

m : size of training set

here :

x : input param
 y : output (given)

h : hypothesis function

h is used to predict ' y ', given ' x ', after learning from m data sets.

$$h(x) \rightarrow \text{straight line} \Rightarrow h(x_i) = \theta_0 + \theta_1 x_i$$

~~needs~~

$x_i \rightarrow i^{\text{th}}$ training set input

$y_i \rightarrow i^{\text{th}}$ training set output (given)

$h(x_i) \rightarrow i^{\text{th}}$ training set (predicted)

h should be as close

h should have minimal error, wrt. y .

\Rightarrow

let

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

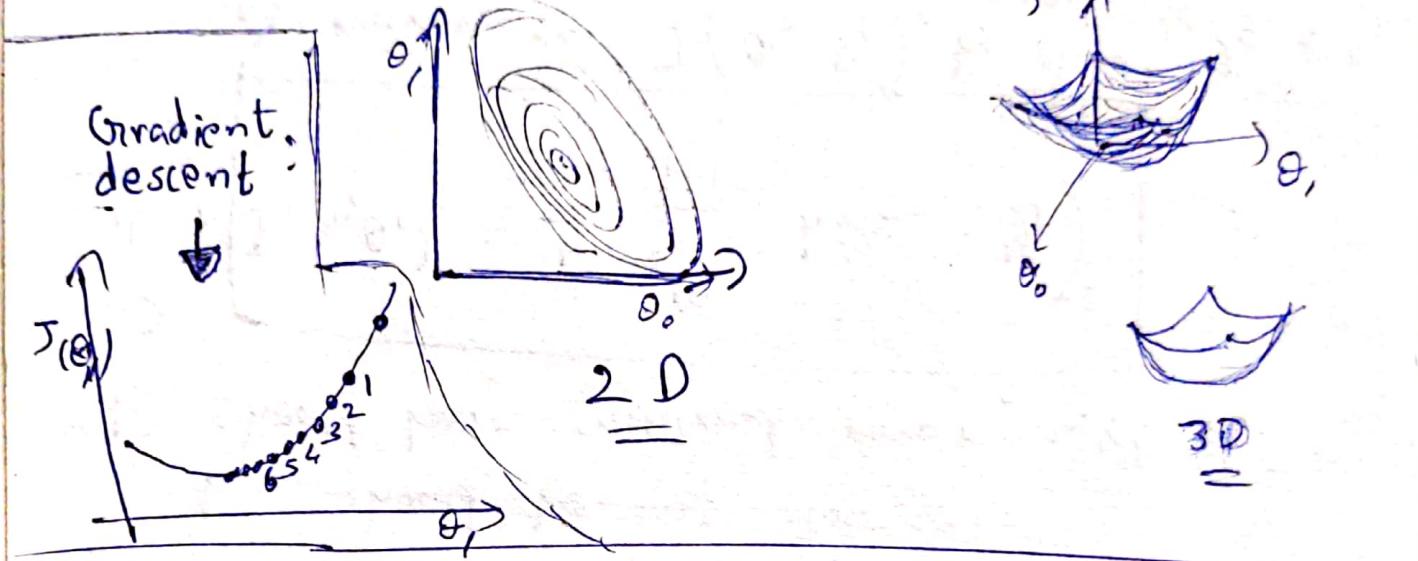
\Rightarrow minimise $J(\theta_0, \theta_1)$, by choosing appropriate θ_0 & θ_1 .

Q: Why is $\frac{1}{2}$ there?

Answer: In order to cancel with square term, upon differentiation.

Contour Plot

when $J(\theta_0, \theta_1)$ is plotted as a fn of θ_0 & θ_1
 the plot is a 'bowl' shaped in 3D,
 and 'contours' shaped in 2D--



How to find $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$? by using

Gradient Descent

Algo :-

- 1) Start with some θ_0, θ_1 ,
- 2) Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until minima is reached.

mathematically speaking :-

repeat until convergence:-

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}, \quad \text{for } j = 0, 1$$

here, $\alpha \Rightarrow$ learning rate

$\hat{=}$ sign \Rightarrow ~~continuous~~ simultaneous update.

