# Machine Learning (ML)
# with Python

## Linear Model

Dr. Aeshah Alsughayyir

Collage of Computer Science and Engineering

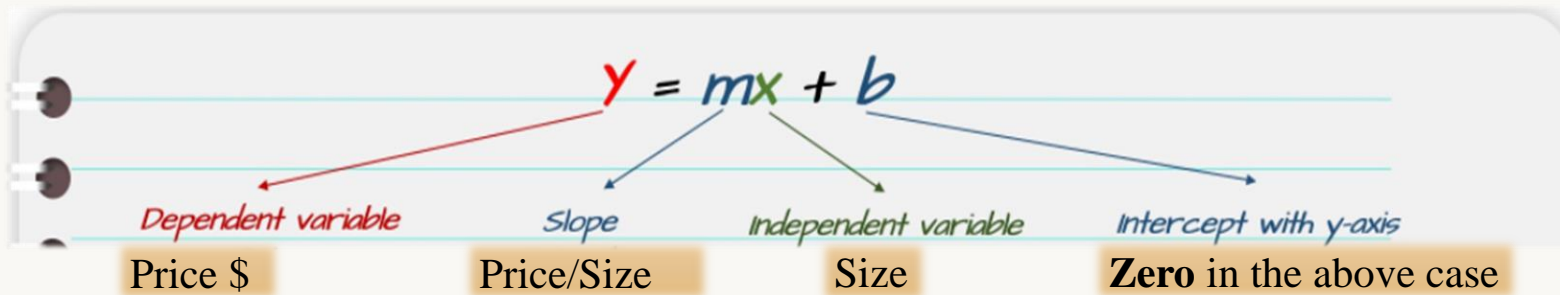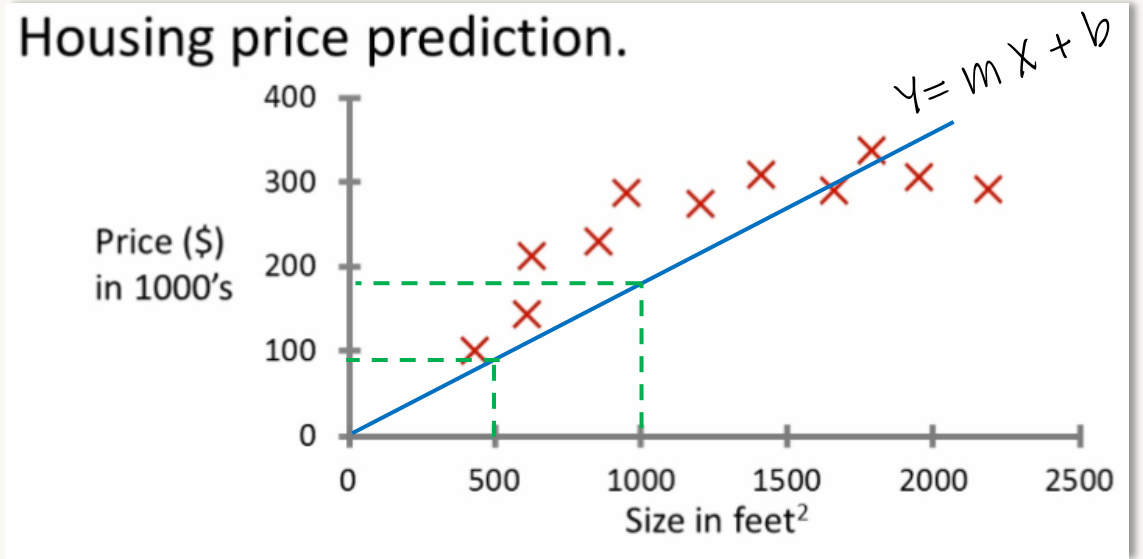Taibah University

2021-2022

# Linear Model

- Model Representation and Hypothesis
- Cost Function of Linear Regression
- Gradient Descent

Dr. Aeshah Alsughayyir

# Motivation Example

| Living area (feet$^2$) | Price (1000$s) |
|:---:|:---:|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| $\vdots$ | $\vdots$ |

## Housing price prediction.

$y = mx + b$

Price ($) in 1000's

Size in feet$^2$

$Y = mx + b$

- **Y** → Dependent variable → Price $
- **m** → Slope → Price/Size
- **x** → Independent variable → Size
- **b** → Intercept with y-axis → **Zero** in the above case

$$Slop(m) = \frac{Price}{Size}$$
$$= \frac{y2 - y1}{x2 - x1}$$
$$= \frac{190 - 90}{1000 - 500} = 0.2$$

Dr. Aeshah Alsughayyir

# Model Representation and Hypothesis

**Notations**

- $m$: Number of training examples

- $x$'s: Input variables / features

- $y$'s: Output / Target variables

- $(x, y)$: One training example

- $(x^{(i)}, y^{(i)})$: $i^{th}$ training example

$x$'s

$y$'s

| Living area (feet$^2$) | Price (1000\$s) |
|---|---|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| $\vdots$ | $\vdots$ |

$(x,y)$

$m$

$(x^{(i)},y^{(i)})$

# Model Representation and Hypothesis (Cont.)

- If we could find the equation of line **y = mx+b** that we use to fit the data represented by the <span style="color:blue">blue</span> inclined line (see the motivation example), <u>then we can easily find the model that can predict the housing prices for any given area.</u>

  For example: **What would be the best-estimated price for area 2450 feet square?**

- In machine learning lingo, function **y = mx+b** is also called a **hypothesis function** where **m** and **b** can be represented by $theta_0$ and $theta_1$ respectively.  $theta_0$ is called a ***bias term or intercept,***  and $theta_0$, $theta_1$,.. are also called ***coefficient or weights***.

> The job of the learning algorithm would be to produce a <span style="color:red">hypothesis (**h**),</span> such that it takes the size of house as input and predicts its price
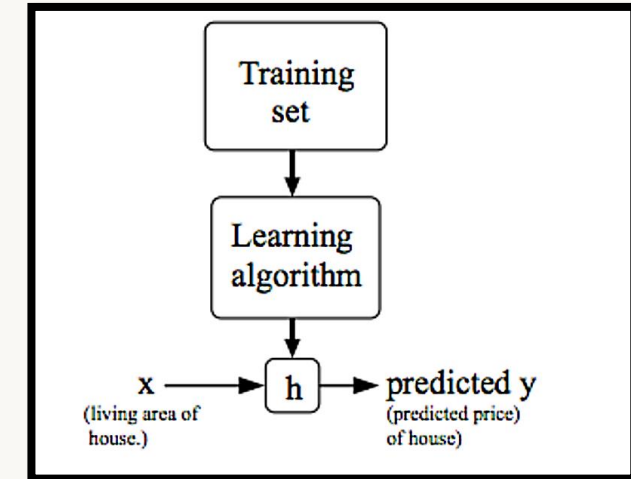
# Model Representation and Hypothesis (Cont.)

**Hypothesis** is the function that is to be learnt by the learning algorithm by the training process for making the **predictions** about the **unseen** data.

