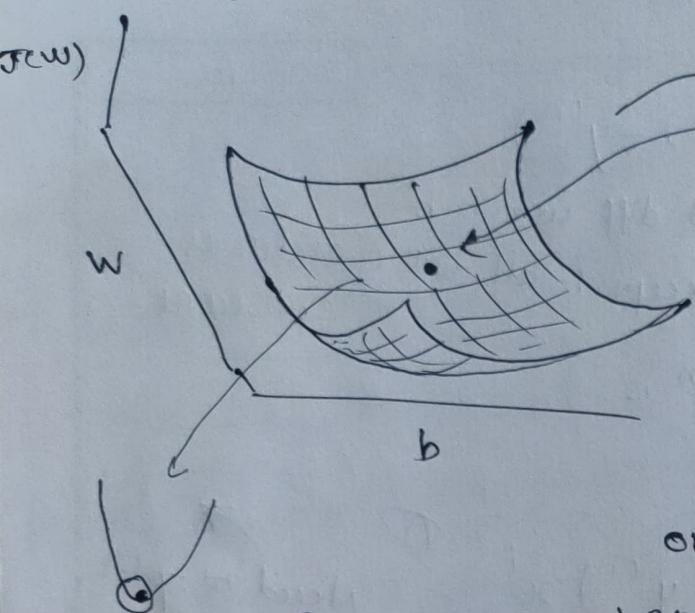


If we go for squared error cost



For this we have only one Global minima ~~at~~
We don't have local minima here.
(Technically, this cost function is convex function)

a bowl shaped function and it don't have local minima & it have only one global minimum

So, when you implement gradient descent on convex function, one nice property is that so long as your learning rate is chosen appropriately, it will always converge to global minimum. ~~(*)~~

Batch Gradient descent

Batch: Each step of gradient descent uses all training examples

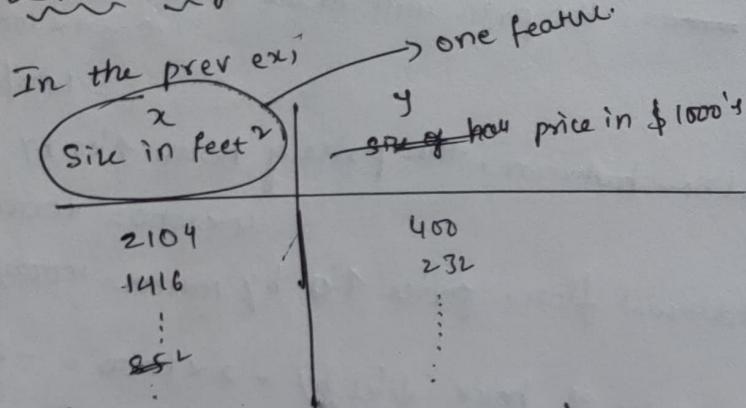
	x size in feet ²	y price in £1000's
1	2104	400
2	1592	232
3	1	?
:	1	?
47	3210	870

$m=47$
 $\sum_{j=1}^{47} (f_{w,b}(x^{(j)}) - y^{(j)})^2$

Week 2 :

MULTIPLE FEATURES: (Variables)

Linear Regⁿ With Multiple Features:



Here we use only one feature to predict price in 1000\$ so it is Linear Regⁿ with one feature Variable predict

Multiple feature means here we use multiple features to find the price of the house

size in feet x_1	No. of bedrooms x_2	No. of floors x_3	Age of home in year x_4	Price (\$) in \$1000's (y)
2104	5	1	45	460
1416	3	2	40	232
;	;	;	30	315
;	;	;	;	;

features of independent feature

Label \$1000's
Dependent feature

$x_j = j^{\text{th}}$ feature (x_1, x_2, x_3, \dots) , $j=1 \text{ to } 4$

$n = \text{no. of features}$ ($n=4$ we have 4 features & 1 o/p)

$x_i^j = \text{feature of } i^{\text{th}}$ example

$x^j_{\text{vector}} = [x_1^j \ x_2^j \ \dots \ x_n^j] \rightarrow \text{row vector}$

$x^j = [1416 \ 3 \ 2 \ 40] \rightarrow \text{row vector}$

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

$\text{ex: } x_3^{(2)} \Rightarrow 2$
In 3rd feature, 2nd element

Previously, $f_{w,b}(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$

But Now, $f_{w,b}(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$

$$f_{w,b}(x) = 0.1 x_1 + 4 x_2 + 10 x_3 + (-2) x_4 + 80$$

Let assume
price of house
is 80000\$

without
Bedroom, floor area, size

By every square foot ~~x~~, the price will increase by 10%. i.e. $0.1 \times 1000 = 100\$$

for every addition bathroom, the price of house increases by $4 \times 1000 = 4000\$$

for each additional floor, price increases by $10 \times 1000 = 10000\$$

For each every year, the price of house will decrease by $-2 \times 1000 = -2000\$$

If we have 'n' features then our model will looks like

$$f_{w,b}(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b$$

vector $\vec{w} = [w_1, w_2, w_3, \dots, w_n]$ → w ' vector
 b is a number } parameters of model
 $\vec{x} = [x_1, x_2, x_3, \dots, x_n]$

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b \Rightarrow [w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b]$$

dot product
of
two vectors

THIS IS MULTIPLE LINEAR REGRESSION

To Implement Multiple linear Reg'n we need to know neat trick called Vectirization.

Used for other algo as well.

