

Addressing overfitting: Overfitting'i address etmek için 2 seçenek var:

Options:

1. Reduce number of features.

→ — Manually select which features to keep. : Gereksizlerini atarız.

→ — Model selection algorithm (later in course). : Algorithms automatically decide which feature to keep, which to throw out.

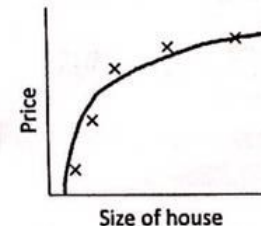
2. Regularization. →  $\theta$ 'ların değerini düşürürs. Böylece atmadan overfitting'den kurtulabiliriz.

→ — Keep all the features, but reduce magnitude/values of parameters  $\theta_j$

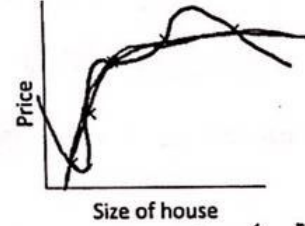
→ — Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

→ Overfitting'e sebep olan terimlerin  $\theta_3 x^3 + \theta_4 x^4$  ise ben bunları minimize ederseniz overfitting'i çözerim, yani  $\theta_3 \approx 0$  ve  $\theta_4 \approx 0$  yaparsam problem çözülür! Çünkü hipotez yine 2. derece gibi olacak gelecektir.

Intuition



$$\theta_0 + \theta_1 x + \theta_2 x^2$$



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

Suppose we penalize and make  $\theta_3, \theta_4$  really small.

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

$$\theta_3 \approx 0 \quad \theta_4 \approx 0$$

$\theta_3$  ve  $\theta_4$  parametrelerini (overfitting'in sebebi) cost function'a ekliyoruz böylece artık bu parametreleri de minimize edip hipotez olan etkisini düşürürs.



Machine Learning

## Regularization

## Cost function



Gözetiminde sadece  $\theta_2$  ve  $\theta_4 \approx 0$  yaparsak simpler hypothesis elde ettilik.  
More broadly:  $\theta$ 'lar küçülürse simpler hypothesis ve less change for overfitting.  
Regularization.

Small values for parameters  $\theta_0, \theta_1, \dots, \theta_n$   
 - "Simpler" hypothesis  
 - Less prone to overfitting

Housing:

- Features:  $x_1, x_2, \dots, x_{100}$
- Parameters:  $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

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

$\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$

100 tane  $\theta$  var hangi parametre geneliştir, hangisi higher order bilmiyoruz. Regularization ile cost function for linear regression'ı alıp bütün  $\theta$  parametrelerini shrink etmek için modify ederim. Reg. Term bütün  $\theta$ 'ları shrink eder,  $\theta_0$  hariçinde!  $\theta_0 = 1$  hep! ...

In regularized linear regression, we choose  $\theta$  to minimize

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

What if  $\lambda$  is set to an extremely large value (perhaps for too large for our problem, say  $\lambda = 10^{10}$ )?

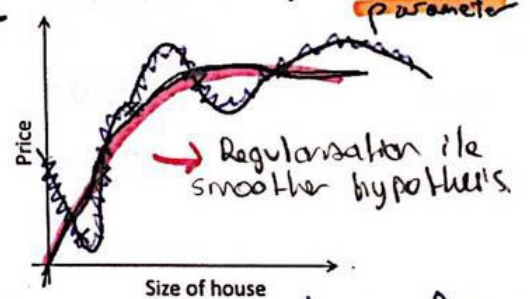
- ✗ Algorithm works fine; setting  $\lambda$  to be very large can't hurt it
- ✗ Algorithm fails to eliminate overfitting.
- ✗ Algorithm results in underfitting. (Fails to fit even training data well).
- ✗ Gradient descent will fail to converge.