- For example, for **Linear Regression of One Variable** or **Univariate Linear Regression**, the hypothesis, $h$ is given by:

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

- Where
  - $h_\theta(x)$ is the hypothesis function, also denoted as $h(x)$ sometimes
  - $x$ is the independent variable
  - $\theta_0$ and $\theta_1$ are the parameters of the linear regression that need to be learnt

Different values for these parameters will give different hypothesis function based on the values of slope and intercepts.



Training set

Learning algorithm

x → h → predicted y
(living area of house.) (predicted price) of house)

Dr. Aeshah Alsughayyir

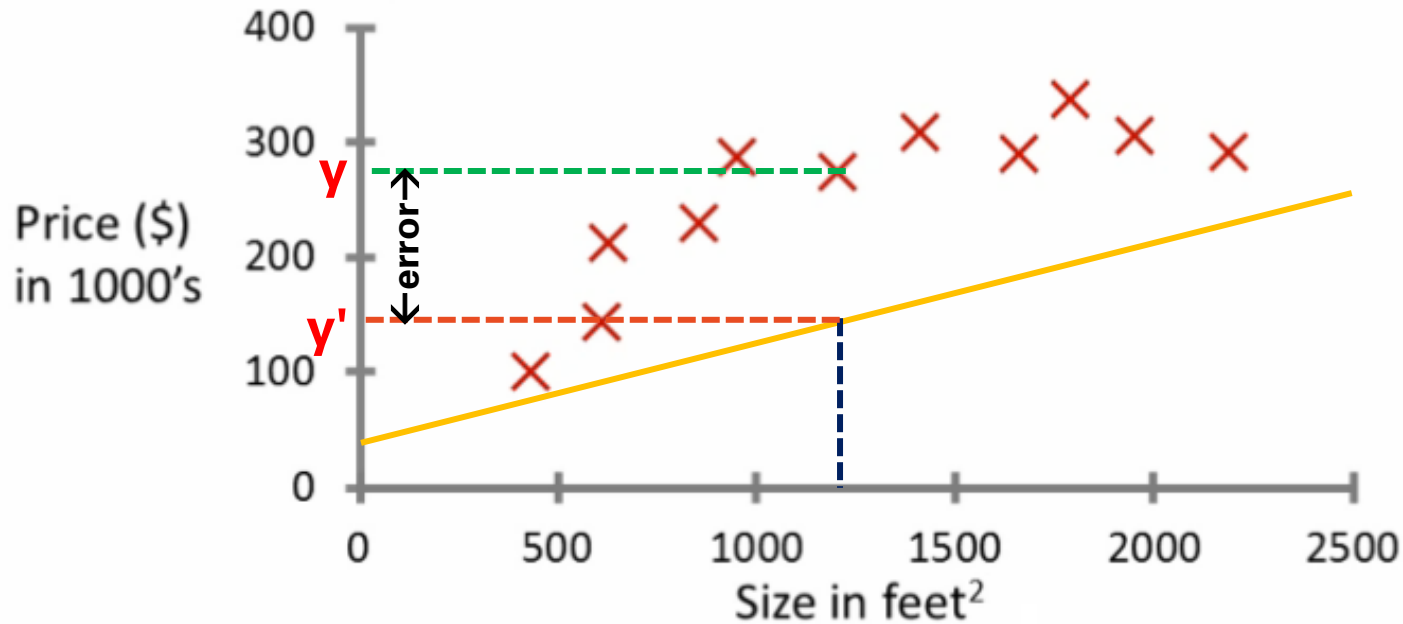# Cost Function of Linear Regression

- Once the parameter values, i.e., **intercept** and **$theta_1$** are randomly initialized:

  - the hypothesis function is ready for prediction,

  - then the **error** (|**predicted value** – **actual value**|) is calculated to check whether the randomly initialized parameter is giving the right prediction or not.

  - Repeat{

    - If (the error is too high) then

      - the algorithm updates the parameters with new values

    }

The algorithm continues this process **until the error is <u>minimized</u>**

To **minimize the error,** we have a special algorithm called *Gradient Descent* but before that, we are going to understand <u>what **Cost Function** is and how it works?</u>

# Cost Function of Linear Regression (Cont.)

## Housing price prediction.



$$\text{cost} = \frac{1}{m}\sum_{i=1}^{m}(y' - y)$$

**Mean Square Error**

$$\text{MSE} = \frac{1}{2m}\sum_{i=1}^{m}(y' - y)^2$$

**Root Mean Square Error**

$$\text{RMSE} = \sqrt{\frac{1}{m}\sum_{i=1}^{m}(y' - y)^2}$$

**Mean Absolute Error**

$$\text{MAE} = \frac{1}{m}\sum_{i=1}^{m}|y' - y|$$

Dr. Aeshah Alsughayyir

# Cost Function of Linear Regression (Cont.)

So, **cost function** is defined as follows,

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

Predicted by **h**          Actual true values

And **learning objective is to minimize the cost function** i.e.

$$minimize_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

Dr. Aeshah Alsughayyir

# Cost Function of Linear Regression (Cont.)

- **Linear regression** finds the parameters *(weights)* that minimize the _mean squared error_ between predictions and the true regression targets *(y)*, on the training set.
    - i.e., choose $\theta_0$ and $\theta_1$ so that $h(x)$ is close to $y$ for the training examples *(x, y)*.

- The _mean squared error_ is the sum of the squared differences between the predictions and the true values.

- This can be mathematically represented as,

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

- Where
    - $h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$
    - $(x^{(i)}, y^{(i)})$ is the $i^{th}$ training data
    - m is the number of training example



**Fit model by *minimizing mean of squared errors***

Dr. Aeshah Alsughayyir

# Cost Function of Linear Regression (Cont.)

Terminologies:



Training example: $x^{(i)}$

Error or loss: $h_\theta(x^{(i)}) - y^{(i)}$

Hypothesis Function: $h_\theta(x) = \theta_0 + \theta_1 x$

**Parameters:** $\theta_0, \theta_1$

Actual Value: $y^{(i)}$

Predicted Value: $h_\theta(x^{(i)})$

**Cost Function:** $J(\theta, \theta_1) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$

*Note: (i) in the equation represents the ith training example, not the power.*

Dr. Aeshah Alsughayyir

# Cost Function of Linear Regression (Cont.)

One Feature



two Features



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

# Cost function visualization

- Consider a simple case of hypothesis **by setting $\theta_0=0$**, then $h$ becomes : $h_\theta(x)=\theta_1 x$
- Each value of $\theta_1$ corresponds to a different hypothesis as it is the **slope** of the line.

which corresponds to different lines passing through the **origin** as shown in plots below as **y-intercept** i.e., $\theta_0$ is nulled out.

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

Calculating each value of $\theta_1$ corresponds to a different hypothesizes:

At $\theta_1=2$, $\quad J(2) = \dfrac{1}{2*3}(1^2 + 2^2 + 3^2) = \dfrac{14}{6} = 2.33$

At $\theta_1=1$, $\quad J(1) = \dfrac{1}{2*3}(0^2 + 0^2 + 0^2) = 0$

At $\theta_1=0.5$, $\quad J(0.5) = \dfrac{1}{2*3}(0.5^2 + 1^2 + 1.5^2) = 0.58$



Simple Hypothesis

# Cost function visualization (Cont.)

What is the optimal value of $\theta_1$ that minimizes $J(\theta_1)$ ?

