

Linear Regression

Multiple features (variables)

Single feature		
Size (ft ²)	Price (k)	
2104	460	Use only x ₁
1416	232	1 feature to predict output y

Size (ft ²) x ₁	Num of bedrooms x ₂	Num of floors x ₃	Age of home (years) x ₄	Price (k) y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178

m=4

Notation: n=4

n = num. of features n=4

$x^{(i)}$ = input (features) of i th training example

$x_j^{(i)}$ = value of feature j in i th training example $x_3^{(2)} = 2$

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

Previously: $h_\theta(x) = \theta_0 + \theta_1 x$ (single)

Multiple:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

Eg: $h_\theta(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$

Annotations:
 - 80: normal house price is 80k
 - 0.1: up abit with the num of bedrooms
 - 0.01: up abit more with num of floors
 - 3: up abit more with the age of home
 - -2: down abit more with the age of home

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

For convenience of notation, define $x_0 = 1$ ($x_0^{(i)} = 1$) - zero feature

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \theta^T x$$

(1 x (n+1) matrix) (row vector) x

Multivariate linear regression

Note: θ_0 as basic price of house, θ_1 as price per ft, θ_2 as price per floor, etc.
 x_1 be the number of square meters in house, x_2 be num of floor, etc.

Can use - For convenience purpose, assume $x_0^{(i)} = 1$ for $i \in \{1, \dots, m\}$. This allows us to do matrix operations with theta and x. Hence making 2 vectors θ and $x^{(i)}$ match each other element-wise (that is, having same num of elements: n+1).

Reason behind setting $x_0^{(i)} = 1$:

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \end{bmatrix}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

As a result, can calculate hypothesis as a vector with:

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

Gradient Descent for Multiple Variables

$$\text{Hypothesis: } h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ θ [n+1-dimensional vector]

Cost function:

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

Gradient descent:

Repeat {

$$n \geq 1 \quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

Let's say 3 features:

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

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

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

MATLAB/OCTAVE

vectorized form [cost function]

$$J(\theta) = \frac{1}{2m} (X\theta - \vec{y})^T (X\theta - \vec{y})$$

$$X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(m)})^T - \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

Gradient Descent in Practice I - Feature Scaling

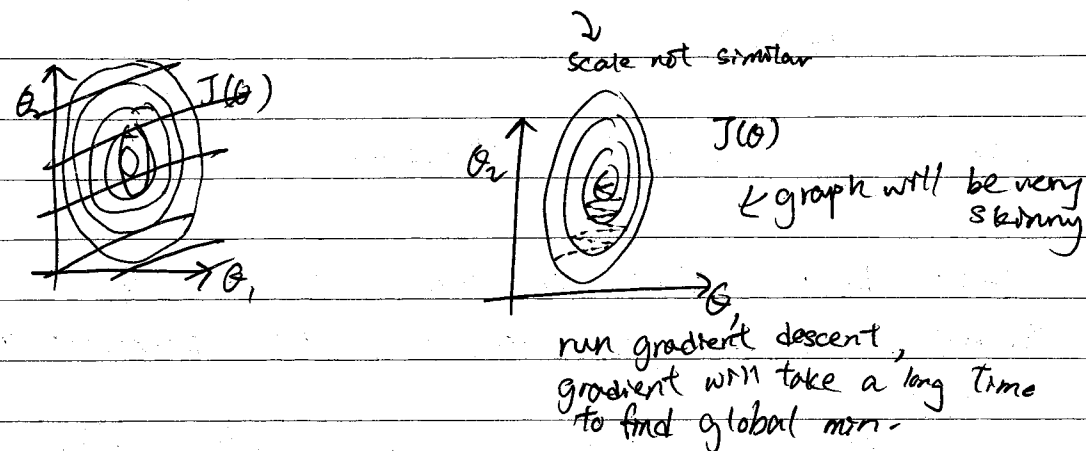
No. _____
Date _____

Feature Scaling

Idea: Make sure features are on a similar scale.

