

Week 2

multiple features

$n \rightarrow$ # of features

$x^{(i)} \rightarrow$ é o i -ésimo vetor de features

$x_j^{(i)} \rightarrow$ é o j -ésimo elemento desse vetor.

Hypothesis function \rightarrow tenta prever, por ex., o preço de uma casa.

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

for convenience of notation, define $x_0^{(i)}$

\rightarrow pense nisso como definir uma zero feature adicional.

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

$n \times 1$

Feature
vector

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

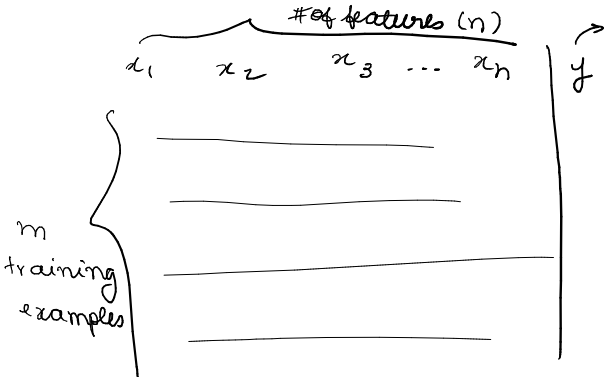
$n \times 1$

Parameters
vector

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

\hookrightarrow inner product

multivariate linear regression



Hypothesis function:

$$h_0(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

How do fit the parameters?

Como usar GD for linear regression w/
multiple features?

Cost function

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_0(x^{(i)}) - y^{(i)})^2$$

GD Descent:

Repeat {

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

} (simultaneously update $\forall j = 0, \dots, n$)

Gradient Descent

Previously ($n=1$):

Repeat {

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

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \underbrace{x_1^{(i)}}_{x_1}$$

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$):

Repeat {

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

(simultaneously update θ_j for $j = 0, \dots, n$)

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

$x_0^{(i)} = 1$

}

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

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

Andrew Ng

São coisas idênticas! $x_0 = 1$

2 features, x_1 e x_2

$x_0 = 1$

todos os
vetores $x^{(i)}$
da i^{a}
feature

Practical tricks for making GD work well

→ Feature Scaling — divide by the range of the input values

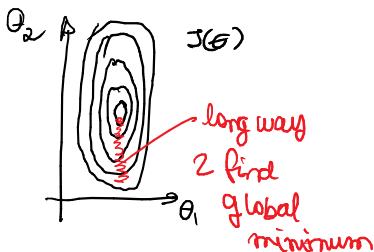
↳ se ve tem muitas features — certifique-se de que elas estão numa escala similar (assume ranges de valores similares)

↳ assim, GD converge + rápido

Ex: x_1 — assume valores de 0 — 2.000

x_2 — assume valores de 1 a 5

Al plots cost function $J(\theta)$:



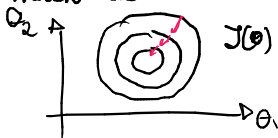
SOLUÇÃO

scaling the features

$$x_1 = \frac{\text{feature 1}}{2.000}$$

$$x_2 = \frac{\text{feature 2}}{5}$$

Contours much less skewed!



Get every feature into approx- $-1 \leq x_i \leq 1$ range.
 → Quase + ou - entre -1 e 1 , mas tudo
 bem se passar um pouco.

$$0 \leq x_1 \leq 3$$

$$-2 \leq x_2 \leq 0,5$$

$$\begin{cases} -100 \leq x_3 \leq 100 \rightarrow \text{VERY different values!} \\ -9,0001 \leq x_4 \leq 9,0001 \rightarrow \text{very different!} \end{cases}$$

features poorly scaled!

ñ precisa estar exatamente na mesma
 escala!

→ Mean normalization (in addition to
 feature scaling)

É bom p/ ter suas features perto da média
 zero.