It is clear that best value for $\theta_1$ =1 as $J(\theta_1) = 0$, which is the minimum.

How to find the best value for $\theta_1$ ?

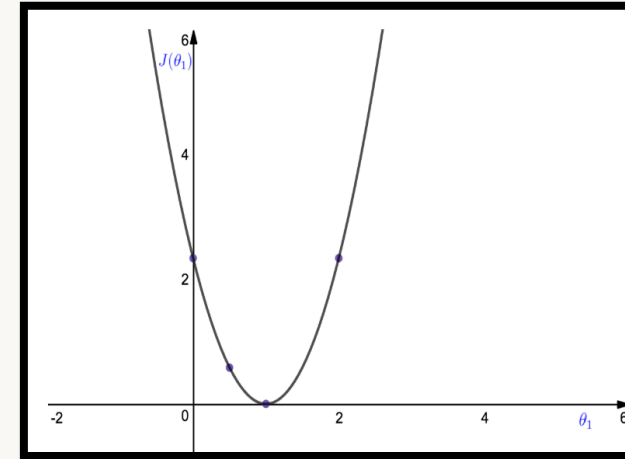Plotting ?? Not practical specially in high dimensions!

**The solution :**

1. Analytical solution: Not applicable for large datasets
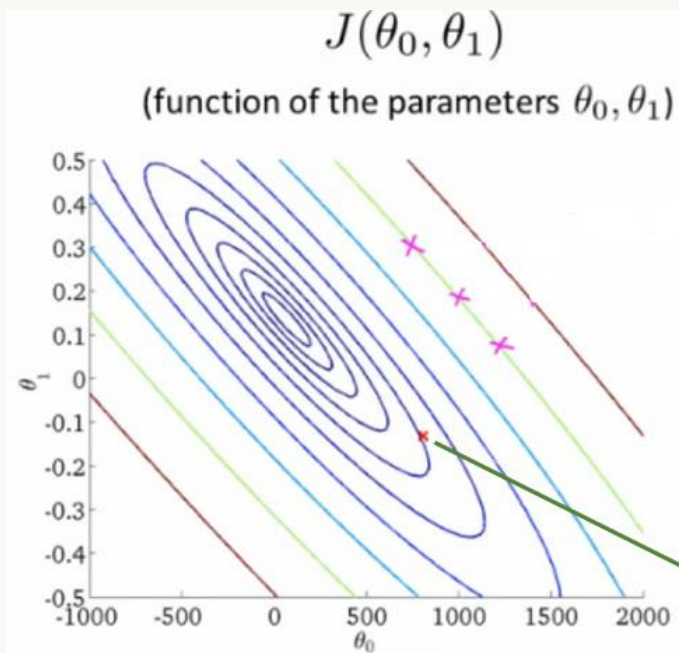2. Numerical solution: e.g., Gradient descent.

# Cost function visualization (Cont.)

- What if we plot our cost function with respect to the value of $\theta_1$ :
  - $\theta_1$ vs $J(\theta_1)$



- What if we plot our cost function with respect to the values of both $\theta_1$ and $\theta_0$ :
  - *Plot becomes a bit more complicated*
  - *Generates a 3D surface plot where axis are:*
    - $X = \theta_1$
    - $Z = \theta_0$
    - $Y = J(\theta_0, \theta_1)$

# Cost function visualization (Cont.)

- The following is a contour plot *of the cost function* where:
  - Ellipses in different color
  - Each color shows the same value of $J(\theta_0, \theta_1)$
  - Imagine a bowl-shape function coming out of the screen, so the middle is the concentric circle.



$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

Each point (like the red one above) represents a pair of parameter values for Θ0 and Θ1

# Cost function visualization (Cont.)

○ Our example here put the values at
  ■ $\theta_0 = {\sim}800$
  ■ $\theta_1 = {\sim}{-}0.15$
○ Not a good fit
  ■ i.e. these parameters give a value on our contour plot far from the center

# Cost function visualization (Cont.)

○ If we have
- $\theta_0 = \sim 360$
- $\theta_1 = 0$
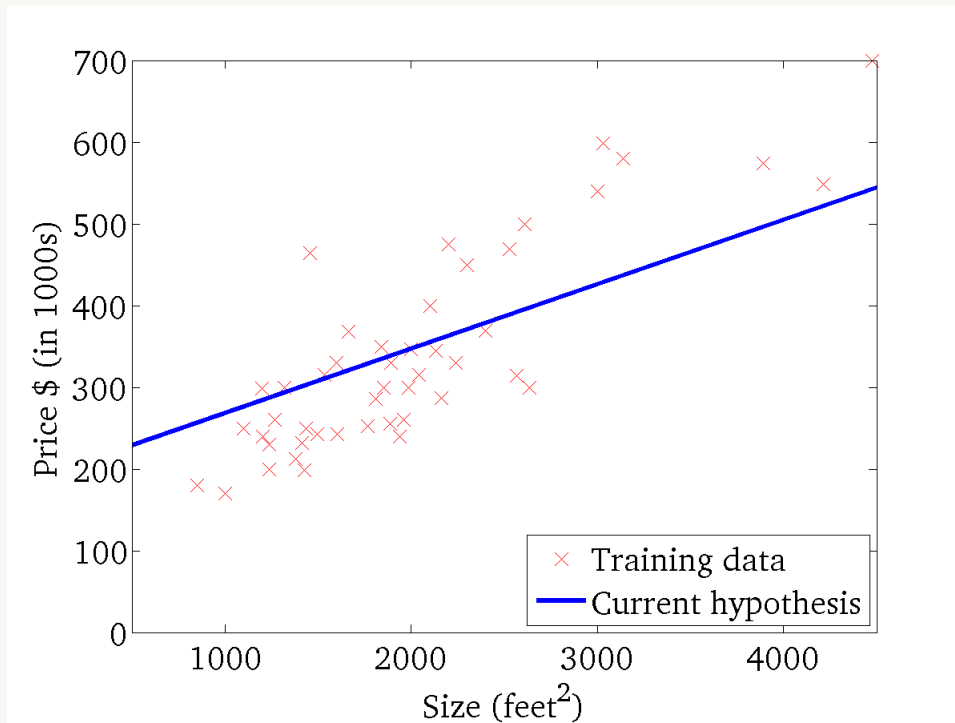- This gives a better hypothesis, but still not great - not in the center of the countour plot
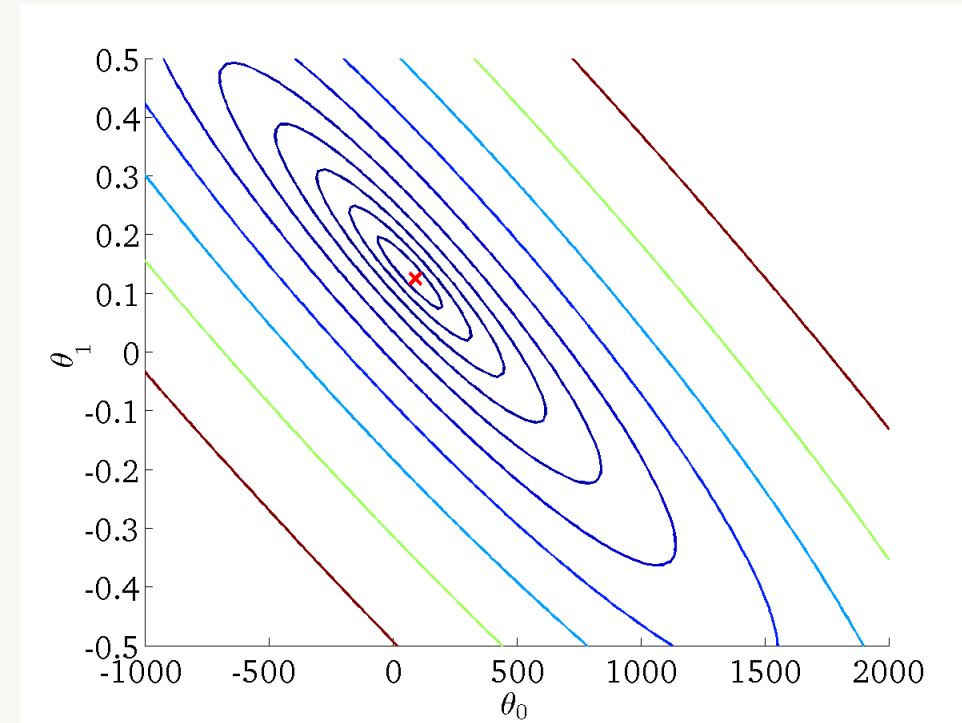
# Cost function visualization (Cont.)