$\lambda$ : Controls the tradeoff b/w two goals: 1) Fitting the training set well. 2) Keeping the parameters small: Avoid overfitting.  
 $\lambda$ : Çok büyükse, underfitting olur!  
 Çok küçülürse, overfitting!  
 Regularization.

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

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

Goal 1



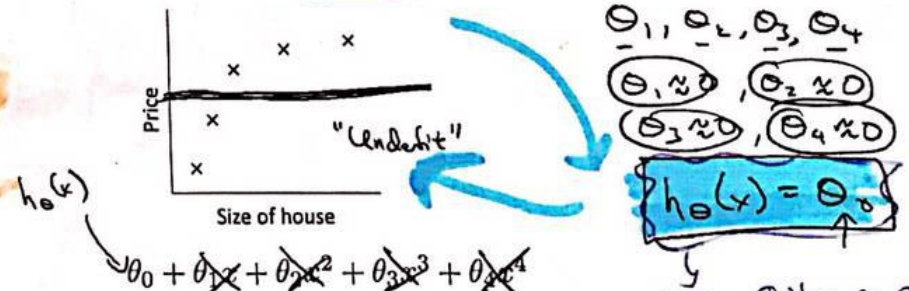
Nedenini anlamadan normal diyen. Ben de gör çalışın, bunu böyle bil diyen!

$\lambda$  büyükse  $\theta$ 'lar  $\approx 0$  olacak underfitting occurs

In regularized linear regression, we choose  $\theta$  to minimize

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

What if  $\lambda$  is set to an extremely large value (perhaps for too large for our problem, say  $\lambda = 10^{10}$ )?



Yani  $\lambda$  doğru seçilmeli!

diğer  $\theta$ 'lar  $\approx 0$  olursa sonuç horizontal flat line olur!



→ Linear Regression için iki farklı algoritma öğrenmiştik, 1 based on gradient descent, 1 based on normal equation

→ Şimdi ikisini de alıp, generalise ederek, Regularized Linear Regression'i oluşturacağız



Machine Learning

# Regularization

## Regularized linear regression

Gradient Descent den  $\Theta_0$ 'yu aynı şekilde minimize ettiğimizde gördüğümüz üzere, regularization term  $\lambda = 1$  den büyükler  $\Theta_0$ 'i etkilemez.

### Regularized linear regression

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

regularization term

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

→ Amacımız regularized cost function'ı minimize etmek!

→ Bunun için Gradient Descent kullanacağız

**Gradient descent**

Repeat {

$\Theta_0$  için aynı yazdık kollarlar için aynı!

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

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

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

$1 - \alpha \frac{\lambda}{m} < 1$   $0.99$   $\theta_j \times 0.99$

Regularization etkisi ile  $\theta_j$  biraz küçükler  $1 - \alpha \frac{\lambda}{m}$  genelde 1 den biraz küçükler 0.99 gibi,  $\alpha$  ve  $\lambda > 0$  ve  $m > 0$ .

son form gösteriyor ki regularization ile yaptığımız aslında henüz önceki  $\theta_j$ 'yi 0.99 gibi bir rakamla çarpıp shrink etmek. Gerisi tamamen aynı!

$\frac{\partial}{\partial \theta_0} J(\theta)$   $\frac{\partial}{\partial \theta_j} J(\theta)$   $J(\theta)$   $\theta_j^2$



\* Min  $\theta$ 'ları bulmak için 2. algoritma Normal Equation idi.  $X^T$  ve  $y$  yi bildükten sonra  $\min \theta = (X^T X)^{-1} X^T y$  olarak bulunabilir

## Normal equation

$$\underline{X} = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \leftarrow \begin{matrix} m \times (n+1) \end{matrix}$$

$$\underline{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \leftarrow \mathbb{R}^m$$

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

$$\frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set}}{=} 0 \quad m \rightarrow$$

$$\Rightarrow \underline{\theta} = (X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & 0 & \ddots & 0 & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}}_{(n+1) \times (n+1)})^{-1} X^T y$$

e.g.  $n=2$   $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

- $\rightarrow$  Kırımlı önceden uygulanan algoritma idi!
- $\rightarrow$  Hani yeni eklenen term. Regularization ile gelen term. (Yeni cost function'dan bulunabilir)
- $\frac{\partial}{\partial \theta_j} J(\theta) = 0$  için  $\theta$  dâğılır!

