

Model and Cost function

Not's :

m : no. of training ex.s

x 's : input var.s / features

y 's : output var.s / ~~features~~ target var. that we are trying to predict

(x, y) : one training ex. / observⁿ

$(x^{(i)}, y^{(i)})$: i th training ex. / observⁿ

$(x^{(i)}, y^{(i)})$; $i = 1, 2, \dots, m$: training set

i : index into the training set

X : space of I/P val.s

Y : space of O/P val.s

n : no. of features

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

θ 's : para. of model

h maps from x 's to y 's

Gradient Descent

now we have a hypoth. funⁿ and a way of measuring how well it fits into the data. Now we need to estimate the params in the hypoth. funⁿ. That's where gradient descent comes in.
(θ_0, θ_1)

Gradient Descent is a first-order iterative optimizⁿ algo. for finding a local min. of a differentiable funⁿ.

it tweaks the funⁿ's params iteratively to minimize a given funⁿ to its local min.

grad. desc. is a gen. algo. i.e., it applies not only just to cost funⁿ

Outline :

- start w. some θ_0, θ_1 (any val. say, $\theta_0 = 0, \theta_1 = 0.5$)
- keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until we hopefully end up at a min.

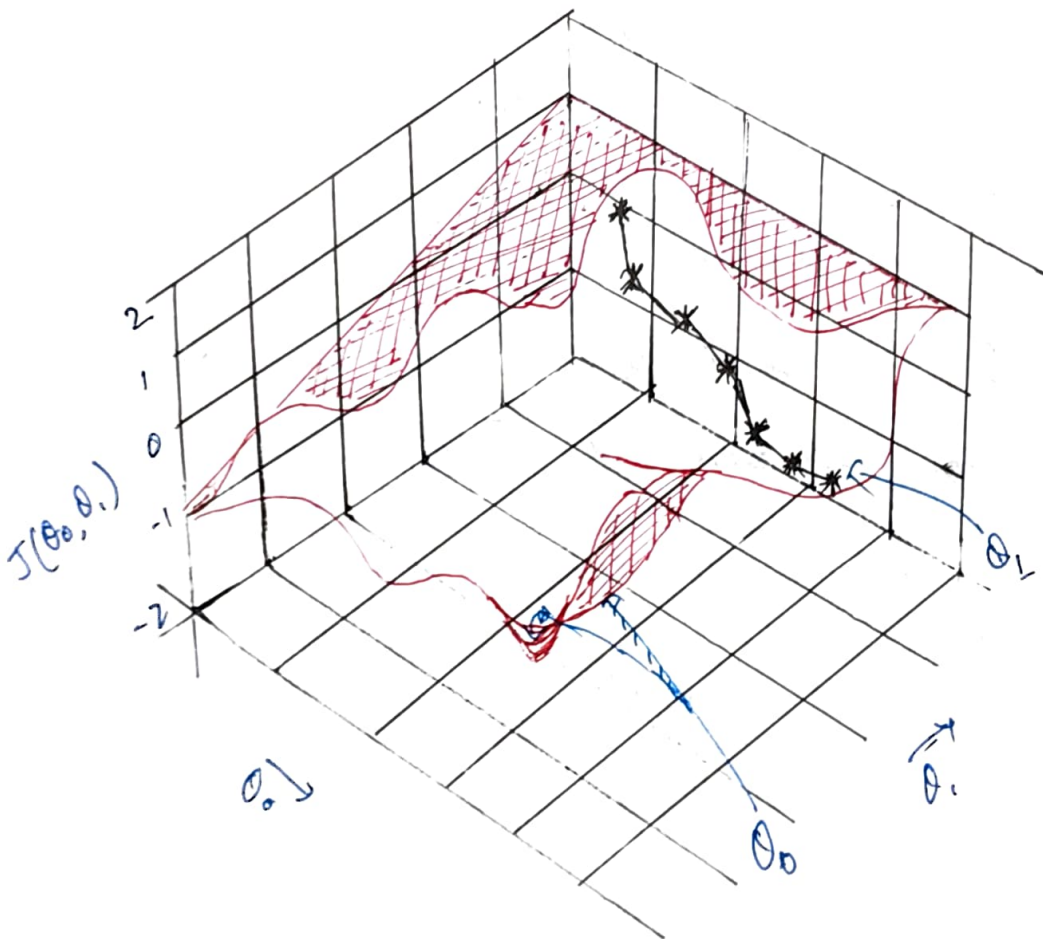
(context : a pt. on a slopy park w. hills like windows xp wallpaper) In grad. desc. what happens is that you look around yourself 360° and ask "if you were to take a step in some direcⁿ and get downhill asap, what direcⁿ do you take that step in?" You take that step and repeat until you've reached the (lowest) local optimum.

so, 2 things to notice are : direcⁿ of your step, size of your step.

