

Notes - Week 2

Multiple features

house example continued:

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$/1000)	
x_1	x_2	x_3	x_4	y_1	output
2106	3	2	40		
1416	3	2	30		

Φ features.

n = number of features

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

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

index into training set.

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix} \in \mathbb{R}^4 = 4\text{th dimensional vector}$$

\mathbb{R}^n

$$x_3^{(2)} = 2$$

hypothesis in linear regression.

now with example (x)

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

$$= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$(x_0^{(i)} = 1)$ $x_0 = 1$

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

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

notes week 2 (multiple features intro)

$$h_0(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

\uparrow
= 1

$$= \Theta^T x$$

$$\Theta^T = [\theta_0 \ \theta_1 \ \dots \ \theta_n]$$

Θ^T $(n+1) \times 1$ matrix

or row vector...

multivariate linear regression

↳ (many variables) for which we try to predict the value Y .

Notes - Week 2

Gradient descent for Multiple Variables

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: θ $n+1$ dimensional vector

$$\hookrightarrow \theta_0, \theta_1, \dots, \theta_n$$

Cost function:

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

\therefore Gradient descent:

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

new algorithm for when features $(n) > 1$

Repeat

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

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

$$\boxed{x_0^{(i)} = 1} \quad \begin{matrix} ||| \\ \dots \end{matrix}$$

Notes - Week 2

Gradient Descent in Practice I

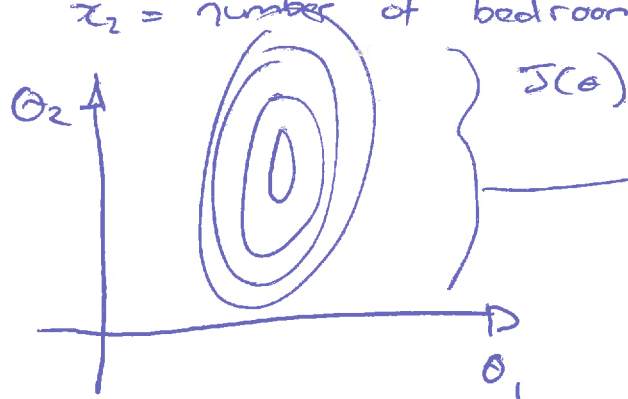
feature Scaling

Feature scaling

Idea: Make sure features are on a similar scale

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

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



skewed due to the difference in scale between Δ & \square

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

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

then we scale $J(\theta)$ so it looks like so



Feature scaling

to get every feature into approximately a $-1 \leq x_i \leq 1$ range.

Mean normalization

Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean (DO NOT Apply to $x_0 = 1$)

e.g. $x_1 = \frac{\text{size} - 1000}{2000}$
 $-0.5 \leq x_1 \leq 0.5$

$x_2 = \frac{\# \text{bedrooms} - 2}{5}$
 $-0.5 \leq x_2 \leq 0.5$

$$X_i = \frac{\mu}{s_i}$$

average value
of X_i in training set

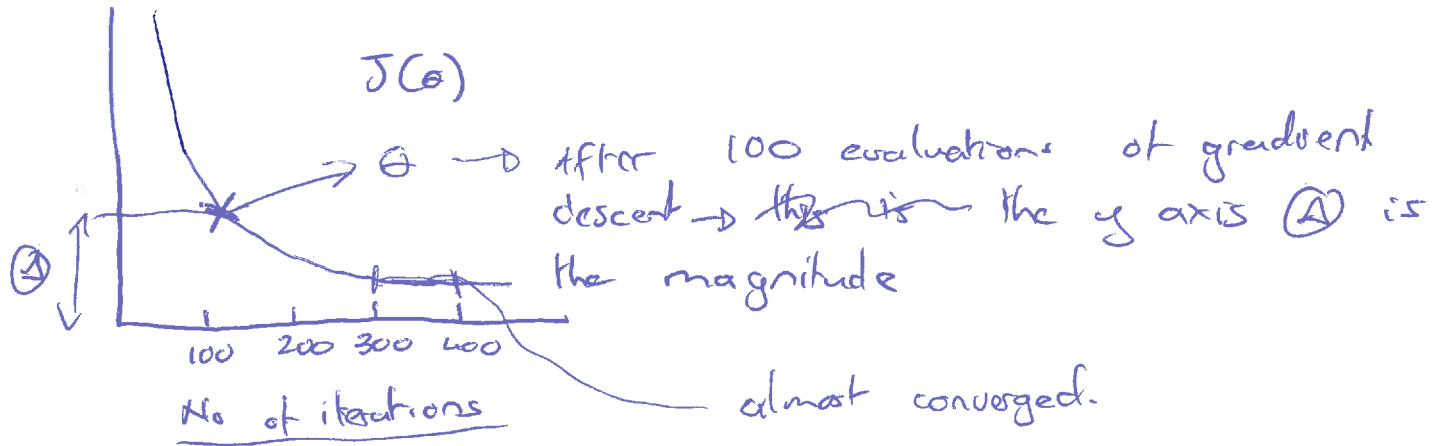
(max-min)
range of X_i (or standard
deviation.)

