

## Logistic Regression

Logistik regresyon, verileri ayırık sonuçlara ayırmak için kullanılan bir yöntemdir. Örneğin; bir e-postayı spam olarak veya spam değil olarak sınıflandırmak için lojistik regresyonu kullanırız. Bir başka örnek ise bir ticaret sitesinde kullanıcıların kalıntı kredi kartı kullanıp kullanmadığını onaylayan bir sistem.

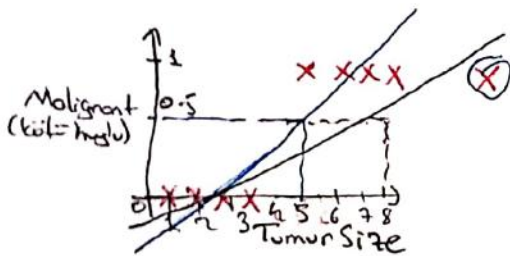
### Classification

İşminde de onbölücüye göre belli verileri kümellemek (sınıflandırmak) Tümenin iyi huylu olup olmadığı bir sınıflandırma örneğidir.

Basit sınıflandırmada  $y$  değişkenimiz 0 ya da 1'dir. Mesela tümör iyi huyludur veya kötü huyludur. Mail spamdir ya da değildir. Tabi gelişmiş sınıflandırma problemlerinde 0, 1, 2, 3, 4, ... olabilir.

$$y \in \{0, 1\} \quad \begin{array}{l} 0: \text{"Negative Class"} \\ 1: \text{"Positive Class"} \end{array} \quad \text{olarkta adlandırılır}$$

0 ve 1'e değer atamak (spam veya değil) genelde 0 olumlu durumken 1 olumsuz durumdur. Ama bunun bir önemi yok size kalmış bir durum.



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

Threshold 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 " $y=0$ "

Yukarıda görüldüğü gibi sınıflandırma için lineer regresyon kullanabiliriz. 0.5'te üstü kötü huylu, altı iyi huylu diye sınıflandırabiliriz. Ancak bu yöntem iyi çalışmaz. Çünkü yukarıda iki doğru ver gördüğünüz gibi mori doğru yörünceğine alınmış çarpı yok sayılarak çizilmiştir. Budurunda 5 tümör boyutunun sol tarafı iyi huylu iken sağ tarafı kötü huyluyu gösteriyor. Sağ üst tarafa bir veri koyduğumuzda doğrunuz değişiyor ve yanlış doğru oluyor. Bu durumda 8 tümör boyutundan sol tarafı iyi huylu, sağ tarafı kötü huylu oluyor. Bu bizim asıl kötü huylu tümörlerimizin seplenememesine neden olur. Bu niye oldu tet bir verinin modeli değiştirmesi yüzünden oldu bu yüzden lineer regresyon kullanılmamalı. Başka bir nedenele  $y$  çıktıların  $h_{\theta}(x)$  arasında değer vermemeye ihtimali. Çok büyük değerler verebilir ve bunu normalleştirebiliriz.

## Logistik Regresyon Hipotez Tanımı

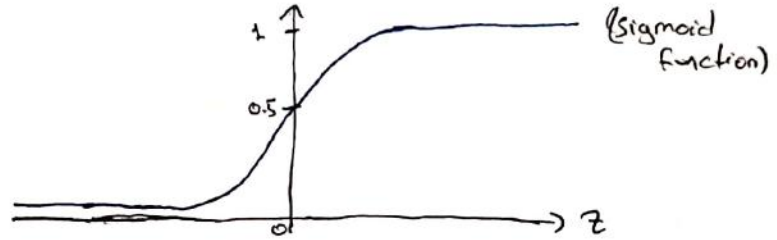
Hipotez olarak kullanacağımız fonksiyon **sigmoid function** veya **logistic function** denot geçmektedir.

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

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

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

$$z \in \mathbb{R}$$



Grafığı yorumlarsak  $z$  değeri  $-\infty$  yaklaştıkça sigmoid fonksiyonu 0'a yaklaşıyor.  $z$  değeri  $+\infty$  yaklaştıkça 1'e yaklaşıyor. Bu sayede  $h_{\theta}(x)$  çıktılarını  $0 \leq h_{\theta}(x) \leq 1$  arasında tutmuş oluyoruz.

