

Where to use?

- Database mining
Large datasets from growth of automation/web.
E.g. web click data, medical records, biology, engineering.
- Applications can't program by hand
E.g. Autonomous helicopter, handwriting recognition, Natural Language Processing, Computer Vision.
- Self-optimising programs
E.g. Amazon, Netflix product recommendations
- Understanding human learning (brain, neural AI).

What is machine learning?

- (1959): Field of study that gives computers the ability to learn without being explicitly programmed.
- 1998: Well-posed learning problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E .

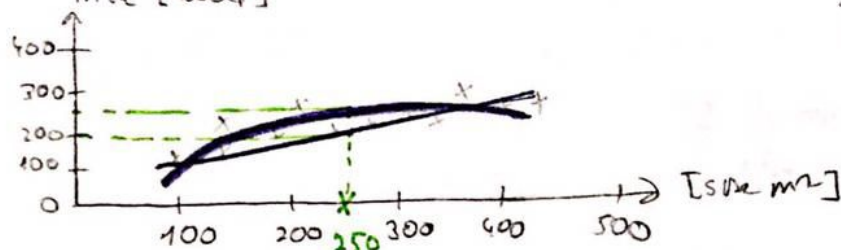
Types of ML algorithms

- Supervised learning
- Unsupervised learning

others: Reinforcement learning, recommender systems

Supervised Learning

Ex: Housing price prediction



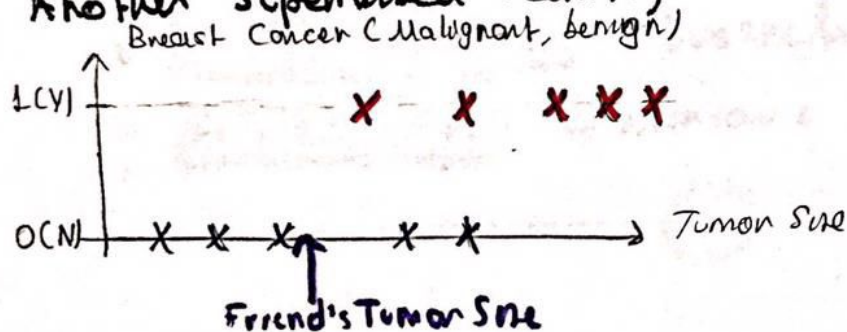
• Böyle bir dataset verdiğimiz, bir ise 250 m²lik bir evin fiyatını tahmin etmeye çalışırız.

• Bunun için datalardan bir straight line veya bir quadratic function öğreniriz.

★ Bu bir **supervised learning** örneğidir. Çünkü **"Right Answers"** is given to the algorithm.

★ This problem is also a **Regression Problem**: Predicting continuous valued output (price)

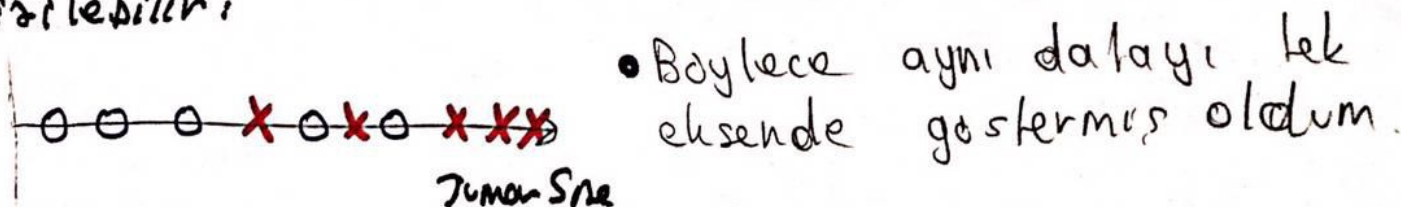
Another supervised learning example:



• Örneğin bir arkadaşın mor noktanın işaret ettiği boyuta bir tımarı var. Bu da bu tımarın malignant olma olasılığını öğrenmek istiyorsanız.