\* Logistic Regression modelinde  $J(\theta)$ 'yi minimize etmek için 2 farklı algorithmler kullanılmaktadır.

- 1) Gradient Descent
- 2) Advanced Optimization Methods

\* Simdi bunları Regularization ile adapt edeceğiz!

## Non-invertibility (optional/advanced).

Suppose  $m \leq n$ ,  $\leftarrow$   
(#examples) (#features)

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

non-invertible / singular

pinv inv

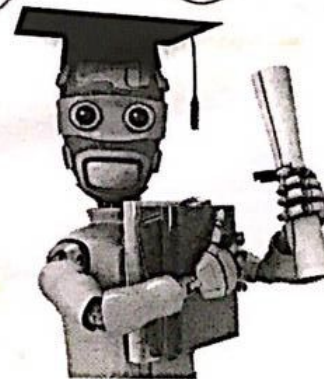
If  $\lambda > 0$ ,

$$\theta = \left( X^T X + \lambda \underbrace{\begin{bmatrix} 0 & 1 & 1 \\ & 1 & \\ & & \ddots \\ & & & 1 \end{bmatrix}}_{\text{invertible}} \right)^{-1} X^T y$$

Andrew Ng

\* Regularization ile non-invertibility problemi çözülür!

Andrew Ng



Machine Learning

## Regularization

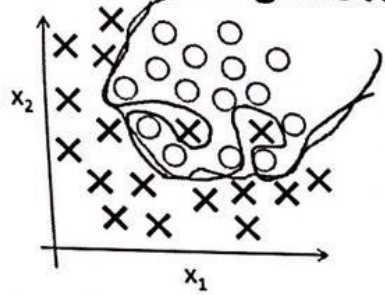
Regularized logistic regression

Scanned with CamScanner

CS



## Regularized logistic regression.



Cost function:

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

$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

By overfittinge neden olabilir.

## Gradient descent

Repeat {

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

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

$j = 1, 2, 3, \dots, n$   
 $\theta_1, \dots, \theta_n$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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

Andrew Ng

cosmetically same Andrew Ng

Birer önceki ile aynı şey çıkar!  
Sadece  $h_{\theta}(x)$  farklı!  
Hypothesis is different!

## Advanced optimization

if minune (e cost function)  $\theta_0, \theta_1, \dots, \theta_n$   $\theta_0, \theta_1, \dots, \theta_n$

function [jVal, gradient] = costFunction(theta)

jVal = [code to compute  $J(\theta)$ ];

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

gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];

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

gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];

$$\left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)} \right) - \frac{\lambda}{m} \theta_1$$

gradient(3) = [code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$ ];

$$\vdots \left( \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)} \right) - \frac{\lambda}{m} \theta_2$$

gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];

$J(\theta)$

Kod aynı yalıtıcı

$J(\theta)$  değişiyor

biraz farklı

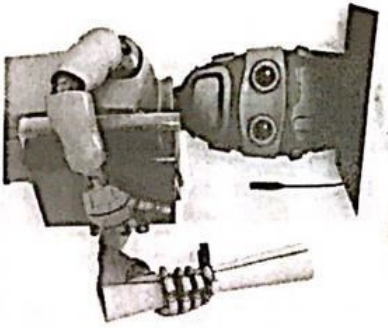
gradients de

değişiyor

sonuçta yine bir algı ile logistic regression "em regülasyon parametresini"  $J(\theta)$ 'yi, min yapan  $\theta$ 'lar elde edilir b



what is overfitting?

