

9/6/22

Machine Learning Course

Week - 2

Multiple Features.

| | Size (feet ²) | No. of bedrooms | No. of floors | Age of home | Price |
|---|---------------------------|-----------------|---------------|-------------|-------|
| 1 | a | 2 | 1 | 10 | 100 |
| 2 | b | 3 | 3 | 15 | 200 |
| 3 | c | 4 | 2 | 20 | 300. |

$$a \xrightarrow{3} \text{size (feet)} \rightarrow 1$$

$$\xrightarrow{\quad c \quad} \rightarrow 3$$

∴ size in feet² of home 2

Cost function -

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

~~Assume~~ $x_0 = 1$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

$\in \mathbb{R}^{n+1}$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Hypothesis function ↑

$$\theta^T \times X \in \mathbb{R}^{(n+1) \times k}$$

$$\theta^T = [\theta_0 \ \theta_1 \ \theta_2 \ \dots \ \theta_n]$$

1 - row, $n+1$ - cols.

$$\theta^T \times X = \begin{matrix} \text{Rows} \\ 1 \end{matrix} \times \begin{matrix} \text{Cols} \\ (n+1) \end{matrix}$$

Vector can be defined as a 1-D data structure.

Eg. Housing dataset with an array containing information -

- No of bedrooms
- sqft of the house

Vectors can be represented as a row as well as a column vector.

A vector can be a matrix but not vice versa.

Hypothesis \Rightarrow

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

$$\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Parameters : $\theta_0, \theta_1, \dots, \theta_n \in \Theta$

Cost Function -

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

$J(\theta) \Rightarrow$ same.

Gradient Descent - For $n = 1$

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

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

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

$$\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta) = \textcircled{1}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x^i \quad \textcircled{11}$$

New algorithm for ($n \geq 1$)

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

θ_0 and θ_1 will remain same as for $n=1$.

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

for second feature

10/6/22

Feature Scaling

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

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

$$\text{Scaling} \rightarrow x_1 = \frac{\text{size}}{2000}$$

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

$$\therefore 0 \leq x_1 \leq 1, 0 \leq x_2 \leq 1$$

Goal of feature scaling is to bring every parameter with a range of $-1 \leq z_i \leq 1$

• Mean normalization. —

Replace x_i with $\frac{x_i - \mu_i}{\sigma_i}$

$$x_1 = \frac{\text{size} - 1500}{2000}$$

$$x_2 = \frac{\# \text{bedrooms} - 2}{5}$$

$$\therefore -0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

Example →

$$x_1 \rightarrow \frac{x_1 - \mu_1}{\sigma_1} \quad \begin{array}{l} \mu_1 = \text{mean} \\ \sigma_1 = \text{range.} \end{array}$$

$$\begin{aligned} a &= 4(3 - 5) + 2 \\ &= 4(-2) + 2 \\ &= -8 + 2 = -6 \end{aligned}$$