Not Formülde gördüğümüz gibi  $\theta$ 'ya ihtiyacımız var. Bunu bulmalıyız.

Örnek

$$x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tümörSiz} \end{bmatrix}$$

$$h_{\theta}(x) = 0.7$$

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

$$P(y=0 | x; \theta) = 1 - h_{\theta}(x)$$

Yukarıdaki örnek belli bir tümör boyutuna göre fonksiyonumuz 0.7 sonucunu veriyor. Bu da bize hastamızın tümörünün %70 olasılıkla kötü huylu olduğunu gösteriyor. İyi huylu olma olasılığı ise  $1 - 0.7 = 0.3$ 'tür. Yani %30 iyi huylu.

## Decision Boundary (Kerem Sınırı)

$$h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

$$h_{\theta}(x) < 0.5 \rightarrow y=0$$

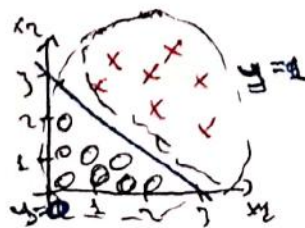
$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

when  $\theta^T x \geq 0$

$$\theta^T x \geq 0 \text{ ise } y=1$$

$$\theta^T x < 0 \text{ ise } y=0$$

$y=0$  ve  $y=1$  sınırlarının nerede olduğunu ayıran çizgidir. Kerem sınırı, veri kümesi ile ilgili değil hipotezin  $\theta$  parametreleri ile ilgilidir.

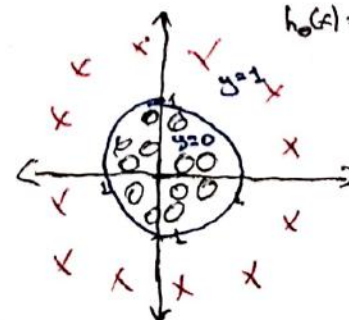


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$-\frac{3}{2} + x_1 + x_2 \geq 0$$

$$x_1 + x_2 \geq \frac{3}{2}$$

$$x_1 + x_2 \geq 3$$



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

$$-1 + x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 1$$

$$x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 0$$

$$x_1^2 + x_2^2 \geq 0$$

daire formülü



## Logistic Regression Model

Cost function

Training set:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_m \end{bmatrix} \in \mathbb{R}^{n+1}$$

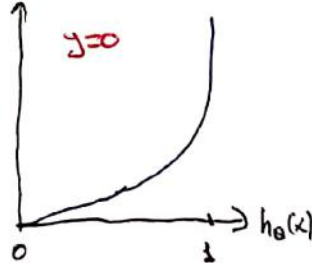
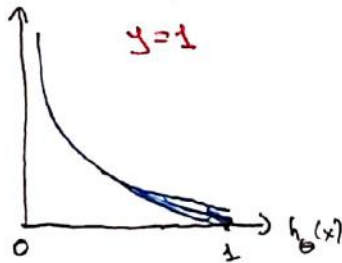
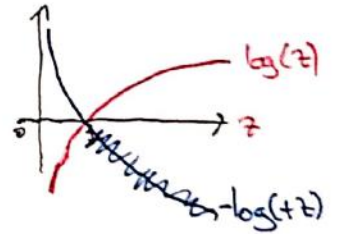
$$x_0 = 1, y \in \{0, 1\}$$

Doğrusal regresyon için kullandığımız aynı maliyet işlevini kullanmayız.  $\left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2\right)$  Çünkü lojistik işlev gücününun dışbükey bir grafik değil, birçok yerel optimuma sahip dalgalı bir grafik verir. Bu durumda yerel optimuma ulaştığımızı hiçbir zaman garanti edemez.

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = -\log(h_\theta(x)) \quad \text{if } y=1$$

$$\text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x)) \quad \text{if } y=0$$



$$\begin{aligned} \text{Cost}(h_\theta(x), y) &= 0 \quad \text{if } h_\theta(x) = y \\ \text{Cost}(h_\theta(x), y) &\rightarrow \infty \quad \text{if } y=0 \text{ ve } h_\theta(x) \rightarrow 1 \\ \text{Cost}(h_\theta(x), y) &\rightarrow \infty \quad \text{if } y=1 \text{ ve } h_\theta(x) \rightarrow 0 \end{aligned}$$

Eğer doğru cevabımız  $y=0$  ise hipotez fonksiyonumuz da 0 verirse maliyet fonksiyonu 0 olur. Hipotezimiz 1'e yaklaşırsa maliyet fonksiyonu sonsuza yaklaşıyor ve bize maliyeti oldukça fazla olur.