Simultaneous Update :- $\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$$\text{temp}_0 = \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}_1 = \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 = \text{temp}_0$$

$$\theta_1 = \text{temp}_1$$

i.e. both θ_s should get updated simultaneously. $\because J$ is a function of two variables (θ_0, θ_1) .

Learning Rate (α) :-

- ① If α is too small, gradient descent can be slow.
- ② If α is too large, gradient descent can overshoot the minima.
It can fail to converge, or even diverge.
- ③ we can keep α fixed, and still reach a local minimum because as we tend converge towards local minimum, gradient descent will automatically take smaller steps.

Now, we know that the gradient descent algorithm is :-

repeat until
(convergence)

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j}, \quad \forall j \in 0, 1$$

↳ ①

We also know :-

Linear regression model: $h(x) = \theta_0 + \theta_1 x$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

from

↳ ②

from ① & ② :-

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \left(\frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2 \right)$$

$$\theta_0 + \theta_1 x_i$$

⇒

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} \left(\frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i - y_i)^2 \right)$$

⇒ ~~for each of~~ upon applying partial differentiation on ~~each of~~ $j=0$ & $j=1$ terms, we gets

we get:-

Gradient descent Algo for Linear Regression:

(in one variable)

repeat until convergence :-

$$\theta_0 := \theta_0 - \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)$$

$$\theta_1 := \theta_1 - \frac{1}{m} \sum_{i=1}^m ((h(x_i) - y_i) * x_i)$$

}

Batch Gradient Descent :-

Each step uses all training set examples.

ex:- the above algorithm is a batch gradient descent.

Linear Regression in Multiple Variables :-

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

where

$n \rightarrow$ no. of features

For uniformity's sake let us introduce

a new variable $x_0 = 1$

$$\Rightarrow h(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Gradient Descent for linear Regression (Multiple variables)

repeat until convergence:-

{

$$\theta_0 := \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m \left[(h(x_i) - y_i) x_{0i} \right]$$

$$\theta_1 := \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m \left[(h(x_i) - y_i) x_{1i} \right]$$

$$\theta_2 := \theta_2 - \frac{\alpha}{m} \sum_{i=1}^m \left[(h(x_i) - y_i) x_{2i} \right]$$

\vdots
 θ_n

}

i Tips

① If gradient descent is not working for certain α , decrease it and try again.

For sufficiently small α , J should decrease with every iteration.

② α too small \Rightarrow slow convergence

α too large \Rightarrow no convergence
(or) even divergence!

③ You can declare convergence if

J decreases by less than ϵ in one iteration. $\epsilon \rightarrow$ user defined (ex: 10^{-3})

Polynomial Regression:-

TIP :- we can combine multiple features into one.
ex:- $x_3 = x_1 \cdot x_2$

Polynomial Regression:- our hypothesis function need not be linear, if it doesn't fit the data well.

func

Linear :- $\theta_0 + \theta_1 x_1$

Quadratic :- $\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$

Cubic :- $\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$

SquareRoot :- $\theta_0 + \theta_1 x_1 + \theta_2 \sqrt{x_1}$

Let $x_2 = x_1^2$ & $x_3 = x_1^3$,

⇒ we have created 2 new features in cubic function.

④ If we square (or) cube a variable, then its range of values also become squared / cubed.

In order to make plots easier, we Scale the features into appropriate range.

↳ upto you :-

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

$$-0.5 \leq x_i \leq +0.5$$

Feature Scaling & Mean Normalization }

We use these two techniques to reduce the range of each feature.

Why? : → ~~Faster~~ plots made easier
→ Fastly arrive at minima.

$$x_i := \frac{x_i - \mu_i}{\sigma_i}$$

mean normalization

$$x_i - \mu_i$$

mean

Feature Scaling

$$\downarrow \text{dividing by } \sigma_i$$

standard deviation.