(ñ se aplica à feature zero)

$$x_1 = \frac{\text{feature 1} - \text{média das features 1}}{2.000}$$

range -
 máx - min
 val val

$$-0,5 \leq x_1 \leq 0,5$$

$$-0,5 \leq x_2 \leq 0,5$$

x_1 $\xrightarrow[\text{by}]{\text{replaced}}$ $\frac{x_1 - \mu_1}{s_1}$ \rightarrow avg val of x_1 in training set
 \searrow max-min

Códigos do curso avançado std

∞
Std da variável

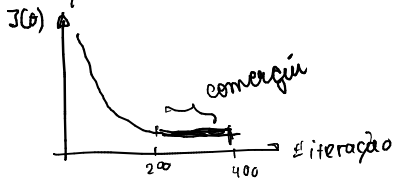
Runs much faster and converges in a lot fewer iterations.

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

debugging - como testar se o gradient descent tá funcionando bem.

- como escolher o α

- a fç do GD é achar o θ q minimiza $J(\theta)$
- Plota $J(\theta)$ de acordo com o # de iterações
- Se estiver funcionando bem, $J(\theta)$ deve cair p/ cada iteração



Automatic Convergence test

→ declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration

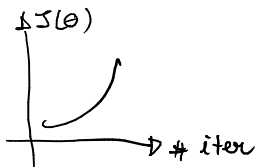
\approx
 ϵ

→ mas é mt difícil achar um ϵ plausível

↓
melhor plotar $J(\theta) \times \# \text{ iter.}$

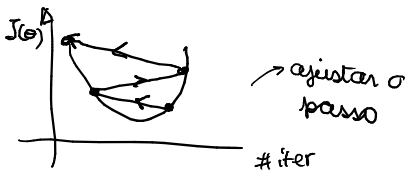
Cases ruins:

①



→ GD not working

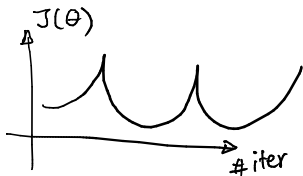
→ use smaller α



I'm getting worse and worse

→ use smaller α

②



α mt pequeno tb
demora um pouco
p/ convergir

→ Sugestão: tentar um range de valores p/α

$\alpha = 0,001$ $\begin{matrix} 1 \\ 0,003 \\ 0,01 \\ 0,03 \\ 0,1 \\ 0,3 \\ \vdots \end{matrix}$ } Plotar $J(\theta)$ $\times \#iter$

Pegue o valor de α
que parece fazer $J(\theta)$
cair + rapidamente.

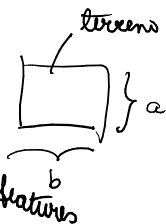
acho um valor bom $\frac{dJ}{d\alpha}$ aí tento

diminuir um pouco p

achar um bom learning rate p /
meu problema.

Choosing appropriate features

Suponha q tenho 2 features:
 a e b



→ posso criar minhas próprias features

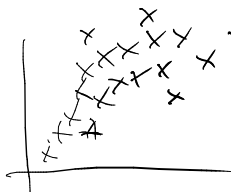
$$\underbrace{area}_{\downarrow}$$

nova feature

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

\downarrow
area

Regressão Polinomial



→ uma linha \bar{n} parece
ser um bom modelo.

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

$$\theta_0 + \theta_1 x_1 + \theta_2 x^2 + \theta_3 x^3$$

forma da h hipótese: (com novas features agora)

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \rightarrow \text{linear}$$

$$\theta_0 + \underbrace{\theta_1 (\text{size})}_{x_1} + \underbrace{\theta_2 (\text{size})^2}_{x_2} + \underbrace{\theta_3 (\text{size})^3}_{x_3}$$

Aplico o método do linear
regression e ajusto um polinômio.

feature scaling becomes increasingly important!

size: 1 - 1000

size²: 1 - 1000 000

size³: 1 - 1000 000 000

features
taking very
different
ranges!

Em vez de ir p/ um modelo cúbico por ex,
posso recombina as features e permaneco

Tem algoritmos q escolhem p/x que features
usar.

Normal Equation

- p/ alguns problemas de reg. linear,
é uma forma melhor de resolver p/o
valor ótimo dos parâmetros θ .
- antes eu tava usando o gradient descent.
minimizava cost function $J(\theta)$, podia
levar mtas iterações.
- normal equation: método p/ resolver p/
 θ analiticamente.

~~Algoritmo iterativo~~ → resolve p/o θ ótimo.

Continua ...