

- Machine learning : field of study that gives computers ability to learn. (1)
 - Grew out of work in AI. We wanted to build machines which can work as humans. But finding an algorithm that works perfectly is very difficult, so we give the task of learning to machines itself, thus emerged Machine learning from AI.

- A computer program is said to learn from experience E with respect to some Task T and some performance measure P if its performance on T measured by P improves with experience E.
 - Eg : An email program watches which email we do or do not mark as spam. here
 - ⊕ T is classifying email as spam or not spam
 - ⊕ E is program watching us label emails
 - ⊕ P is number of mails correctly labeled as spam/not spam by P
 - Eg : Chess engine learning to play chess
 - ⊕ T is the task of playing chess
 - ⊕ E is the experience computer gets from playing chess
 - ⊕ P is the probability of computer winning next game.

- Types of learning :-

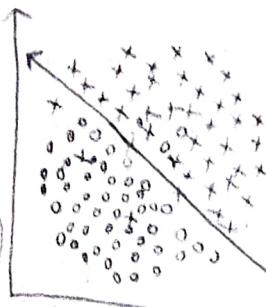
- Supervised learning
- unsupervised learning

- a) Supervised learning :→ (learns from 'right answers')
 - here right answers are already given
 - i.e. we feed a dataset to the computer and according to it we find the answers
 - for eg :→ selling a 750 sq. feet house.
 - > collect dataset i.e. sq. feet and selling value of houses
 - > plot these on graph
 - > choose the best fit { this is an important & tricky task }
 - > according to the fit find value of house.
- Here we had already had dataset i.e. correct answers*

- This type of problem can also be called Regression problem where we predict Continuous valued outputs (Price)
- eg 1 : finding out if Breast cancer is malignant or benign i.e. serious or normal.
 - > This type of problem is different from above Regression Problem
 - > Here there are only 2 possibilities unlike Regression where we have a continuous range of possibilities
 - > Note :> number of possibilities being 2 is not necessary, it can be 3, 4 etc.
 - The only criteria is that possibilities has to be discrete.
 - > Here we use multiple features like Tumor size and age

- > Based on these features plot a graph.
- > now we use our machine learning algorithm to Create a line that partitions between these plot.
- > now we find where does our points lie in this graph. Based on this we predict type of Breast Cancer.

Now $\circ \Rightarrow$ benign cancer datapoint
 $\times \Rightarrow$ malignant cancer datapoint
 • we plot both of these on graph.
 • create the line with algorithm
 • then plot our point on graph & if it is below than Cancer is benign else is malignant.



? In this case we used only 2 features Age & Tumor size Tumah \rightarrow 2d
 but what if we want to use more features, likely infinite features;
 This is a problem coz ∞ features = ∞ datapoints, but we don't have
 ∞ memory to store these points. Here we can use Support Vector mach-
 ine.

Q) Difference between Regression & Classification?

↴ Predicts a number
 ↴ Infinitely many possibilities

↴ Predicts categories
 ↴ Smaller number of Possible outputs

⑤ Unsupervised learning

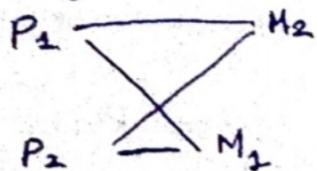
- We have unlabeled dataset
- We plot this unlabeled dataset and try to find clusters in them thus it is also called clustering algorithm.
- Eg: Google news; collects all of the news regarding a particular topic. Here in the beginning there is nothing given to us; but based on the similarities between multiple news we cluster them into some groups. involves clustering
- "I don't know what is the data but can you find structure in this data?"

Cocktail party problem: (A type of unsupervised learning problem)

multiple people talking at the same time and it is difficult to understand anything. taking example with 2 people & 2 Mics

∴ This does not involve clustering.

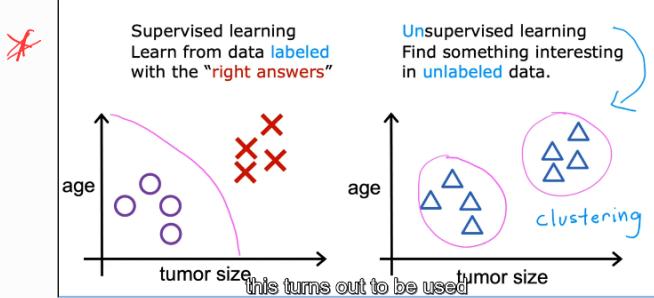
∴ Here we find structure in chaotic environment



These mic M₁ & M₂ will collect voices of both P₁ & P₂. but when both of these audios are used together, and based on the distance we can filter out the voices of P₁ & P₂ individually

- * Data only comes with input x but not output labels y (like in supervised learning)
- * Algorithm has to find **structure** in the data.
- * Clustering → Grouping similar data points together
- * Dimensionality reduction → Compress data using fewer numbers
- * Anomaly Detection → Find unusual data points

Clustering Examples



Clustering: Google news

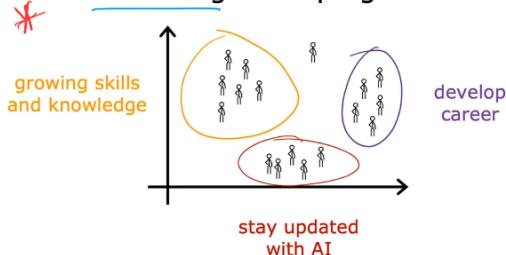
Giant **panda** gives birth to rare **twin** cubs at Japan's oldest **zoo**

USA TODAY · 6 hours ago

- Giant **panda** gives birth to **twin** cubs at Japan's oldest **zoo**
CBS News · 7 hours ago
- Giant **panda** gives birth to **twin** cubs at Tokyo's Ueno **Zoo**
WHBL News · 16 hours ago
- A Joyful Surprise at Japan's Oldest **Zoo**: The Birth of Twin **Pandas**
The New York Times · 1 hour ago
- **Twin** Panda Cubs Born at Tokyo's Ueno **Zoo**
PEOPLE · 6 hours ago

[View Full Coverage](#)

Clustering: Grouping customers



* All the news have the three words common - eg of clustering

* Unsupervised learning Algorithms are used at →

⇒ Given a set of news articles found on the web, group them into sets of articles about the same story.

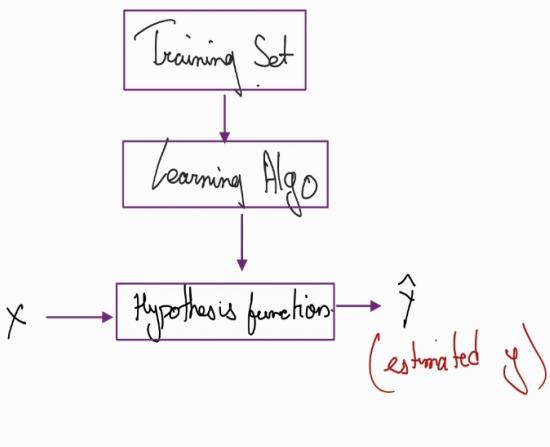
Given a database of customer data, automatically discover market segments and group customers into different market segments.

Linear Regression with One Variable (Univariate Linear Regression) (2)

Model Representation

- It is a type of Supervised learning.
- Regression i.e. predicting Real-valued output (predicts number).
- Here we use Training Set where:
 - m is total number of training example.
 - x is input variable / features.
 - y is output variable / target.

→



→ This training data set is passed through a learning algorithm which gives out h . h is a hypothesis function.

→ For the example of housing if we put in number of square feet of our house in h it will give predicted price of it.

linear equation
 m is slope
 C is constant

→ We represent h as :

$$h(x) = \Theta_0 + \Theta_1 x \quad \{ \text{i.e. } y = mx + c \}$$

→ we are simplifying and sticking to only straight lines for model-fitting. thus a linear equation is used.



Corresponding
Data Set

Data table

size in feet ²	price in \$1000's
2104	400
1416	232
1534	315
852	178
...	...
3210	870

Terminology

Training set:	x → size in feet ²	Data used to train the model
		→ price in \$1000's
(1)	2104	400
(2)	1416	232
(3)	1534	315
(4)	852	178
...
(47)	3210	870

$x^{(1)} = 2104$ $y^{(1)} = 400$
 $(x^{(1)}, y^{(1)}) = (2104, 400)$

$x^{(2)} = 1416$ $x^{(2)} \neq x^2$ not exponent

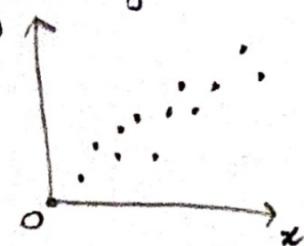
Notation:
 x = "input" variable
 y = "output" variable
 m = number of training examples

(x, y) = single training example

$(x^{(i)}, y^{(i)})$
 $(x^{(i)}, y^{(i)}) = i^{\text{th}}$ training example
 index (1st, 2nd, 3rd ...)

Cost function

→ Suppose we get the plotted graph as below:



→ Here it is impossible to choose a value of Θ_0 & Θ_1 in a way that all points passes perfectly.

→ In this case we construct a line which will create least errors in all cases.

→ Turns out that using squared error function is the best way to deal with Regression problems. thus we get

Cost function
Denoted by J

$$\text{minimize}_{\Theta_0 \Theta_1} \frac{1}{2m} \sum_{i=1}^m (h_{\Theta}(x^i) - y^i)^2$$

→ $h_{\Theta}(x)$ is the value we get from our hypothesis function

∴ Thus if we minimize this value and get a value of Θ_0 & Θ_1 we get a hypothesis function that creates least average error. → $\frac{1}{2m}$ is for average error.