If calculating σ_i is difficult,

~~get~~ replace it by range of that x_i

(max. value of x_i - min. value of x_i)

Classification :-

Output values are discrete and finite. (and small number)

Linear Regression doesn't work well in classification problems, since classification is not a linear function.

Binary classification :- output takes 2 values
(0 (or) 1)

In order to make Linear Regression work for Binary classification, we use a different form of hypothesis function, called:-
Logistic (or) Sigmoid Function.

defined as : $g(z) = \frac{1}{1+e^{-z}}$

we can plug $\theta^T x$ into z , to get our new hypothesis function :-

$$\Rightarrow h_{\theta}(x) = g(\theta^T x)$$
$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

where $\theta^T \rightarrow$ Transpose of vector :-
i.e. $\underline{\theta^T x} = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

we can also say that :-

$$h_{\theta}(x) = P(y=1 | x; \theta) = \underline{1 - P(y=0 | x; \theta)}$$

read as :- Probability that $y=1$, given x , parameterized by θ .

meaning :-

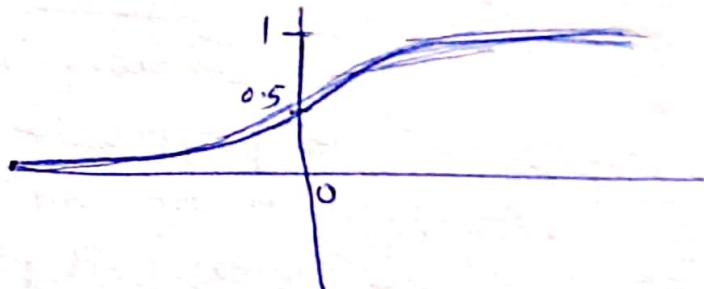
$h_{\theta}(x)$ will give probability that o/p is 1.

ex:- if $h_{\theta}(x) = 0.7 \Rightarrow$ Probability of $y=1$ is 0.7
(70%)

Logistic function (or)

Sigmoid function & its behaviour :-

$$g(z) = \frac{1}{1+e^{-z}}$$



From the figure, we can say that:-

$$\begin{aligned} g(z) > 0.5 &\Rightarrow y=1 \\ g(z) < 0.5 &\Rightarrow y=0 \end{aligned}$$

$$g(z) \geq 0.5 \text{ when } z \geq 0 \rightarrow ①$$

when:-

$$\text{also, } z=0, g(z) = \frac{1}{2}$$

$$z=+\infty, g(z) = 1$$

$$z=-\infty, g(z) = 0$$

In order to fit this function into our Binary classification, we can translate our o/p of hypothesis 'function':-

$$\left\{ \begin{array}{l} \text{If } h_0(x) \geq 0.5 \Rightarrow y = 1 \\ \text{If } h_0(x) < 0.5 \Rightarrow y = 0 \end{array} \right. \rightarrow ②$$

now, replacing z with $\theta^T x$ in ①, we get!

$$\theta^T x \geq 0.5 \quad \text{when } \theta^T x \geq 0.5$$

③

from ② + ③, we get:-

$$\left\{ \begin{array}{l} \theta^T x \geq 0 \Rightarrow y = 1 \\ \theta^T x < 0 \Rightarrow y = 0 \end{array} \right.$$

example :-

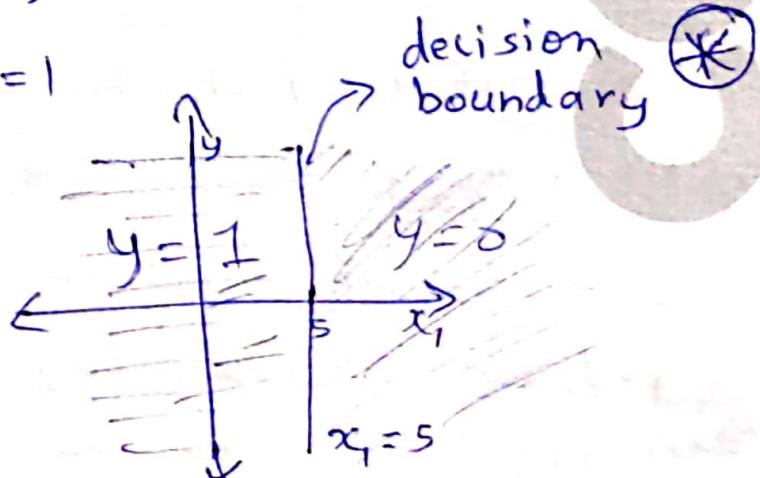
$$\theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}$$