Vectorization part 1:

When we use vectorization, the code will be shorter & efficient to run.

parameters and features

$$\vec{w} = [w_1, w_2, w_3] \quad n=3$$

b is a Number

$$\vec{x} = [x_1, x_2, x_3]$$

In Linear algebra : count starts from 1

$$w = \text{np.array } ([w[0], w[1], w[2]]) \quad \left. \right\} \text{In Python index starts from '0'}$$

$$b = 4$$

$$x = \text{np.array } ([10, 20, 30])$$

Without Vectorization:

$$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$\underbrace{1, 2, 3}_{n=3}$ No. of features

$$f = w[0] * x[0] + w[1] * x[1] + w[2] * x[2] + b$$

But if $n=100$ then it is hard to write code so, here we need to use
if we write $f=100$ times then
it is inefficient to run &
Bad for compute to compute

Vectorization to
run the code efficiently

$$\sum f_{\vec{w}, b}(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \rightarrow \text{without Vectorization}$$

$$f = 0$$

for j in range (0, ~~b~~ n): \rightarrow it is better than but it
is also inefficient

$$f = f + w[j] * x[j]$$

$$f = f + b$$

Using Vectirization:

$$f_{\vec{w}, b}(\vec{x}) = \underbrace{\vec{w} \cdot \vec{x}}_{\text{dot prod}} + b$$

To know about dot product
Refer that ^{vector} calculate both
once &

----> Using Vectirization.
(faster, efficient)

Vectirization Part 2:

without Vectirization the computer takes lot of time to execute stmts if n is high
But if we use " " the code will run efficiently.

without Vectirization

for j in range(0,16):

$$f = f + w[j] * x[j]$$

At t₀ to

$$f + w[0] * x[0]$$

at t₁

$$f + w[1] * x[1]$$

:

at t₁₅

$$f + w[15] * x[15]$$

Vectirization

$$\text{np.dot}(w, x)$$

to

$$\boxed{w[0] \boxed{w[1]} \dots \boxed{w[15]}}$$

$$\begin{matrix} * & * & * \\ \boxed{x[0]} & \boxed{x[1]} & \dots & \boxed{x[15]} \end{matrix}$$

At t₁

$$w[0] * x[0] + w[1] * x[1] + \dots + w[15] * x[15]$$

Codes with Vectirization will do calculation in less time as compare to without Vectirization.
This matters more when you running algorithms on large datasets or trying to train large models.

Gradient Descent

Suppose, $\vec{w} = (w_1, w_2, \dots, w_{16})$

(ignoring) for this case
b

$$\vec{d} = (d_1, d_2, \dots, d_{16}) \rightarrow \text{derivative}$$

$$\left\{ \begin{array}{l} w = np.array([0.5, 1.3, \dots, 3.4]) \\ d = np.array([0.3, 0.2, \dots, 0.4]) \end{array} \right.$$

$$\text{Compute } w_j = w_j - 0.1 d_j \text{ for } j = 1, 2, \dots, 16$$

↓
Assumed $\alpha = 0.1$

Without Vectorization

$$\begin{aligned} w_1 &= w_1 - 0.1 d_1 \\ w_2 &= w_2 - 0.1 d_2 \\ &\vdots \\ w_{16} &= w_{16} - 0.1 d_{16} \end{aligned}$$

python:

```
for j in range(0, 16):
    w[j] = w[j] - 0.1 * d[j]
```

With Vectorization:

$$\vec{w} = \vec{w} - 0.1 \vec{d}$$

~~$w = w - 0.1 * d$~~

$w = [0.5, 1.3, \dots, 3.4] \quad 0.1 \times 16$
 $0.1 * d = [0.3, 0.2, \dots, 0.4] \quad 0.1 \times 16$
 $= [0.2, 1.1, \dots, 2.4] \quad 0.1 \times 16$

GRADIENT DESCENT FOR MULTIPLE LINEAR REGRESSION:

previous Notation

parameters

$w_1, w_2, w_3, \dots, w_n$
 b

Model $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$

cost function $J(w_1, \dots, w_n, b)$

Vector Notation

88

→ vector of length n

$\vec{w} = [w_1, \dots, w_n]$

$b \rightarrow$ number

$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
↑ dot product

Gradient

Descent

repeat until convergence {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_1, \dots, w_n, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w_1, \dots, w_n, b)$$

repeat until convergence {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

Gradient Descent

One feature

(we saw about this already)

repeat {

$$w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w, b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$\hookrightarrow \frac{\partial}{\partial w} J(w, b)$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w, b}(x^{(i)}) - y^{(i)})$$

$\hookrightarrow \frac{\partial}{\partial b} J(w, b)$

Simultaneously update w, b

* *

n features ($n \geq 2$)

repeat {

$$w_j = w_j - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}$$

$\hookrightarrow \frac{\partial}{\partial w_j} J(\vec{w}, b)$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$$

Simultaneously update

w_j (for $j=1, \dots, n$) and b

here $x^{(i)}, x_1^{(i)}, \dots, x_n^{(i)}$ as vectors

{ THIS IS THE GRADIENT DESCENT }

→ Next we will see

"An alternative to gradient descent".

An alternative to Gradient Descent:

An alternate way to find w and b for Linear regression. This method is called Normal equation.

Normal equation: We know gradient descent is best algo for minimizing cost function $J(\theta)$ to find w and b there is one other algorithm that works only for Linear regression and pretty much none of other algorithms you see in this specialization for solving w & b . and this other method does not need any iterative gradient descent algorithm. called Normal equation method, it turns out to be possible to use an advanced linear algebra library to just solve for w and b all in one go! without iterations

Some Disadvantages of normal equation method are

- First unlike gradient descent, this is not generalized to other learning algorithms such as logistic reg'n algorithm
- Slow when the no. of features are large ($> 10,000$)

(Just know some of machine learning libraries that implement linear reg'n by using Normal equation method).

(But Gradient descent is the recommended method for finding parameters w & b)

Gradient Descent in Practice:

Feature Scaling Part 1:

Feature scaling is enable gradient descent to run much faster.

Let's start by taking a look at the relationship b/w the size of a feature that is how big are the numbers for that feature and the size of its associated parameter.

As a concrete example, let's predict the price of house using two features x_1 and size of the house and x_2 the no. of bedrooms. Let say x_1 range is (300 - 2000 ft²) & x_2 range is (0 - 5)

$$\text{price} = w_1 x_1 + w_2 x_2 + b$$

↓ ↓
 estimated size # of bedrooms

x_1 : size (feet²)
range: 300 - 2000

x_2 : # bedrooms
range: 0 - 5

for this example x_1 take very large range of values & x_2 takes on relatively small range of values.

(x_1 takes large range of values 300 - 2000 ft² whereas x_2 takes on small range of values 0 - 5).

Now let's take an example, size of house = 2000 feet², No. of bedrooms = 5

& price = \$500k

What is the reasonable w_1, w_2 values for this ex?

Let's look at one possible set of parameters

say $w_1 = 50, w_2 = 0.1, b = 50$

$$\begin{aligned} \text{estimated price} &= 50 * 2000 + 0.1 * 5 + 50 \\ &= 100050.5 \text{K} \end{aligned}$$

→ it is very far away from \$ 500k.

