

Stanford ML

WEEK 2 :

Linear Regression with Multiple Variables

- Let $n \equiv \# \text{ of features}$
- $x^{(i)} \equiv \text{input features of } i\text{th training example}$
- $x_j^{(i)} \equiv \text{value of feature } j \text{ in } i\text{th training example.}$
- Now for linear regression with n features, our hypothesis becomes :

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

↳ The coefficients can be +VE or -VE and indicates how sensitive our hypothesis is to each feature x_i .

↳ Let $x_0 = 1$ ALWAYS $\therefore x_0^{(i)} = 1$ } ALLOWS COMPACT NOTATION

$$\Rightarrow x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}; \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

↳ Representing the input features for a particular example as an $n+1$ vector and doing same for the parameters we see that

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

MULTIVARIATE
LIN. REG.

Multivariate Gradient Descent

- Remember the cost function is defined as follows:
- $$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]^2$$
- ↳ where $\theta \in \mathbb{R}^{n+1}$

which is equivalent to:

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

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left[\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right]^2$$

Gradient descent:

Repeat {

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

}

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

New Algorithm

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

Repeat {

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

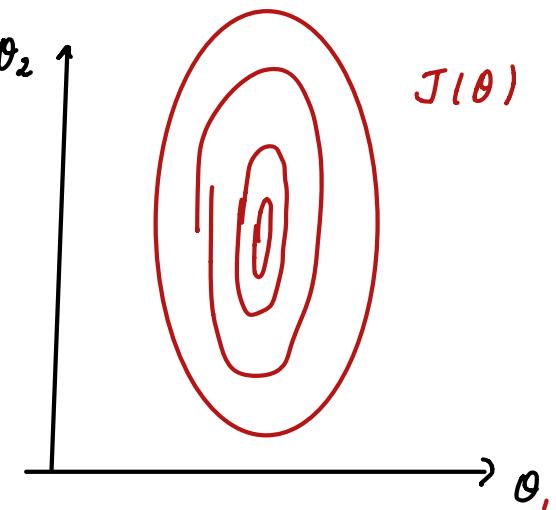
} ↳ simultaneously update θ_j for
 $j = 0, \dots, n$

Feature Scaling

- Trick used to ensure all features are on a similar scale to ensure optimal performance of the gradient descent algorithm

E.g. $x_1 = \text{size (0-2000 feet)}$

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



$J(\theta)$

- Having one feature taking on much larger creates LONG, SKINNY contours of $J(\theta)$
- \therefore Re-scale to ensure fast CONVERGENCE

$$\Rightarrow x_1 = \frac{\text{size (feet}^2)}{2000} \quad x_2 = \frac{\#\text{ of bedrooms}}{5}$$

- Using the new $x_1, x_2 \rightarrow$ more circular contours \rightarrow faster convergence
- Try to get every feature:

$$-1 \leq x_i \leq 1$$

note you can poorly scale features both too SMALL and too LARGE

OR

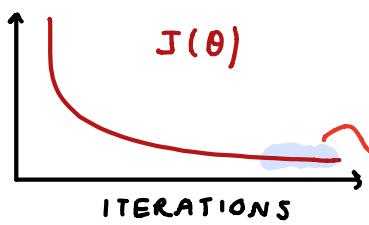
$$x_i' = \frac{x_i - \mu_i}{\sigma_i} \quad \left\{ \begin{array}{l} x_1' = \frac{\text{size} - 1000}{2000} \\ x_2' = \frac{\#\text{bedrooms} - 2}{5} \end{array} \right.$$

STDEV.

MEAN NORMALIZATION

Setting α and Debugging

- To ensure gradient descent we may wish to plot $J(\theta)$ as a function of time



We expect $J(\theta)$ to DECREASE

judging convergence

EVERY ITERATION

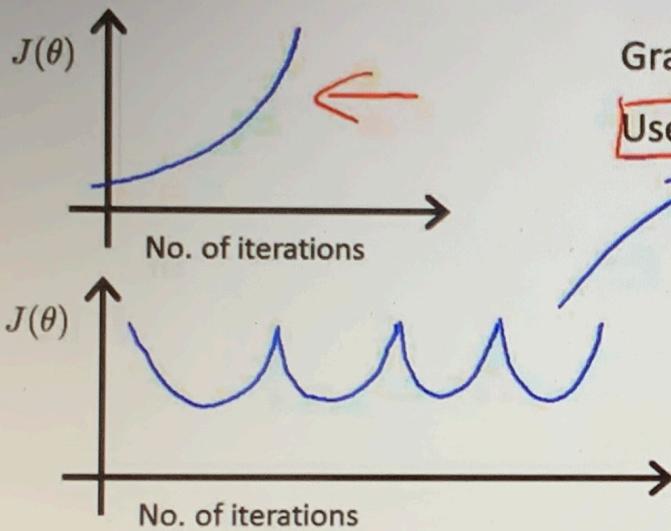
MATH PROVEN.