$$\Rightarrow \theta^T x = 5 - x_1 \quad \theta^T x \geq 0 \Rightarrow y = 1$$

$$\Rightarrow 5 - x_1 \geq 0 \Rightarrow y = 1$$

$$\Rightarrow x_1 \leq 5 \Rightarrow y = 1$$

plotting this, we get:-



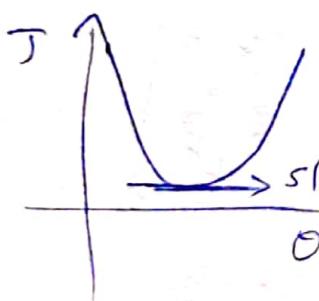
Decision boundary need not be a straight line
 it was a straight line because we put
 $z = \theta_0 + \theta_1 x_1$, . we can put anything in z ,
 as long as we get a decision boundary,
that separates $y=0$ & $y=1$ regions.

ex: $z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2 \Rightarrow$ circle shape

Alternative method to get minima :- Normal Equation.

Normal Equation:

Alternative to gradient descent.



$$\text{slope} = 0 \text{ at } \underline{\text{minima}} \Rightarrow \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} = 0$$

so, minima is directly found by calculating derivative of J , and equating it to 0.

Thus, we get required θ_j values.

Normal Equation
Formula :

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

, where

where $X = \begin{bmatrix} x_0 & x_1 & \dots & x_n \\ x_0 & x_1 & \dots & x_n \\ x_0 & x_1 & \dots & x_n \\ x_0 & x_1 & \dots & x_n \end{bmatrix}$ \rightarrow ~~m rows~~ m rows
~~n+1 columns~~ n+1 columns

and $x_0 = 1$

Gradient Descent

- * feature scaling sometimes necessary
- * no need to choose α
- * many iterations needed.
- * $O(Kn^2)$ Time complexity
Preferred method
- * for large n ($\approx 10^4$)

Normal Equation

no need of feature scaling

no need to choose α
no iterations

$O(n^3)$, need to calculate
 $(X^T X)^{-1}$

Preferred if n is small ($\approx 10^4$)

Note :- Sometimes $X^T X$ may be noninvertible

Solution :- 1) delete redundant features.
i.e. closely related (Linearly dependent),
ex :- X_1 : feet
 X_2 : metre

2) Too many features ($m \leq n$) :-

delete some features
(or)

use Regularization

Logistic Regression Model:

Cost Function:

Linear Regression : $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_\theta(x_i) - y_i)^2$

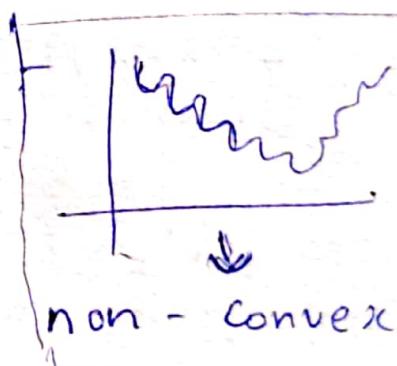
Let this be
 $\text{Cost}(h_\theta(x_i), y_i)$

so, for linear regression, cost function is:

$$\text{Cost}(h_\theta(x), y) = \frac{1}{2} (h_\theta(x) - y)^2$$

but, for logistic regression, this square-error formula cannot be used, because it has a lot of local minima, like this:

plot of $\text{Cost}(h_\theta(x), y)$, where
 $h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$



So gradient descent won't converge to global minimum.