{ So, these w_1, w_2 are not a good values (prediction)

Let's look for another possibility

$$w_1 = 0.1 \text{ (small)}, w_2 = 50, b = 50$$

$$\hat{\text{price}} = 0.1 * 2000 + 50 * 5 + 50$$

$\hat{\text{price}} = \$500 \text{K} \rightarrow$ predicted price matching with actual value } More reasonable

How does it relate to gradient descent?

Feature size and parameter size

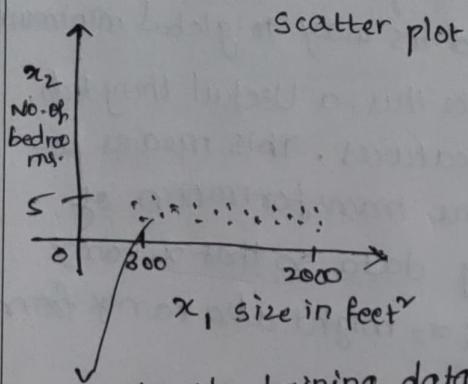
	Size of feature x_j	Size of parameters w_j
Size in feet?	$\xleftarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}}$	$\xleftarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}}$
# bedrooms	\leftrightarrow	$\xleftarrow{\hspace{1cm}} \xrightarrow{\hspace{1cm}}$

Features

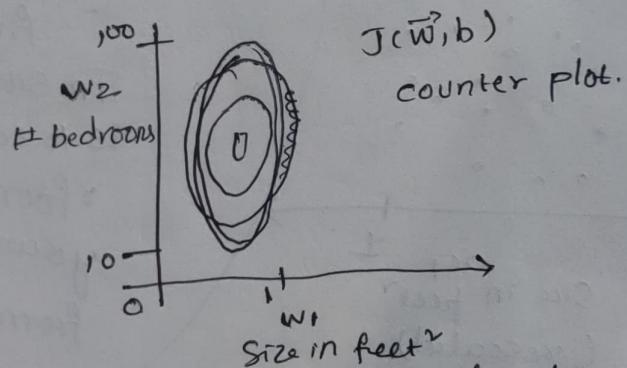
Parameters

Let's take a look at scatter plot of features

Let's see how ~~counter plot~~^{cost function} might look in a counter plot



if we plot the training data, you notice that the horizontal axis is on a much larger scale of much larger range of values compared to vertical axis



You might see a counter plot where the horizontal axis has a much narrower range, say $4 \text{ to } 6$, whereas vertical axis takes on much larger values say 10 & 100.

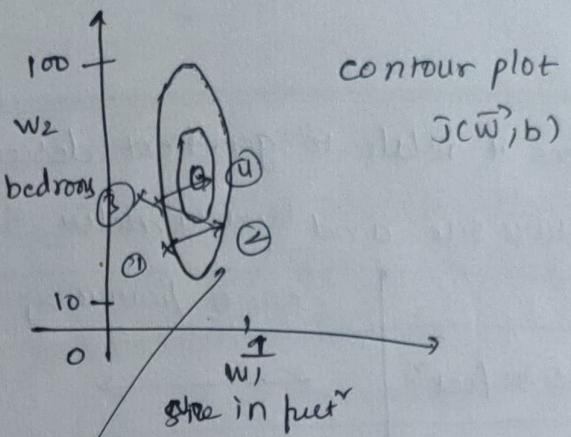
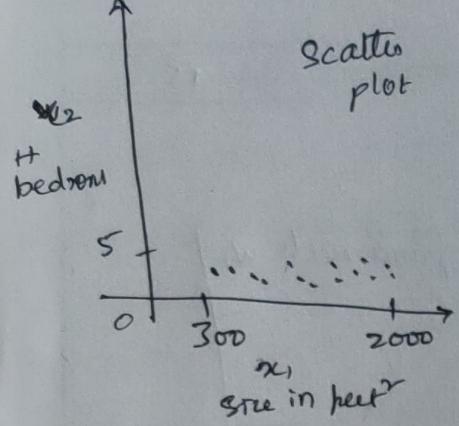
so contours form ovals & ellipses
and they're short on one side
horizontally

is longer on other (vertical)

is longer on other (vertical)
And this because very small change
to w_i can have very large impact
on estimated price & cost $f''(x)$ also.
Because w_i tends to multiplied by very
large number, the size & square feet.

In contrast, it takes a much larger change w_2 in order to change the prediction ~~is~~ much. And thus small changes to w_2 , don't change $\cos r f^n$ nearly as much.

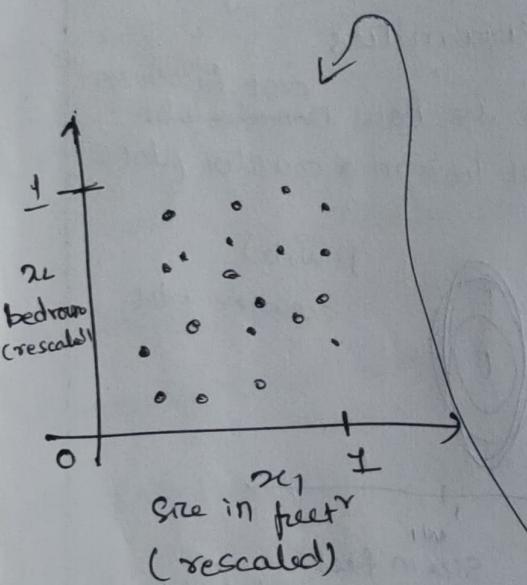
so, where does this leave us?
★ (New page)



This is what we end up happening if you were to run gradient descent, if you were to use your training data as is.

Because the contours are so tall and skinny may end up bouncing back & forth for a long time before it can finally find its way to global minimum.

In situations like this, a useful thing to do is to scale the features. This means performing some transformation of your training data so that x_1 ranges from 0 to 1 & x_2 might also range from 0 to 1.