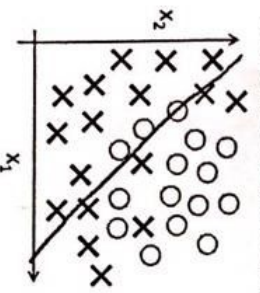


# Regularization

# Machine Learning

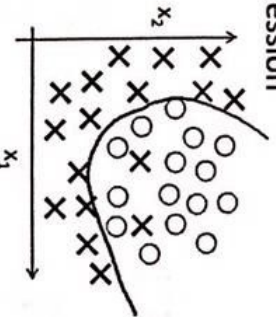
→ Overfitting: logistische Regression kein deterministisch.

### Example: Logistic regression

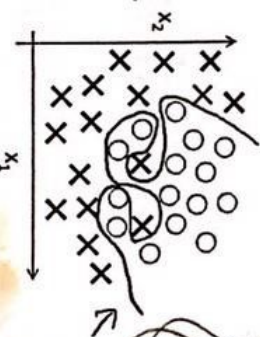


$$\rightarrow h_{\theta}(x) = \underbrace{g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)}_{(g = \text{sigmoid function})}$$

"Undarfit"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 \underline{x_1^2} + \theta_4 x_2^2 + \theta_5 \overline{x_1 x_2}) \quad \swarrow$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots) \leftarrow$$

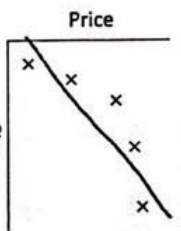
"Ola-fu"

A lot of high order polynomials

⑦ Logische Aggression da bei ihr holo bifa  
sons ut aost getunyer nu? o daman  
grosenst- d'scent T'cal'yi nasil minime  
edereh, ilk ke 2. aqou'thber i'ien  
konusurda 2. yag'timyon Sonur  
o olmasi'ien h'ed'ien 2  
o vermes'tat'm. yene nar b

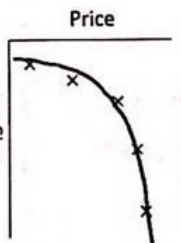
### Example: Linear regression (housing prices)

2. vs 3. Linear Regression always scales?



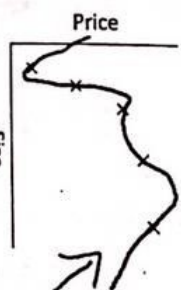
$$\rightarrow \theta_0 + \theta_1 x$$

"Underfit"      "High bias"



$$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$

"Just right"



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

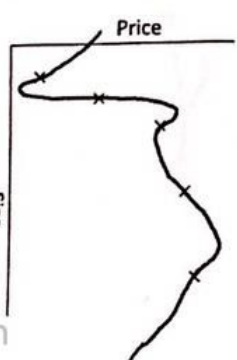
"Overfit" "High variance"

→  $J(\theta) \propto 0$  also dahi, yani previously unseen data  $y_i$  Hypothesis has  
**Overfitting:** If we have too many features, the learned hypothesis  $h_{\text{learned}}$  may fit the training set very well ( $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \approx 0$ ), but fail to generalize to new examples (predict prices on new examples).

② Problem feature sayısından ziyade  $h(x)$ 'in doğrusal mı değil mi?  $x_3, x_4$  lerin içinde girmediği değil mi?

### Addressing overfitting:

- $x_1$  = size of house
- $x_2$  = no. of bedrooms
- $x_3$  = no. of floors
- $x_4$  = age of house
- $x_5$  = average income in neighborhood
- $x_6$  = kitchen size
- $\vdots$
- $x_{100}$



→ Col Post feature varsa ve  
col as training data var  
ise OVERFITTING bir  
problem olabilir.

Baghate me sala yori hamada  
tel bir yosigi horigor cuntu  
overfilad.

Tek bir feature  
kansa, plotting ile  
kacını derece polinom  
kullancagına kadar  
verebiliyim, ama  
gerekli feature sayisi  
çok daha fazla  
olur