In order to solve this problem, we have to make cost fn concave :-

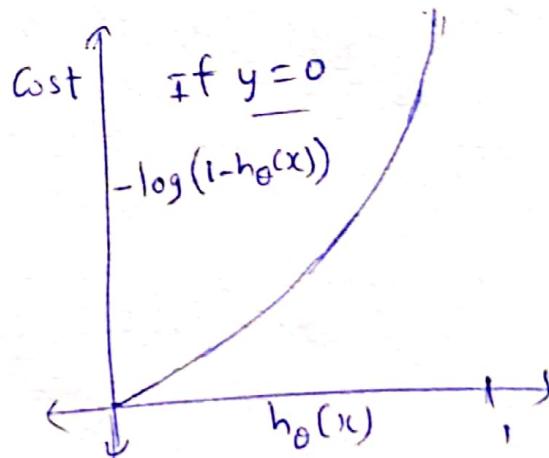
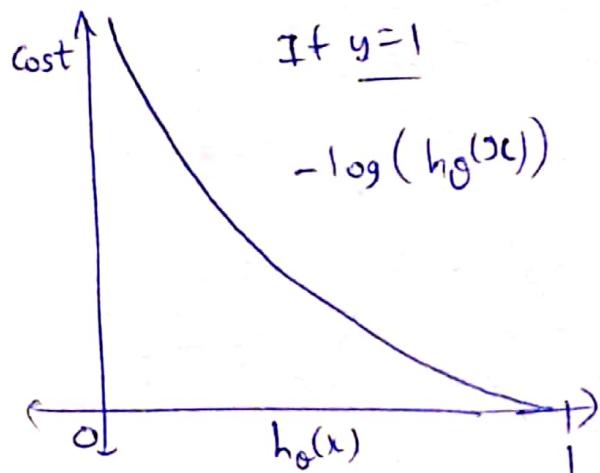
$$\text{cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)), & \text{if } y=1 \\ -\log(1-h_\theta(x)), & \text{if } y=0 \end{cases}$$

logistic Regression's Cost Function :-

↓
one liner
↓

$$\boxed{\text{cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))}$$

This yields a concave function.



Advantages of this kind of cost function:-

~~Penalize the training set example if hypothesis function is~~

cost = 0, if $y=1$ & $h_\theta(x) = 1$
& if $y=0$ & $h_\theta(x) = 0$

But as $h_\theta(x) \rightarrow 0$, cost $\rightarrow \infty$ ($y=1$)
& $h_\theta(x) \rightarrow 1$, cost $\rightarrow \infty$ ($y=0$)

This ensures that we penalize the learning algorithm by a large cost, if we get incorrect prediction.

$$\text{so, } J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[y_i \log(h_\theta(x_i)) + (1-y_i) \log(1-h_\theta(x_i)) \right]$$

for logistic regression

Gradient Descent :-

$$\text{Repeat: } \left\{ \theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j} \right\}$$

3

Partial

upon taking derivative, we get :-

Repeat :-

$$\left. \begin{array}{l} \theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) x_{ij} \end{array} \right\}$$

There are some faster ways of obtaining θ , than gradient descent :-

- 1) Conjugate gradient
- 2) BFGS
- 2) L-BFGS.

Multiclass classification:

classifying data into more than 2 categories.

Classification

Binary Multiclass

$$\text{i.e. } y \in \{0, 1, 2, \dots, n\}$$

set \hookrightarrow n+1 categories

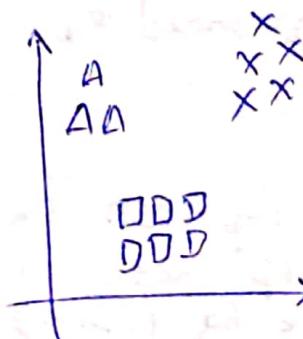
Steps:- how to classify? :-

① Compute $h_\theta^i(x)$, $\forall i \in \{0, 1, 2, \dots, n\}$

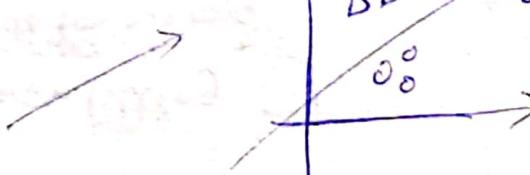
$$\text{so, } h_\theta^i(x) = P(y=i|x; \theta)$$