Cost function: Squared error cost function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

*$\hat{y}^{(i)}$ = predicted value
 $y^{(i)}$ = actual value
error = $\hat{y}^{(i)} - y^{(i)}$*

m = number of training examples

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

Comparison between hypothesis & cost function:

$$\rightarrow h(x) = \Theta_0 + \Theta_1 x$$

i.e. function of x

$$J(\Theta_0, \Theta_1) = \frac{1}{2M} \sum_{i=0}^n (h(x)^i - y^i)^2$$

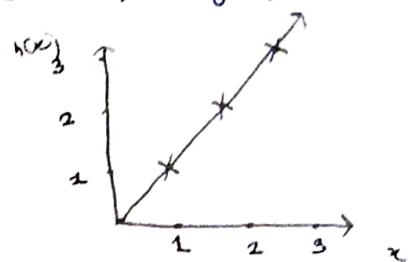
i.e. function of Θ_0 & Θ_1

→ now for each value of Θ_0 & Θ_1 gives us a different fit in the graph. and thus we get a value of J corresponding to this chosen Θ_0 & Θ_1 .

→ for the simplicity of calculation let $\Theta_0 = 0$, i.e. all fits pass through origin.

and the dataset is

x	y
1	1
2	2
3	3



→ for $\Theta_1 = 1$ we get a graph as:

for $\Theta_1 = 1$ & $\Theta_0 = 0$ we get J as

$$J(0, 1) = \frac{1}{2 \cdot 3} \sum_{i=0}^3 (h(x)^i - y^i)^2$$

∴ here for all cases $h(x)^i = y^i \Rightarrow h(x)^i - y^i = 0$

$$\Rightarrow J(0, 1) = 0$$

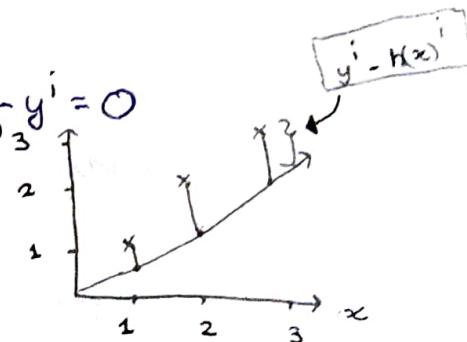
→ for $\Theta_1 = 0.5$ we get a graph as:

$$J(0, 0.5) = \frac{1}{6} [(0.5)^2 + (1)^2 + (1.5)^2]$$

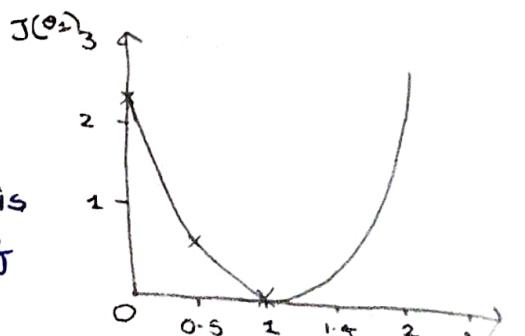
$$= \frac{1}{6} [0.25 + 1 + 2.25]$$

$$= \frac{3.5}{6}$$

$$= 0.58$$



→ similarly if we repeat this process in other values we get the graph of Cost function as :

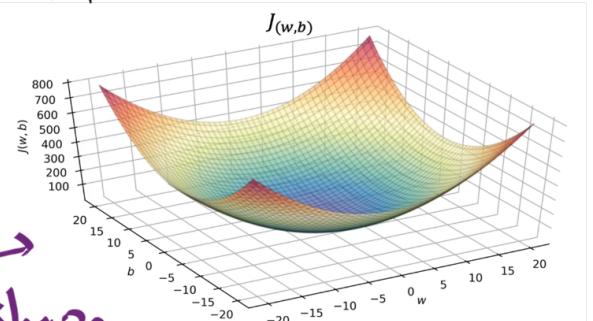


Conclusion: → We get a parabolic graph; by studying this graph we get that at $\Theta_1 = 1$ we get $J(\Theta_1) = 0$ i.e. least possible cost for our scenario & thus we choose $\Theta_1 = 1$.

★ If we do not neglect Θ_0 i.e. $\Theta_0 \neq 0$ then we get a 3-D graph of cost function which is similar to the above graph but just 3-D.

here the bottom point of this curve have least value of J and thus we choose it.

Soup bowl
or Hammock Shape



model:

$$f_{w,b}(x) = wx + b$$

parameters:

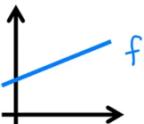
$$w, b$$

cost function:

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

goal:

$$\underset{w,b}{\text{minimize}} J(w, b)$$



simplified

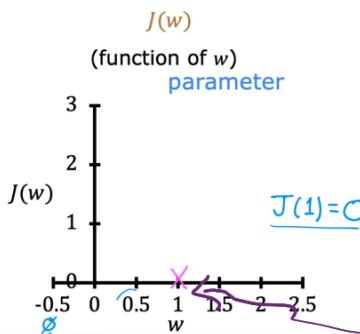
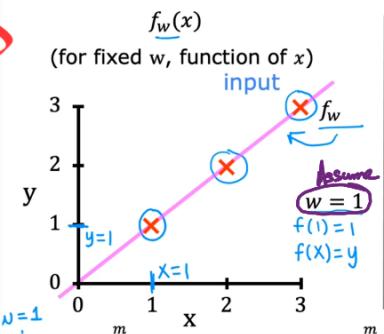
$$f_w(x) = wx \quad b = \emptyset$$



for simplicity of calculation

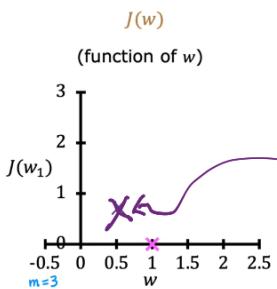
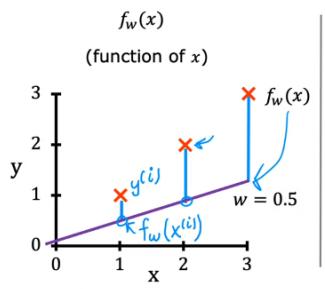
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

$$\underset{w}{\text{minimize}} J(w)$$



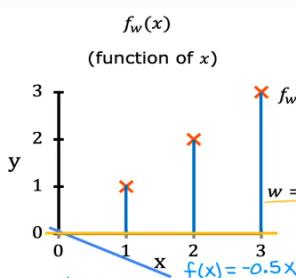
$$\begin{aligned} J(w) &= \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} ((1-1)^2 + (2-2)^2 + (3-3)^2) \\ &= \frac{1}{2m} \times 0 = 0 \end{aligned}$$

Similarly for $w=0.5$,



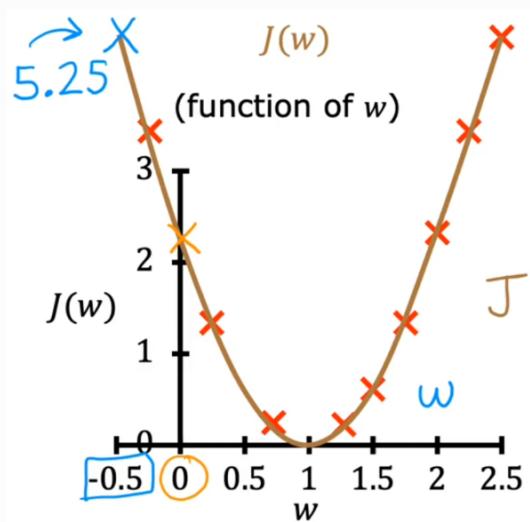
$$\begin{aligned} J(w) &= \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} ((0.5-1)^2 + (1-1)^2 + (1.5-1)^2) \\ &= \frac{1}{2m} \times 3.5 \approx 0.58 \end{aligned}$$

Similarly if we do for $w=0$,



we get $J(w) = 2.3$

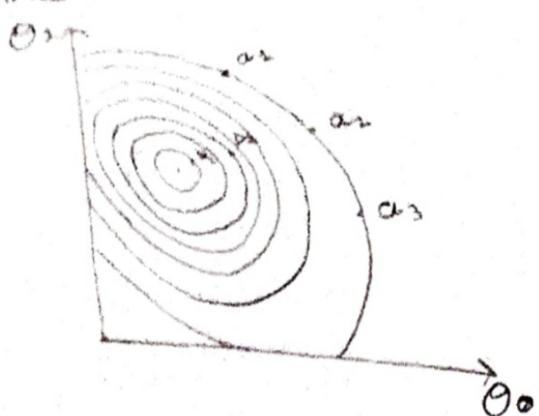
After completely plotting $J(w)$ we get a parabolic curve.



Our final goal is to minimize the difference,
From the above plotted graph we can see that the minimum difference ($J(w)$) is for $w=1$. Hence we choose $w=1$.

Contour Plots : \rightarrow Trying to find $\min J$ in a 3-D graph is diff.
so we use a contour plot rather than the 3D graph. it looks like : \rightarrow

\therefore It is the projection of the ③ 3-D graph on the $\Theta_0-\Theta_1$ plane.
 \therefore for all points a_1, a_2, a_3 value of J is same
 \therefore value of J at $b_2 <$ value of J at b_1
 \therefore Thus more the center we go, better value of $\Theta_0 \& \Theta_1$ is found i.e. values yielding to lower J .



Gradient Descent

→ Up until now all the methods we saw required us to plot the graph for cost function and observe its lowest point for value of θ 's.

→ Gradient Descent does not require construction of graph

→ Procedure :

- Assume a 3-D graph; here choose any random point.
- Height of this graph shows Cost. Thus our objective is to reach to the least height possible.
- At this random point scan 360° around us and choose a direction which leads us to a smaller height. take some steps.
- Repeat till we reach to a position where no direction can decrease height.