Imagine that we graph our hypoth. funⁿ based on its fields θ_0 and θ_1 (actually we are graphing the cost funⁿ as a funⁿ of the para. estimates). We are not graphing x and y itself, but the para. range of our hypoth. funⁿ and the cost resulting from selecting a particular set of para.s.

We put θ_0 on the x -axis and θ_1 on the y -axis, w the cost funⁿ on the z -axis. The pts of our graph will be the result of the cost funⁿ using our cost funⁿ w those specific θ para.s

The graph below depicts such a setup



We will say that we have succeeded when the cost funⁿ is at the very bottom of the pits in our graph i.e., when its val. is the min. The blue arrow shows the min. pt.s in the graph.

The way we do this is by taking the derivative (the tangential line of a funⁿ) of our cost funⁿ. The slope of the ~~derivative~~ tangent is the derivative at that pt. & it will give us a direcⁿ to move towards. We make steps down the cost funⁿ in the direcⁿ of the steepest descent. The size of each step is determined by the para. α , wh. is called the Learning rate.

for ex., the dist. b/w each 'star' in the graph above represents a step determined by our para. ' α '. The smaller α would result in a smaller step and the larger α in a larger step.

Contd.
The direcⁿ in wh. the step is taken is determined by the partial derivative of $J(\theta_0, \theta_1)$.

Depending on where I starts graph, one could end up at diff. pt.s (local optima).

Gradient Descent Algo

repeat until convergence :

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

where

$j = 0, 1$ represents the feature index no.

$:=$ assignment operator

$a := b$ assign val. of b to a

$a := a + 1$ \neq

$=$ truth assertⁿ

$a = b$ assert that
val. of a & b is equal

$a = b + 1$ \times

α is called Learning Rate

it controls how big / small step we take downhill
w creating descent.

$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ is a partial derivative term.

it determines the slope of our descent

also in grad. desc., the parameters $\theta_1, \theta_2, \dots, \theta_n$ are updated simultaneously (& not 1st θ_0 then θ_1)

$$\# \text{ temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \quad \underline{\quad}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp1}$$

X

bc this will change the val. of θ_1 for that respective θ_0

Gradient Descent Intuition

In math, the partial derivative $\frac{\partial}{\partial x}$ of a funⁿ of several var.s is its derivative w.r.t. (only) 1 of those var.s, w the others held const. while in total derivative $\frac{d}{dx}$, all the var.s are allowed to vary.

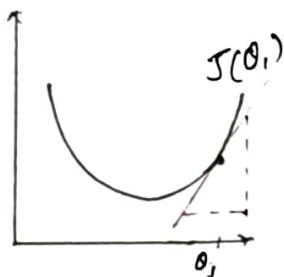
for intuⁿ, let us consider a funⁿ of only 2 para θ_1 .

Cost funⁿ : $\min_{\theta_1} J(\theta_1)$

Grad. desc : repeat until convergence :

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

then,

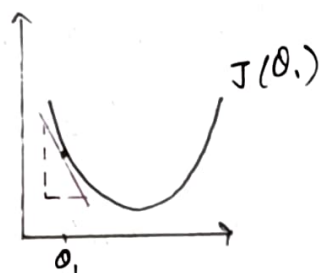


the slope of tangent is +ve

$$\text{so, } \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{some +ve no.})$$

or, θ_1 is moved towards left to obtain local optimal min.



here, the slope of tangent is -ve

$$\text{so, } \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

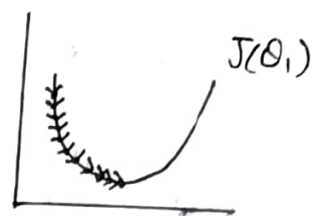
$$\theta_1 := \theta_1 - \alpha \cdot (\text{some -ve no.})$$

$$\theta_1 := \theta_1 + \alpha \cdot (\text{some +ve no.})$$

or, θ_1 moves towards right to obtain local optimal min.