Eğer doğru cevabımız  $y=1$  ise hipotez fonksiyonumuz 1 çıkarsa maliyet fonksiyonu 0 olur. Hipotezimiz 0'a yaklaşırsa maliyet fonksiyonu sonsuza yaklaşıyor.

Not Maliyet fonksiyonunun bu şekilde yapılması dışbükey olduğunu garanti eder.

## Simple Cost Function ve Gradient Descent

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

Cost function'ümüzü yukarıdaki gibi iki satırla ifade ediyoruz ancak gradient descent'te uygulayabilmek için tek bir satıra indirgeyebiliriz.

$$\text{Cost}(h_{\theta}(x), y) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

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

Lineer regresyonun maliyet fonksiyonu yukarıdaki gibidir. Ancak bu cost function yerine Principle of Maximum Likelihood Estimation (Maksimum olasılık tahmini ilkesi) 'de kullanılabilir.

## Gradient Descent

$J(\theta)$ 'ımızı minimize etmemiz gerekiyor ve bunun içinde gradient descent kullanacağız.

$$\min_{\theta} J(\theta):$$

repeat  $\xi$

$$\left. \begin{aligned} \theta_{\hat{j}} &:= \theta_{\hat{j}} - \alpha \frac{\partial}{\partial \theta_{\hat{j}}} J(\theta) \Rightarrow \theta_{\hat{j}} := \theta_{\hat{j}} - \frac{\alpha}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_{\hat{j}}^{(i)} \\ \} \end{aligned} \right\}$$

$J(\theta)$  Vectorized

$$h = g(x\theta)$$

$$\hat{J}(\theta) = \frac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h))$$

Gradient Descent Vectorized

$$\theta := \theta - \frac{\alpha}{m} X^T (g(x\theta) - \vec{y})$$

## Advanced Optimization

### Optimization Algorithms

- Gradient Descent
- Conjugate Gradient
- BFGS
- L-BFGS

- $\alpha$ , öğrenme oranını manuel seçmeniz gerekmez.
- Gradient descent'ten genelde daha hızlıdır.
- Aşırı kompozittir. NE yaptıklarını onarmak oldukça güçtür.

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

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

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

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

$$\text{function } [fVal, \text{gradient}] = \text{costFunction}(\theta)$$

$$fVal = (\theta(1) - 5)^2 + (\theta(2) - 5)^2;$$

$$\text{gradient} = \text{zeros}(2,1);$$

$$\text{gradient}(1) = 2 * (\theta(1) - 5);$$

$$\text{gradient}(2) = 2 * (\theta(2) - 5);$$

$$\text{options} = \text{optimset}('GradObj','on','MaxIter',100);$$

$$\text{initTheta} = \text{zeros}(2,1);$$

$$[\text{optTheta}, \text{functionVal}, \text{exitFlag}] = \text{fminunc}(@\text{costFunction}, \text{initTheta}, \text{options});$$

Not Yukarıdaki optimizasyon algoritmelerini hazır kütüphanelerden kullanarak fminunc'ta Octave'de örnek verilebilecek optimizasyon fonksiyonudur. Daha büyük problemler üzerinde çalışırken gradient descent'ten daha hızlı olan bu algoritmeleri kullanabilirsiniz.

## Logistic Regression with Multi-class Classification

Emcillerinizi iş, orkadeş, arke, hobir olarak 3 sınıfa ayırmak isteyebilirsiniz veya titoli bir burunla hastaneye giden biri hasta olmayabilir, soğukalgınlığı, grip olabilir. Bu iki örnekte multi-class sınıflandırmaya örnektir.