② Prediction = $\max(h_\theta^0(x), h_\theta^1(x), \dots, h_\theta^n(x))$

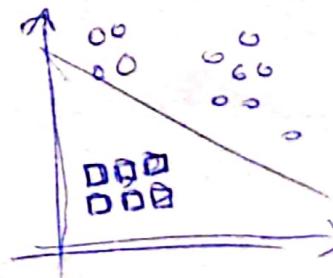
Its corresponding class is the answer.



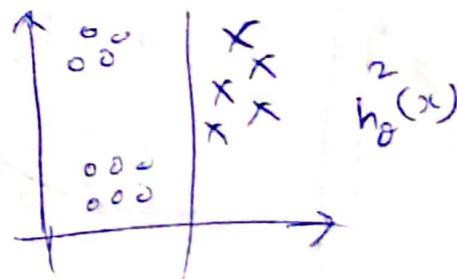
(class 1: $y=0$)
(class 2: $y=1$)
(class 3: $y=2$)



$$h_\theta^0(x)$$

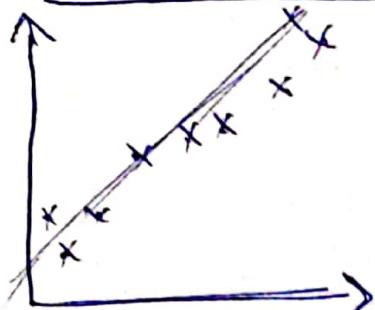


$$h_\theta^1(x)$$

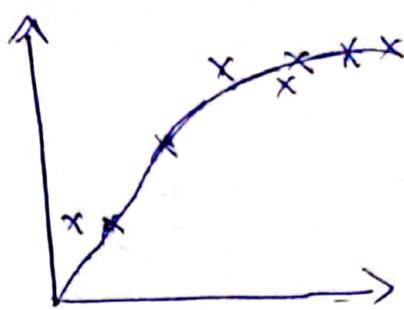


$$h_\theta^2(x)$$

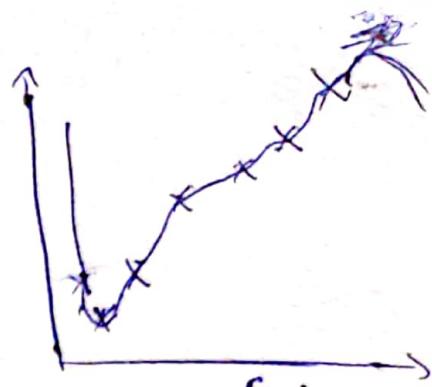
Overfitting :-



under fitting



best fit



overfitting

caused by a fn that
is either too simple,
or uses too few features.

fits the given
data very well,
but does not
generalize well
to predict new
data.

To prevent overfitting :-

1) Reduce #features :-

- manually select which feature to remove
- use model selection algorithm.

2) Regularization :-

→ keep all the features, but reduce the magnitude
of parameters θ .

→ Regularization works well when we have
a lot of slightly useful features.

Regularization :-

This is done in order to prevent
a few parameters dominating the remaining parameters.

How do parameters dominate, in
the first place? :-

At the beginning of gradient
descent, we will assume a hypothesis function &
give some initial values to each θ_j , by using

gaussian distribution. Then we run gradient descent, which will assign new values to each θ_j .

Sometimes, few θ_j will have large values, compared to the remaining θ_j . This means domination.

In order to reduce domination, we penalize all ~~parameters~~ θ_j parameters, by adding a new term, called Regularization Term, into the cost function. Regularization term also solves overfitting by, flattening smoothening the curve, thereby making it more generalized.

When does overfitting occur :-

If we have a small data set ('m' is small), then we can easily come up with a curve to fit all values in the given data set. But it incorrectly predicts the new values.