Eg. $x_1 = \text{size (0-2000 feet}^2\text{)}$

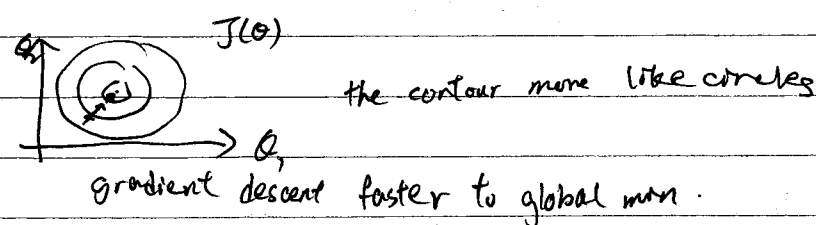
$x_2 = \text{number of bedrooms (1-5)}$



so want to make it better, make it

$$x_1 = \frac{\text{size (feet}^2\text{)}}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$



$$0 \leq x_i \leq 1$$

Feature Scaling

Get every feature into approximately a $-1 \leq x_i \leq 1$ range

$x_0 = 1$

still close to $-1 \leq x_i \leq 1$

$$0 \leq x_1 \leq 3 \checkmark$$

$$-2 \leq x_2 \leq 0.5 \checkmark$$

not close $\rightarrow -100 \leq x_3 \leq 100 \times$

$$-0.0001 \leq x_4 \leq 0.0001 \times$$

$$-3 \leq x \leq 3 \checkmark$$

$$-\frac{1}{3} \leq x \leq \frac{1}{3} \checkmark$$

Mean Normalization

Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$)

Eg. $x_1 = \frac{\text{size} - 1000}{2000} \rightarrow$ find, average size = 1000

$$x_2 = \frac{\text{\#bedrooms} - 2}{5} \rightarrow$$
 average 2 bedrooms

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

$$x_i \leftarrow \frac{x_i - \mu_i}{s_i}$$

μ_i : average value of x_i in training set
 s_i : range of value of x_i (max - min) (or standard deviation)
 \rightarrow become nearer to scale

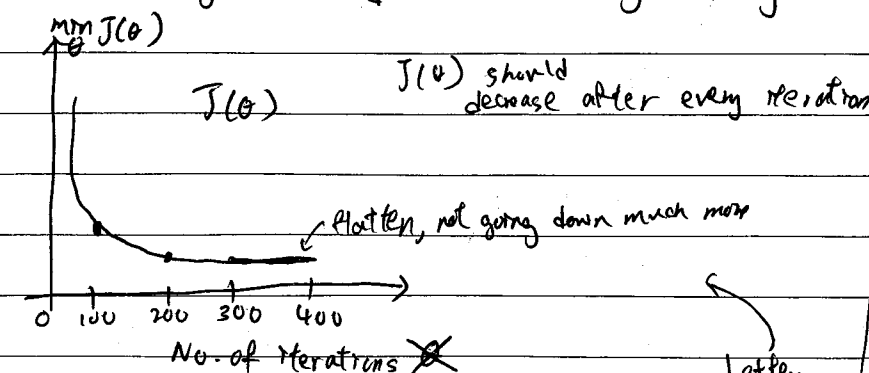
Gradient descent in practice II: Learning rate (α)

Gradient descent update

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

- Debugging: How to make sure gradient descent is working correctly
- How to choose learning rate α .

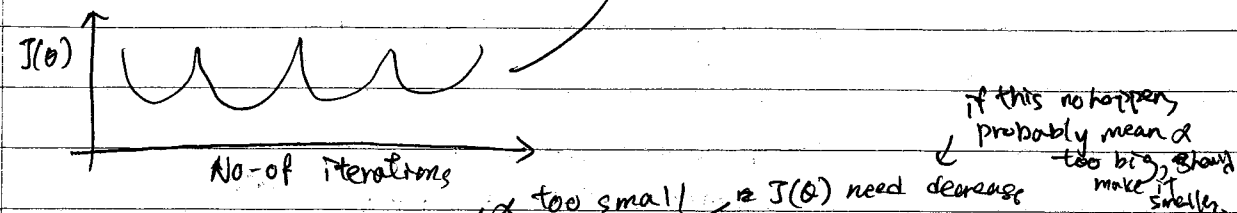
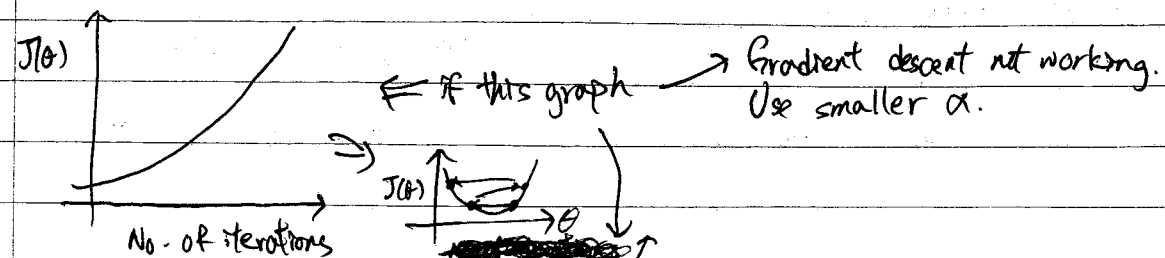
Make sure gradient descent is working correctly



Example automatic convergence test:
Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration.

better to check the gradient descent's converge than this

Making sure gradient descent is working correctly.



