# - Classification Problem -

- ▷ Variable y that we want to predict is discrete valued.
- ⌐ Bu problemin çözümü için **logistic regression** algoritması kulanıbcak.

Exs: • Email: Spam / Not Spam   • Online Transactions: Fraudelant (Yes / No)?
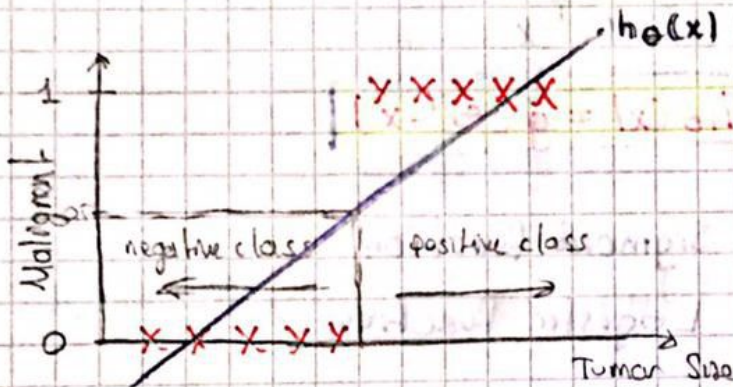• Tumor: Malignant / Benign?

$$y \in \{0,1\}$$

→ absence of something

0: "Negative Class" (e.g. benign tumor)
1: "Positive Class" (e.g. malignant tumor)

→ Şimdilik y sadece 0 ve 1 olabilir, bu durumla ilgileniriz. İleride **Multiclass Classification Problems** göreceğiz, $y \in \{0,1,2,3\}$ olabilecek.
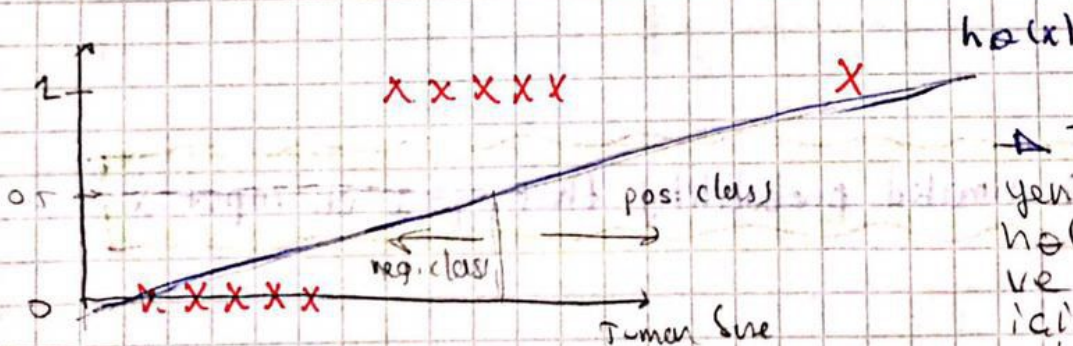
## Binaryclass Classification ( $y \in \{0,1\}$ )



→ Classification problem olmasına rağmen burada linear regression modelini kullanabiliriz.

→ Bir treshold belirleriz ve hypothesis $h_\theta(x) \geq 0.5$ ise $h_\theta(x) = 1$ deriz, küçükse $h_\theta(x) = 0$ deriz.

$$h_\theta(x) = \Theta^T \cdot x$$

Treshold classifier output $h_\theta(x)$ at 0.5:
   If $h_\theta(x) \geq 0.5$, predict "$y=1$"
   If $h_\theta(x) < 0.5$, predict "$y=0$"

⊗ Bu örnek için, linear regression modeli işliyor ama biraz modifiye edersek acaba hala kullanışlı olacak mı?



→ Training set'e yeni eklenen data $h_\theta(x)$'i saptırdı ve aynı treshold için, hypothesis artık doğru sonuç vermiyor.

Tüm linear regression modelini classification için kullanmak yanlışlık örnekte sonsuna çıkıyor.

Classification $\quad y = 0 \quad$ or $\quad 1$

$h_\theta(x)$ can be $> 1$ or $< 0$ ..... far) Linear Regression. model

△ Tahmin edilecek output 0 ve 1 olsa dahi hypothesis
sonuçları 0'dan küçük ve 1 den büyük elabiliyor bu anlamsız.