Mathematically,

$$x = x - \alpha \frac{\partial}{\partial w} J(x, y)$$

Generalizing it for w ,

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

Here
 → $=$ is an assignment operator, i.e. $A := B$ means update value of A such that $A = B$ happens.
 → α is learning rate, i.e. defines how big step are we taking after fixing direction.
 → $\frac{\partial}{\partial w}$ Partial derivative; defines which direction to move.

→ We know we have two variables here and the update for both is necessary

→ We need to keep in mind that before w gets updated we also change b . i.e., simultaneously update w & b

Correct: Simultaneous update

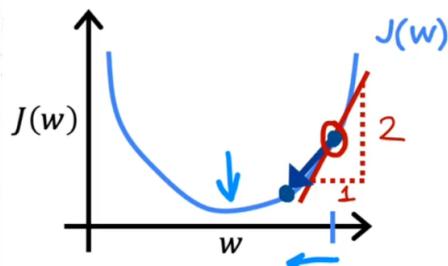
$$\begin{aligned} \text{tmp_w} &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \text{tmp_b} &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w &= \text{tmp_w} \\ b &= \text{tmp_b} \end{aligned} \quad \left. \right\}$$

Incorrect

$$\begin{aligned} \text{tmp_w} &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ w &= \text{tmp_w} \\ \text{tmp_b} &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ b &= \text{tmp_b} \end{aligned} \quad \begin{array}{l} w \text{ is getting} \\ \text{updated} \\ \text{before} \\ J(w, b) \end{array}$$

can be calculated for b

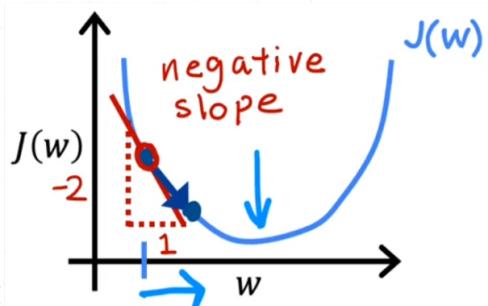
→ As we know $\frac{\partial f(x)}{\partial x}$ will give us slope of x , and if we assume $\theta_0 = 0$, we get a 2-D graph and $\frac{\partial}{\partial \theta_1} = \frac{\partial}{\partial \theta_1}$ i.e. for only 1 variable partial diff. & normal diff. are same thus



$$\therefore \text{here } \theta_1 := \theta_1 - \alpha \frac{\partial J(\theta_1)}{\partial \theta_1}$$

and at this point slope is positive thus $\frac{\partial J(\theta_1)}{\partial \theta_1}$ is a positive number

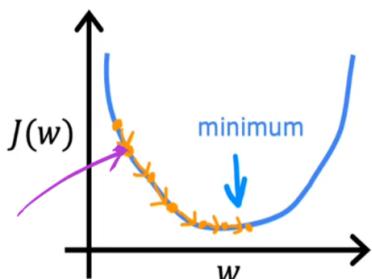
∴ Thus updated value of θ_1 is smaller than previous value and that is what we want.



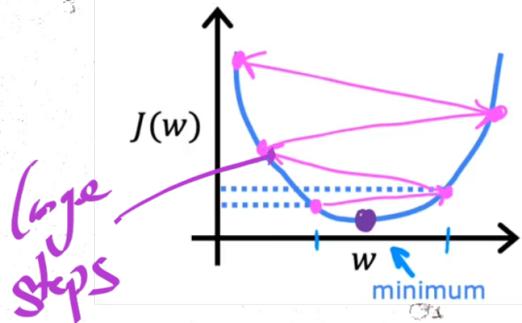
∴ have slope is -ve
 $\Rightarrow \frac{\partial}{\partial \theta_1}$ is -ve ; thus new $\theta_1 >$ prev θ_1
 $(w - \alpha(-ve \Delta)) = w + \alpha(\text{the } \Delta)$
 and that is what we want to reach at θ^*



∴ slope is 0
 thus min already reached and value of θ_1 will not change now.



∴ Very small value of α means we are taking very small steps; thus gradient descent will be slow and take very long to reach to final point.



∴ IF we choose a value of α too large then each step will lead us to a value of J even larger than before. Thus we will never reach to final point.

→ Thus we will want to take larger steps initially and decrease step size as we get closer to the answer. (4)

for this we do not have to change α after every point.

As with each step value of α_1 changes which will lead to a smaller value of $J(\alpha_2)$ and thus automatically decreasing step size in next iteration

∴ so we are good as long as value of α is not too small not too big.

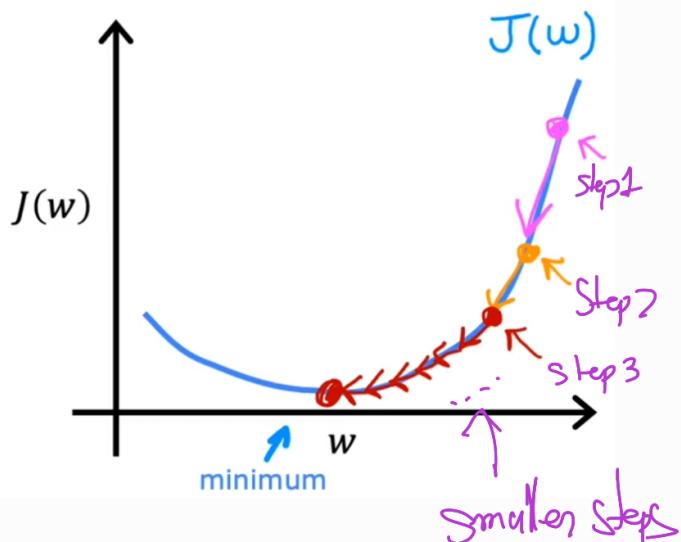
Can reach local minimum with fixed learning rate? $\frac{\alpha}{\alpha}$ Yes !!

$$w = w - \alpha \frac{d}{dw} J(w)$$

smaller
not as large
large

Near a local minimum,
 - Derivative becomes smaller
 - Update steps become smaller

Can reach minimum without decreasing learning rate α



Linear regression model

$$f_{w,b}(x) = wx + b$$

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

Cost function

Gradient descent algorithm

* repeat until convergence {

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

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

}



The derivation of these forms are optional

(Optional)

$$\frac{\partial}{\partial w} J(w, b) = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)}} + b - y^{(i)})^2$$

$$= \cancel{\frac{1}{2m} \sum_{i=1}^m} (\underline{wx^{(i)}} + b - y^{(i)}) \cancel{2x^{(i)}} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}}$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)}} + b - y^{(i)})^2$$

$$= \cancel{\frac{1}{2m} \sum_{i=1}^m} (\underline{wx^{(i)}} + b - y^{(i)}) \cancel{2} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})}$$

no $x^{(i)}$

Gradient descent algorithm

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

repeat until convergence {

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

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

} (update w & b simultaneously)

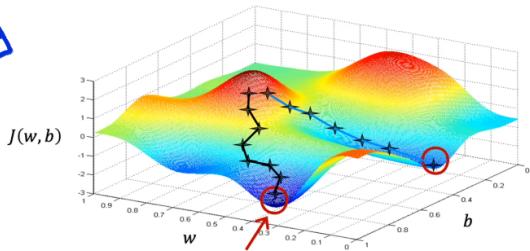
final Algorithm

- This type of gradients descent is called **batch gradient**
descent coz entire range of $i=0$ to $i=m$ is considered single batch and processed.
- Disadvantage : → For a 3-D graph with multiple local minima & a global minima ; this process will give the local minima nearest to starting point ; not the global minima

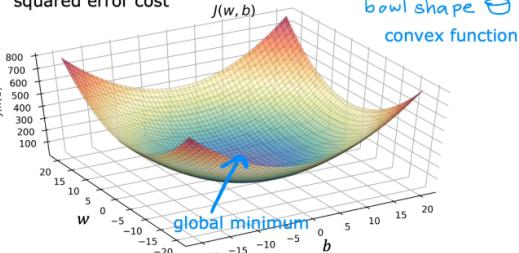
Each step of gradient descend uses all the training examples.

But this is not a problem in linear regression Coz graphs are Convex in nature i.e only 1 global minima.

More than one local minimum



squared error cost



Matrices and Vectors

→ Matrix : $\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$

size : Row x Column
i.e. 3×3

→ If Rows = Column \Rightarrow square matrix

→ If only 1 column i.e. $(n \times 1)$ then it is called vector.

→ $\begin{bmatrix} a & b \\ c & d \end{bmatrix} + \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} a+e & b+f \\ c+g & d+h \end{bmatrix}$ i.e if size same then normal addition possible otherwise invalid.

→ $4 \times \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} 4a & 4b \\ 4c & 4d \end{bmatrix}$ i.e irrespective of size multiplication with scalar is possible

→ $\begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix} \times \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} ax + by \\ cx + dy \\ ex + fy \end{bmatrix}$ i.e multiplication possible only if no. of columns in 1st = no. of rows in 2nd.

$(3 \times 2) \quad (2 \times 2) \quad (3 \times 1)$

→ $\begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix} \times \begin{bmatrix} w & x \\ y & z \end{bmatrix} = \begin{bmatrix} aw + by & ax + bz \\ cw + dy & cx + dz \\ ew + fz & ex + fz \end{bmatrix}$

$(3 \times 2) \quad (2 \times 2) \quad (3 \times 2)$

→ $A * B \neq B * A$ {i.e. not commutative}

→ $A * (B * C) * = (A * B) * C$ {i.e associative}

→ Identity matrix : square matrix with all elements 0 except diagonal which is 1 i.e.

a 4 dimension identity matrix: $\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$

→ Identity matrix denoted by I. i.e.: $\rightarrow A * I = A$

→ If $A * x = I$ we say x is the inverse of A i.e $x = A^{-1}$

→ Octave has 2 commands inv & pinv to find inverses. inv finds inverses but gives error if the matrix is irreversible. pinv means Pseudo-inverse. It gives us values despite the fact that matrix is irreversible.