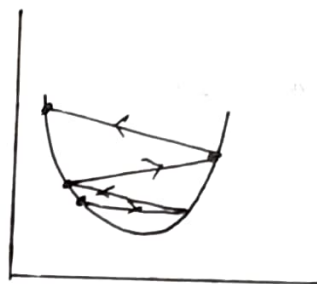
we see that regardless of the slope's sign for $\frac{d}{d\theta_1} J(\theta_1)$, θ_1 eventually converges to its min. val.

On a side note, we should ~~ensure~~ adjust 'a' to ensure that the grad. desc. converges in a reasonable time.



If α is too small, grad. desc. can be slow.

If α is too large, grad. desc. can overshoot the min. It may fail to converge or even diverge.



by def. local min. is when

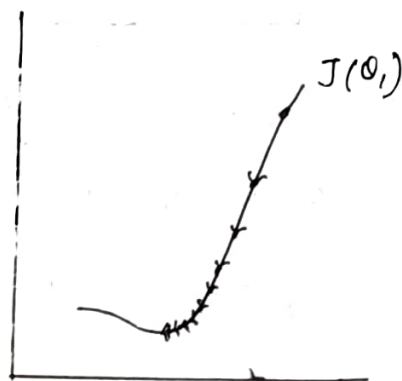
$$\frac{d}{d\theta_1} J(\theta_1) = 0 \text{ (slope} = 0\text{)}.$$

grad. desc. will converge even when α is fixed.

this is bc as we approach the local min, the slope gets steeper i.e., derivative automatically becomes smaller.

so, $\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ is smaller and

hence we do not need to dec. 'a'



Gradient Descent for Linear Regression

here we will put together grad. desc. w our cost funⁿ and that will give us an algo for Linear Regression, or putting a straight line to our data.

grad. desc. algo :

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j=0$ & $j=1$)

}

Linear Regression Model:

$$h_\theta(x) = \theta_0 + \theta_1 x \quad \text{Linear hypoth.}$$

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

squ. error cost funⁿ

Now, we need to calc. $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$

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

$$= \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^i - y^i)^2$$

so, for

$$j=0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \quad \text{(partial diff. wrt } \theta_0)$$

$$j=1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \cdot x^i$$

putting them in grad. desc. algo

repeat until convergence {

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

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

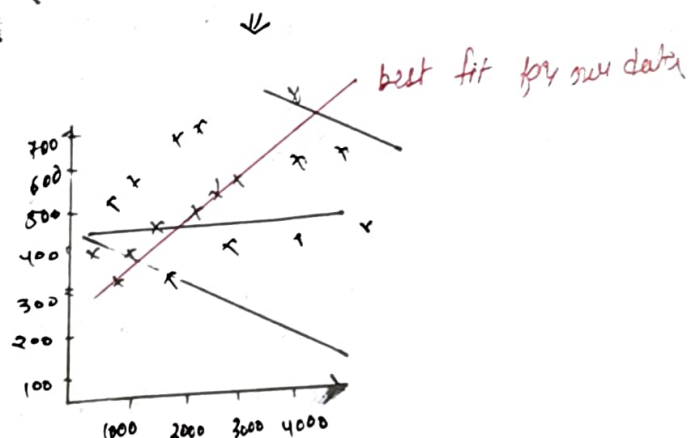
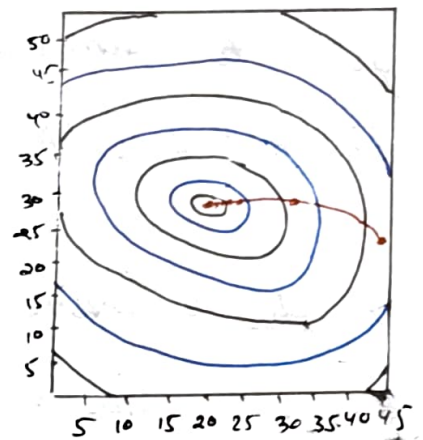
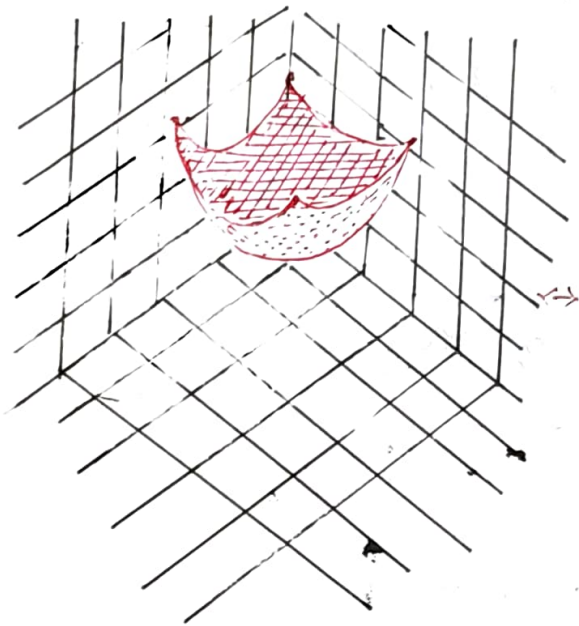
≠ This algo is also k/a Batch grad. desc. ^(fig) where Batch means that every step of grad. desc. uses all training ex.s $\sum_{i=1}^m (h_{\theta}(x^i) - y^i)$