△ ileriki bölümlerde `Logistic Regression` modeline girecez.

Logistic Regression için $\quad 0 \leq h_\theta(x) \leq 1$

↳ İsmin içinde regression olması kafa karıştırmasın
bu bir classification algorithm.

~~~~~~~~~~~~~~~~~~~~~~~~~~

**Logistic Regression** $\quad \to 0 \leq h_\theta(x) \leq 1 \quad$ istiyoruz.

△ Linear regression için $\quad h_\theta(x) = \theta^T \cdot x \quad$ idi.

△ Logistic Regression için $\quad \boxed{h_\theta(x) = g(\theta^T \cdot x)}$

where
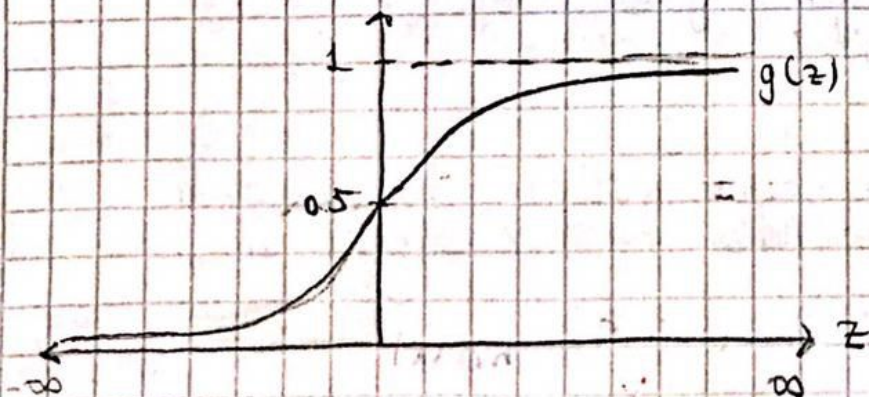
$$g(z) = \frac{1}{1 + e^{-z}}$$

$\quad \to$ Sigmoid function
$\quad \to$ Logistic function $\quad \}$ Synonyms

$z$ is real number

Ⓧ Thus $\quad h_\theta(x) = g(\theta^T \cdot x) = \dfrac{1}{1 + e^{-\theta^T \cdot x}}$



$g(z)$

△ $\theta$ parametrelerini
fit etmeyi daha sonra
göreceğiz.

$h_\theta(x) = $ Estimated probability that $y = 1$ on input $x$

**Ex:** If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{Tumor Size} \end{bmatrix}$ and $h_\theta(x) = 0.7$

→ Tell patient that 70% change of tumor being malignant.

→ Verilen tumor size için $y=1$ olma olasılığı 70%

→ $\boxed{h_\theta(x) = p(y=1 \mid x; \theta)}$ " probability that $y=1$ given $x$, parametrized by $\theta$ "

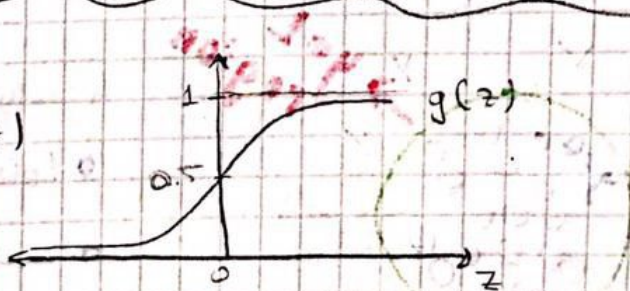→ $p(y=0 \mid x; \theta) = 1 - p(y=1 \mid x; \theta) = 1 - 0.7 = 0.3$

---

**Ekstra açıklama:** Regression problem için, linear regression modelini kullandık ve bir $h_\theta(x)$ atadık. Daha sonra training set ile hypothesis arası farkları minimize ettik bunu da cost function $J$'i, Gradient Descent algorithm ile veya analitik olarak minimize ederek yaptık. Sonuç olarak elde edilen $h_\theta(x)$'e girilen rasgele bir $x$ bize en gerçekçi output $y$'yi verdi.

Classification dan öte, Logistic Regression modeli kullanarak hypothesis oluşturuldu $h_\theta(x)$ ve bu cost function minimize edilerek en iyi $h_\theta(x)$'e ulaşıldı. Fakat burada $h_\theta(x)$ in sonucu $y=1$ olma olasılığını veriyor.

---

**Logistic Regression**

○ $h_\theta(x) = g(\theta^T x) = p(y=1 \mid x; \theta)$

○ $g(z) = \dfrac{1}{1+e^{-z}}$



Suppose predict "$y=1$" if $h_\theta(x) \geq 0.5$

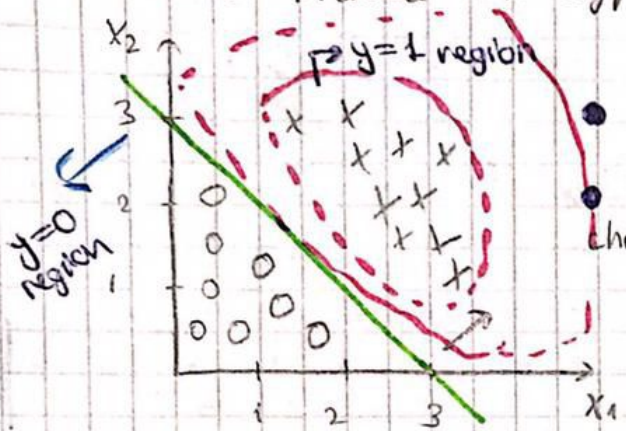predict "$y=0$" if $h_\theta(x) < 0.5$

$g(z) \geq 0.5$ when $z \geq 0$

since

$h_\theta(x) = g(\theta^T x) \geq 0.5$ when $\theta^T x \geq 0$

Olay şu $\theta^T x \geq 0$ için $g(\theta^T x) = h_\theta(x) \geq 0.5$ oluyor. Yani $y=1$ olma olasılığı $\geq 0.5$ oluyor o zaman $y=1$ prediction'ı yapabiliriz.