→ Transpose function: If $A = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix}$

then $A^T = \begin{bmatrix} a & d \\ b & e \\ c & f \end{bmatrix}$

Multiple Features

→ Till now we used only 2 features i.e $h(x) = \theta_0 + \theta_1 x$; here $x_0 & x_1$ are 2 features where $x_0 = 1$

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

$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n] \Rightarrow$ Row Vector (parameters of the model)

b is a number.

$\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n] \Rightarrow$ Row Vector of multiple features

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

↑
(dot product) ↑
Multiple Linear Regression

Multiple features (variables)

Size in feet ² x_1	Number of bedrooms x_2	Number of floors x_3	Age of home in years x_4	Price (\$) in \$1000's
2104	5	1	45	460
i=2 1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$x_j = j^{\text{th}}$ feature

$n = \text{number of features}$

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

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

$j = 1 \dots 4$

$n = 4$

$$\vec{x}^{(2)} = [1416 \ 3 \ 2 \ 40]$$

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

Vectorization :-

Parameters:-

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

b is a number

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

Python Implementation :-

$$w = np.array([1.0, -2.5, -3.3])$$

$$b = 4$$

$$x = np.array([10, 20, 30])$$

Without Vectorization :-

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

$$\begin{aligned} f &= w[0] * x[0] + \\ &\quad w[1] * x[1] + \\ &\quad w[2] * x[2] + b \end{aligned}$$

Without Vectorization :- (loop)

$$f_{\vec{w}, b}(\vec{x}) = \left(\sum_{j=1}^m w_j x_j \right) + b$$

\therefore for j in range (m)

$$\begin{aligned} f &= f + (w[j] * x[j]) \\ f &= f + b \end{aligned}$$

Using Vectorization :-

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

$$f = np.dot(w, x) + b$$

Why do Vectorized Code faster than loops ?

→ Vectorized Code usually use parallel processing in the hardware making the calculations faster. whereas loops are a serial implementation and thus slower relatively.

Gradient Descent For multiple Variables

→ now we have $\theta_0, \theta_1, \theta_2, \dots, \theta_n$; update all of these simultaneously till convergence.

Gradient descent

repeat { One feature }

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

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

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

simultaneously update $\underline{\underline{w}}, \underline{b}$

}

n features ($n \geq 2$)

repeat {

$j=1$ $\underline{\underline{w}}_1 = \underline{\underline{w}}_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\underline{\underline{w}}, b}(\bar{\underline{x}}^{(i)}) - \underline{y}^{(i)}) \underline{x}_1^{(i)}$

⋮

$j=n$ $\hookrightarrow \frac{\partial}{\partial \underline{\underline{w}}_1} J(\underline{\underline{w}}, b)$

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

}

Feature & Parameter Scaling :-

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

\uparrow \uparrow
 Size #bedrooms

$$x_1 \text{ range} \rightarrow 300 - 2000 \text{ ft}^2$$

$$x_2 \text{ range} \rightarrow 0 - 5$$

House: $x_1 = 2000, x_2 = 5, \text{ price} = 500K \leftarrow \text{taking one training data for example}$

According to the above what should be the value of w_1, w_2 ?

$\checkmark \quad w_1 = 50, w_2 = 0.1, b = 50$

$$\widehat{\text{price}} = \frac{50 * 2000}{100,000K} + \frac{0.1 * 5}{0.5K} + \frac{50}{50K}$$

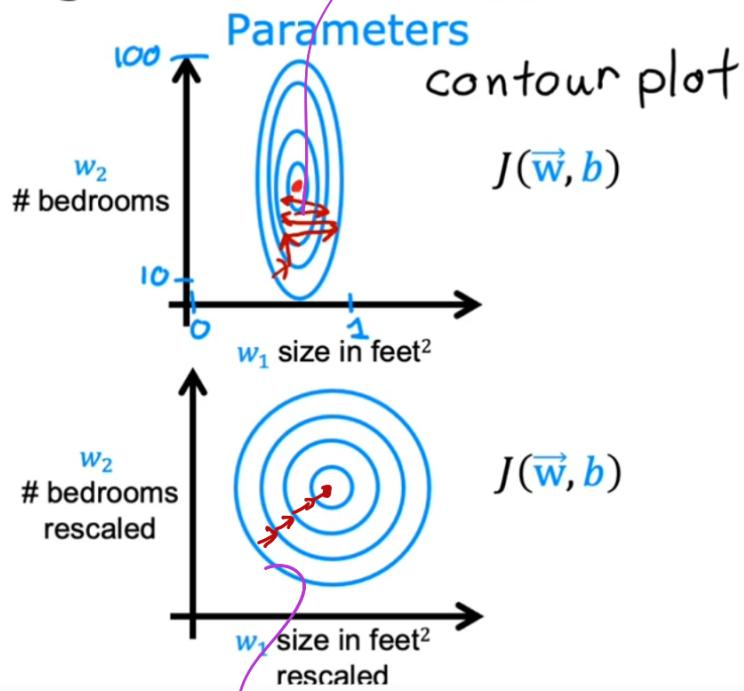
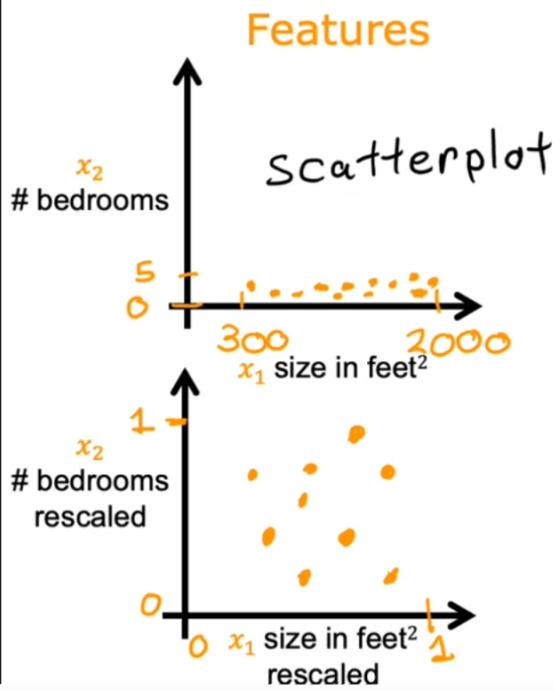
$$\widehat{\text{price}} = \$100,050.5K = \$100,050,500$$

$w_1 = 0.1, w_2 = 50, b = 50$
small large

$$\widehat{\text{price}} = \frac{0.1 * 2000K}{200K} + \frac{50 * 5}{250K} + \frac{50}{50K}$$

$$\widehat{\text{price}} = \$500K \text{ more reasonable}$$

Feature size and gradient descent



The values keeps jumping back before reaching the global minimum

The value directly reaches the global minimum.

When Feature Scaling is Used :-

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

$-3 \leq x_j \leq 3$	}	acceptable ranges
$-0.3 \leq x_j \leq 0.3$		

$$0 \leq x_1 \leq 3$$

okay, no rescaling

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

okay, no rescaling

$$-100 \leq x_3 \leq 100$$

too large \rightarrow rescale

* Anything close to -1 to 1 is acceptable else rescale it!

How to choose feature Scaling :-

1) Divide by Maximum :-

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

$$x_{1\text{scaled}} = \frac{x_1}{2000} \quad x_{2\text{scaled}} = \frac{x_2}{5}$$

$$0.15 \leq x_{1\text{scaled}} \leq 1 \quad 0 \leq x_{2\text{scaled}} \leq 1$$

2) Mean Normalization :-

Usually scaled answer ranges from -1 to 1. We center everything around zero.

e.g. Assuming μ_1 (mean x_1) = 600, μ_2 (mean x_2) = 2.3

$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$x_1 = \frac{x_1 - \mu_1}{\text{Max-Min}}$$

$$x_2 = \frac{x_2 - \mu_2}{\text{Max-Min}}$$

$$\therefore -0.18 \leq x_1 \leq 0.82$$

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

3) Z-Score Normalization :-

We calculate standard deviation for this

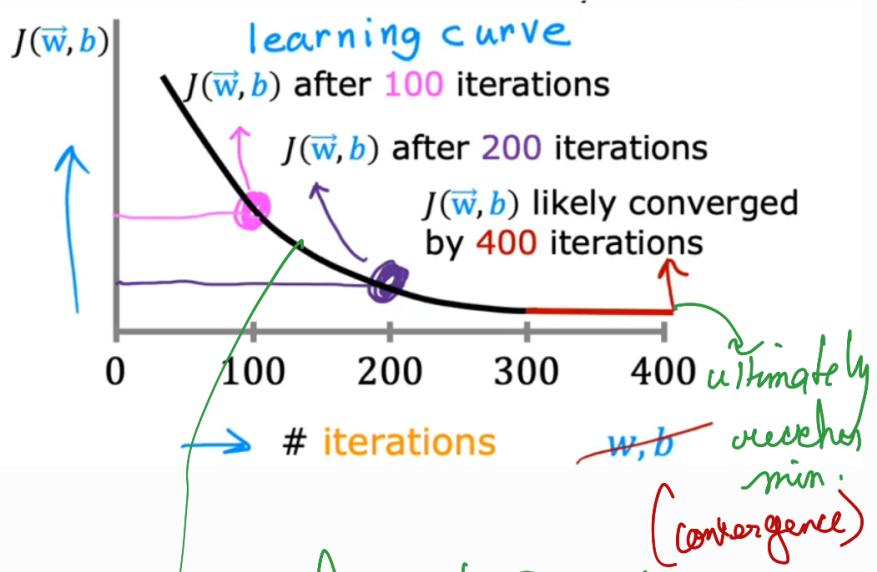
$$x = \frac{x - \mu}{\sigma}$$

$$-0.67 \leq x_1 \leq 3.1$$

$$-1.6 \leq x_2 \leq 1.9$$

Making Sure Gradient Descent is Working:-

Objective: $\min_{\vec{w}, b} J(\vec{w}, b)$



According to Gradient Descent

$J(\vec{w}, b)$ should decrease after every iteration. Hence the downward curve.