★ **Classification Problem**: Discrete valued output (0 on 1)
Burada 2 den fazla gruba da classify edebiliriz. Output 0, 1, 2, 3 ... şeklinde gidebilir.

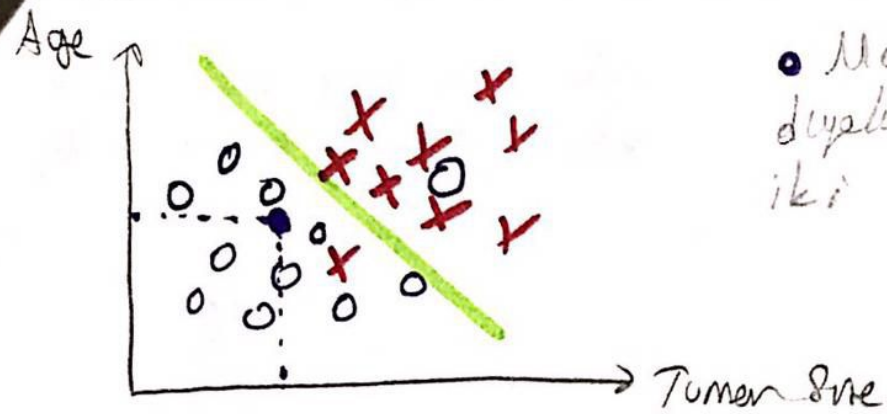
★ classification probleminde data farklı şekilde de çizilebilir:



• Böylece aynı datayı tek eksenle göstermiş olduk.

There is only one **Feature** or **Attribute** (Tumor Size)

Multifeatured Classification Problem



• Mon noktaya denk duzen bir tumör var dıyelim. Learning algoritminin yapacağı şey iki gruba bir arada ayırmak

* daha fazla feature de kullanılabilir.

Sonuç olarak:

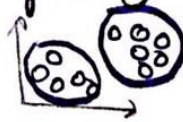
- Supervised ian doğru cevapla veriluyen (Ev fiyatı veya tumörün B/N olması)
- Classification: Discrete output
- Regression: Continuous output

Unsupervised Learning

o Supervised için data plot şuna benziyor:



o Unsupervised için ise data set için label olmaksızın her datanın neye dahil geldiği söylenmes.



o İstenilen şey şu: Here is the data-set can you find some structure in the data. Learning algorithm bu durumda iki cluster oluşturabilir. Buna **clustering** denir.

Cocktail Party Algorithm

o İki kişi 1 den 10'a sayıyor. İki müzikten olsun birini ing. sayan diğer italyanca sayan olsun. Sonuçta ikisinde de iki farklı sayım duyulur birinin sesi daha yüksek.

* Take this recordings give to an unsupervised learning algorithm and tell the algorithm find structures.

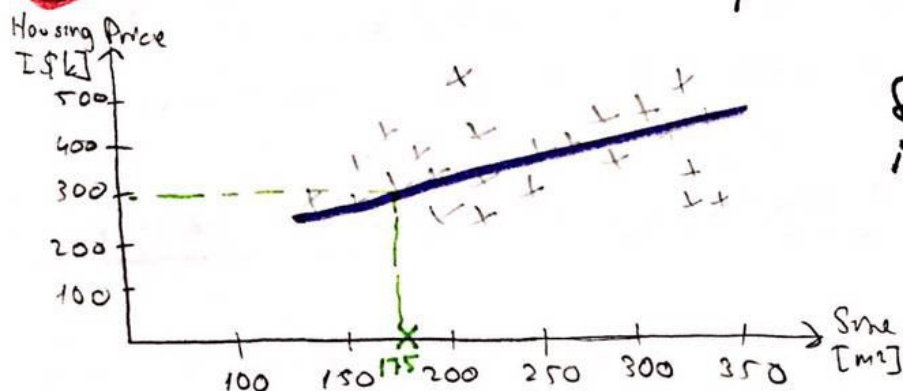
* What algorithm will do is, listen them and say: sound like two recordings are added together, moreover cocktail party algorithm will separate them.