However, some algo.s also use subsets of training ex.s and not all of them.

If we start w a guess for our hypoth. & then repeatedly apply grad. desc. eq's, our hypoth. will become more accurate.

* the cost funⁿ for Linear Regression is always going to be a convex funⁿ (bow-shaped funⁿ) and so, it doesn't have the local optimum except the one global optimum. so, the grad. desc. will always converge to the 1 local optimum.

$J(\theta_0, \theta_1)$ funⁿ of params θ_0, θ_1



Gradient Descent for Multiple Variables

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

Params : $\theta_0, \theta_1, \theta_2, \dots, \theta_n = \theta$

$\theta = [\theta_0, \dots, \theta_n]$ is an $n+1$ dimension vec.

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

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

Gradient Descent

for single feature / var,

for $n=1$

Repeat

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

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

(simultaneously update θ_0, θ_1)

}

The gradient descent equation itself is generally the same form, we just have to repeat it for our 'n' features.

$n \geq 1$.

repeat until convergence

{

$$\theta_j := \theta_j^0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(simultaneously update θ_j for $j = 0, 1, \dots, n$)

}

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

sim. to prev θ_0 (for $n=1$) as

for convenience, we have assumed $x_0 = 1$

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

sim to prev θ_1 (for $n=1$)

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

Gradient Descent in Practice 1: Feature Scaling

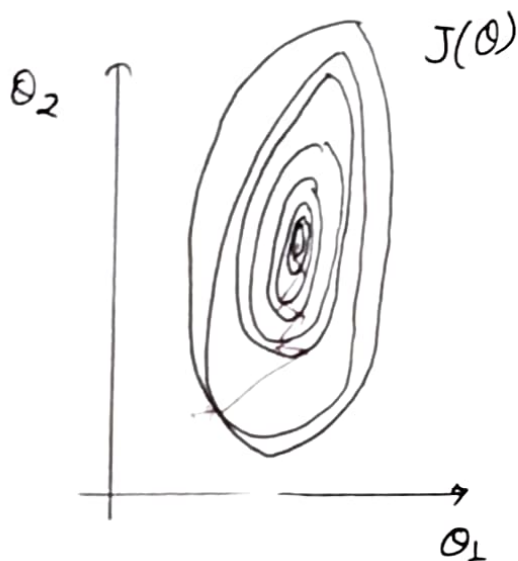
* ~~total~~ Feature scaling can make grad. desc. run much faster and converge in a lot fewer iter's.