Regularization steps :-

Once you have fixed the parameters (area, length, weight, etc) and finalized some rough form of hypothesis function, you can do regularization by using the following cost function :-

$$\min_{\theta} \frac{1}{2m} \sum_{j=1}^m (h_{\theta}(x_j) - y_j)^2 + \lambda \sum_{j=1}^n \theta_j^2$$

$\lambda \rightarrow$ regularization parameter.

↳ usually 0-1

↳ if it is too large, under-fitting may happen.

note that θ_0 is not optimized penalized.

Regularization Formulae :-

Linear Regression :-

Gradient Descent :-

Repeat:-

$$\theta_j := \theta_j \left(1 - \frac{\alpha \lambda}{m}\right) - \frac{\alpha}{m} \sum_{i=1}^m ((h_\theta(x_i) - y_i) x_{ij})$$

3, $j \in \{0, 1, 2, \dots, n\}$

Note:- ① $1 - \frac{\alpha \lambda}{m} < 1$, always.

$\Rightarrow \theta_j$ is reduced by small amount on every update.

② remaining term \rightarrow no change.

Normal Equation :-

$$\theta = (X^T X + \lambda L)^{-1} X^T y, \text{ where } L = \begin{bmatrix} 0 & & \\ & \ddots & \\ & & 0 \end{bmatrix}$$

$(n+1) \times (n+1)$ matrix

Note:-

If $m < n$, $(X^T X)^{-1} \rightarrow$ non-invertible.

But, if we add λL term to it,

$(X^T X + \lambda L)^{-1} \rightarrow$ invertible

(i.e. inverse is possible)



Logistic Regression :-

cost function :-

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left(y_i \log(h_\theta(x_i)) + (1-y_i) \log(1-h_\theta(x_i)) \right) \\ + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Gradient Descent :-

Repeat {

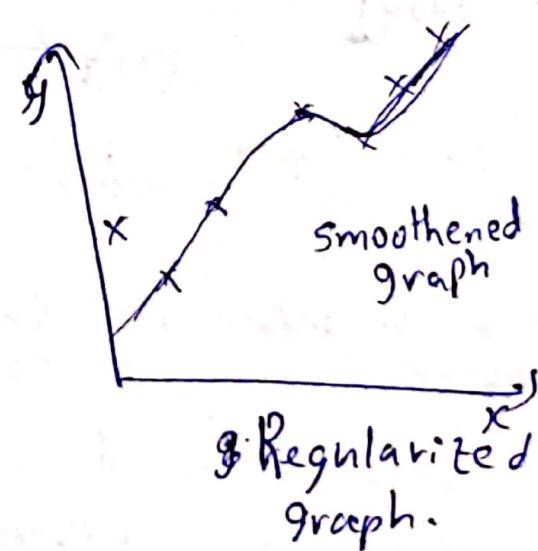
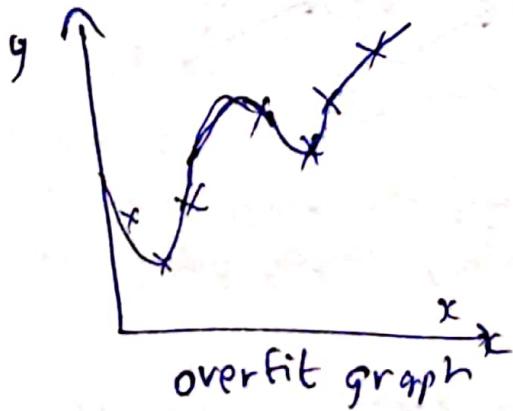
$$\theta_0 := \theta_0 - \alpha \sum_{i=1}^m (h_\theta(x_i) - y_i) x_{0i}$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) x_{ji} + \frac{\lambda}{m} \theta_j \right]$$

~~+ θ_j~~

}

where $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$



• Regularized graph.