The best fit…

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$ , this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$ )

# Next…

Gradient descent

Dr. Aeshah Alsughayyir

# Motivation

**To reach the lowest point**

*We should move the ball **step** by **step** in the gradient direction*

Dr. Aeshah Alsughayyir

# The idea of gradient descent.

- Imagine that this is a landscape of grassy park, and you want to go to the lowest point in the park as rapidly as possible …



Starting point

local minimum

Red: means high
blue: means low

Dr. Aeshah Alsughayyir

# The idea of gradient descent (Cont.)

- With different starting point

Red: means high
blue: means low

New Starting point

New local minimum

# What is Gradient Descent?

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

- The objective is $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

> *Gradient descent* is an **iterative optimization algorithm** for finding the minimum of a function.

There are **two basic steps** involved in this algorithm:

- Start with some random values of $\theta_0, \theta_1$.

- Keep changing the values of $\theta_0, \theta_1$ in order to reduce $J(\theta_0, \theta_1)$, until we hopefully end up at a minimum.

Dr. Aeshah Alsughayyir

# Gradient Descent (Cont.)

**Why do we need a Gradient Descent?**

In short **to minimize the cost function**, But How? **Let's see**

- The *cost function* only works when it knows the <u>parameters' values</u>
  - Note: in the previous example we manually choose the parameters' value each time
- During the algorithmic calculation, once the parameters' values are randomly initialized it's the *gradient descent* who must decide what params value to choose.
- In the next iteration, in order to minimize the error, it's the gradient descent who decide by how much to increase or decrease the params values.

**Rule:**

- How far you can jump each time is given by this algorithm

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

New position | Update to new position | Previous position | Minimize the function | Learning rate | Derivative of cost function

# Gradient Descent (Cont.)

$$\text{repeat until convergence}\{\theta_j := \theta_j - \alpha\frac{\partial}{\partial\theta_j}J(\theta_0, \theta_1)\ \forall j \in \{0, 1\}\}$$

- Where
  - := is the assignment operator

  - $\alpha$ is the **learning rate** which basically defines how big the steps are during the descent

  - $\frac{\partial}{\partial\theta_j}J(\theta_0, \theta_1)$ is the **partial derivative** term

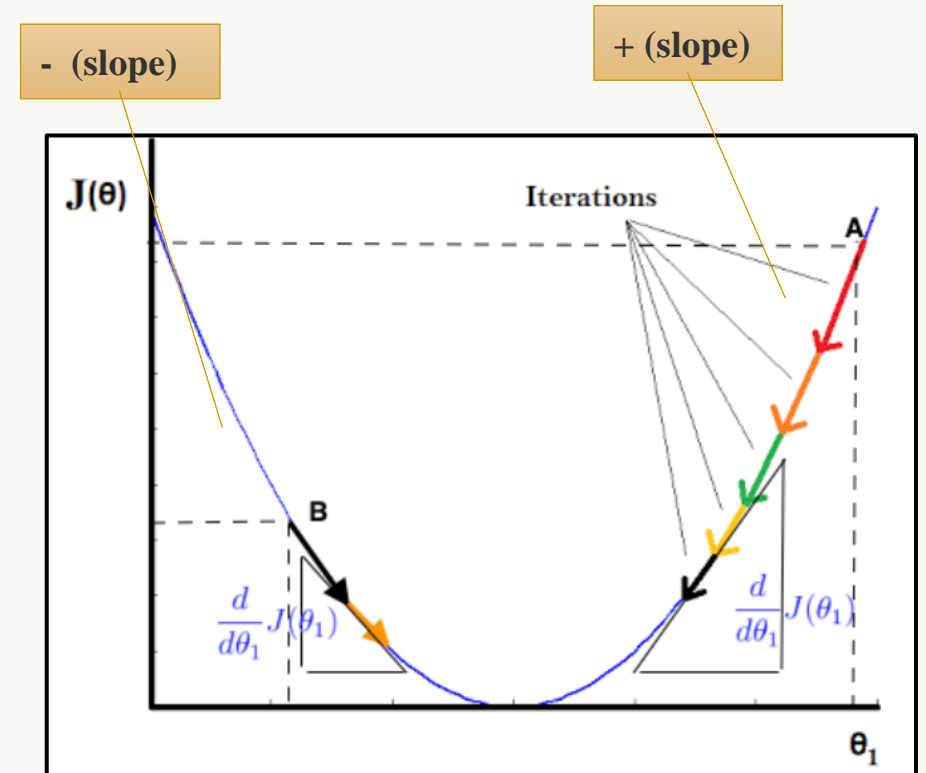  - j = 0, 1 represents the **feature index number**

Also the parameters should be **updated simultenously**, i.e. ,

$$temp_0 := \theta_0 - \alpha\frac{\partial}{\partial\theta_0}J(\theta_0, \theta_1)$$

$$temp_1 := \theta_1 - \alpha\frac{\partial}{\partial\theta_1}J(\theta_0, \theta_1)$$

$$\theta_0 := temp_0$$

$$\theta_1 := temp_1$$

# Understanding Gradient Descent

Consider a simpler cost function $J(\theta_1)$ and the objective: $\min_{\theta_1} J(\theta_1)$

Repeat{
$$\theta_1 = \theta_1 - \alpha \, \frac{d}{d\theta_1} J(\theta_1)$$
}

- If the slope is **positive** and since α is a positive real number, then overall term $-\alpha \, \frac{d}{d\theta_1} J(\theta_1)$ would be negative. This means the value of $\theta_1$ will be decreased in magnitude as shown by the arrows.