$$\begin{aligned} b &= -4(5 - 3) \\ &= -4(+2) \\ &= -8 \end{aligned}$$

$$\begin{aligned} 2b &= -16 \\ |2b| &= +16 \end{aligned}$$

$$2|a| - |2b|$$

$$\Rightarrow 2 \times 6 - +16$$

$$\Rightarrow 12 - 16$$

$$\Rightarrow -4$$

$$3T: 20$$

$$1 < 2a < 18$$

$$0 < 2a < 17$$

$$0 < a < 8.5$$

$$x + y = 180$$

$$(x - 70) + 2y - (x + y) = 20$$

$$(x - 70) + 2y - 180 = 20$$

$$(x - 70) + 2y = 170$$

$$x + 2y = 240$$

$$x + y + y = 240$$

$$y = \frac{240 - 180}{2}$$

$$y = 60$$

$$(x + y) - [(x - 70) + 2y] = 20$$

$$180 - 20 = x - 70 + 2y$$

$$130 = x + y + 70$$

$$200 = 180 + y$$

$$y = 20$$

$$x = 100$$

$$y' = 100$$

$$x' = 30$$

Q

$$\frac{14}{x} \rightarrow \frac{14}{x-y}$$

$$\frac{14}{x-y} - \frac{14}{x}$$

$$\frac{14x - 14(x-y)}{x(x-y)}$$

$$\frac{14y}{x^2 - xy}$$

$$t = 5$$

$$u = -2$$

$$7, -7$$

$$\sqrt{2} + 0.22 \frac{2}{9} y - \frac{1}{2}$$

$$\frac{2}{9} + \sqrt{2} = 0.3$$

$$0.22 + 1.732 = 0.3$$

$$1.9 - 0.3$$

$$\Rightarrow -1.1$$

$$1.414 - 0.5 +$$

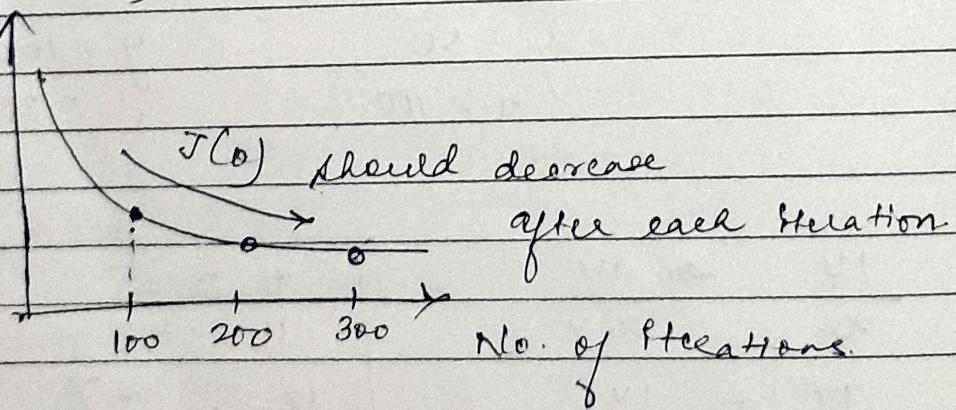
$$110 \rightarrow 70 \rightarrow 56 \\ M \qquad \qquad T$$

$$a \left[\frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} \right]$$

$$a \left[\frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \frac{1}{6} \right]$$

Learning Rate

$\min(J(\theta))$



After 300th iteration we see that Gradient Descent has converged as there is no further reduction in the cost function.

Declare convergence if $J(\theta)$ decrease by less than 10^{-3} in one iteration.

Acceptable learning rates \rightarrow

$$\alpha \in 0.001, 0.01, 0.1, 1, \dots$$

Normal Equation

$$X = \begin{bmatrix} X_0 \\ X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

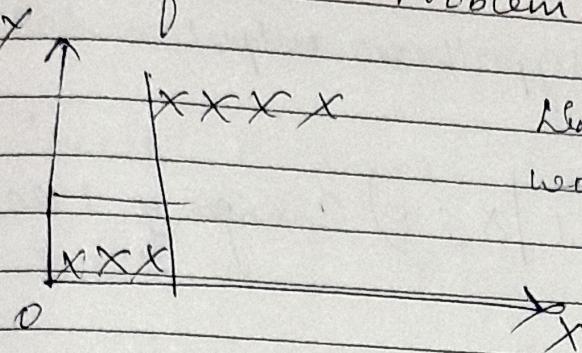
design matrix

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

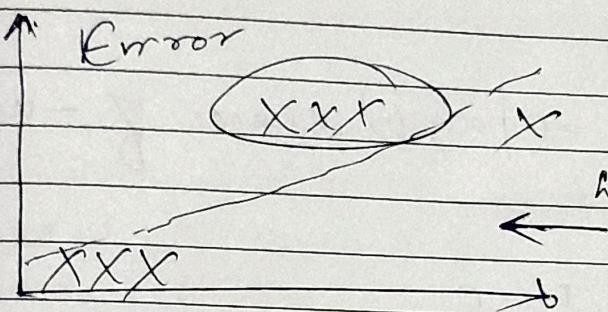
| X_1 | X_2 | Price |
|-------|-------|-------|
| 4 | 89 | 100 |
| 5 | 81 | 150 |
| 6 | 60 | 200 |

$$X = \begin{bmatrix} 4 & 89 & 1 \\ 5 & 81 & 1 \\ 6 & 60 & 1 \end{bmatrix} \quad Y = \begin{bmatrix} 100 \\ 150 \\ 200 \end{bmatrix}$$

Classification Problem.



Linear reg may be working here



But adding an extreme point makes it fail.
here.

Use a classifier for a classification problem.

~~Logistic~~ Logistic Regression is a classification algorithm.

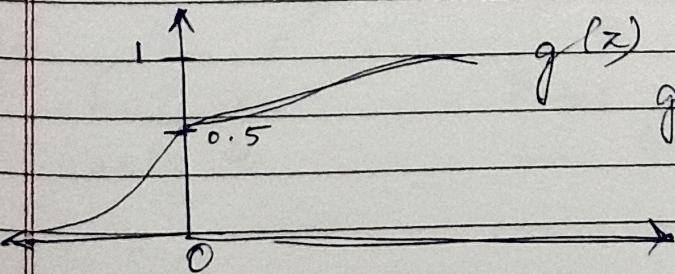
Logistic regression Model →

$$0 \leq h_0(z) \leq 1$$

Sigmoid function.

$$h_0(z) = g(z^T \alpha) \quad z = \gamma_1 + e^{-z}$$

$$g(z) = \frac{1}{1 + e^{-z}} \quad \text{Sigmoid function}$$



$g(z)$ approaches 1 as it tends to $+\infty$

and 0 as it approaches to $-\infty$

14/6/22

Interpretation of Hypothesis output.