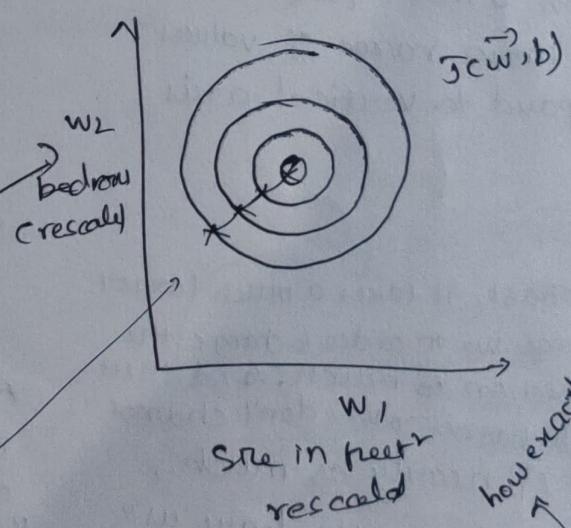


rescaling x_1 and x_2 both now taking comparable range of values to each other

if you run gradient descent on a cost J to find on this rescaled x_1 & x_2 using the transformed data, then the contours will look like circles & less tall & skinny.

And gradient descent can find a much more direct path to the global minimum.

RECAP: When you have different features that take on very different range of values, it can cause gradient descent to run slowly but rescaling diff features so they all take on comparable range of values. because speed, upgrade & descent significantly.



how exactly it done?
see next
:-)

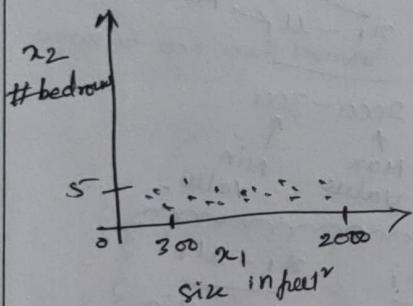
Feature Scaling part 2:

Implementation of feature scaling:

To take features that take on very different range of values and scale them to have comparable ranges of values to each other.

How do you actually scale features?

If x_1 ranges from 300 to 2000



Every value from x_1 feature will divide by Max value

$$300 \leq x_1 \leq 2000$$

$$x_1 \text{ scaled} = \frac{x_1}{2000} \quad \text{Max value}$$

Then scaled x features ranges from

$$\frac{300}{2000} \leq x_1 \leq \frac{2000}{2000}$$

$$0.15 \leq x_1 \leq 1$$

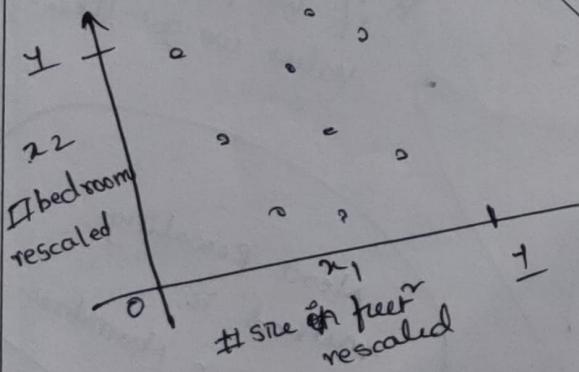
x_2 ranges from 0 $\leq x_2 \leq 5$ min max

$$x_2 \text{ scaled} = \frac{x_2}{5} \quad \text{Max value}$$

$$\frac{0}{5} \leq x_2 \leq \frac{5}{5}$$

$$0 \leq x_2 \leq 1$$

if we plot their values in graph (x_1 vs x_2)



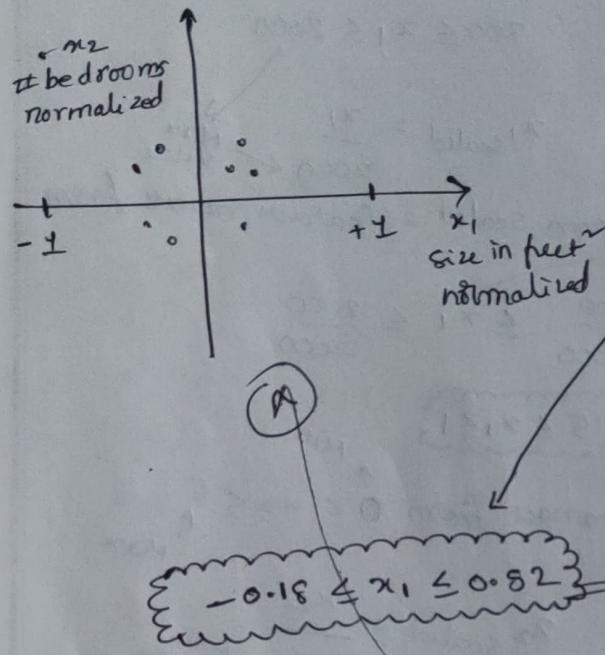
In addition to dividing by Maximum value, You can also do what's called "mean normalization"

Mean Normalization:

First we started with original features and we rescale them so that both of them are centred around zero.

whereas before (earlier scaling method) they only had values (>0) greater than zero but now they have both +ve and -ve values (usually $b/w -1 \text{ to } +1$)

after mean normalization:



To calculate Mean Normalization of x_1 ,

i) Find Mean of $x_1 (\mu_1)$

$$\Sigma: \text{Mean} = \mu_1 = 600$$

$$\text{Then } x_1 = \frac{x_1 - \mu_1}{\text{range from } 300 \text{ to } 2000}$$

$$= \frac{2000 - 300}{\text{Max value} - \text{Min value}}$$

$$x_1 \text{ (rescaled)} = \frac{x_1 - \mu_1}{2000 - 300} = \frac{x_1 - 600}{2000 - 300}$$

$$\frac{300 - 600}{2000 - 300} \leq x_1 \leq \frac{2000 - 600}{2000 - 300}$$

Similarly,

for x_2 , also, x_2 ranges: $0 \leq x \leq 5$

$$x_2 = \frac{x_2 - \mu_2}{\text{Max-Min}} \quad \text{Let's say } \mu_2 = 2.3$$

Then

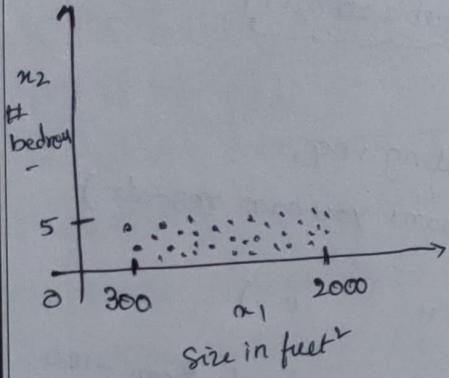
if we plot these values we got like

$$\frac{0 - 2.3}{5 - 0} \leq x_2 \leq \frac{5 - 2.3}{5 - 0}$$

$$-0.46 \leq x_2 \leq 0.54$$

Next Rescaling method is Z-score Normalization.