- If the initialization was at <u>point B</u>, then the **slope** of the line **would be negative** and then *update term would be positive* leading to increase in the value of $\theta_1$ as shown in the plot.



So, no matter where the $\theta_1$ is initialized, the algorithm ensures that parameter is updated in the **right direction** towards the minima, **given proper chosen value for α** .

Dr. Aeshah Alsughayyir

**How can we choose the learning Rate: α?**

# Understanding Gradient Descent(Cont.)

## Choosing learning rate:



| Too low | Just right | Too high |
|---|---|---|
| A small learning rate requires many updates before reaching the minimum point | The optimal learning rate swiftly reaches the minimum point | Too large of a learning rate causes drastic updates which lead to divergent behaviors |

•The **yellow** plot shows the **divergence** of the algorithm when the learning rate is really high wherein the learning steps overshoot.

•The **green** plot shows the case where learning rate is not as large as the previous case but is high enough that the steps keep **oscillating** at a point which is not the minima.

•The **red** plot would be the **optimum curve** for the cost drop as it drops steeply initially and then saturates very close to the optimum value.

•The **blue** plot is the least value of αα and **converges very slowly** as the steps taken by the algorithm during update steps are very small.



Learning Rate — Set the learning rate carefully

If there are more than three parameters, you cannot visualize this.

Very Large
small
Large
Just make

Loss
No. of parameters updates

Dr. Aeshah Alsughayyir

# Understanding Gradient Descent(Cont.)

**Practical tips:**

- Gradient Descent Stopping condition

    **Automatic Convergence Test**: Gradient descent can be considered to be converged if the drop in cost function is not more than a preset threshold say $10^{-3}$.

- How to choose the learning rate?

    In order to choose optimum value of ( $\alpha$ ), run the algorithm with different values like: *1, 0.3, 0.1, 0.03, 0.01, etc.,* and plot the learning curve to understand whether the value should be increased or decreased.

# Understanding Gradient Descent(Cont.)

- **Is it required to decrease the value of $\alpha$ by time?**

Gradient descent can converge to a local minimum, even with the learning rate $\alpha$ fixed.

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease $\alpha$ over time.

Slope1> slope2> slope3>  …..
By time the slope is getting smaller, consequently the value of the step is automatically is getting smaller by time, so no need to change the value of alpha.

# Applying Gradient Decent To The Linear Model

Dr. Aeshah Alsughayyir

# Applying Gradient Descent

**Model against dataset**
the housing price problem

|   | area | price |
|---|------|-------|
| 1 |      |       |
| 2 |      |       |
| ... |    |       |
| m |      |       |

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

|   | area | age | price |
|---|------|-----|-------|
| 1 |      |     |       |
| 2 |      |     |       |
| ... |    |     |       |
| m |      |     |       |

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$n$

|   | area | .... | #Floors | price |
|---|------|------|---------|-------|
| 1 |      |      |         |       |
| 2 |      |      |         |       |
| ... |    |      |         |       |
| m |      |      |         |       |

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

Dr. Aeshah Alsughayyir

# Applying Gradient Descent (Cont.)

## Linear Regression Model

1

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

2

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

## Gradient descent algorithm

3

repeat until convergence {

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

(for $j = 1$ and $j = 0$)

}

In order to apply gradient descent, the derivative term $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ needs to be calculated.

Dr. Aeshah Alsughayyir

# Applying Gradient Descent (Cont.)

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

$$= \frac{1}{2m} \left( \frac{\partial}{\partial \theta_j} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 \right)$$

$$= \frac{1}{2m} \left( \frac{\partial}{\partial \theta_j} \sum_{i=1}^{m} (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2 \right)$$

$$= \begin{cases} \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) & \text{for } j = 0 \\ \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} & \text{for } j = 1 \end{cases}$$

$$\text{repeat until convergence} \begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h(x^{(i)} - y^{(i)})) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h(x^{(i)} - y^{(i)})) \cdot x^{(i)} \end{cases}$$

# Applying Gradient Descent (Cont.)

**Gradient descent algorithm**

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

repeat until convergence {

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

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

}

update $\theta_0$ and $\theta_1$ simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

Dr. Aeshah Alsughayyir

# Applying Gradient Descent (Cont.)

## Risk of meeting different local optimum

The problem of existing more than one local optima, may make the gradient decent converge to the non global optima.

# Applying Gradient Descent (Cont.)

- Cost/error function is always convex:



o The linear regression cost function is always a **convex function** - always has a single minimum
- Bowl shaped
- One global optima
- So gradient descent will always converge to global optima

Dr. Aeshah Alsughayyir

# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

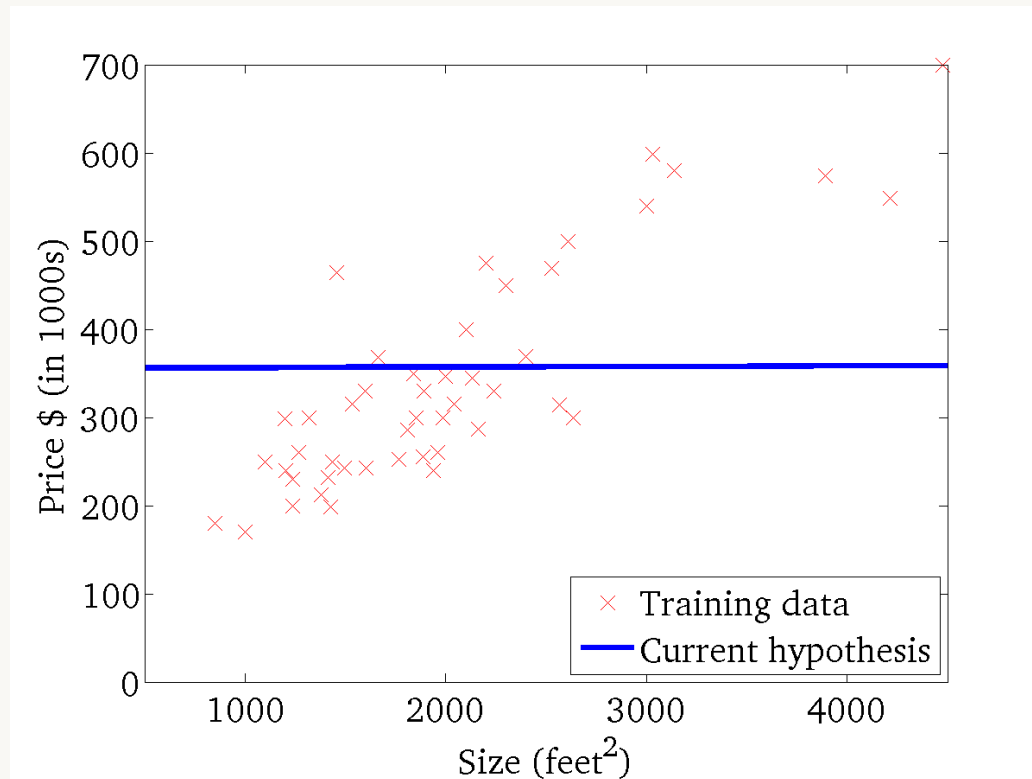(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Applying Gradient Descent (Cont.)