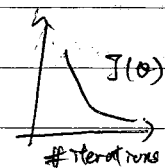
- For sufficiently small α , $J(\theta)$ should decrease on every iteration
- But if α is too small, gradient descent can be slow to converge.

Summary:

- if α is too small cause \rightarrow slow convergence (take a lot of iterations to reach global min)
- if α is too large: $J(\theta)$ may not decrease on every iteration, may not even converge. (slow convergence also possible)

plot $J(\theta)$ vs # iterations over many runs of gradient descent. If $J(\theta)$ ever increases, need to decrease α .

In order to debug all of these things, often plotting that $J(\theta)$ as a function of number of iterations can help you figure out what's going on.



To choose α , try running

0.001, 0.01, 0.1, 1
0.003 0.03 0.3

keep trying the values until found one value that's too small and make sure found one value that's too large

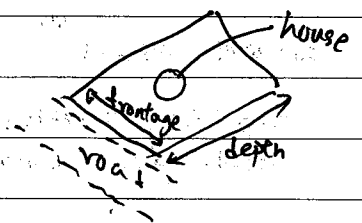
Linear Regression with multiple variables

Features and polynomial regression

Housing prices prediction 2 features

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$$

how wide the area of house



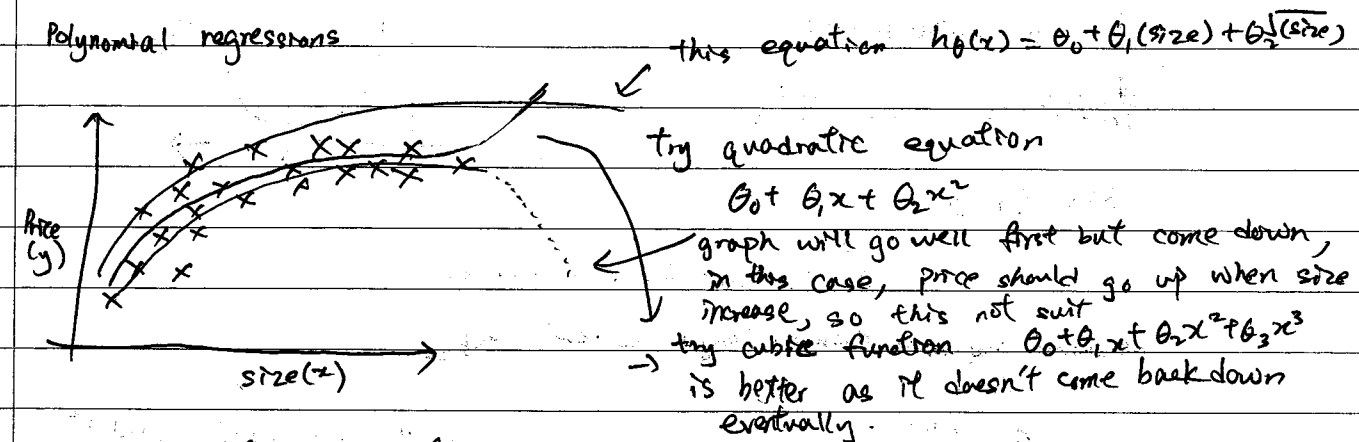
Area

new feature $\rightarrow x = \text{frontage} \times \text{depth}$

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

land area

Polynomial regressions



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$$

$$x_1 = \text{size}$$

$$x_2 = (\text{size})^2$$

$$x_3 = (\text{size})^3$$

$$\begin{array}{l} \text{size} = 1 - 1000 \\ \text{size}^2 = 1 - 1000,000 \\ \text{size}^3 = 1 - 10^9 \end{array}$$

Note: we can improve our features and form of our hypothesis function in a couple different ways.

- combine features into one. For example, combine x_1 and x_2 into a new feature x_3 by taking $x_1 \cdot x_2$.

Polynomial Regression

- our hypothesis function need not be linear (a straight line) if it does not fit the data well.
- change the behavior or curve of our hypothesis function by making it a quadratic, cubic or square root function (or any other form)
- important thing, if you choose features this way then feature scaling becomes very important.