* Bu algorithm 1 satırlık bir algoritmadır!

* Octave ve Matlab ile machine learning alg. yapmaları daha kolay olur.

Model Representation

Let's look at a housing price prediction application.



175 m²'lik evimin fiyatını tahmin etmek istiyoruz.

One thing we could do is fit a **MODEL**. Maybe fit a straight line to this data.

Remember, this is an example of supervised learning algorithm. Moreover this is an example of regression problem.

For supervised learning we have a training set and our job is to learn from this data to predict prices.

Here are some notations

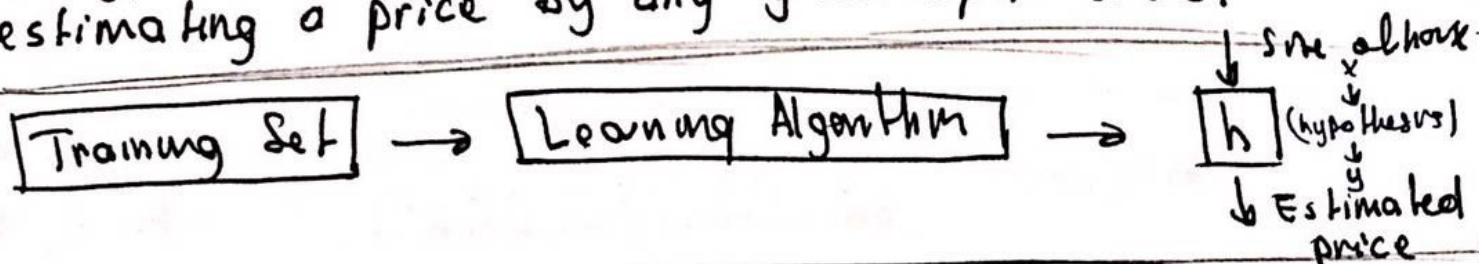
Size in m² (x)	Price in \$K (y)
205	307
257	335
487	405
323	370
...	...

} $m = 47$

- m : # of training examples
- x 's: "input" variable / features
- y 's: "output" variable / "target" variable
- (x, y) : one training example
- $(x^{(i)}, y^{(i)})$: i th training example

$$\begin{aligned} x^{(1)} &= 205 \\ y^{(1)} &= 307 \\ &\vdots \end{aligned}$$

- We feed our training to our learning algorithm. The job of the learning algorithm is to output a function h (hypothesis function). The job the hypothesis is estimating a price by any given input size.



- Hypothesis h is a function maps x 's to y 's

- How do we represent h ?

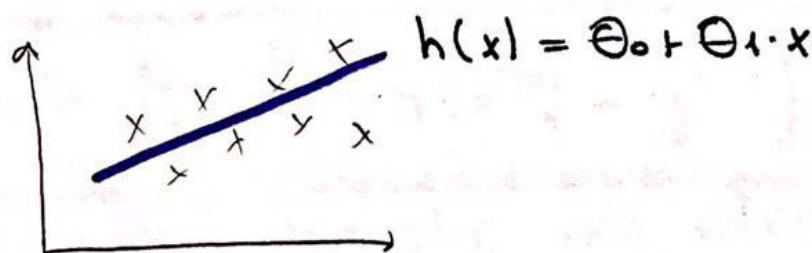
- Our initial choice will be

$$h_0(x) = \theta_0 + \theta_1 \cdot x$$

shorthand $h(x)$ instead of $h_0(x)$

- Note that this h is a straight line function.

- Zamanla daha complex nonlinear hypotheses-ler de kullanacağız. Başlangıç için bunu kullacağız!