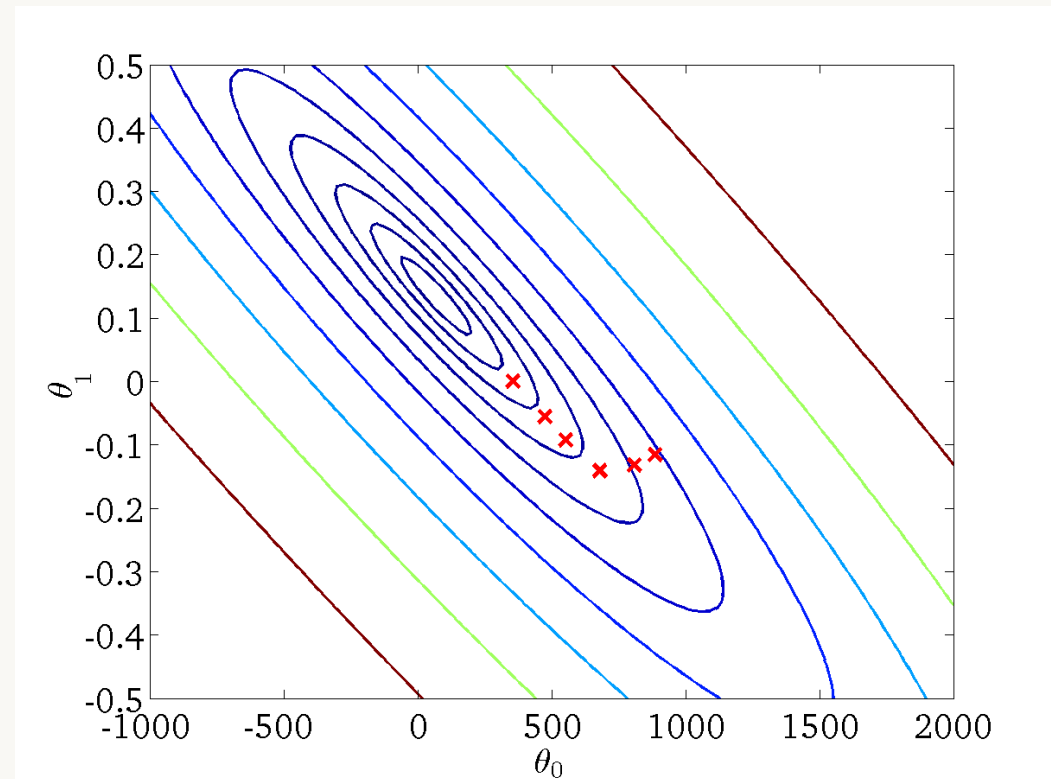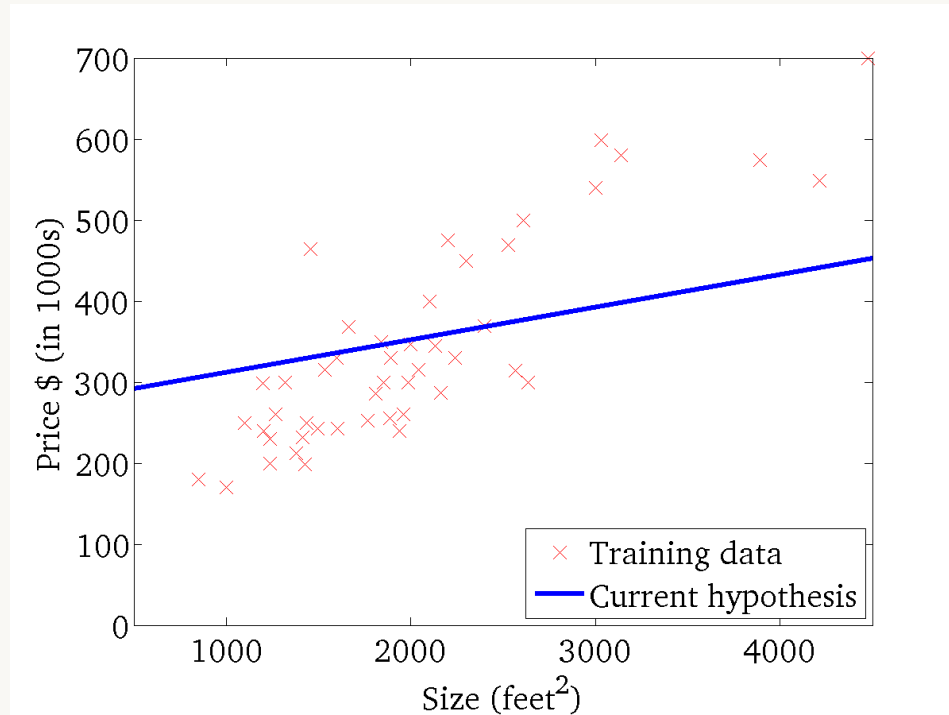$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Applying Gradient Descent (Cont.)

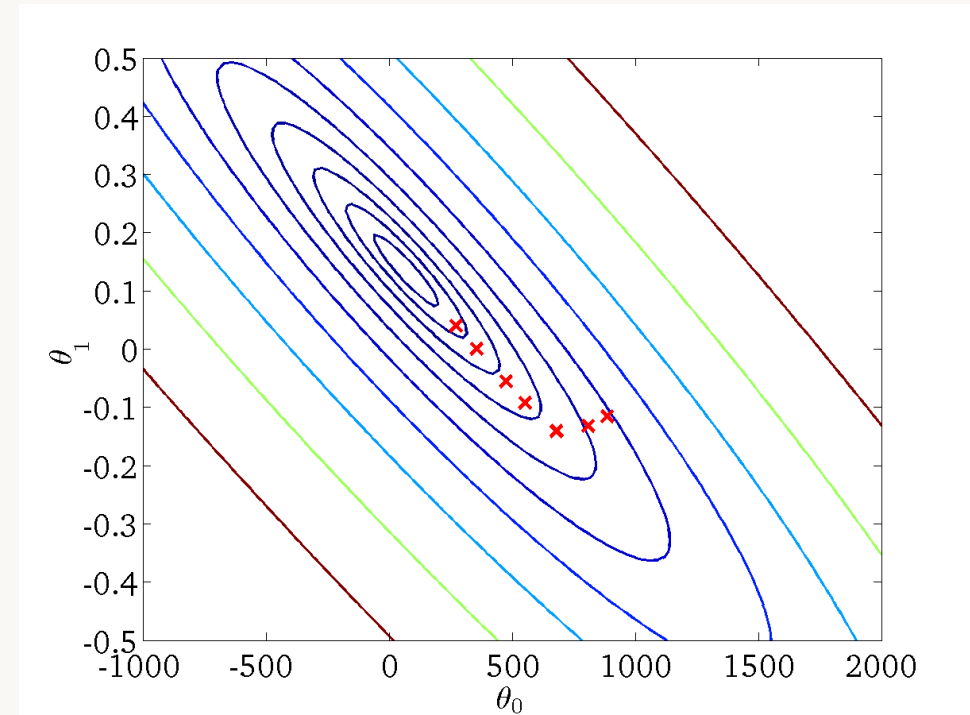$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