Idea: Put diff. features on a similar scale of val's

ex:

$$x_1 = \text{size (0-2000 ft}^2\text{)}$$

$$x_2 = \text{no. of bedrooms (1-5)}$$

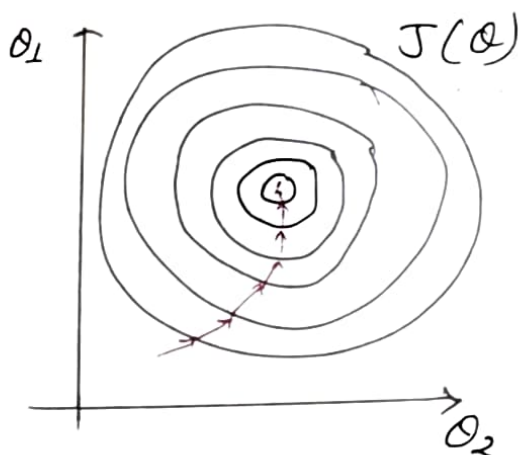


Here, contours of $J(\theta)$ are tall/skiny/skewed ellipses. On running grad. desc., grad. desc. will oscillate back & forth until it gets to the local min. thus taking long time. to reach the local min.

in these settings, a useful thing to do is scale the features

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

$$x_2 = \frac{\text{no. of bedrooms}}{5}$$



on scaling, contours of $J(\theta)$ become much less skewed looking more like circles. On running grad desc., grad desc. takes a much more directed path rather than a convoluted path to get to local min. so, by scaling the features so that there are consistent range of val's (here $0 \leq x_1/x_2 \leq 1$), implementⁿ of grad. desc. converges much faster.

We can speed up grad. desc. by having each of our 1/p vars in roughly same range.

This is bc θ will descend quickly on small ranges and slowly on large ranges, and so will oscillate inefficiently down to the optimum when the vars are very uneven. (as shown in ~~fig~~ left fig)

The way to prevent this is to modify the ranges of our 1/p vars so that they are all roughly the same. Ideally:

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

\nrightarrow not necessarily this range, but a range in gen. wh is scalable / compatible / not too big or small to plot

ex:

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

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

$$-100 \leq x_3 \leq 100 \quad \times$$

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

\nrightarrow It's fine if features are not exactly in the same range (or scale) as long as they are close enough (to the grad. desc.)

The goal is to get all 1/p vars into roughly 1 of these ranges, give or take a few.

Feature Scaling involves dividing the $1/p$ vars by the range (max. val - min val) of the $1/p$ var, resulting in a new range of just 1.

Mean Normalization.

Replace x_i w $x_i - \mu_i$ to make features have approx. ly zero mean

Implementing both feature scaling & mean normalization:

$$x_i^o := \frac{x_i - \mu_i}{s_i}$$

where, μ_i = avg. val of feature (i)

s_i = range of vals (max - min)

or standard deviation for feature (i)

so, here,

$$x_1 \leftarrow \frac{x_1 - \mu_1}{s_1} = \frac{\text{size} - 1000}{2000} \quad (\text{say}) \rightarrow -0.5 \leq x_1 \leq 0.5$$

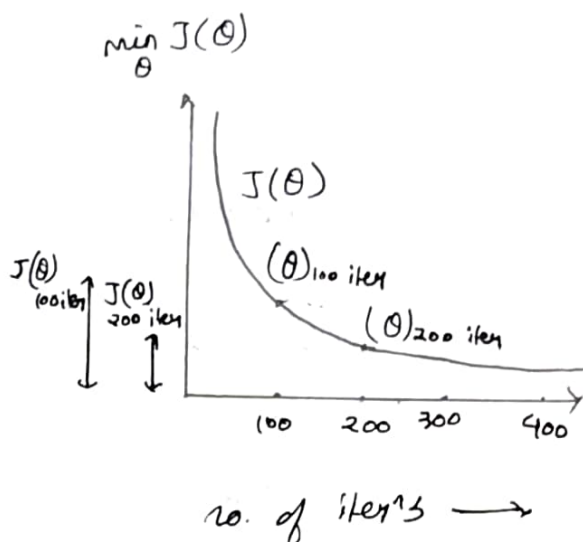
$$x_2 \leftarrow \frac{\text{no. of bedrooms} - 2}{54} \quad (\text{say}) \rightarrow -0.5 \leq x_2 \leq 0.75$$

Gradient Descent in Practice 2: Learning Rate (α)

terminology: Debugging + how to make sure grad. desc. is working correctly.

→ how to choose learning rate α ?

→ If grad. desc. is working correctly, $J(\theta)$ should dec. after every iterⁿ.



→ When $J(\theta)$'s val. doesn't go down much after iter's (i.e. when the curves start flattening), it is safe to assume that grad. desc. has converged bc the cost funⁿ isn't going down much.