Automatic Convergence test:-

Let ϵ (epsilon) be 10^{-3}

If $J(\vec{w}, b)$ decreases $< \epsilon$
declare Convergence.

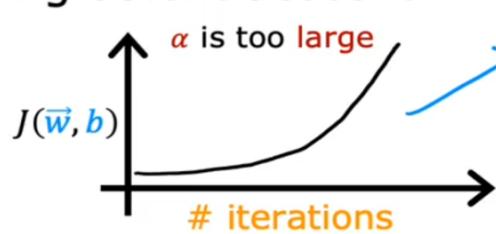
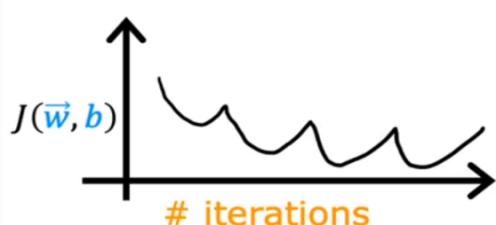
i.e., $J(\vec{w}, b)$ has reached the flat line part (red part) of the graph.

∴ We have found parameters \vec{w}, b to reach global minimum.

Learning Rate

- To confirm that we have chosen a perfect α ; plot $J(\theta)$ vs no of iteration
- if everything is right value of $J(\theta)$ decreases after every iteration
- else value of α is too large
- $J(\theta)$ can be declared as converged if the decrease in $J(\theta)$ after every iteration is less than 10^{-3} .

Identify problem with gradient descent



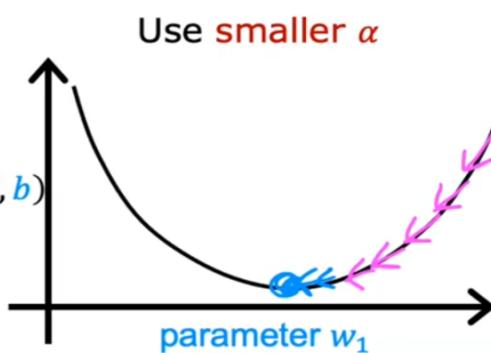
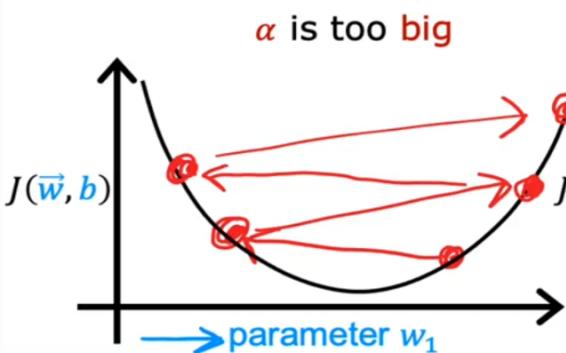
or learning rate is too large

$$w_1 = w_1 + \alpha d_1 \quad ||$$

use a minus sign

$$w_1 = w_1 - \alpha d_1 \quad ||$$

Adjust learning rate



With a small enough α , $J(\vec{w}, b)$ should decrease on every iteration

If α is too small, gradient descent takes a lot more iterations to converge

Normal Equation → Only to solve linear Regression (Usually not recommended)

→ another method to minimize J . unlike gradient descent it is a non-iterative approach.

$$\text{usually used in ML & Deep learning} \quad \theta = (X^T X)^{-1} X^T Y$$

→ eg: →

x_0	size x_1	no of bedroom x_2	no. of floors x_3	Age of house x_4	Price y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
2	852	2	1	36	178

here $X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 2 & 852 & 2 & 1 & 36 \end{bmatrix}$ $y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$

→ Here the data set is small so matrix multiplication & inverse is computationally inexpensive ; but for larger dataset gradient descent is preferable

→ feature scaling does not affect the Normal Equation procedure

-e

∴ Gradient descent vs Normal equation :

Gradient descent	Normal equation.
• Need to choose α	• no need of α
• Iterative process	• no iteration.
• Complexity : $O(kn^2)$ suited when n is large	• $O(n^3)$ { coz not to find inverse } suited when n is small preferably $< 10,000$

→ in Octave Use `pinv` for inverse so that we get value of θ if $(X^T X)$ is non-invertible.

→ Reasons due to which $X^T X$ can become non-invertible : →

- Redundant features. { i.e. 2 or more feature linearly related }
- Too many features { $m \leq n$ }

Feature engineering

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

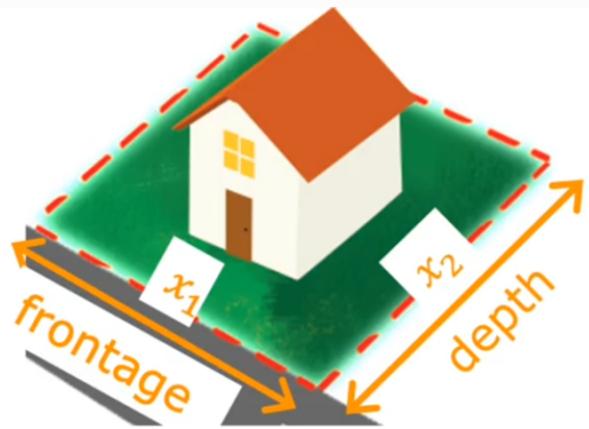
frontage depth

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

$$x_3 = x_1 x_2$$

new feature

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



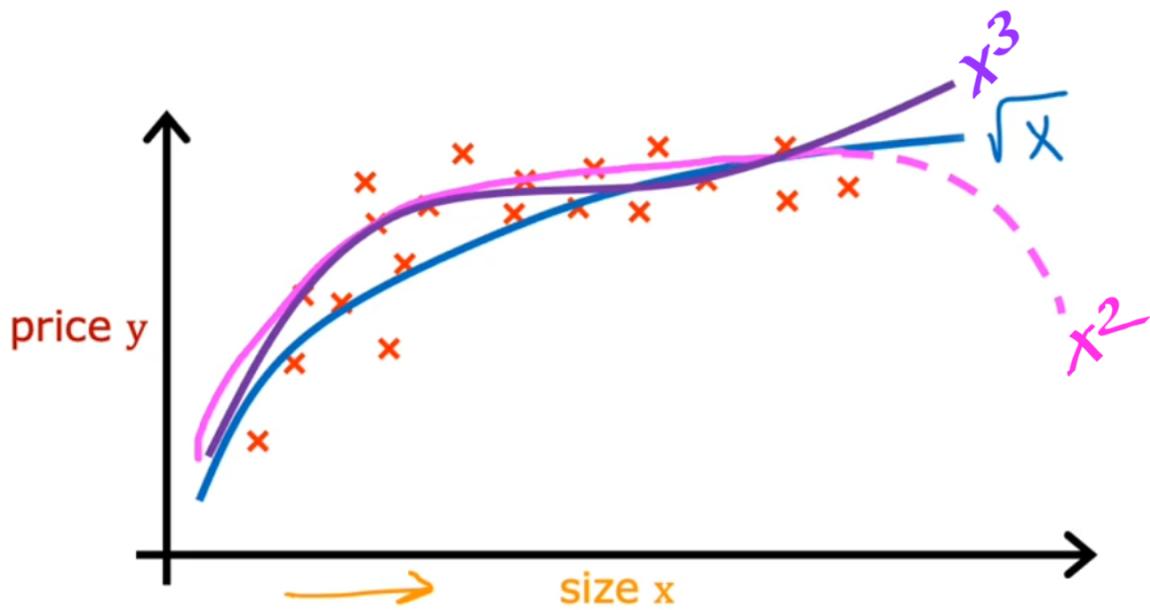
Feature engineering:
Using **intuition** to design
new features, by
transforming or combining
original features.

- * Sometimes the features we start with doesn't help us predict correctly. Hence we need to take certain new features while working. Here x_3 .
- * Feature engineering can be used to fit not just straight lines but also curves to a set of data.

Polynomial

Regression :-

-> Helps us fit curve, non linear function to the data.



⇒ $f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + b$ ⇒ We can see this as x^2 comes down in the graph plotting ↗

$\uparrow \text{size}$ $\uparrow \text{Size}^2$

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

$\uparrow \text{size}$ $\uparrow \text{size}^2$ $\uparrow \text{size}^3$

We are taking different values of $x \rightarrow$ Feature Scaling.

⇒ $f_{\vec{w}, b}(x) = w_1 x + w_2 \sqrt{x} + b$

$\uparrow \text{size}$ $\uparrow \sqrt{\text{size}}$

CLASSIFICATION

⇒ Output is discrete

⇒ The answer can only have 2 values → Yes or No.