Notes - Week 2

Gradient descent in Practice II - learning rate

"Debugging" gradient descent

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

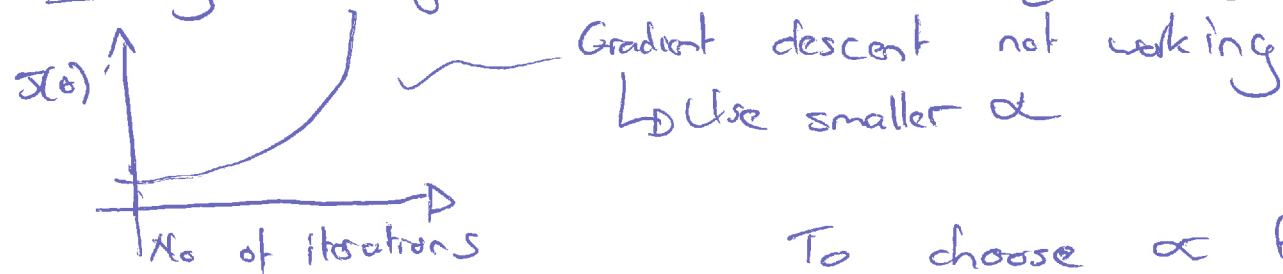


$\rightarrow J(\theta)$ should decrease after every iteration.

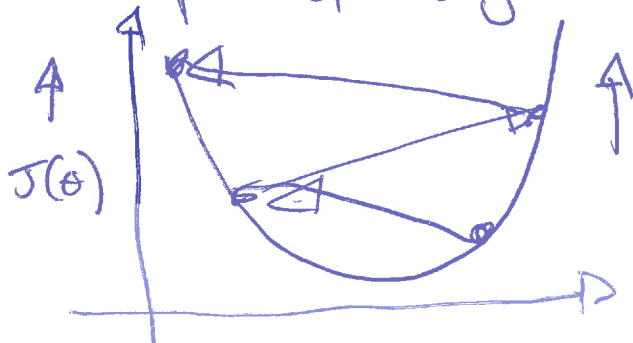
an example of an automatic convergence test:

\rightarrow "declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration" \checkmark good idea.

Making sure gradient descent is working correctly



example of big α .



To choose α try:

0.001, 0.01, 0.1, 1

$\downarrow 3\times$ $\downarrow 3\times$ \downarrow \downarrow

0.003 \downarrow 0.03 \downarrow 0.3 \downarrow 3.

Week 2 - Features and polynomial regression

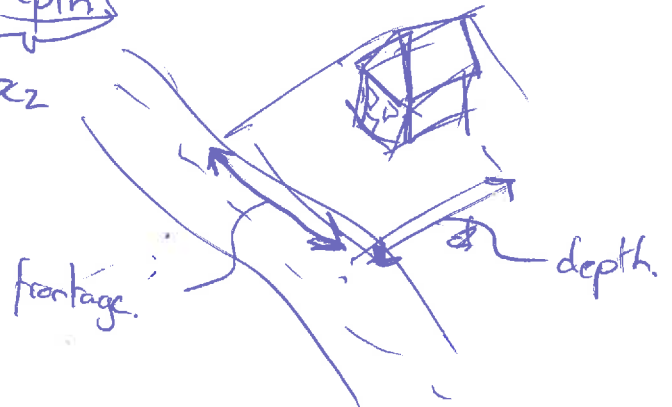
$$h_0(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

Area

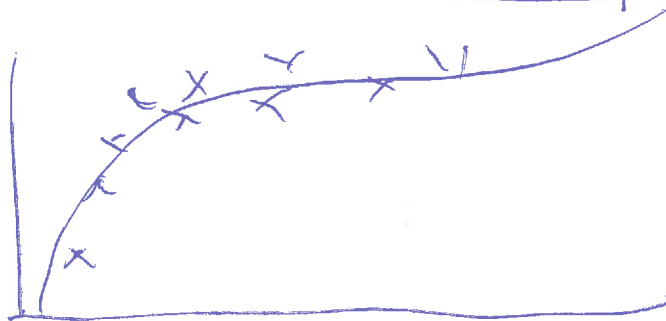
$$x = \text{frontage} \times \text{depth}$$

$$h_0(x) = \theta_0 + \theta_1 x$$

new feature. \uparrow Land area.



Polynomial regression



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

$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

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

$$x_1 = (\text{size})$$

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

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

$$\text{size} : \begin{array}{c} \text{house} \\ 1 - 1000 \end{array}$$

$$\text{size}^2 : 1 - 1000,000$$

$$\text{size}^3 : 1 - 10^9$$

Normal equation

examples: $m=4$

x_0	Size x_1	Number of bedrooms x_2	number of floors x_3	Age of home 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

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

$m \times (n+1)$
matrix

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

m -dimensional vector.

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

↑ value of θ that minimizes cost function

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

e.g. If $x^{(1)} = \begin{bmatrix} 1 \\ x_1^{(1)} \end{bmatrix}$

constructed
matrix
 ~~X~~

$$X = \begin{bmatrix} 1 & x_1^{(1)} \end{bmatrix}$$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

week 2
notes

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

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

or have $\text{pinv}(X^T X) \cdot X^T \cdot y$

$$\left| \begin{array}{l} \text{set } A = X^T X \\ (X^T X)^{-1} = A^{-1} \end{array} \right.$$

m training examples, n features

Gradient descent

- need to choose α
- needs many iterations
- Works well even when n is large.

Normal equation

- no need to choose α
- don't need to iterate.

- Need to compute

$$\left| (X^T X)^{-1} \right|_{n \times n}$$

= slow if n is very large.

$$O(n^3)$$

↑
cost to invert matrix.

$$n = 100 \quad \text{NE}$$

$$n = 1000 \quad \text{NG}$$

$$n = 10000 \quad \text{HG}$$

$$n > 10000 \quad \text{NE}$$

Week 2
notes

Normal equation non-invertibility

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

Redundant features (linearly dependent)

e.g. $x_1 = \text{size in feet}^2$

$x_2 = \text{size in m}^2$

$\therefore X^T X$ becomes non-invertible.

* too many features (e.g. $m \leq n$)

- delete some features, or use regularization.

do this

↓
1st
and
this

& 2nd