- Linear Regression with one variable (x)

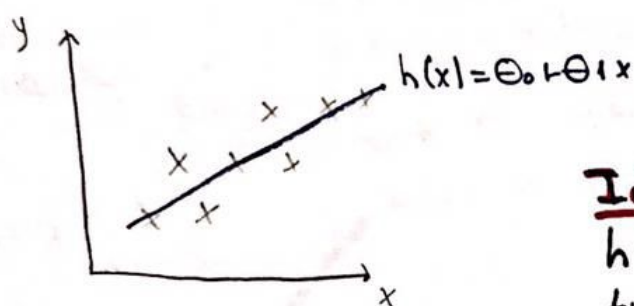
Cost Function

⊛ Cost function will let us figure out how to fit the best possible straight line to our data.

⊛ m : # of training examples olduğumuzu unutma.

Hypothesis: $h(x) = \theta_0 + \theta_1 x$ demektir.

How to choose **PARAMETERS** ($\theta_0, \theta_1, \dots$)?



• How do we come up with θ_0 and θ_1 ?

Idea: Choose θ_0, θ_1 so that $h(x)$ is close to y for our training examples.

⊛ I want to minimize $\sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$ over θ_0, θ_1 .

Result: If I can minimize the **COST FUNCTION**

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

over θ_0 and θ_1 , sonuçta elde edilen hipotezim başarılı bir hipotez olacaktır.

⊛ Burada cost function Squared Error Function olarak seçildi. Linear Regression problemlerinde bu popülerdir ve işe yarar.

⊛ kısaca cost function tüm training examples için hata karesinin toplamı ile olur.

- Cost Function Intuition I -

In order to visualize cost function we will use a simplified hypothesis:

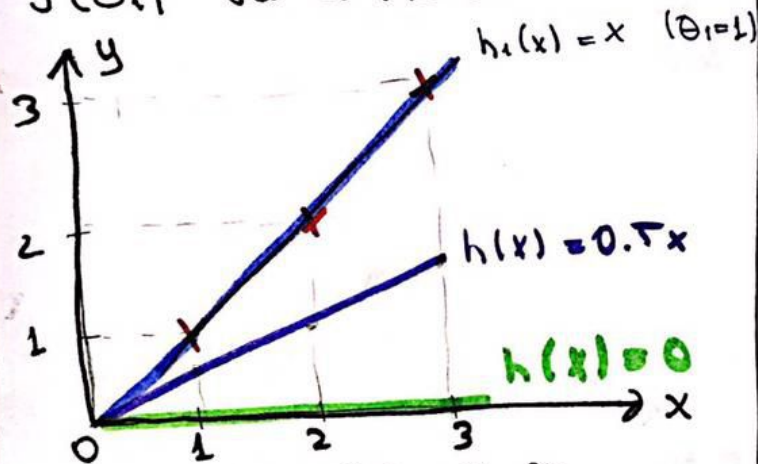
$$h(x) = \theta_1 x$$

assume $\theta_0 = 0$

Hypothesis $h(x)$

For fixed θ_1 , it is a function of x .

Now let's say we have a training set with $m=3$ as below. We will find the value of $J(\theta_1)$ for different θ_1 values.