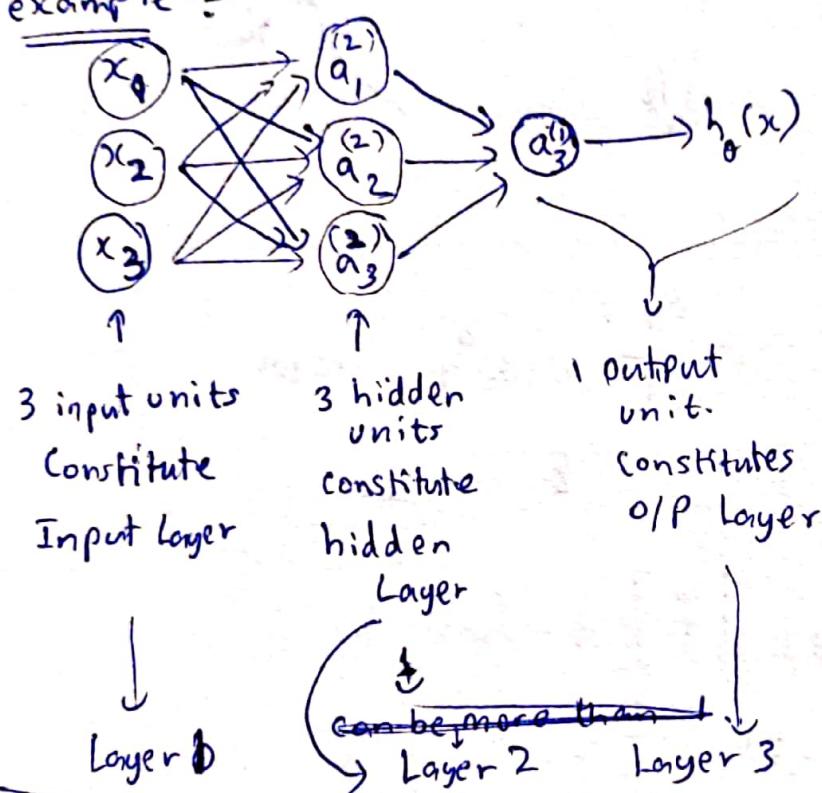
Neural Networks

Mimic the brain's learning ability.
electric

In each neuron, we have dendrites, which take inputs and axons, where electric pulses are sent & output).

In this model, dendrites are input features (x_1, x_2, \dots, x_n), and output is the result of our hypothesis function.

example:



3 input units
constitute
Input Layer

3 hidden
units
constitute
hidden
Layer

1 output
unit.
constitutes
O/P Layer

There can be
more than 1
hidden/Intermediate
Layer.

Hidden Layer

↓
a.k.a

↓
Intermediate Layer.

hidden layer units \rightarrow denoted by $a_i^{(j)}$
(activation of unit i in layer j)

$\theta^{(j)}$ is the matrix of weights, controlling the function mapping from layer j to $j+1$.

$\theta_{ab}^{(j)}$ \rightarrow ath row, bth column of $\theta^{(j)}$

$\theta^{(j)}$ \rightarrow It has dimension: $(s_{j+1}) \times (s_j + 1)$, where ~~s_j~~ = number of units in



where $s_j \rightarrow$ # units in layer j
 $s_{j+1} \rightarrow$ # units in layer $j+1$

$x_0 \rightarrow$ "bias unit"
 \rightarrow always = 1

$\theta_{qb}^j \rightarrow$ "weights"

How to compute $h_\theta(x)$ for above example?

$$a_1^{(2)} = g(\theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \cancel{\theta_{13}^{(1)}x_3})$$

$$a_2^{(2)} = g(\theta_{20}^{(1)}x_0 + \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \theta_{23}^{(1)}x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)}x_0 + \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3)$$

and, $a_1^{(3)} = g(\theta_{10}^{(2)}a_0^{(2)} + \theta_{11}^{(2)}a_1^{(2)} + \theta_{12}^{(2)}a_2^{(2)} + \theta_{13}^{(2)}a_3^{(2)})$

and, $h_\theta(x) = a_1^{(3)}$
=.

Here, $g \rightarrow$ sigmoid
logistic function
activation ↴
 $\frac{1}{1+e^{-z}}$

$$x_0 = 1, a_0^{(2)} = 1$$