$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \begin{bmatrix} \theta_0 & \theta_1 & \theta_1 & \dots & \theta_n \end{bmatrix} \times \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

Gradient Descent (linear reg. multiple variables)



Hypothesis: $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_n$

Cost function:

$$J(\theta_0, \theta_1, \theta_2, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Gradient descent:

repeat {

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$$

}

Simultaneously update for every j (0 to n)

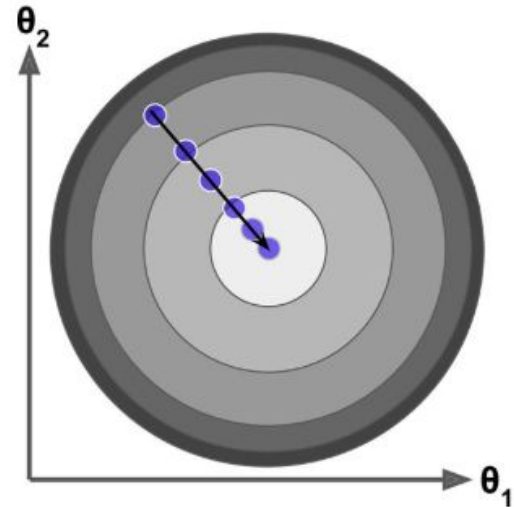
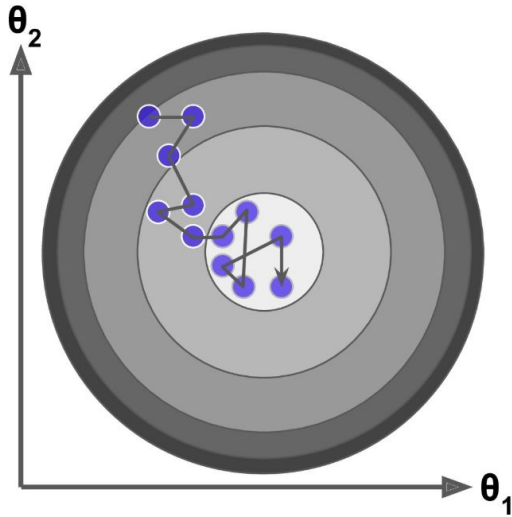
Gradient descent:

repeat until the convergence {

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad \text{for } j = 0, \dots, n$$

}

Gradient Descent: **trick #1** - Feature Scaling



Gradient Descent: **trick #1** - Feature Scaling

Z-Score or Standardization

$$z = \frac{x - \mu}{\sigma}$$

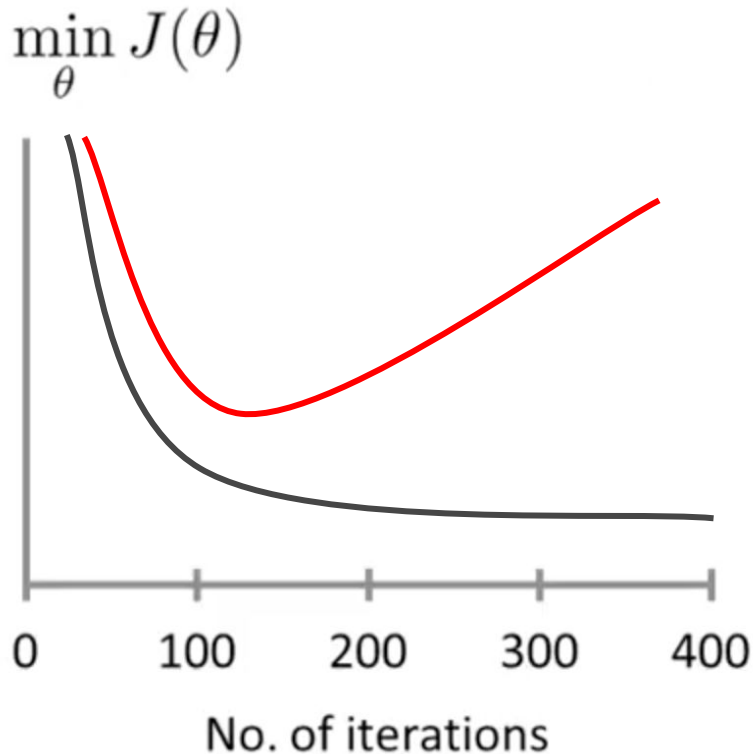
$$\begin{aligned}\mu &= 0 \\ \sigma &= 1\end{aligned}$$

Min-Max Scaling

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

range (0,1)

Gradient Descent: **trick #2** - Debugging α



1. Make a plot with number of iterations on the x-axis.
2. Now plot the cost function, $J(\theta)$ over the number of iterations of gradient descent.
3. If $J(\theta)$ ever increases, then you probably need to decrease α .
4. It has been proven that if learning rate α is sufficiently small, then $J(\theta)$ will decrease on every iteration.
5. Automatic convergence test

Gradient Descent: **trick #2** - Debugging α

- If α is too small:
 - Slow convergence
- If α is too large:
 - $J(\theta)$ may not decrease on every iteration;
 - $J(\theta)$ may not converge.

To choose α :

..., 0.0001, ..., 0.001, ..., 0.01, ..., 0.1, ..., 1, ..., 10, ...

normal equation: method to
solve for θ analytically

Normal Equation

Lot Area	Overall Qual	Year Built	Yr Sold	SalePrice
31770	6	1960	2010	215000
11622	5	1961	2010	105000
14267	6	1958	2010	172000
11160	7	1968	2010	244000
13830	5	1997	2010	189900

$$X\theta = y$$

$$X^T X \theta = X^T y$$

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

$$\begin{bmatrix} 1 & 31770 & 6 & 1960 & 2010 \\ 1 & 11622 & 5 & 1961 & 2010 \\ 1 & 14267 & 6 & 1958 & 2010 \\ 1 & 11160 & 7 & 1968 & 2010 \\ 1 & 13830 & 5 & 1997 & 2010 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \\ \theta_5 \end{bmatrix} = \begin{bmatrix} \theta_0 + 31770\theta_1 + 6\theta_2 + 1960\theta_3 + 2010\theta_4 \\ \theta_0 + 11622\theta_1 + 5\theta_2 + 1961\theta_3 + 2010\theta_4 \\ \theta_0 + 14267\theta_1 + 6\theta_2 + 1958\theta_3 + 2010\theta_4 \\ \theta_0 + 11160\theta_1 + 7\theta_2 + 1968\theta_3 + 2010\theta_4 \\ \theta_0 + 13830\theta_1 + 5\theta_2 + 1997\theta_3 + 2010\theta_4 \end{bmatrix}$$

There is **no need** to do feature scaling with the normal equation.

The following is a comparison of gradient descent and the normal equation:

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
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

With the normal equation, computing the inversion has complexity $\mathcal{O}(n^3)$. So if we have a very large number of features, the normal equation will be slow. In practice, when n exceeds 10,000 it might be a good time to go from a normal solution to an iterative process.

If $X^T X$ is **noninvertible**, the common causes might be having :

- Redundant features, where two features are very closely related (i.e. they are linearly dependent)
- Too many features (e.g. $m \leq n$). In this case, delete some features or use "regularization" (to be explained in a later lesson).