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

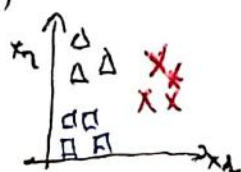
$$h_{\theta}^{(0)}(x) = P(y=0|x;\theta)$$

$$h_{\theta}^{(1)}(x) = P(y=1|x;\theta)$$

$$h_{\theta}^{(n)}(x) = P(y=n|x;\theta)$$

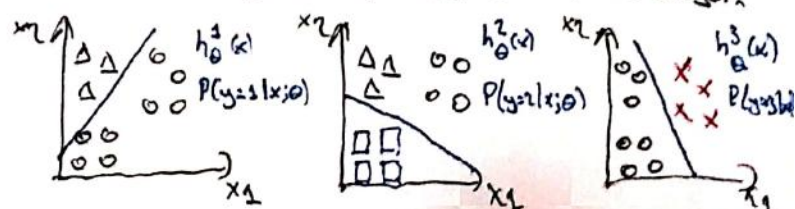
$$\text{prediction} = \max h_{\theta}^{(i)}(x)$$

class 1:  $\Delta$   
class 2:  $\square$   
class 3:  $\times$



One vs All (one vs rest)

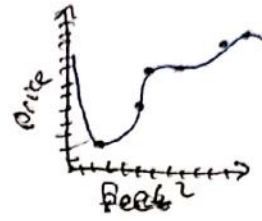
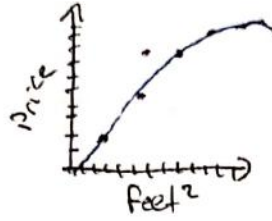
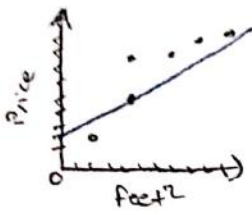
Bu yöntem tüm sınıfların içinden birini seçip diğerlerini bir gibi düşünüp logistic regresyon uyguluyoruz ve bu yöntemi tüm sınıflar için yapıyoruz. Sağda sınıf sayısı kadar hipotez çıkarır ve en çok olasılığı veren, örneğimizin hipotezi olur.





## Regularization (Düzenleme)

### Over fitting Problemi



En fiyat tahmini üzerinden devam edelim. En soldaki şekilde  $y = \theta_0 + \theta_1 x_1$  modelinde verilerin düz bir çizgi üzerinde olmadığını ve uyumun (fit) çok iyi olmadığını görüyoruz.

Orta şekilde  $y = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$  modeli verilerimize daha iyi uyum bir model sergiliyor. Ancak daha fazla özellik eklemişle beraber verilere daha iyi uyum gösteren model çıkmasına karşın yeni örnekleri iyi bir tahminde bulunma olasılığı düşecektir.

En soldaki şekilde model verileri iyi bir şekilde yakalayamıyor buna **under fitting** deniyor. En sağdaki şekil ise verilere çok fazla uyum sağlamış buna da **overfitting** deniyor.

Underfitting veya yüksek bias, hipotez fonksiyonumuz ile verilerin zayıf eşleşmesidir. Genellikle çok basit veya çok az özellik kullanan bir işlevden kaynaklanır.

Overfitting veya düşük varyans, mevcut verilere (training set) çok iyi uyum ancak iyi bir genelleştirme yapamayan hipotez işlevinden kaynaklanır. Genellikle verilerle ilgisi olmayan çok sayıda gereksiz eğri ve ağı oluşturan karmaşık bir fonksiyondan kaynaklanır.

Örnek olarak lineer regresyon örneği verdik ama logistic regresyonda da meydana gelir.

Over fittingin özelliğini giderebilmek için iki ana yol vardır:

1. Özellik sayısını azaltın

- Seçilecek özellikleri manuel seç.
- Bir model seçim algoritması kullan.

2. Regularization (Düzenleme)

- Tüm özellikleri koruyun, ancak parametrelerin ( $\theta_j$ ) büyüklüğünü azaltın.
- Düzenleme, birçok kullanışlı özelliklere sahip olduğumuzda iyi çalışır.

## Regularization - Cost Function

Eğer fonksiyonumuzda overfitting varsa, neden olan parametrelerin maliyetlerini artırarak daha küçük değer almasını sağlayabiliriz.