For $\theta_1 = 1$, $h(x) = x$ so

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

$$= \frac{1}{6} \cdot [0^2 + 0^2 + 0^2] = 0 //$$

For $\theta_1 = 0.5$

$$J(\theta_1) = \frac{1}{6} \cdot [0.5^2 + 1^2 + 1.5^2]$$

$$\approx 0.58 //$$

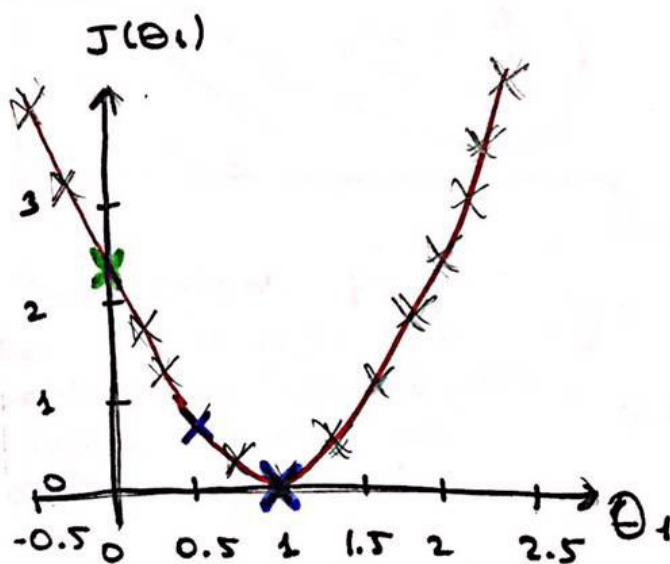
For $\theta_1 = 0$

$$J(\theta_1) = \frac{1}{6} \cdot [1^2 + 2^2 + 3^2]$$

$$\approx 2.3 //$$

Cost Function $J(\theta_1)$

It is a function of parameter θ_1 . (For fixed training set)



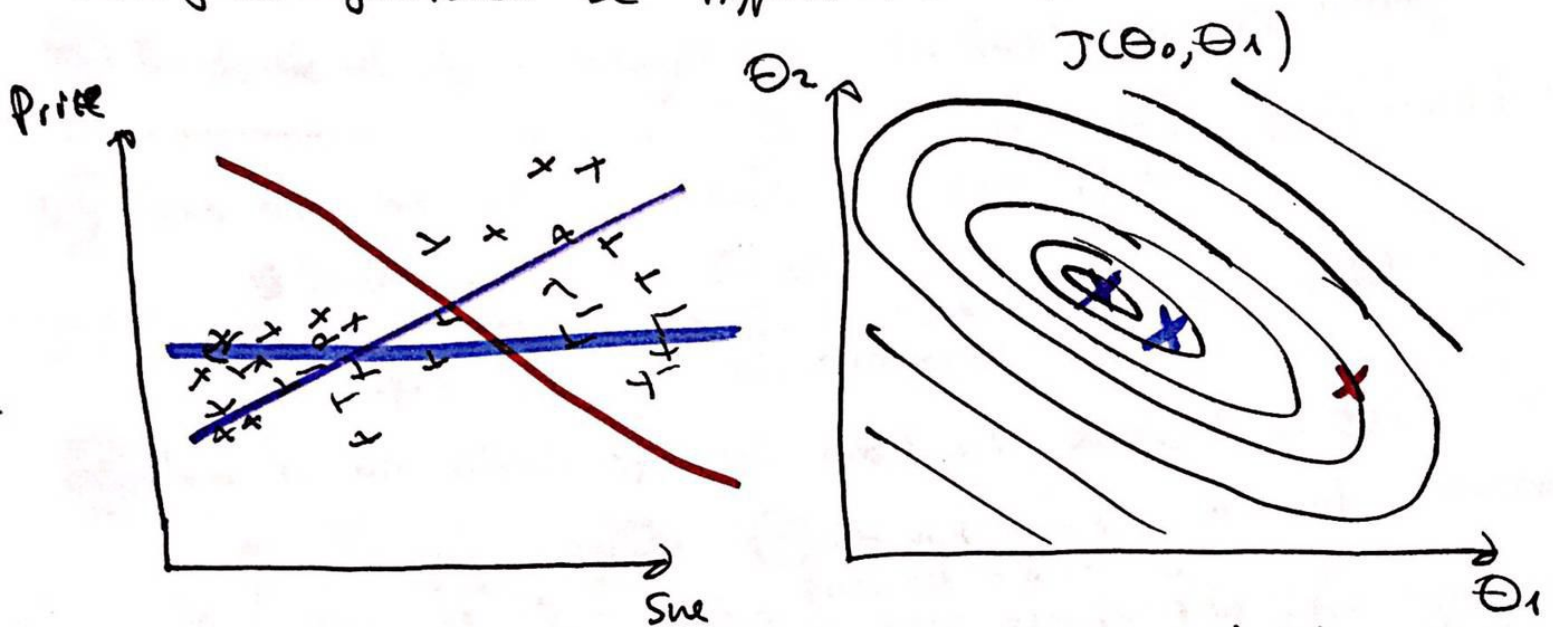
Üç noktada hesaplandı, daha fazla hesaplanırsa da minimum değenden uzaklaştıkça büyüyen bir fonksiyon şekli ortaya çıkar.