Z-Score Normalization:



$$300 \leq x_1 \leq 2000 \quad 0 \leq x_2 \leq 5$$

To implement Z-score normalization we need to calculate standard deviation of each feature

$$Z = \frac{x - \mu}{\sigma}$$

For instance if x_1 feature has $\sigma_1 = 450$

$\& \mu_1 = 600$
Then $\text{z-score Normalization of } x_1 \text{ is}$ \rightarrow each value of feature

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

likewise for x_2

$$x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

for instance std-dev
of x_2 feature is $\sigma_2 = 1.4$
 $\mu_2 = 2.3$

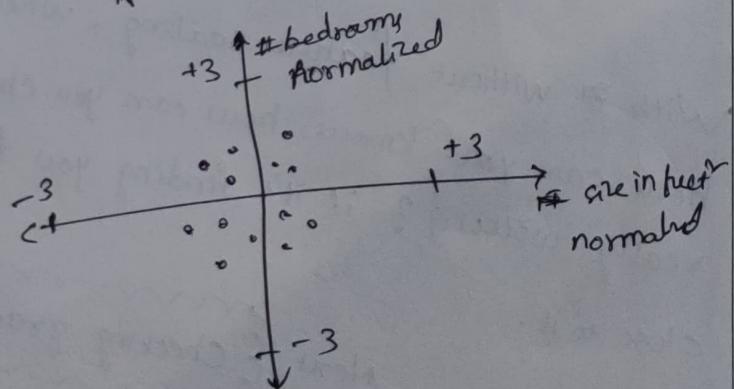
$$\frac{300 - 600}{450} \leq x_1 \leq \frac{2000 - 600}{450}$$

$$-0.666 \leq x_1 \leq 3.111$$

$$\frac{0 - 2.3}{1.4} \leq x_2 \leq \frac{5 - 2.3}{1.4}$$

$$-1.64 \leq x_2 \leq 1.928$$

if we plot $x_1 \& x_2$



Next we will see about
Thumb rule of Normalization

As a thumb rule, when you perform feature scaling, you might want aim for getting the features to range from: ~~-1 to 1~~

aim for about $-1 \leq x_j \leq 1$ for each feature x_j

$$\left. \begin{array}{l} -3 \leq x_j \leq 3 \\ -0.3 \leq x_j \leq 0.3 \end{array} \right\} \text{Acceptable ranges}$$

$0 \leq x_j \leq 3$ (Okay, no rescaling req,
But if you want you can rescale)

$$-2 \leq x_2 \leq 0.5 \quad (", ", ")$$

Too large car

$-100 \leq x_3 \leq 100$ (But here range is high from -100 to +100 i.e. 200 types of values ($\approx -100 \text{ to } 100, 0.05/100$))
(Too large) So Not acceptable, Need to rescale

Too small car

$-0.001 \leq x_4 \leq +0.001$ (These values are too small)
(Need to rescale)

$98.6 \leq x_5 \leq 105$ (Values are around 100, i.e. pretty large.
compared to other scale features so it will cause gradient descent to run more slowly.
So we need to rescale the feature to make fast the gradient descent)

With or without feature scaling, when you run gradient descent how can you know, how can you check if gradient descent is really working? if it's finding you the global minimum or something close to it.

Next { Checking gradient descent for Convergence }

Checking Gradient Descent for Convergence :-

When running gradient descent, how can you tell whether it is converging? That is whether it's helping you to find parameters close to the global minimum of cost function.

By learning to recognize what a well-running implementation of gradient descent looks like, we will also look how to choose better learning rate (α). Let's take a look.

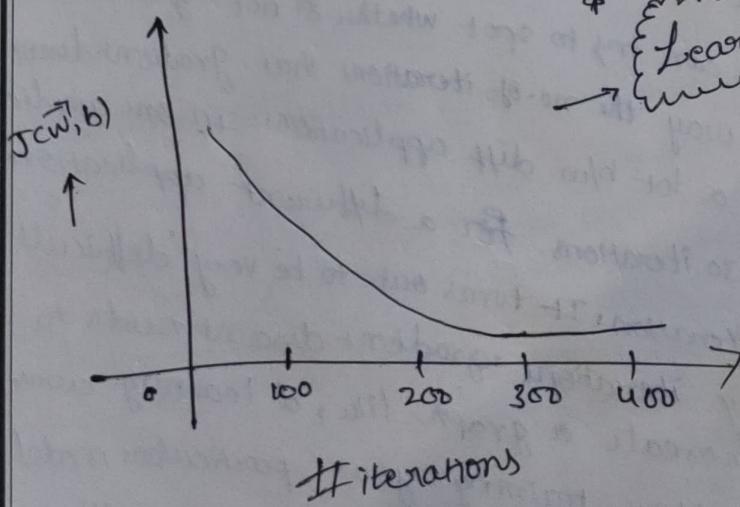
$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$