Normal Equation

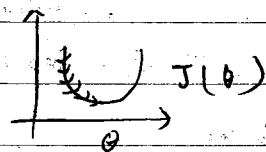
No. _____
Date _____

x_i = row
 y_j = column

No. _____
Date _____

Gradient Descent

- keep taking iterations to reach global minimum

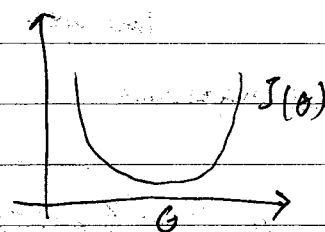


Normal equation: Method to solve for θ analytically.

Intuition: If $\theta \in \mathbb{R}$

$$J(\theta) = a\theta^2 + b\theta + c$$

To minimize $J(\theta)$ $\frac{d}{d\theta} J(\theta) = \dots = 0$
Solve for θ



↑
This is just for real number

we should deal with theta no longer a real number

$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$\frac{d}{d\theta} J(\theta) = \dots = 0$ (for every j)
Solve for $\theta_0, \theta_1, \dots, \theta_n$

How to minimize $J(\theta)$
- take partial derivative of J , with respect to every parameter of θ in turn, and then, to set all of these to 0

(7)

Examples: $m=4$

x_0	Size (ft ²) x_1	Num. of bedrooms x_2	Num. of floors x_3	Age of home (years) x_4	Price (\$1000) y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

Add 1 column and name as x_0

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$m \times (n+1)$ dimensional matrix

m -dimensional vector

$$\theta = (X^T X)^{-1} X^T y \quad (\text{minimize theta that minimize cost function})$$

Let's say

m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n features.

$$x^{(1)} = \begin{bmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \\ x_n^{(1)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$X = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix}$$

(Design matrix)

Eg. If $x^{(1)} = \begin{bmatrix} 1 \\ x_1^{(1)} \\ x_2^{(1)} \end{bmatrix}$ $x_0^{(1)} = 1$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$m \times 2$ vector y

Only 1 feature

(8)

Normal Equation and Noninvertibility

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

$(X^T X)^{-1}$ is inverse of matrix $X^T X$

let $A = X^T X$
 $(X^T X)^{-1} = A^{-1}$

← using this, it's ok not to use this (not necessary)

Octave: $\text{pinv}(X' * X) * X' * y$

$X' = X^T$ in Octave

min J(θ)
Feature scaling (gradient descent)
 $0 \leq x_1 \leq 1$
 $0 \leq x_2 \leq 1000$
 $0 \leq x_3 \leq 10^{-5}$

m training examples, n features.

Gradient Descent

- Need to choose α
- Needs many iterations

Disadvantage

→ you need to choose α , and keep running diff. α to see the best, and so that is sort of extra work and extra hassle
→ needs many more iterations, may make it slower

advantage → even have millions of features, you can run gradient descent and it still efficient, it will do something reasonable.

• Work well even when n is large

Normal Equation

- No need to choose α (convenient, simple to implement)
- Don't need to iterate
→ no need to plot J of θ or check convergence or take all those extra steps.

Disadvantage

- Need to compute $(X^T X)^{-1}$ $\frac{m \times (n+1) \times m}{\text{matrix}}$
- Slow if n is very large

n large
 $n = 10^6$

n small
 $n = 100$ ✓
 $n = 1000$ ✓
 $n = 10000$ actually a modern computer will work not that fast.

$$O(kn^2)$$

$$O(n^3)$$

need to calculate inverse of $X^T X$

With normal eqn, computing inversion has complexity $O(n^3)$. So if we have very large num of features, normal eqn will be slow. In practice, when n exceeds 10k, it might be a good time to go from a normal solution to an iterative process.

Normal equation

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

- What if $X^T X$ is non-invertible? (singular / degenerate)

- Octave: $\text{pinv}(X' * X) * X' * y$

non-invertible matrices

pinv (pseudo-inverse)
inv (inverse)

What if $X^T X$ is non-invertible?

2 common causes?

first • Redundant features (linearly dependent).

E.g. $x_1 = \text{size in feet}^2$, $x_2 = \text{size in m}^2$
 $x_1 = (3.28)^2 x_2$
related, 其中一个是另一个的倍数

second • Too many features (e.g. $m \leq n$).

- Delete some features or use regularization (later topic in this course)

example: $m = 10$
 $n = 100$

$$Q = R^{101}$$

sometimes work but not always be a good idea

Solutions to above problems include
deleting a feature that is linearly dependent with another or deleting one or more features when there are too many features.

(10)