$J(\theta)$ üzerindeki her noktanın bir hipotezin hata değerini temsil ettiğini unutmamasak ise en az hataya sahip cost function'ı bulmak.

Cost Function her zaman can seklinde mi? Kısıtları neler? Linear regression can mı bityle? Nonlinear olsa nasıl olur? Classification can mı bityle?

- Cost Function Intuition II -

⊗ Bu videoda $J(\theta_0, \theta_1)$ tan plot'un nasıl olduğunu gösterdi ve hypothesis bağlantılarını gösterdi



Bu contour plot her circle aynı yüksekliğe denk gelir. Yani bu daire boyunca cost function $J(\theta_0, \theta_1)$ değeri aynıdır!

⊗ Sıradaki adım şu olmalı, bir efficient algorithm ile cost function'ı minimize etmek! Böylece en iyi hipotezi buluruz!

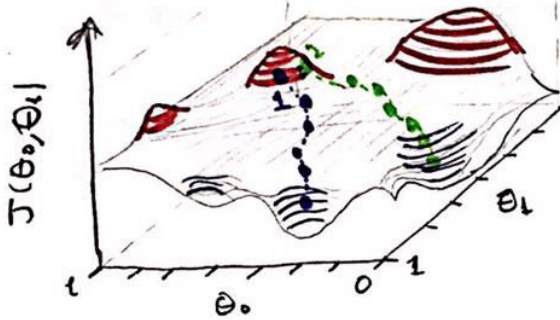
Gradient Descent

- "Gradient Descent" is an algorithm for minimizing the cost function.
- We use GD algorithm all over the machine learning not just for linear regression.
- Bu denkle GD alg. ile rastgele bir J function'ı minimize edeceğiz.

Assume we want to minimize $J(\theta_0, \theta_1)$

- 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 minimum.

Here is an intuition; what gradient descent does:



θ_0 ve θ_1 için 1 numaralı noktadan başladık diyelim. GD ile yapılan şey şu kendi etrafımızda 360° döneceğiz ve aşağı doğru inmek için hangi yöne doğru bir baby step atmamız gerekliliğini buluyoruz.

İlk adımları attıktan sonra aynı süreci tekrarlarız. Böylece monoton yolu izleyerek bir local minimum'a ulaşırız. Burada önemli bir nokta şudur:

Eğer initial θ_0 ve θ_1 farklı seçilseydi bambaşka bir local minimum'a ulaşabilirdik. Bu GDA'nın bir özelliği bunu daha sonra tartışacağız.

That's the intuition in pictures

let's look at the math!

ML W L =

②

Gradient Descent Algorithm

repeat until convergence {

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

assignment learning rate

We have to make a simultaneous update:

$$\begin{aligned} \text{temp0} &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp1} &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 &= \text{temp0} \\ \theta_1 &= \text{temp1} \end{aligned}$$

$$\begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \theta_1 &:= \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \end{aligned}$$

X WRONG implementation!

Scanned with CamScanner
Simultaneous olmayan da muhtemelen çalışır ama gradient descent ilkiden ilkinisinin kendine has başka özellikleri vardır.

Gradient Descent Intuition

GD Algorithm was:

repeat until convergence

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

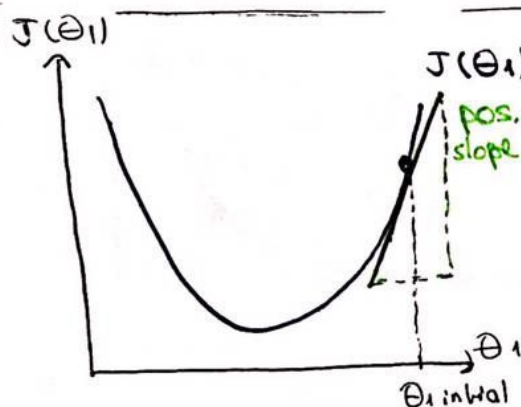
learning rate

Derivative term

(Simultaneously by update $J=0$ and $J=1$)

To understand what GDA is doing let's assume we have a hypothesis as $h(x) = \theta_1 \cdot x$ meaning $\theta_0 = 0$. Thus the cost function will be: $J(\theta_1)$ where $\theta_1 \in \mathbb{R}$.