$\theta^T x = 0$ bir çizgi verir bu çizgi bize $h_\theta(x) = 0.5$'i veren $x$ ve $\theta$ değerlerini gösterir.

Yani tahmini predict $y=1$ if $\theta^T x \geq 0$

$y=0$ if $\theta^T x < 0$

Ex: Let's suppose we have a training set as below and we have obtained an hypothesis as below:



$\rightarrow$ y=1 region

$\leftarrow$ y=0 region

$\rightarrow$ Decision Boundary

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

- Let's say by using an algorithm we chose: $\theta_0 = -3 \quad \theta_1 = 1 \quad \theta_2 = 1$

- $\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$

$\rightarrow$ Verilen hypothesis'i kullanarak, bu $h_\theta(x)$'in nerelerde y=1 nerelerde y=0 prediction'ı yaptığını bulalım!
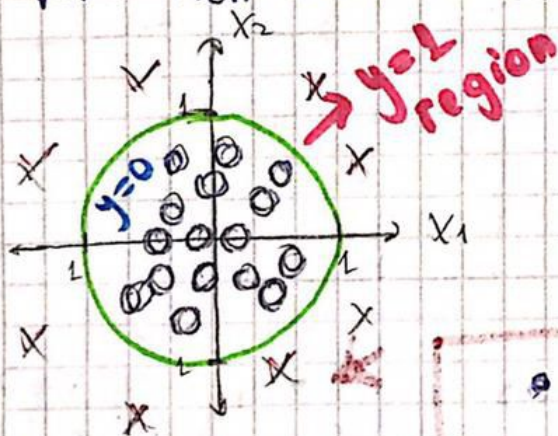
✗ Yada nerelerde $h_\theta(x) = 0.5$ olduğunu bulalım!

We know: • Predict "y=1" if $\overbrace{\boxed{-3 + x_1 + x_2 \geq 0}}^{\theta^T x}$

$\rightarrow x_1 + x_2 \geq 3$ için $h_\theta(x) \geq 0.5$, predict y=1.

✗ Decision Boundary yalnızca $h_\theta(x)$'e bağlıdır, training set'e bağlı değil.

---

Ex: Non-Linear Decision Boundaries



X ➔ y=1 region

y=0

Assume our hypothesis:

• $h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$

• Let say $\theta = [-1 \ 0 \ 0 \ 1 \ 1]^T$

• Decision boundary'i bulalım:

$\rightarrow$ Predict y=1 when $-1 + x_1^2 + x_2^2 \geq 0 \Rightarrow \boxed{x_1^2 + x_2^2 \geq 1}$

---



• $h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \cdots)$

✗ Farklı $h_\theta(x)$'ler için daha farklı ve complex Decision Boundaries elde edilebilir.

# ~ Logistic Regression - Cost Function ~

- Training Set: $\{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots (x^{(m)}, y^{(m)}) \}$ ..examples

- $x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$   - $x_0 = 1$   - $y \in \{0,1\}$

- $h_\theta(x) = g(\theta^T x) = \dfrac{1}{1 + e^{-\theta^T x}}$
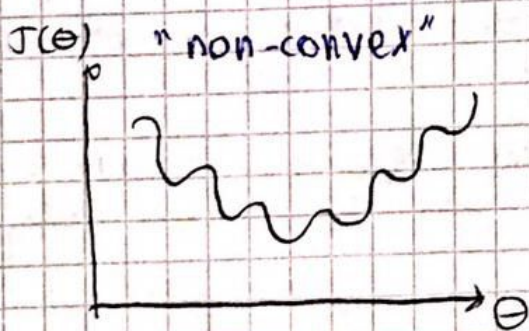
Question: How to choose parameters $\theta$?

- Linear Regression için  $\boxed{J(\theta) = \dfrac{1}{m} \sum_{i=1}^{m} \dfrac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2}$

- Let's say $\boxed{\text{cost}(h_\theta(x^{(i)}), y^{(i)}) = \dfrac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2}$
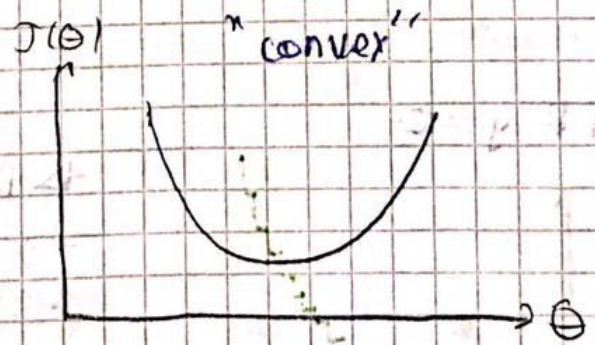
  $\boxed{\text{cost}(h(x), y) = \dfrac{1}{2}(h_\theta(x) - y)^2}$ (Simplified)

- Logistic Regression için bu cost function'ı kullanırsak J non-convex bi form olur yani bir çok local optima olur.