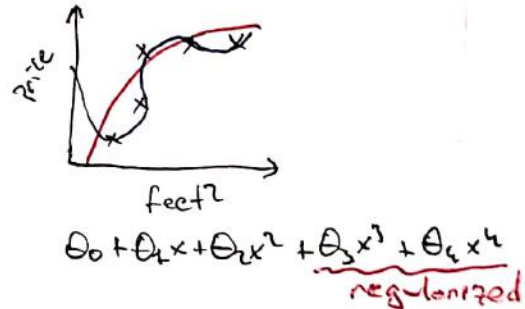
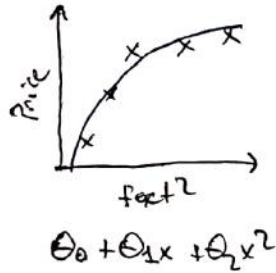
Diyeelim ki aşağıdaki fonksiyonu karesel yapmak istiyoruz:

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

$\theta_3 x^3$  ve  $\theta_4 x^4$  parametrelerinin etkisini kaldırmamız gerekecektir. Bu özelliklerden kurtulmadan ya da hipotezimizin biçimini değiştirmeden, **cost function**'imizi değiştirebiliriz.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

Maliyet fonksiyonumuzun sonuna maliyetleri artırmak için iki terim ekledik. Şimdi maliyet fonksiyonun sifıra yaklaşması için  $\theta_3$  ve  $\theta_4$  sifıra yakın olması gerekecek.



Gök fazla özelliği den bir problemde hangi özelliği overfittinge neden oluyor anlayabiliriz. Bunun için tüm parametreleri düzenleyen genel bir formül kullanırız. Overfitting olmasa dahi parametreleri regülarize ederek küçültmek modeli daha pürüzsüz hale getiririz.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2$$

**Not**  $\theta_0$  regülarize edilmez. Eğer regularization parametresi ( $\lambda$ ) çok büyük seçilirse hipotez  $\theta_0$ 'a yakınsayan bir fonksiyona dönüşür,  $\lambda$  seçimine dikkat edilmelidir.

## Regularization - Linear Regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

### Gradient Descent

Cost function'ımız değiştiğine göre tekrar kısmi türev alıyoruz ve aşağıdaki sonucu çıkar:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j = (1, 2, \dots, n)$$

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

Bu formülde  $1 - \alpha \frac{\lambda}{m}$  'den biraz küçük bir sayı çıkacaktır. Çünkü  $\alpha$  öğrenme oranımız küçük eğitim veri sayımız büyük olduğu için buradan çok küçük bir pozitif sayı çıkacaktır. Burada 1'den çıkardığımızda 0,80-0,99 arasında bir değer elde ederiz ve bu sayede  $\theta$  güncellemelerinde  $\theta$ 'ümüzü daha küçük bir şekilde güncellemiş oluruz.

### Normal Equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

$m \times (n+1)$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$m \times 1$

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

$$\theta = (X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 1 & 0 \end{bmatrix}}_{(n+1) \times (n+1)})^{-1} X^T y \quad \Rightarrow \text{Laplace regularization formül}$$

boğutlarındadır

Not Regularization' normal equation'deki non-invertibility (tersine çevrilemezlik) sorunu da ortadan kaldırır. Tek şart  $\lambda$ 'nın 0'dan büyük olmasıdır.

Not Düzenleme yaparken  $\theta_0$ 'i düzenlemediğimizden  $\theta_0$ 'i formülde ayrı hesaplıyoruz.



## Regularization - Logistic Regression

$$J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Gradient Descent  
repeat  $\xi$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j = (1, 2, 3, \dots, n)$$

Aslında gradient descent linear regresyonunki ile aynı yazım tarzında tek farkı hipotezlerinin farklı olmasıdır.

Not Logistic regresyonun doğru çalışıp çalışmadığını onlamak için  $J(\theta)$  güncelleme geçmişini ve iterasyon sayısını içeren grafik çizmek <sup>bu işin sonucunda</sup> ve azalan bir grafik görmek gerekir.

## Regularization - Advanced Optimization Algorithm (fminunc)

function [f\_val, gradient] = cost\_function(theta)

f\_val = [code to compute  $J(\theta)$ ]  $\Rightarrow J(\theta) = - \left[ \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$

gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ]  $\Rightarrow \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)} \right]$

gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ]  $\Rightarrow \left[ \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_n^{(i)} \right] + \frac{\lambda}{m} \theta_n$

Not Regularization her zaman eğitim setinde en iyi sonucu verir diyemeyiz.

Not Yeni bir özellik eklemek modeli daha iyi bir hale getirebilir.