We know $J(\theta)$ was looking like below graph. What we want to do is understand how GDA works.



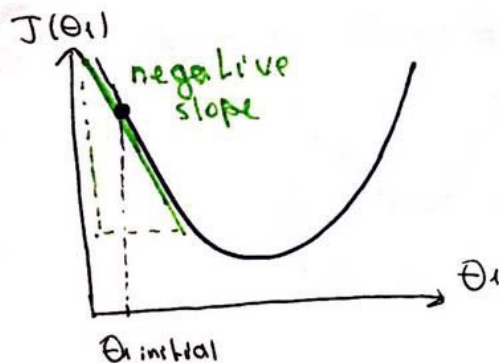
$J(\theta_1)$

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

→ ilgili θ_1 noktasında $J(\theta_1)$ 'in eğimi!

Bu durumda derivative term ≥ 0 dir!
Bu sebeple θ_1 'in yeni değeri atılır
yani minimum'a doğru slope ve α
ile bağlantılı bir adım atan!

α is always a positive number!



$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1) \rightarrow \leq 0!$$

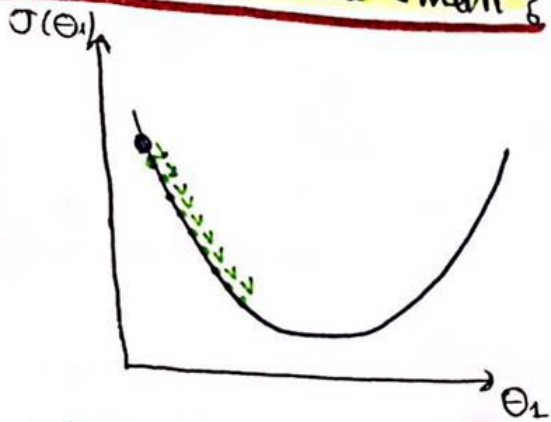
Bu durumda yeni θ_1 değeri artacaktır
böylece local minimum'a bir adım atılmış
olur.

Sonuç olarak GDA ile her durumda GDA ile θ parametreleri kendini yokuş aşağı bırakmaya yönelimlidir! Bu şekilde minimum'a ulaşılır!

Now let's try to understand what learning rate α is doing...

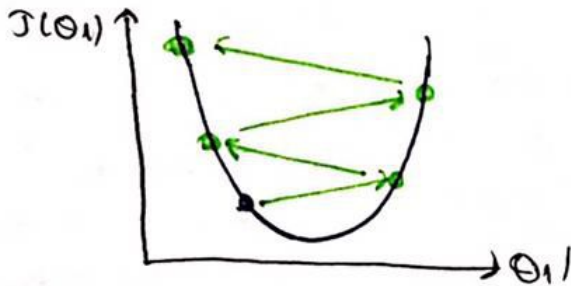
Learning Rate, α Intuition

1) If α is too small? \Rightarrow Gradient Descent can be slow.

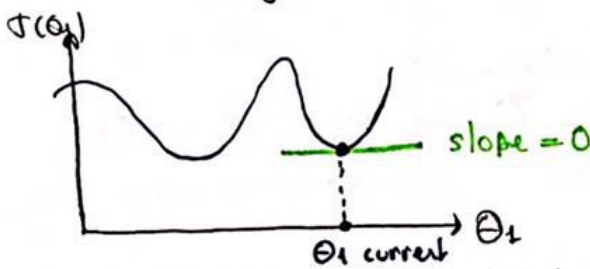


• α çok küçük olursa her update (θ) ile atılan adım boyutu çok küçük olacağından bu yüzden minimum'a ulaşmak gerektiren fazla zaman alabilir!