  $\rightarrow$ Yani  $h_\theta(x) = \dfrac{1}{1+e^{-\theta^T x}}$'i  $J(\theta) = \dfrac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$  de yerine koyarsak J linear regression'da olduğu gibi convex olmaz şöyle olur.

$J(\theta)$ "non-convex"

$J(\theta)$ "convex"

$\rightarrow$ Böyle bir $J(\theta)$ için gradient descent global minimum'u bulamayabilir.

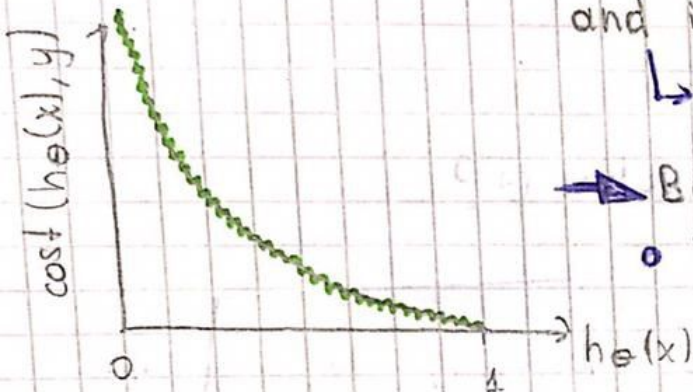$\not\otimes$  $h_\theta(x) = \dfrac{1}{1+e^{-\theta^T x}}$  logistic function very nonlinear olduğu için eski cost function tanımımız non-convex oluyor. Yeni bir cost function tanımlayarak $J(\theta)$'yı convex yapacağız.

# Logistic Regression Cost Function:

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

If y=1



cost (h₀(x),y) vs h₀(x)

➡ If hypothesis $h_\theta(x)$ predicts $y=1$ and indeed $y=1$ then cost = 0
  ↳ Cost = 0 if $y=1$, $h_\theta(x)=1$
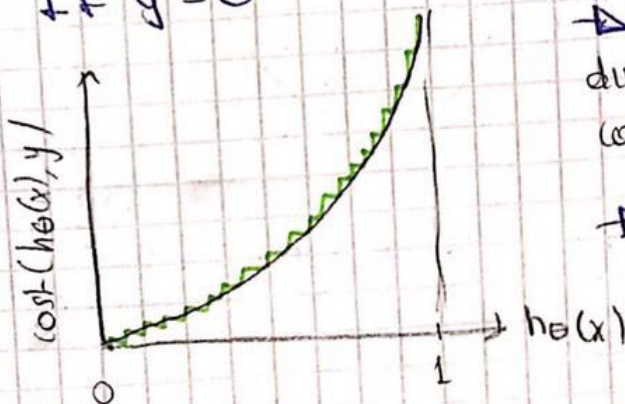
➡ But as $h_\theta(x) \to 0$, cost $\to \infty$
  • It means if $h_\theta(x)=0$ (predict $P(y=1|x;\theta)=0$), but $y=1$ we will penalize learning algorithm by a very large cost

➡ Logistic regression için, $h_\theta(x)$'in 0 ile 1 arasında olduğunu unutma çünkü $h_\theta(x)$ aslında $y=1$ olma olasılığını veriyor.

➡ Olay şu, cost'u hesaplarken önceden y değeri fark etmeksizin kare fark alıyorduk. Şimdi $y=1$ ise farklı cost, $y=0$ ise farklı.

➡ $y=1$ iken $h_\theta(x)$ bize 1 veriyorsa süper çalışıyor bu $x^{(i)}$ için cost 0! Fakat $y=1$ iken $h_\theta(x)=0$ veriyorsa, yani tümör malignant olduğunu vermemne rağmen hypothesis bunun olasılığı 0 diyorsa cost $\to \infty$ tamamen yanlış çalışıyor bu $h_\theta(x)$, $\theta$'ların değişmesi lazım.

If y=0



cost (h₀(x),y) vs h₀(x)

➡ Demin'kinin reversi, $h_\theta(x)=0$ diyorsa ve $y=0$ ise sıkıntı yok, cost = 0

➡ Fakat $y=0$ iken $h_\theta(x)=1$ verirse cost $\to \infty$. $\theta$'ların değişmesi şart.

Bu tek bir data like için yeni bir $x^{(i)}$ için ve
bir $y = 0$ a $y = 1$ değeri için cost hesaplar.
$J(\theta)$'yı bilmak için tüm bunları toplam ve $\frac{1}{m}$'le böler

## – Simplified Cost Function and Gradient Descent –

Cost Function for Logistic Regression

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

$$cost(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases}$$

ikisi
aynı
şey!

Note: $y = 1$ or $0$ Always

✱ Simdi $y = 0$ ve $1$ için aynı ayrı iki cost tanımlamalı yerine
bunları birleştirerek simplify ederogr:

$$cost(h_\theta(x), y) = -y \cdot (\log(h_\theta(x)) - (1-y)\log(1 - h_\theta(x))$$

why
this
cost function

it can be derived from statistics using the principle of