Recalling the objective of gradient descent : minimize $J(\vec{w}, b)$

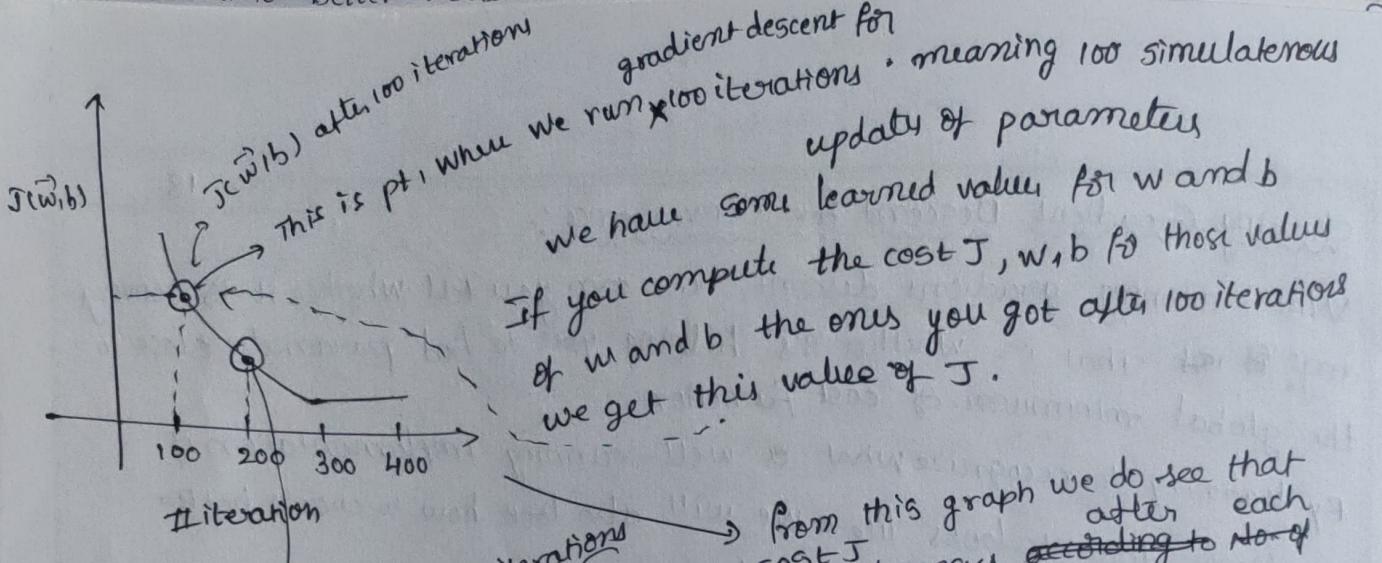
↳ is finding parameters w, b that hopefully minimize cost function.

plot cost f^n vs J (which is calculated on training set) & I plot the J at each iteration. Each iteration means each simultaneous update of parameters w and b . In this plot the horizontal axis is no. of iterations of gradient descent you've run so far. Then we get a curve like below & vertical axis is cost f^n .



↑
Learning Curve
→ cost f^n vs no. of iterations

There are some more types of learning curves in this machine learning.
We will see later in the course



If you compute the cost J , w , b for those values of w and b the ones you got after 100 iterations we get this value of J .

From this graph we do see that how cost J changes after each iteration of gradient descent.

If gradient descent is working properly, then cost J should decrease after every single iteration.

If J ever increases after one iteration that means either α is chosen poorly (and it usually means α is too large) or there could be a bug in the code.

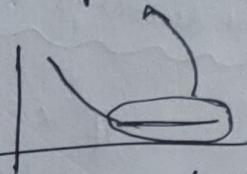
Another useful thing that this part can tell you is that if you look at this curve by the time we reach 300 iterations also, the cost J is levelling off and no longer decreasing much. By 400 iterations, it looks like curve has flattened out. This means that gradient descent has more or less converged because the curve is no longer decreasing.

By looking into the learning curve, you can try to spot whether or not gradient descent is converging. By the way the no. of iterations that gradient descent takes a convergence can vary a lot b/w diff applications. In one application, it may converge after 30 iterations. For a different application it could take 1,000 or 10,000 iterations. It turns out to be very difficult to tell in advance how many iterations gradient descent needs to converge, which is why we can create a graph like, a learning curve. Try to find out when you can start training your particular model. Another way to decide when your model is done training, is with an automatic convergence test.

Automatic convergence test :

Let ' ϵ ' be 10^{-3} i.e. 0.001 (a small number)

If $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration, declare "convergence"
 i.e. you're likely on flattened part of curve



Remember convergence, hopefully in the case that you found parameters w and b that are close to minimum possible value of J . I usually find that choosing the right threshold ' ϵ ' is difficult. I actually tend to look at graphs & rather than rely on automatic convergence tests.

You've now seen what the learning curve should look like when gradient descent is running well.

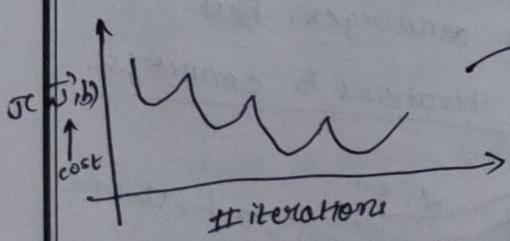
CHOOSING LEARNING RATE: (α) We need to choose best α .

If α is small, it will run slowly & if α is high (too large) it may not converge.

Let's take a look how can choose a good learning rate for your model?

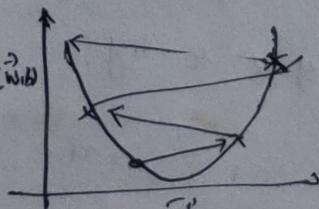
(Q: 17)

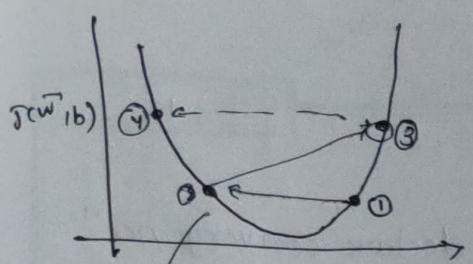
Identify the problem with gradient descent



when we plot cost vs no. of iterations, we notice that cost sometimes up & sometimes down. You should take that clear sign that gradient descent is not working properly. This could mean that there is a bug in code & sometimes it could mean that your learning rate is too large. (not converge)

so, here is illustration of what might happen

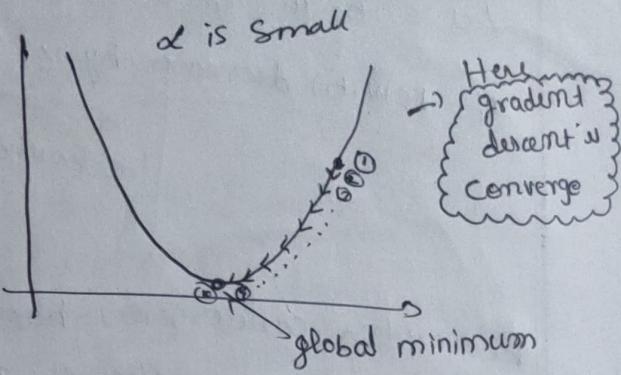




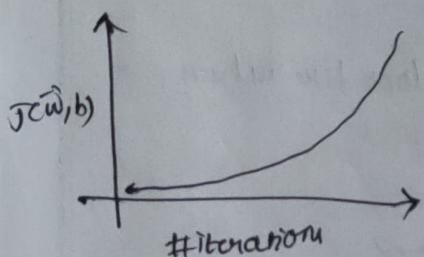
If α is too big
The cost will increase instead of decreasing.
To fix this we need to take α value too small to limit this. Then updates may start smaller

Here gradient descent is diverge

if α is small,
the cost decrease until it reaches global minimum.