2) If α is too large? \Rightarrow Gradient descent can overshoot the minimum. It may fail to converge or even diverge



• If θ_1 is already at a local minimum:



$$\begin{aligned} \theta_1 &:= \theta_1 - \alpha \cdot \frac{d}{d\theta_1} J(\theta_1) \\ &= \theta_1 - 0 = \theta_1 \end{aligned}$$

θ_1 stays unchanged !!!

• As we approached a local minimum, gradient descent will automatically take smaller steps (even if the α is fixed). So no need to decrease α over time



• Local minimuma yaklaştıkça eğim azalır bu yüzden update ile atılan adım α sabit olsa dahi küçülür!

Gradient Descent for Linear Regression with one Variable

Gradient Descent Algorithm

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

Linear Regression Model & Cost Function

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

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

Apply gradient descent algorithm to this cost function!

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

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

Sonuçta burada J θ_0 ve θ_1 'in bir fonksiyonu her bir farklı θ_0 ve θ_1 parametresi için ayrı bir J değeri elde ediliyor! Şu şekilde bir J fonksiyonu var orada ve partial derivative alınabilir.

As a result GDA becomes:

Simultaneous update / olmalı

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

Cost function for Linear Regression is always going to be a "Convex Function". Meaning it doesn't have a local optima except for global optima! (Bowl Shape Cost Function)



Cost function şekli hipotez ve cost function formuna bağlı Linear Regression için sürekli lineardır ve Error Devire!

④ Daha önce de bahsedildiği gibi her update ile θ_0 ve θ_1 değerleri sürekli olarak hipotetik kendini minimum hataya doğru adım adım yaklaşımlar olacaktır.

"Batch" Gradient Descent

④ Batch gradient descent algorithm ile her update iterasyonu için tüm training examples kullanılır.

④ Böyle olmayan GDA'lar da vardır. Daha sonra bunlardan da bahsedeceğiz.

Normal Equations Method

Bu method ile linear algebra kullanılarak $J(\theta_0, \theta_1)$ 'in minimum yapan θ_0 θ_1 değerleri iknatif eşitime gerek kalmadan bulunabilir. Bunu göreceğiz.

Longer data sets için GDA NEM'e göre daha iyi çalışıyor!

Matrix: Rectangular array of numbers.
 A_{ij} : "i, j entry" in the ith row, jth column
 Vector: An $n \times 1$ matrix.

Matrix Addition: Aynı boyuttaki matrisleri eleman elemana toplarız.

$$\begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix}$$

Scalar Multiplication:

$$3 \times \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \times 3$$

Matrix Vector Multiplication: $(A_{m \times n}) \times (x_{n \times 1}) = y_{m \times 1}$

$$\begin{bmatrix} 1 & 3 \\ 4 & 0 \\ 2 & 1 \end{bmatrix}_{3 \times 2} \times \begin{bmatrix} 1 \\ 5 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 16 \\ 4 \\ 7 \end{bmatrix}_{3 \times 1}$$

$$\begin{aligned} 1 \times 1 + 3 \times 5 &= 16 \\ 4 \times 1 + 0 \times 5 &= 4 \\ 2 \times 1 + 1 \times 5 &= 7 \end{aligned}$$

Matrix Matrix Multiplication: $(A_{m \times n}) \times (B_{n \times o}) = C_{m \times o}$

Matrix Multiplication Properties:

- ⊗ $A \times B \neq B \times A$ (not commutative)
- ⊗ $A \times (B \times C) = (A \times B) \times C$ (Associative)
- ⊗ $A \cdot I = I \cdot A = A$

Matrix Inverse: If A is an $n \times n$ matrix, and if it has an inverse:

$$A \cdot (A^{-1}) = (A^{-1}) \cdot A = I$$

Matrix Transpose:
 If $B = A^T$ and $A_{m \times n}$
 then $B_{n \times m}$ and

$$B_{ij} = A_{ji}$$