maximum likelihood estimation. And it is convex!
Thus
Burada $\log \rightarrow \ln$ ;

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^{m} y^{(i)} \cdot \log(h_\theta(x^{(i)})) + (1-y) \log(1 - h_\theta(x^{(i)})) \right]$$

↳ Bu fonksiyon logistic regression için convex olacaktır.

✱ Cost function $J(\theta)$'yı yazdık, artık amacımız bu $J(\theta)$'yı
minimize edecek $\theta$ parametrelerini bulmak!

↳ $\min\limits_{\theta} J(\theta)$ : Get $\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}$

✱ After obtaining $\theta$, then we have $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$
so we can make a prediction for given new $x^{(i)}$.

$$h_\theta(x) = P(y = 1 | x; \theta) \text{ olduğunu unutma}$$

✱ Cost function $J(\theta)$'yı minimize etmek için yine
gradient descent kullanırız sonuçta bu yine convex ban
function. Mantık aynı!

3. videoda Advanced Optimazation tekniklerini
gösteriyor α'yı seçmeden daha hızlı converge eden
complex methodlar.

## Gradient Descent

Want min'e $J(\theta)$:

Repeat $\{$

$$\theta_J := \theta_J - \alpha \boxed{\frac{\partial}{\partial \theta_J} J(\theta)} \qquad \text{Aynı} \, b$$

→ Simultaneously update all $\theta_J$.

$$\frac{\partial}{\partial \theta_J} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_J^{(i)} \qquad \begin{array}{l}\text{Öncekı ile} \\ \text{birebir aynı!}\end{array}$$

$$\theta_J := \theta_J - \alpha \sum_{i=1}^{m} (h(x^{(i)}) - y^{(i)}) \cdot x_J^{(i)} \qquad \begin{array}{l}\text{Simultaneously update} \\ \text{all } \theta_J\end{array}$$

→ $\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}$ hepsi için simultaneously update

※ Sonuç olarak gradient descent algorithm eskisi
ile birebir aynı fakat $h_\theta(x)$ değişti!

※ Önceden gradient descent'i monitor ederek converge
ettiğini görmekten bahsetmiştik aynısı geçerli. $J(\theta)$'yı
iteration'a göre plot edersek hep azalması lazım.

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

vector formda

※ Feature scaling'ide logistic regression'da daha
hızlı bir convergence için kullanabılırız.

# — Multiclass Classification —   $y=1$   $y=2$   $y=3$   $y=4$
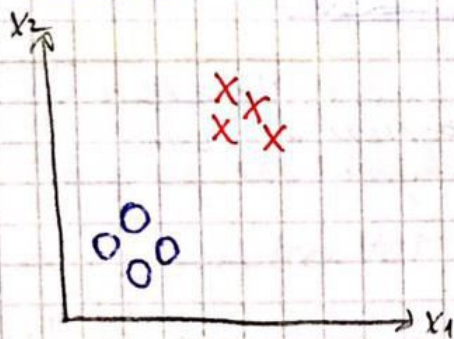
**Example:** Email foldering / tagging: Work, Friends, Family, Hobby

✗ Şimdiye kadar binaryclass classification gördük, $y=1$ yada $y=0$ olmak zorundaydı. Burada $y=0,1,2,3,\dots$ olabilir.

$y=1$   $y=2$   $y=3$

**Example:** Medical Diagnosis: Not ill, Cold, Flu
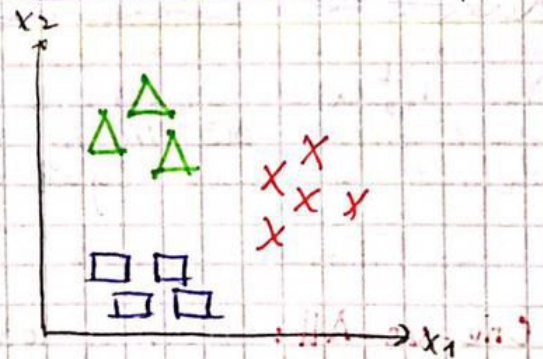
$y=1$   $y=2$   $y=3$   $y=4$