Email is spam → Y or N

Tumour is Malignant → Y or N

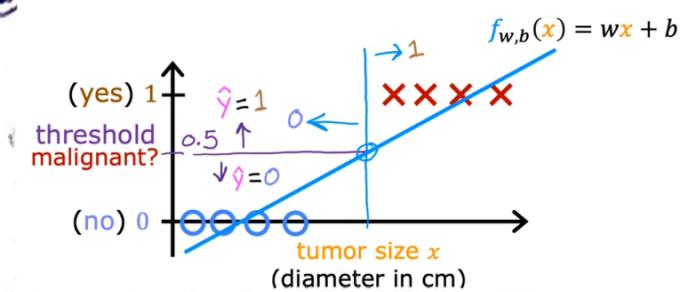
⇒ Also known as binary classification

⇒

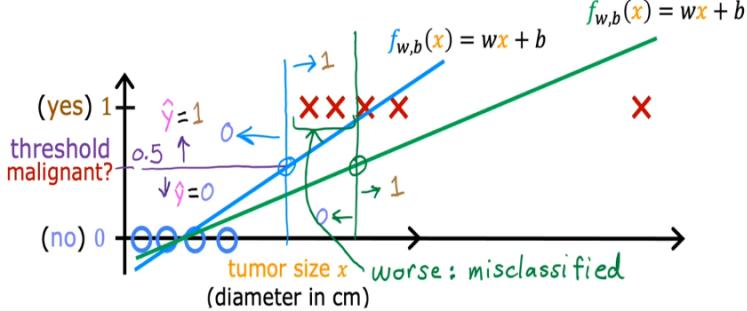
False 0 Negative Class	True 1 Positive Class
------------------------------	-----------------------------

⇒ If we try using linear Regression, we fail miserably

(ex 1)



- Any value > 0.5 is positive (1) & others are negative (0).
- Seems to work fine with this example.



- It is obvious that adding a training data set with very large Tumor size should be a Malignant tumor.
- But here our hypothesis does not work fine, starting results does not match !!

Sigmoid Function :-

- ⇒ It is one of the most used functions in Machine/Deep learning.
- ⇒ Also known as squashing function.
- ⇒ Takes input and returns values between 0 and 1.
- ⇒ S-shaped Curve (Sigmoid Curve)

$$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}, \quad \begin{array}{l} \text{Domain} \rightarrow \mathbb{R} \\ \text{Range} \rightarrow (0,1) \end{array}$$

⇒ As x approaches $+ve \infty$, $\sigma(x)$ approaches 1
 " " -ve ∞ , $\sigma(x)$ " 0

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

* Let $y = \sigma(x) = \frac{1}{1+e^{-x}}$

Let $u = 1 + e^{-x}$. Thus, $y = 1/u$.

First, find the derivative of u with respect to x :

$$du/dx = -e^{-x}.$$

Then, find the derivative of y with respect to u :

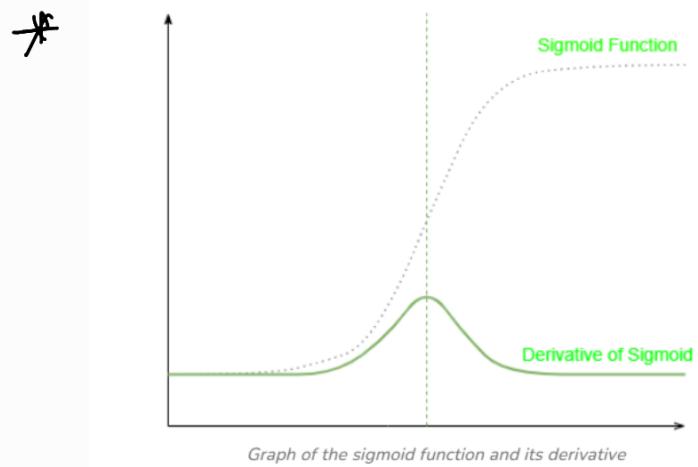
$$\frac{dy}{du} = -\frac{1}{u^2}$$

Apply the chain rule:

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx} = -\frac{1}{u^2} \cdot (-e^{-x}) = \frac{e^{-x}}{(1+e^{-x})^2}$$

$$\sigma(x) = \frac{1}{1+e^{-x}}, 1-\sigma(x) = \frac{e^{-x}}{1+e^{-x}}$$

$$\text{Thus, } \sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \left(\frac{1}{1+e^{-x}}\right) \left(\frac{e^{-x}}{1+e^{-x}}\right) = \sigma(x)(1-\sigma(x))$$



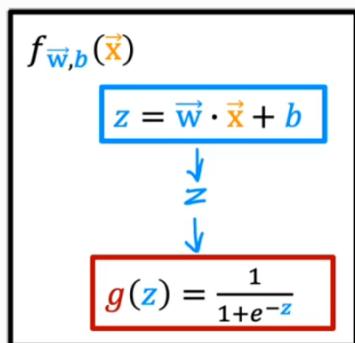
Using Sigmoid for logistic Regression :-

We know $\sigma(z) = \frac{1}{1+e^{-z}}$

We have seen in linear Regression

$$z = f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

Now we put this in the sigmoid function,



$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

∴ Logistic Regression,

$$g(z) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

Interpretation of logistic regression output :-

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"probability" that class is 1

$$f_{\vec{w}, b}(\vec{x}) = P(y = 1 | \vec{x}; \vec{w}, b)$$

Probability that y is 1,
given input \vec{x} , parameters \vec{w}, b

Example:

x is "tumor size"

y is 0 (not malignant)
or 1 (malignant)

$$P(y = 0) + P(y = 1) = 1$$

$$f_{\vec{w}, b}(\vec{x}) = 0.7$$

70% chance that y is 1

Decision Boundary :-

⇒ Decision Boundary is the line on the map which helps us in classification.

⇒ H is the threshold line (0.5)

$$f_{\vec{w}, b}(\vec{x}) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

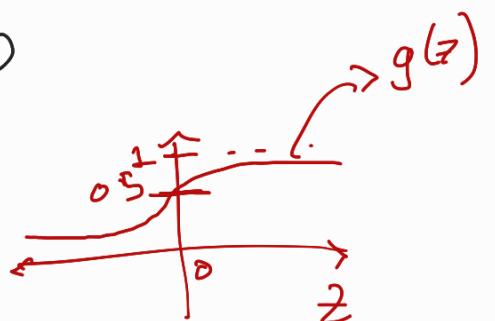
$$= P(y=1 | \vec{x}; \vec{w}, b)$$

If $f_{\vec{w}, b}(\vec{x}) > 0.5$,

Yes : $\hat{y} = 1$ No : $\hat{y} = 0$

when $z > 0.5$

$g(z) > 0.5$



when $z > 0$

We can see from the graph

$$\vec{w} \cdot \vec{x} + b > 0$$

$$\therefore \hat{y} = 1$$

$$\vec{w} \cdot \vec{x} + b < 0$$

$$\hat{y} = 0$$

Different example for Decision Boundary :-

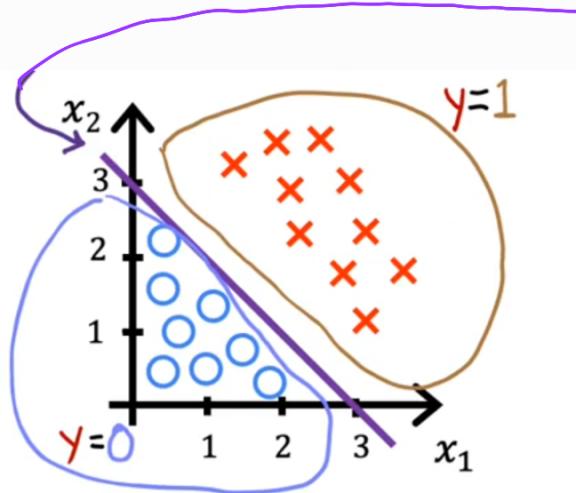
$$\Rightarrow f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1x_1 + w_2x_2 + b)$$

Assume

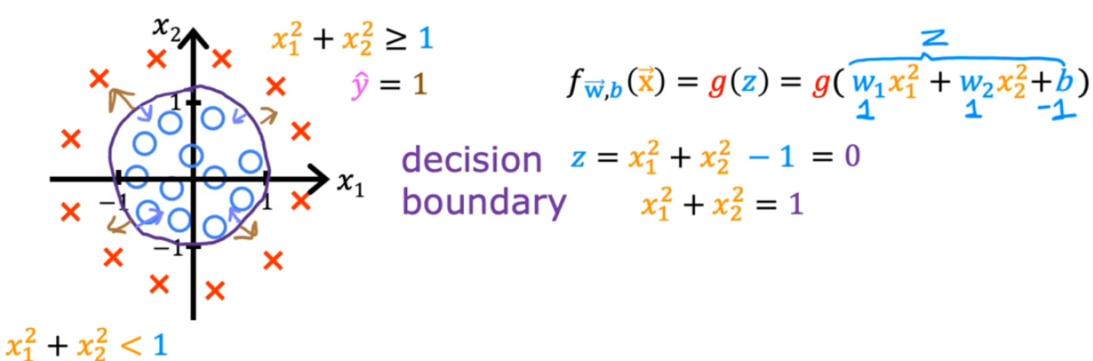
Decision Boundary. $z = \vec{w} \cdot \vec{x} + b = 0$

$$z = x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$

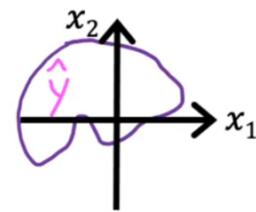


Non Linear Decision Boundary —



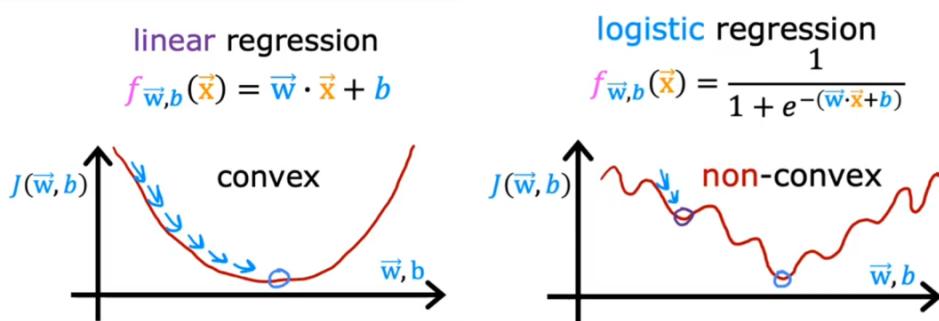
\Rightarrow Decision Boundary can be a even more complex shape.
It depends on the no. of features and their degree.