Sometimes we may see that cost increases w.r.t iteration



→ This is also due to α is too large or sometimes due to bug in the code
so to fix this we need to take smaller learning rate(α)

one debugging tip
for correct implementation
of gradient descent is
1) using smaller ' α '
2) cost $J(w)$ should decrease on
every iteration. (if we follow
steps, 2nd one will automatically done)

incorrect → like if we write code like

$$w_1 = w_1 + \alpha \frac{\partial J(w, b)}{\partial w_1}$$

usually here we need
to put (-ve) sign but we place
it's there so that is called
'bug in code'

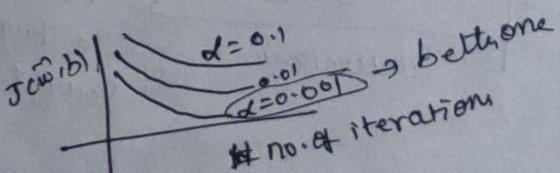
$$w_1 = w_1 - \alpha \frac{\partial J(w, b)}{\partial w_1}$$

correct

Note: If α is too small, the gradient descent converges. But α is too small gradient descent takes a lot more iterations to converge.

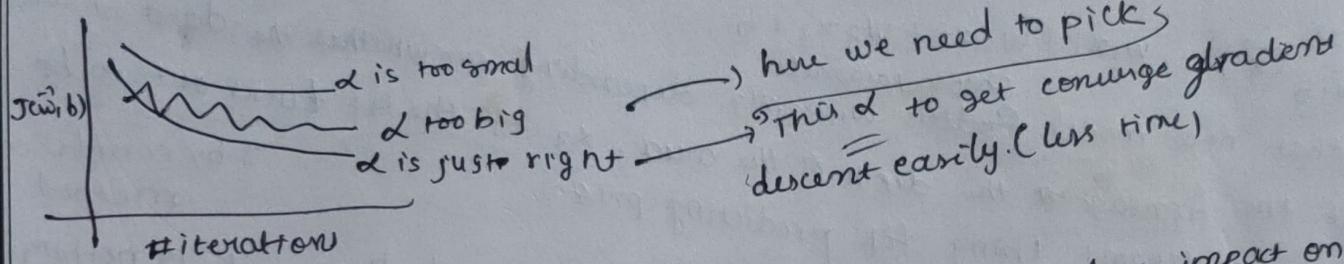
Values of α to try:

... 0.001 0.01 0.1 1 ...



Try range of values & pick the 'best' value which giving best loss
time.

Suppose if we got curves like



here we need to pick
This α to get converge gradient
descent easily (less time)

FEATURE ENGINEERING: The choice of features can have huge impact on your learning algorithm's performance. In fact, for many practical applications choosing a suitable right features is a critical step to making the algorithm work well.

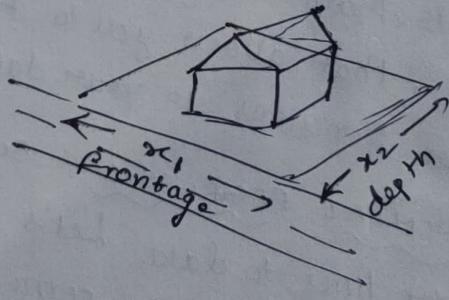
Let's see how you can choose or engineer the most appropriate features for the learning algorithm

Let's take example of predicting price of house.

Let's say we have two features, $x_1 \rightarrow$ width of the lot side of the plots of land that house is built on & In real estate, we called frontage of the lot) $x_2 \rightarrow$ (second feature) depth of the lot size (Assume the house built on rectangular plot)

By using two feature we built a model to predict price is like.

$$\begin{cases} f_{w,b}(x) = w_1 x_1 + w_2 x_2 + b \\ \text{Frontage} \quad \text{Depth} \end{cases}$$



This model is okay, but here is another option for how you might choose a different way to use these features in the model that could be more effective. You might notice that area of land calculated by frontage & width times the depth

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

We have an intuition that area of land is more predictive of the price, than the frontage & depth as separate features. So you might define a new feature, x_3 as $x_1 * x_2$, $x_3 \rightarrow$ area of plot of land with the feature x_3 we have a model, for

$$f_{\tilde{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

choose params
so that new model ~~predicts~~ w_1, w_2, w_3 depending on whether the data shows that the frontage & the depth or the area, x_3 of the lot turns out to be most important thing for predicting price \rightarrow by correlation method.

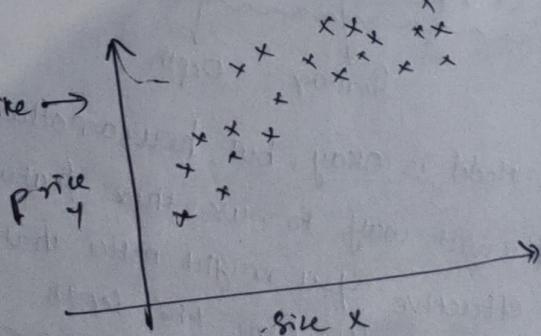
We have created new feature for predicting price of house, it is nothing but "feature engineering". \rightarrow Using intuition to derive new features by transforming & combining original features ^{problem} to make it easier for learning algorithm to make accurate predictions.