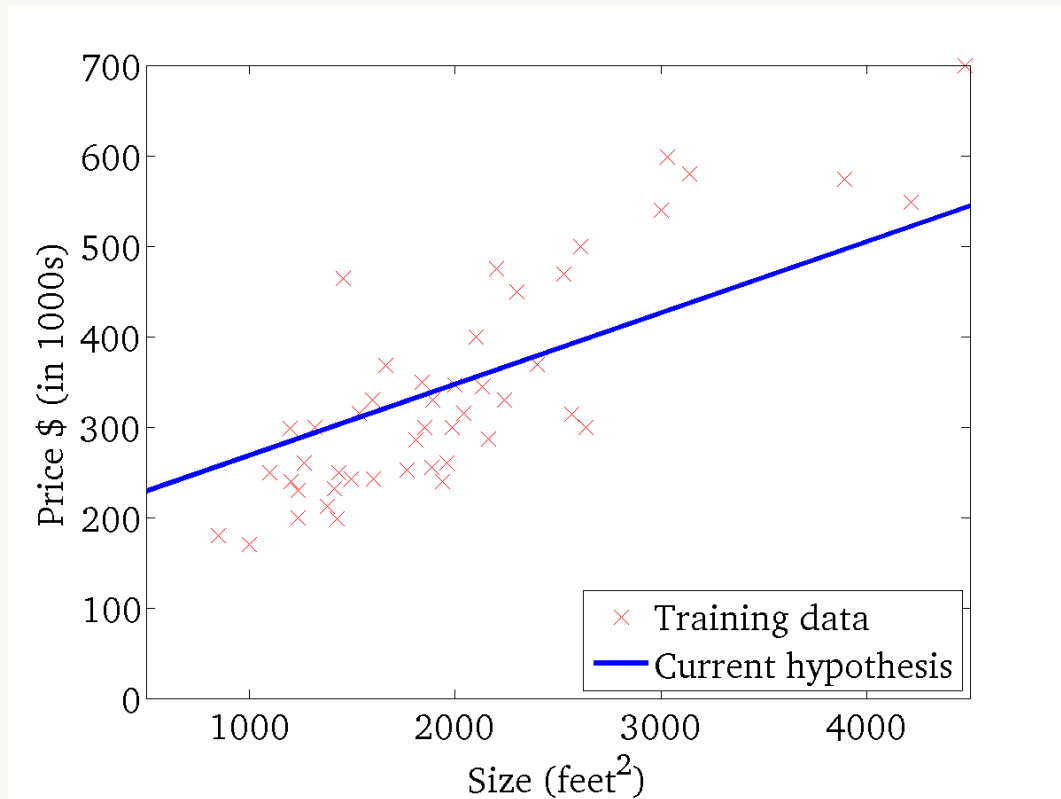(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)
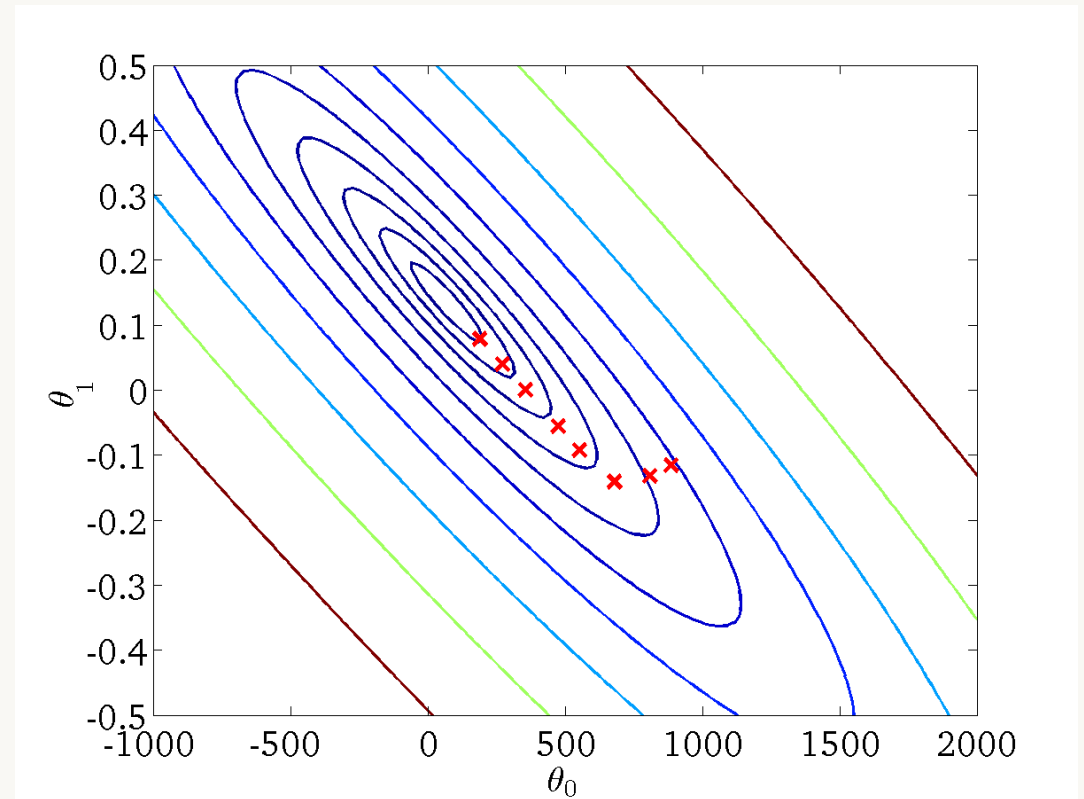
# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