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_1x_2 + w_5x_2^2 + w_6x_1^3 + \dots + b)$$



Cost Function:-

- * The cost function that we used in linear regression \rightarrow Squared Error Cost cannot be used for Regression -
- * If we use the Squared Error Cost we get a non-convex curve that has a lot of local minima and we keep getting stucked in it.



Squared Error Cost $\rightarrow \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{2} (f_{\vec{w}, b}(\vec{x}_i) - y_i)^2 \right]$

we take this as a new function

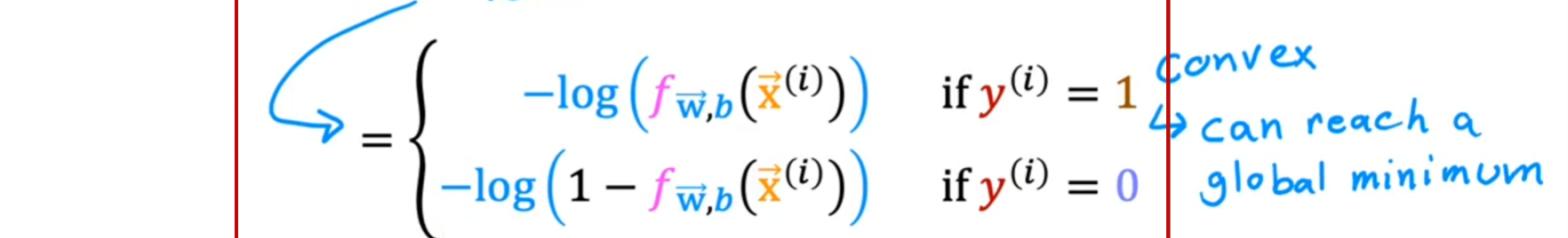
$$\text{Loss} \rightarrow \begin{cases} f_{\vec{w}, b}(\vec{x}_i), y(i) \end{cases}$$

Logistic loss function:-

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

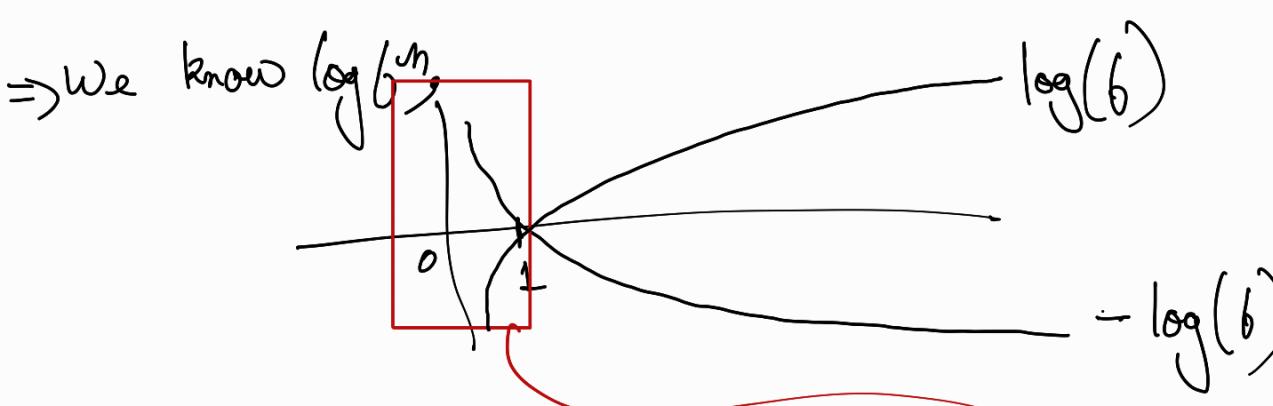
cost

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



$$= \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

convex
 ↳ can reach a global minimum

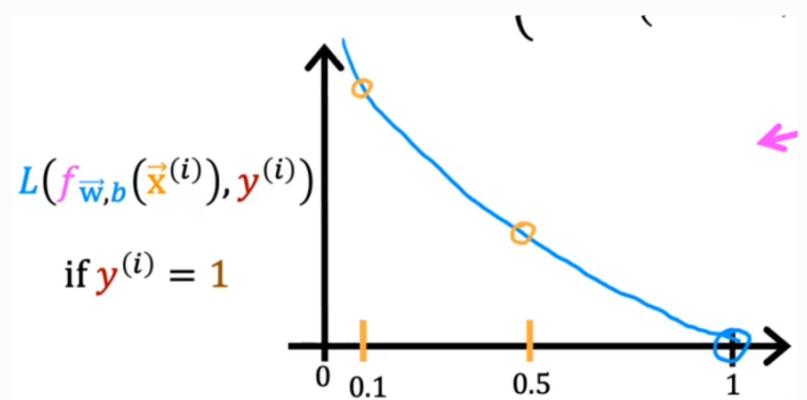


\Rightarrow In logistic regression answers can be only between 0 and 1 and hence we will only focus on that part.

\Rightarrow If $y^{(i)} = 1$,

If $f_{\vec{w}, b}(\vec{x}_i) \rightarrow 1$ then,
loss $\rightarrow 0$

If $f_{\vec{w}, b}(\vec{x}_i) \rightarrow 0$ then
loss $\rightarrow \infty$.

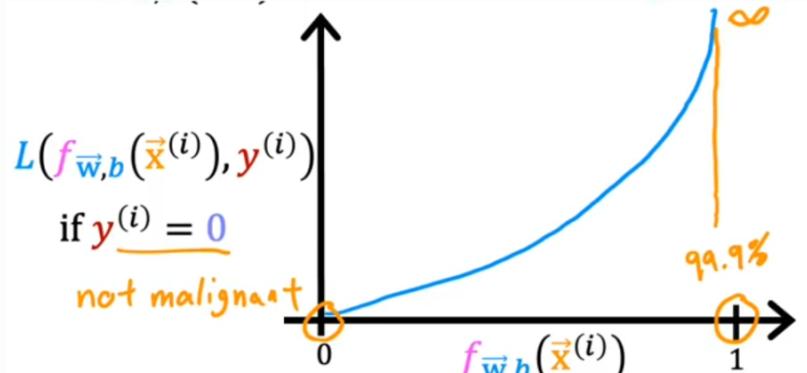


\therefore Loss is lowest when $f_{\vec{w}, b}(\vec{x}_i)$ predicts closest to the true value of $y^{(i)}$

\Rightarrow If $y^{(i)} = 0$,

If $f_{\vec{w}, b}(\vec{x}_i) \rightarrow 1$ then,
loss $\rightarrow \infty$

If $f_{\vec{w}, b}(\vec{x}_i) \rightarrow 0$ then
loss $\rightarrow 1$



\therefore The further the prediction $f_{\vec{w}, b}(\vec{x}^{(i)})$ is from target $y^{(i)}$, the higher the loss.

Simplified loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$\Rightarrow L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

⇒ if $y^{(i)} = 1$:

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -1 \log(f_{\vec{w}, b}(\vec{x}))$$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

⇒ if $y^{(i)} = 0$:

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -(1 - 0) \log(1 - f_{\vec{w}, b}(\vec{x}))$$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

Simplified cost function

$$loss \quad L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

$$cost \quad J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})]$$

convex

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

Gradient descent

cost

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right]$$

repeat {
 $j = 1 \dots n$
 $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$
 $b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$
} simultaneous updates

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

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

Gradient descent for logistic regression

repeat {

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

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

} simultaneous updates

↑

We might think that this looks similar to linear regression but the $f_{\vec{w}, b}(\vec{x}^{(i)})$ is different for both.

Linear regression $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

Logistic regression $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

Same concepts:

- Monitor gradient descent (learning curve)
- Vectorized implementation
- Feature scaling

These are same as linear regression

Multiclass Classification

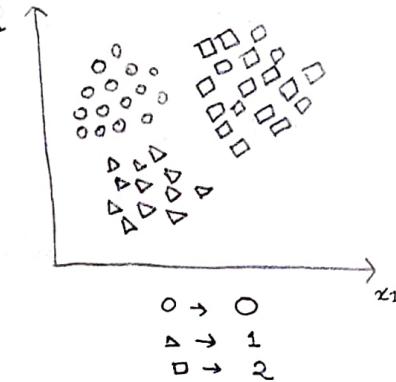
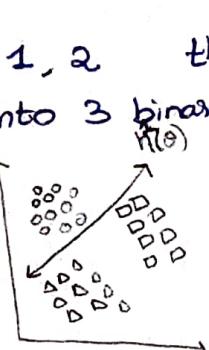
→ Approach for multiple classes ; i.e $y = 0, 1, 2, 3 \dots n$ {multiple classes but still discrete}

→ For $y = 0, 1, 2 \dots n$ we divide the problem into smaller $n+1$ binary classifications.

eg: → If $y = 0, 1, 2$ then :→
we divide this into 3 binary classifications

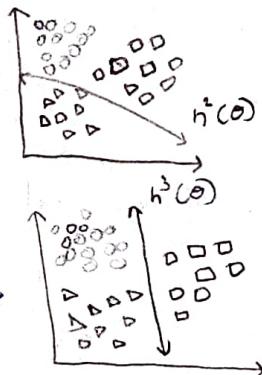
① for class 0 :

from this we get
 $h^1(\theta)$



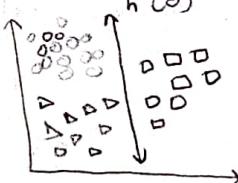
② for class 1 :

from this we get
 $h^2(\theta)$



③ for class 2 :