- Or we can simply use some tolerance (10^{-3}) and if $\Delta J(\theta) \leq 10^{-3} \Rightarrow \text{CONVERGED!}$

↳ TOL is HARD TO CHOOSE
 \therefore PLOT OFTEN BETTER

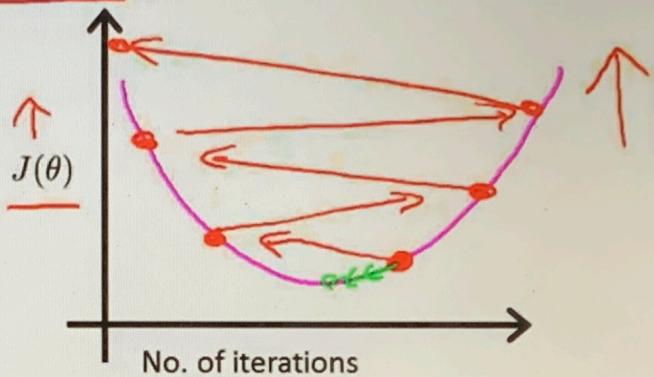
Common Errors

Making sure gradient descent is working correctly.



Gradient descent not working.

Use smaller α .



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

- α takes values

$$\alpha \in [10^{-2}, 1] \text{ in ...}$$

in steps of $\times 3$

Features and Polynomial Regression

Ex: $h_\theta(x) = \theta_0 + \theta_1 \times \text{frontage}(x_1) + \theta_2 \times \text{depth}(x_2)$

- Can make new feature:

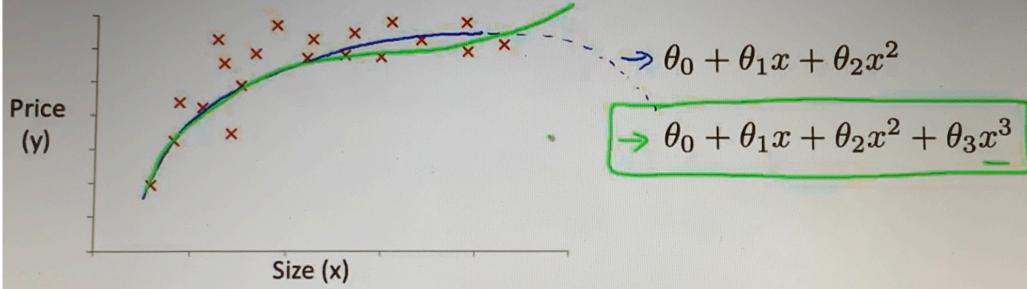
Area : $x = \text{frontage} \times \text{depth}$

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

↳ AREA



Polynomial regression



How do we fit higher order polynomials to data?

- We could

$$\begin{aligned} h_{\theta} &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \end{aligned}$$

↳ define features as so, then simply use linear regression

TRICK: HIDE THE POWERS IN THE DEFNS!

NOTE: Feature scaling is now CRITICAL for fast convergence

- Another reasonable choice:

$$h_{\theta}(x) = \theta_0 + \theta_1 (\text{size}) + \theta_2 (-\sqrt{\text{size}})$$



Normal Equation

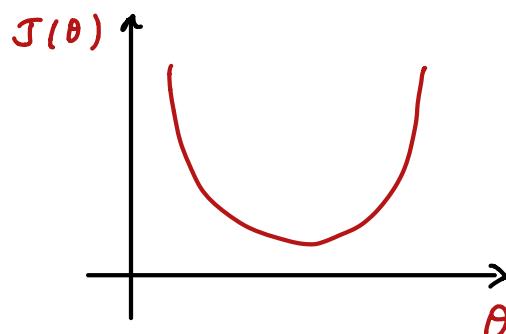
* Method to solve for optimal θ value ANALYTICALLY

INTUITION: Consider 1D case where $\theta \in \mathbb{R}$

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

$$\frac{dJ(\theta)}{d\theta} = 2a\theta + b = 0$$

$$\Rightarrow \theta = -\frac{b}{2a}$$



REALITY:

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

① $\Rightarrow \frac{\partial}{\partial \theta_j} J(\theta) \quad \forall j = 0, \dots, n \stackrel{\text{set}}{=} 0$

② \Rightarrow value for $\theta_0, \theta_1, \dots, \theta_n$

Ex:

$m \times (n+1)$

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

x_0	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	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178



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

m -dimensional vector

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

General Case:

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

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

X =
(design matrix)

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

• If using this method

$[m \times (n+1)]$

FEATURE SCALING NOT NEEDED

• If $n \geq 10^5$ maybe switch to GRADIENT DESCENT

m training examples, n features.

Gradient Descent

- • Need to choose α .
- • Needs many iterations.
- Works well even when n is large.

$$\nearrow n = 10^6$$

← -

Normal Equation

- • No need to choose α .
- • Don't need to iterate.
- Need to compute $(X^T X)^{-1}$ $\underset{n \times n}{\text{---}}$ $O(n^3)$
- Slow if n is very large.

$$n = 100$$

$$n = 1000$$

$$\dots n = 10000$$

Normal Equation Non-Invertibility

- What if :

$X^T X$ is non-invertible ?

↳ MATLAB has special methods to handle this so it always works (PSEUDO INVERSE)

CAUSES

- ① Redundant features $\begin{cases} \text{size in feet}^2 \\ \text{size in m}^2 \end{cases}$ } LINEARLY DEPENDENT
- ② Too many features e.g. ($m \leq n$)
 \therefore ~~use~~ some features, or use REGULARIZATION (easier)