**Example:** Weather: Sunny, Cloudy, Rain, Snow

✗ y can take small number of discrete values. $y=0$'dan da başlayabilir, önemi yok.

Binary Classification              Multi-class classification



? How to get a learning algorithm?

→ Binary için, logistic regression modelini kullanırım ve belki bir çizgi çekerim, decision boundary olarak. Bir tarafı positive class $y=1$, bir tarafı negatif class $y=0$ olur. Bu line'da $h_\theta(x)=0.5$'i yeni $\theta^T \cdot x = 0$'ı temsil eder.

✗ "One vs All Classification" adında bir idea kullanarak, Eski modeli alarak multi-class classification için de çalışır hale getiririz.
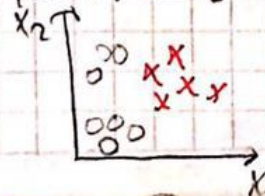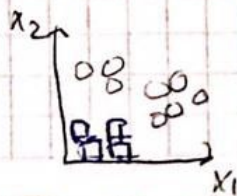
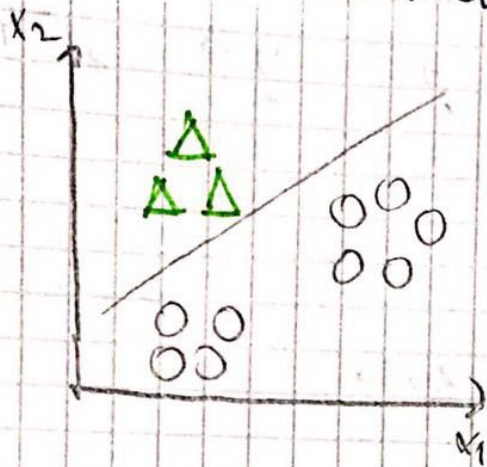✗ Nasıl Çalışır: Yukarıdaki gibi bir training set olsun.

Class 1: △ ← $y=1$  
Class 2: □ ← $y=2$  
Class 3: ✗ ← $y=3$  

} Yapılacak şey, training set'i alıp, 3 ayrı binary classification problemine dönüştürmek!

Önce Class 1'i ele alalım  <u>class 1 vs all</u>:



❌ Class 1 positive class oldu, diğerleri negative class,

→ classifier: Bunu training set ile eğitecez.

❌ $h_\theta^{(1)}(x)$ assign edicek ve heo yaptığımızı yapıcak ve mesela şekildeki gibi bir decision boundary elde edeceğiz.

Sonra Class 2 için: Yeni bir logistic regressien classifier $h_\theta^{(2)}(x)$ fit ederiz.

Sonra Class 3 için: $h_\theta^{(3)}(x)$

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

❌ We have fit three classifiers for $i = 1, 2, 3$.

$$h_\theta^{(i)}(x) = P(y = i \mid x; \theta) \qquad (i = 1, 2, 3)$$

Yani $h_\theta^{1}(x)$: girilen x datası için ve $\theta$'lar için $y = 1$ olma olasılığını veriyor. Diğerleri de aynı şekilde $y = 2$ ve $y = 3$ olma olasılığını veriyor.

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

**One-vs-All**: Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $y = i$.

On a new input $x$, to make a prediction, pick the class $p$ that maximizes

$$\max_p h_\theta^{(i)}(x)$$

Yani en yüksek olasılıklıyı seçeriz.

# W3 - Advanced Optimization Algorithms

→ Bunları kullanarak logistic regression'ın GDA ile çok daha hızlı çalışmasını
ve large scale problemlere rahat uygulanabilmesini sağlar.

→ GDA'nın dışında başka optimization algorithms de var:

- Conjugate gradient
- BFGS
- L-BFGS

+ No need to manually pick α
+ often faster than gradient descent
− More complex

## Example:

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

Let's apply an advanced opt. algorithm
to minimize $J(\theta)$:

```
function [JVal, gradient] = costFunction(theta)

JVal = (theta(1)-5)^2 + (theta(2)-5)^2;
gradient = zeros(2,1);
gradient(1) = 2*(theta(1)-5);
gradient(2) = 2*(theta(2)-5);
```

Now we can call adv. opt. function:
fminunc:

```
options = optimset('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros(2,1);
[optTheta, functionVal, exitFlag] = fminunc(@costFunction,
                 initialTheta, options);
```

Önce üst taraftaki function'a end ekle  costFunction ismi ile kaydet!
Sonra aşağıdaki kısmı başka mfile'a yaz '100' yerine inL istiyor
direkt 100 yaz.      help fminunc da yardımcı olabilir.

## Sonuçta Logistic Regression için:

$$theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \begin{matrix} - theta(1) \\ - theta(2) \\ \\ - theta(n+1) \end{matrix}$$

```
function [JVal, gradient] = costFunction(theta)
    JVal = [code to compute J(θ)]
    gradient(1) = [code to compute ∂/∂θ₀ J(θ)]
    ⋮
    gradient(n+1) = [code to compute ∂/∂θₙ J(θ)]
```

Sonra fminunc'ı
kullanır ve
optimized theta
values will be
obtained.

**Ideal for large
scale (too many
features) problems**