~~start from scratch~~ Depending on what insights you may have into the application, rather than just taking the features that you happen to have started off with sometimes by defining new features, you might be able to get a much better model \heartsuit

This is feature engg, It turns out that this one flavor of feature engg, that allow you to fit not just straight lines, but curves, non-linear functions to your data. Let's see that

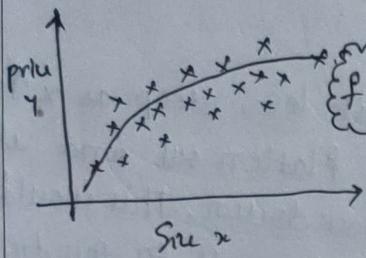
POLYNOMIAL REGRESSION: so far we have seen that, we are fitting straight lines to data. Let's take the idea of multiple Linear Regression and feature engineering to come up with a new algorithm called "polynomial regression". Which you let you fit curves with non-linear functions to your data

Let's say you have having data like \rightarrow
feature \rightarrow size in feet
feature \rightarrow price



It looks like ~~not~~ straight line fit this dataset very well. Maybe you want to fit a curve may be a quadratic function to the data like

$$\{ f_{\tilde{w}, b}(\vec{x}) = w_1 x^2 + w_2 x^v + b \}$$

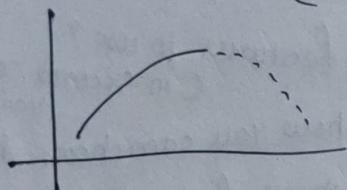


$$f_{\tilde{w}, b}(x) = w_1 x + w_2 \tilde{x} + b$$

size \tilde{x}

May be this give better fit to data.

But then you may decide that your quadratic model doesn't really make sense because a quadratic function, eventually comes back down. Well we wouldn't really expect housing price to go down when size 1's Big houses seems like they should usually cost more. Then you may choose cubic function where we now have not only x squared, but x cubed.



$$f_{\tilde{w}, b}(x) = w_1 x + w_2 \tilde{x}^2 + w_3 \tilde{x}^3 + b$$

size \tilde{x} size \tilde{x}^2 size \tilde{x}^3

May be this Model produces this curve like, which is somewhat better fit to data because the size does eventually come back up as the size increases.

These are both examples of polynomial Regression. Because you took optional feature x & raised it power to 2 & 3...

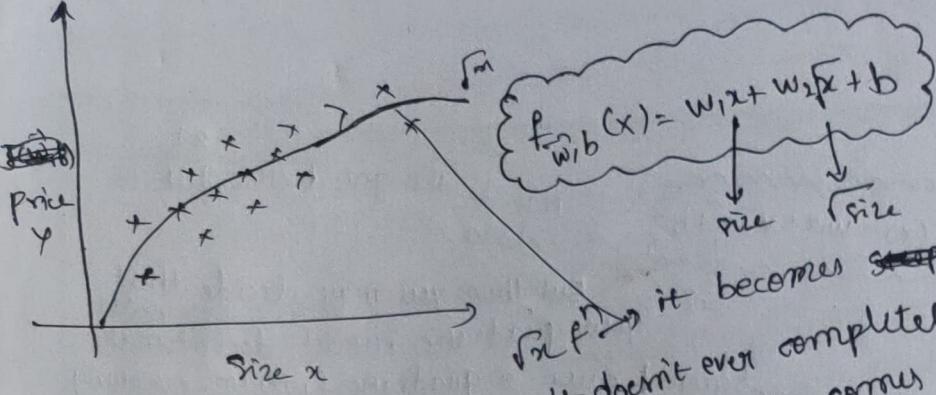
I just want to point out one thing, which is that if you create features that are these powers like square of original feature like this, then feature scaling becomes increasingly important.

if size of house ranges $1 - 10^3$ then $\tilde{x} \rightarrow 1^2 - (10^3)^2 \Rightarrow 1 - 10^6 = 1 \text{ million}$
 then $\tilde{x}^3 \rightarrow 1^3 - 10^9 \Rightarrow 1 \text{ billion}$

These two features (\tilde{x}^2, \tilde{x}^3) take on very different ranges of values compared to the original feature x . If you are using gradient descent, it's important to apply feature scaling to get your features into comparable range of values. Finally, we have one last example of how you really have a wide range of choice of features to use. Another reasonable alternative to

Smt. B. SEETHA POLYTECHNIC, VISHNUPUR, BHIMAVARAM

taking the \tilde{x}^2 & \tilde{x}^3 is to say we square root of x then Model look like



be another choice of features that might work well for this data set as well

$f_w,b(x) = w_1x + w_2x^2 + b$

but it doesn't ever completely flatten out and it certainly never ever comes back down. This would

You may ask yourself, how do I decide what features to use?
 In second course (advance learning algos) you see how you can choose diff features and diff models that include or don't include these features. & you have a process for measuring how well these diff models perform to help you to decide which features to include or not include. For now, you need to aware that you have a choice in what features you have.
 By using feature engg, and polynomial f^n , you can potentially get a much better model for your data.

LABS:

computing cost: 1st Lab.

$$\text{cost} = J(w, b) = \frac{1}{2m} \sum_{i=0}^{m-1} (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

$f_{w,b}(x^{(i)}) = w x^{(i)} + b$ → it is the prediction (\hat{y})

squared diff b/w target value & prediction

>> import numpy as np

>> import matplotlib.pyplot as plt

>> x-train = np.array([1.0, 2.0]) → size in feet

y-train = np.array([300.0, 500.0]) → price in \$

compute cost method: $\sum_{i=0}^{m-1} (y^{(i)} - \hat{y}^{(i)})^2$ → w & b value

def compute_cost(x, y, w, b):

$m = x.\text{shape}[0]$ → no. of training ex's ($\text{len}(x)$) also works

cost = sum = 0

for i in range(m):

$f_{w,b} = w * x[i] + b$

for i in range(m):

$$f_{-wb} = w * x[i] + b \longrightarrow \cancel{w}$$

$$\text{cost} = (f_{-wb} - y[i])^{**2}$$

$$\text{cost_sum} = \text{cost_sum} + \text{cost}$$

$$\text{total_cost} = (1/(2*m)) * \text{cost_sum}$$

return total_cost

LAB 2: Implementing gradient descent.

def compute_gradient(x^{x-train}, y^{y-train}, w, b):

$$m = x.shape[0]$$

$$dj_dw = 0$$

$$dj_db = 0$$

for i in range(m):

$$f_{-wb} = w * x[i] + b$$

$$dj_dw_i = (f_{-wb} - y[i]) * x[i]$$

$$dj_db_i = f_{-wb} - y[i]$$

$$dj_dw += dj_dw_i$$

$$dj_db += dj_db_i$$

$$dj_dw = dj_dw/m$$

$$dj_db = dj_db/m$$

return dj(dw), dj(db)

[Refer page NO: 10]

(Just to know concept).