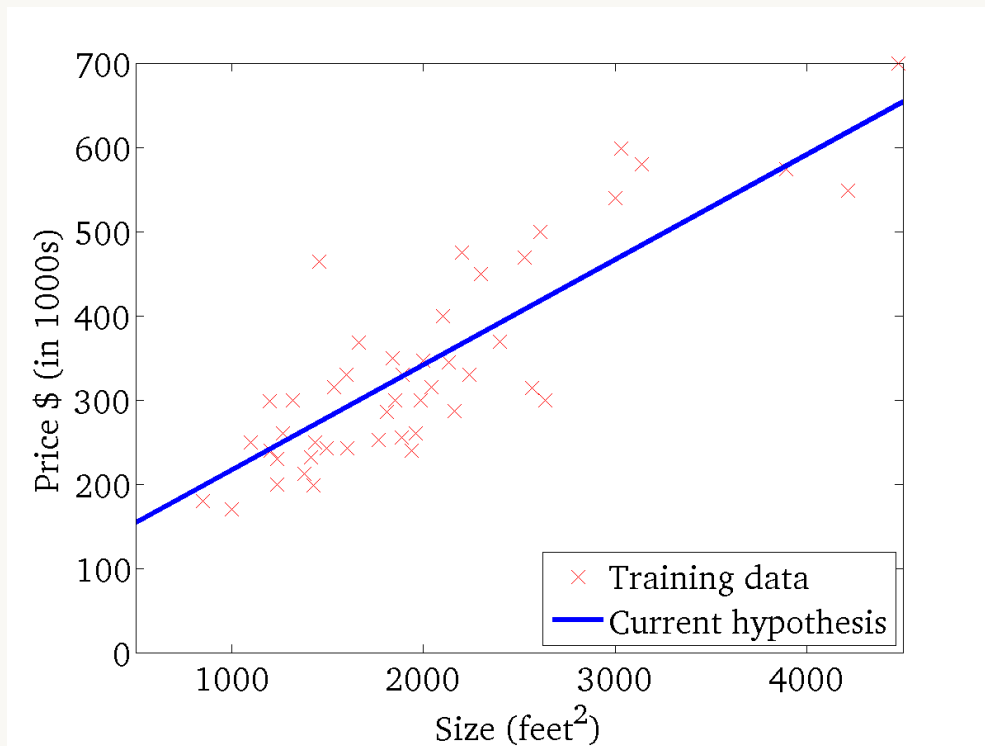(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)
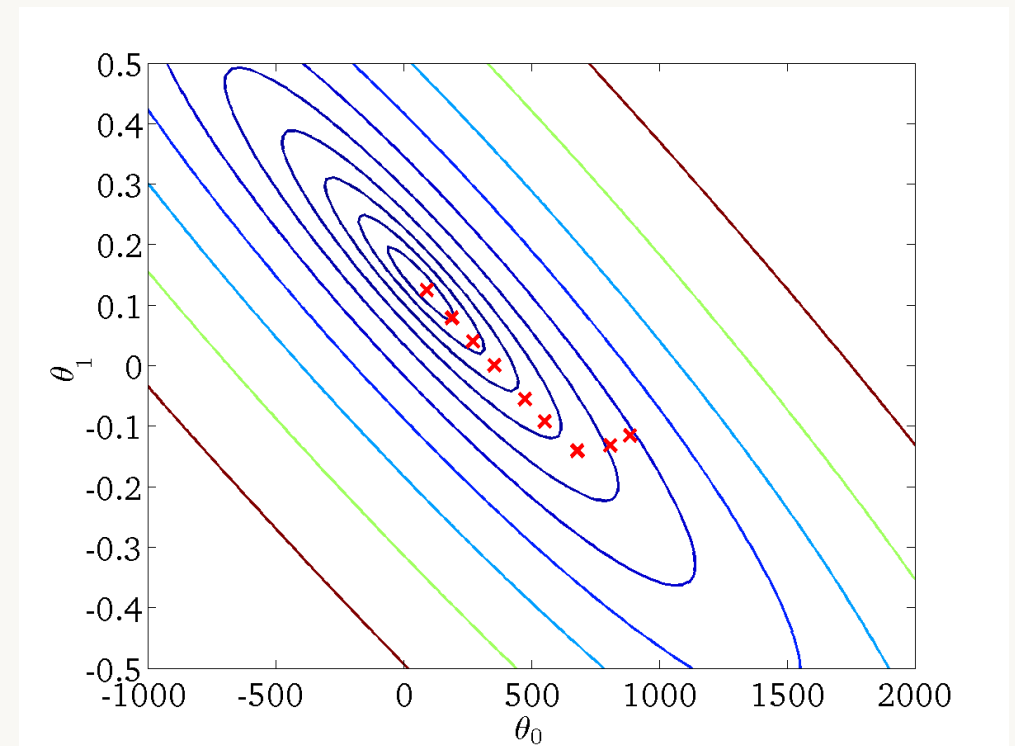
# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

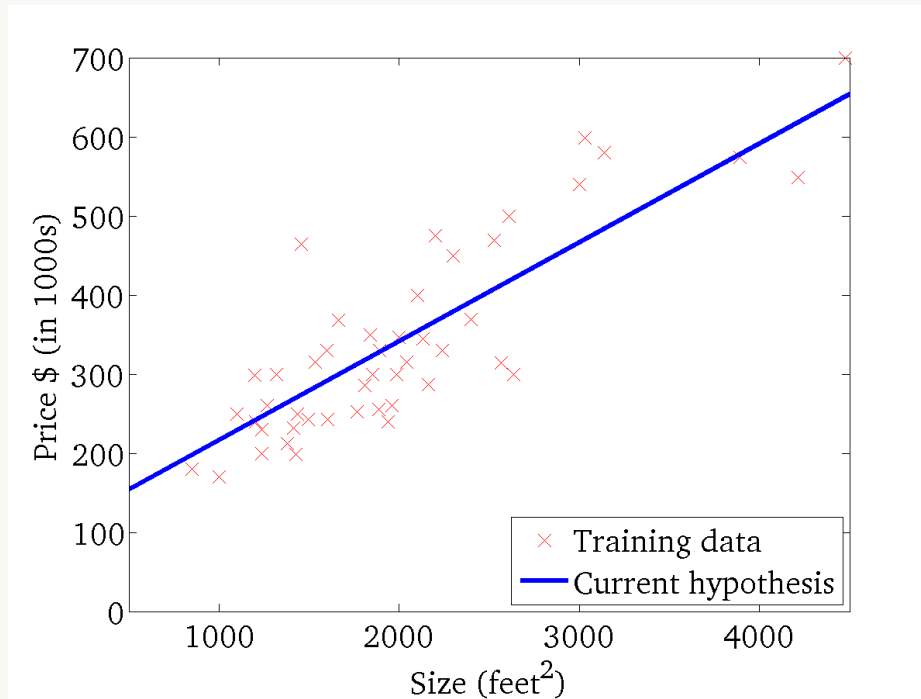(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)



Dr. Aeshah Alsughayyir

# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# Applying Gradient Descent (Cont.)

$$h_\theta(x)$$

(for fixed $\theta_0$, $\theta_1$, this is a function of x)



Now you can do prediction, given a house size!

Dr. Aeshah Alsughayyir

# Batch Gradient Descent

- The gradient descent technique that uses all the training step is called **Batch Gradient Descent**. This is basically the calculation of the derivative term over all the training examples as it can be seen it the equation above.

**Gradient descent algorithm**

repeat until convergence {

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

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

}

*Note*:
The equation of linear regression can also be solved using **Normal Equations** method, but it poses a disadvantage that it does not scale very well on larger data while gradient descent does.

# More about Gradient Descent

➢ There are plenty of optimizers

- Adagrad
- Adadelta
- <u>Adam</u>
- Conjugate Gradients
- BFGS
- Momentum
- Nesterov Momentum
- Newton's Method
- RMSProp
- SGD

# Summary

- With the cost function the learning algorithm determine whither the current parameters' values are the best or not (by calculating the error).

- Using the gradient descent algorithm is mainly to automate the process of updating the values of the parameters towards the global minimum.

- There are **two basic steps** involved in gradient descent algorithm:
  - ❑ Start with some random value of $\theta_0$, $\theta_1$.
  - ❑ Keep updating the value of $\theta_0$, $\theta_1$ to reduce $J(\theta_0, \theta_1)$ until minimum is reached.

## *Any Question?*

# Appendix

Dr. Aeshah Alsughayyir