→ The no. of iter's grad. desc. takes to converge for a phys. al applicⁿ can vary a lot (maybe 30 or 30000 ~~times~~ iter's) & it's difficult to predict how many iter's grad. desc. would need to converge.

So it is often done w plotting such a graph b/w cost funⁿ as we inc. no. of iter's.

→ Debugging gradient descent:

Make a plot w no. of iter's on x-axis. Now plot the cost function, $J(\theta)$ over the no. of iter's of gradient descent.

If $J(\theta)$ ever inc.s, then you probably need to dec. α .

→ Automatic Convergence Test:

Declare convergence if $J(\theta)$ decreases by less than E in 1 iterⁿ, where E is some small val. (such as 10^{-3} or ϵ).

However in practice it's difficult to choose this threshold val.

this algo tells us if grad. desc. has converged.

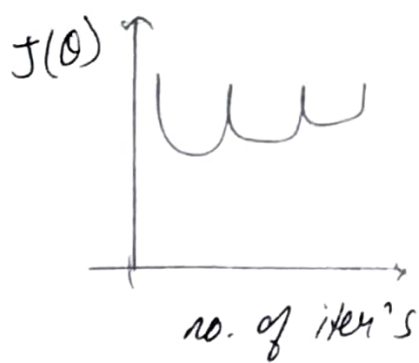
→ Making sure grad. desc. is working correctly:



a plot like this clearly suggests grad. desc. isn't working correctly.

This is bc α is too large & it overshoots the min.

(dec. val. of α)

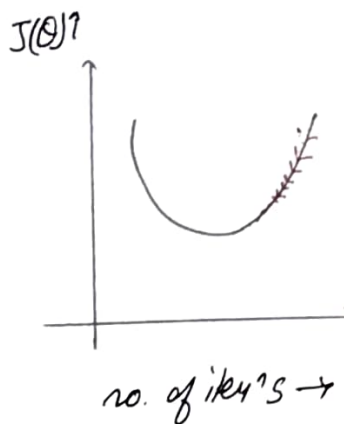


a plot like this is also bc α is too large.

(dec. val. of α)

• for sufficiently small α , $J(\theta)$ should dec w/ "iter".

but if α is too small, grad-des. can be too slow to converge.



→ Summary

- If α is too small : slow convergence
 - If α is too large : $J(\theta)$ may not dec. on "iter"; it may not even converge.
- (slow converge is also possible but usually not the case).

To choose α , try

(plotting for following val.s of α until best fit is found).

..., 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ...

Vectorization

Vectorization is basically the art of getting rid of explicit 'for loops' from code.

a simple ex. would be multiplying 2 matrices using $*$ operator instead of looping through all el's

$$\text{ex: } h_0(x) = \sum_{j=0}^1 \theta_j^0 x_j^0 = \theta^T x$$

$$\theta = \begin{bmatrix} \theta_0^0 \\ \theta_1^0 \\ \theta_2^0 \end{bmatrix}; x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Unvectorized Implementation

prediction = 0.0;

for j = 1:n+1,

prediction = prediction

+ theta(j) * x(j);

end;

Vectorized Implementation

prediction = theta' * x;

Grad. Desc.

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

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

⋮

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

vectorized Implementⁿ:

$$\begin{array}{ccccc} \theta & = & \theta & - & \alpha \delta \\ \uparrow & & \uparrow & & \uparrow \\ \in \mathbb{R}^{n+1} & & \in \mathbb{R}^{n+1} & & \in \mathbb{R}^{n+1} \\ & & \text{(is a no.)} & & \end{array}$$

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \end{bmatrix}$$

$$\text{where, } \delta = \frac{1}{m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}) - y^{(i)})}_{\substack{\in \mathbb{R}^n \\ \text{no.}}} \underbrace{x^{(i)}}_{\substack{\text{vec.} \\ \in \mathbb{R}^{n+1}}}$$

$$\delta = \begin{bmatrix} \delta_0 \\ \delta_1 \\ \delta_2 \end{bmatrix}$$

$$\delta_0 = \frac{1}{m} \sum (h(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)} \quad \left| \quad \delta = (h(x^{(1)}) - y^{(1)}) x^{(1)} \right.$$

$$+ (h(x^{(2)}) - y^{(2)}) x^{(2)}$$

$$+ (h(x^{(3)}) - y^{(3)}) x^{(3)}$$