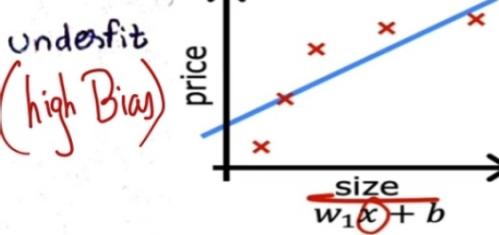
from this we get
 $h^3(\theta)$



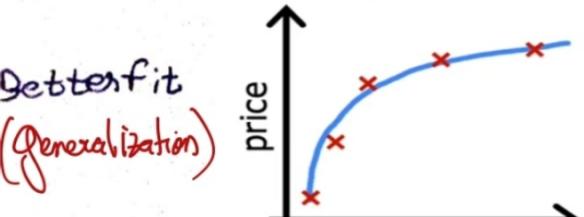
By combining $h^1(\theta), h^2(\theta)$ and $h^3(\theta)$ we can classify this multi class classification

$$\text{Prediction} = \max_i (h_i(\theta))$$

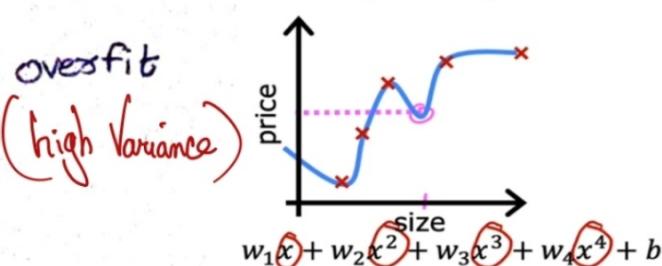
Overfitting :-



using $y = \theta_0 + \theta_1 x$
Does not pass through any points. of training set
errors on training set
errors in real time values

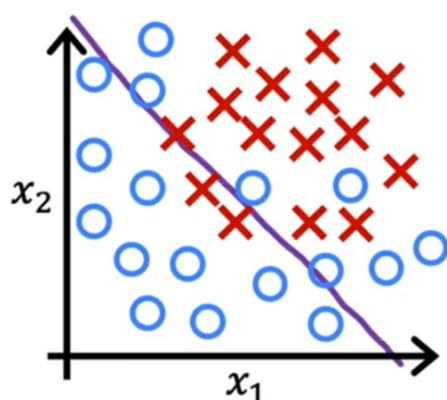


adding an extra x^2 feature
 $y = \theta_0 + \theta_1 x + \theta_2 x^2$
Decent fit; Passes through almost all points.
Decent on training set
Decent on real time values



adding too many features
 $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_1 x_2^2 + \theta_4 x_1^2 x_2$
Perfect fit; cost = 0.
Perfect on training set
errors in real time values

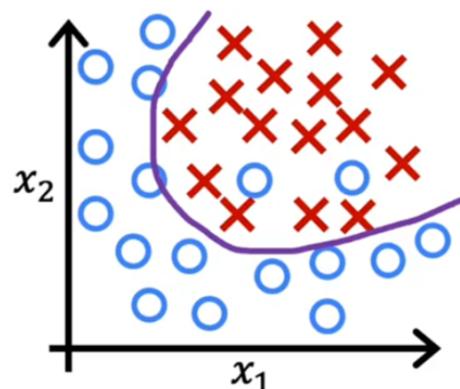
Classification



$$z = w_1 x_1 + w_2 x_2 + b$$

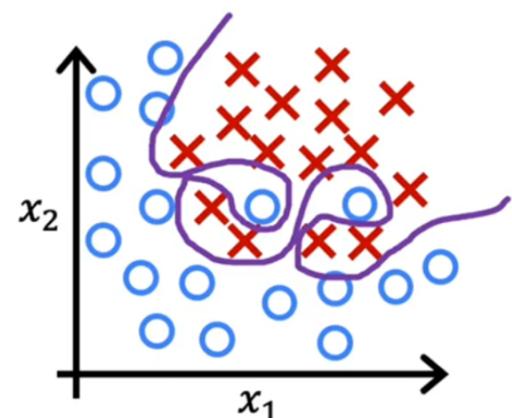
$$f_{\vec{w}, b}(\vec{x}) = g(z)$$

g is the sigmoid function
underfit high bias



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + b$$

just right



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + w_5 x_1^2 x_2^3 + w_6 x_1^3 x_2 + \dots + b$$

overfit

The problem with overfitting is we try hard to fit the training dataset due to which we loose the generalized behaviour of curve. So we have cost 0 for training set but while calculating real time values it shows very high cost.

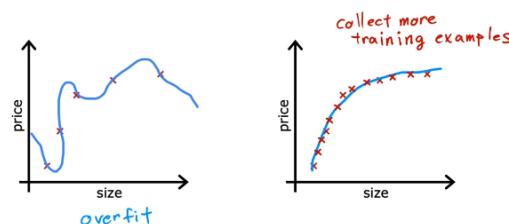
Loose generalized behavior

Reducing Overfitting :-

- 1) Collect more Data.
- 2) Reduce the no. of features.
- 3) Regularize the features.

Collect More Data:-

- We can always collect more data to convert an overfit data to a normal Data.
- The only problem is that sometimes more data is just not available.

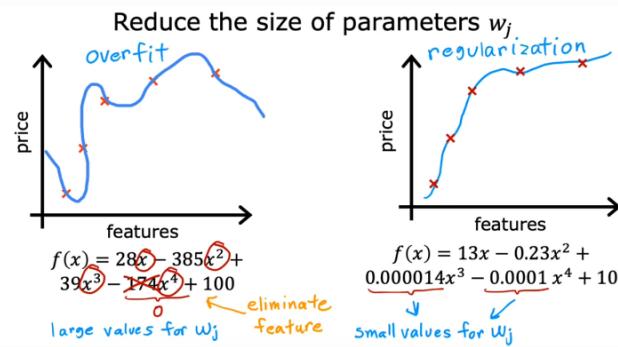


Reduce the no. of features:-

- If there are too many features and we have insufficient data then it will result in Overfitting.
- We can reduce it by taking a subset of the features.
- Disadv → As we are taking a limited no. of features and not all of them, we may lose out on important features.

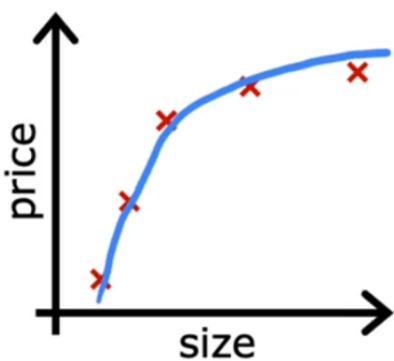
Regularization :-

- Instead of just simply omitting a feature we can take a smaller (Regularized value) of the parameter.

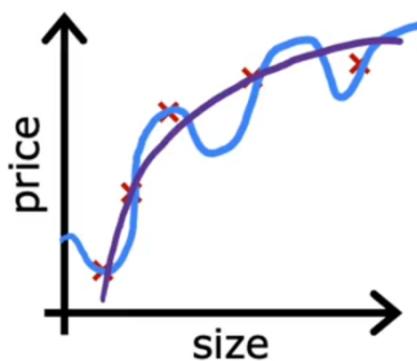


Regularization :-

Intuition



$$w_1x + w_2x^2 + b$$



$$w_1x + w_2x^2 + \cancel{w_3x^3} + \cancel{w_4x^4} + b$$

≈ 0 ≈ 0

make w_3, w_4 really small (≈ 0)

$$\min_{\vec{w}, b} \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + 1000 \frac{w_3^2}{0.001} + 1000 \frac{w_4^2}{0.002}$$

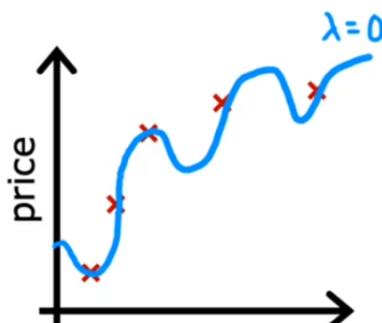
$$J(\vec{w}, b) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \underbrace{\lambda \sum_{j=1}^n w_j^2}_{\text{"lambda"} \text{ regularization parameter}} + \underbrace{\frac{\lambda}{2m} b^2}_{\text{regularization term}} \right]$$

can include or exclude

Choosing λ :-

If $\lambda=0$, the regularization term is zero and we end up with the same overfitted data.

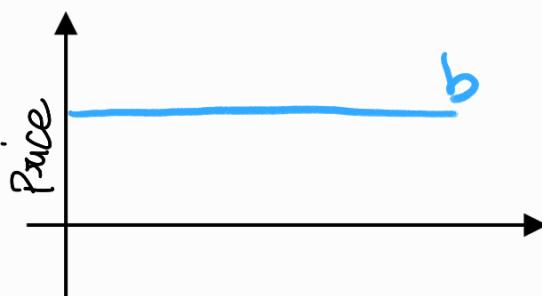
$$f_{\vec{w}, b}(\vec{x}) = w_1x + w_2x^2 + w_3x^3 + w_4x^4 + b$$



If $\lambda=\text{Very high } (10^{10})$, then the parameters have to be very very small and tends to 0.

$$f_{\vec{w}, b}(\vec{x}) = \cancel{w_1x} + \cancel{w_2x^2} + \cancel{w_3x^3} + \cancel{w_4x^4} + b$$

$$\therefore f(x) = b$$



Regularized linear regression

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right]$$

Gradient descent

repeat {

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

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

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

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

} simultaneous update

don't have to regularize b

(as if w_0 optional)

Implementing gradient descent

repeat {

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

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

} simultaneous update

Regularized logistic regression

$$\min_{\vec{w}, b} J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Gradient descent

repeat {

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

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

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

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

Looks same as for linear regression!

logistic regression

don't have to regularize