$$h_0(x) = P(y=1|x; \theta) \text{ chance of +ve.}$$

$$y = 0 \text{ or } 1$$

$$h_0(x) = P(y=0|x; \theta) \text{ chance of -ve}$$

Q

$$\theta_0 = 5, \theta_1 = -1, \theta_2 = 0$$

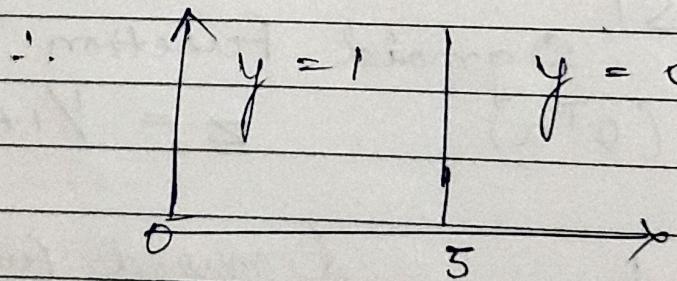
$$\begin{aligned} h_0(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ &= 5 - x_1 + 0 \\ &= 5 - x_1 \end{aligned}$$

$$5 - x_1 > 0$$

$$x_1 < 5$$

$$x_1 < 5$$

$h_0(x)$ stays greater than 0



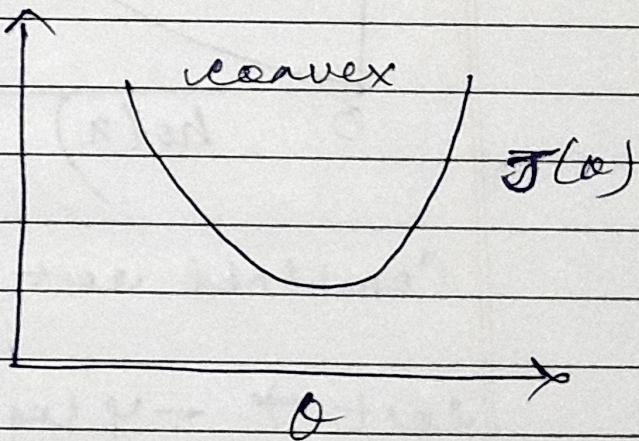
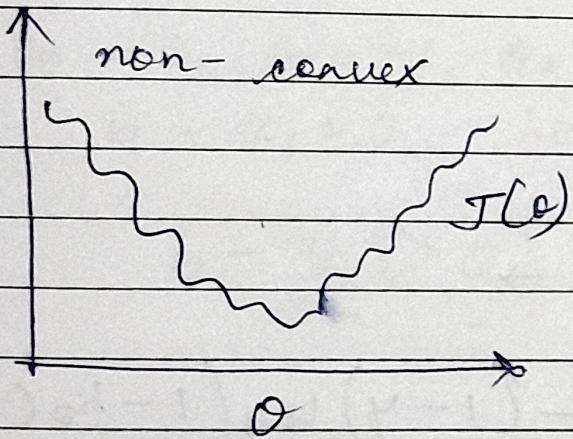
equation of a circle.

7/22

Training set = $\{(x^1, y^1), (x^2, y^2), (x^m, y^m)\}$

$$x \in [x_0 \\ x_1 \\ \vdots \\ x_n]$$

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

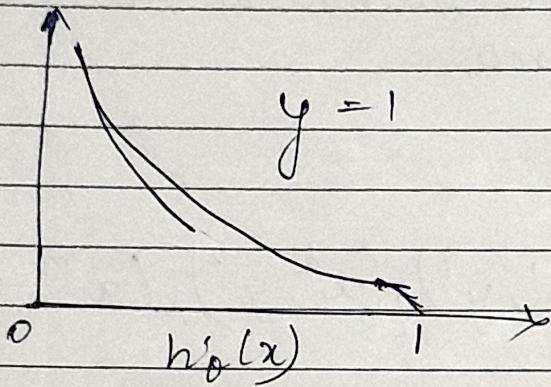


Logistic regression cost function -

$$\text{cost}[h_\theta(x), y] \Rightarrow$$

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

a)



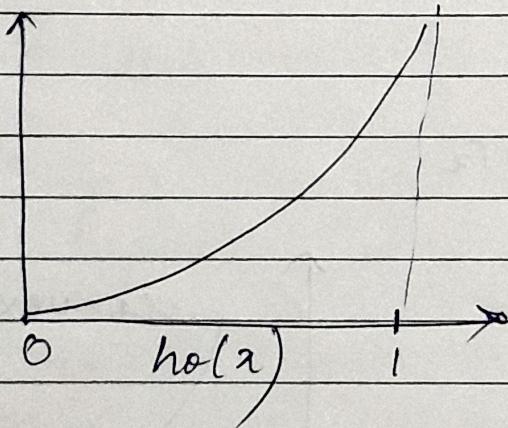
If $y = 1$,

$$h_\theta(x) = 1 \Rightarrow \text{Cost} = 0$$

But as $h_\theta(x) \rightarrow 0$
Cost $\rightarrow \infty$

b)

$$y = 0$$



Combined cost function \rightarrow

$$\text{cost} \rightarrow -y \log(h_\theta(x)) - (1-y) \log(1 - h_\theta(x))$$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^i \log(h_\theta(x)) + (1-y^i) \log(1 - h_\theta(x)) \right]$$

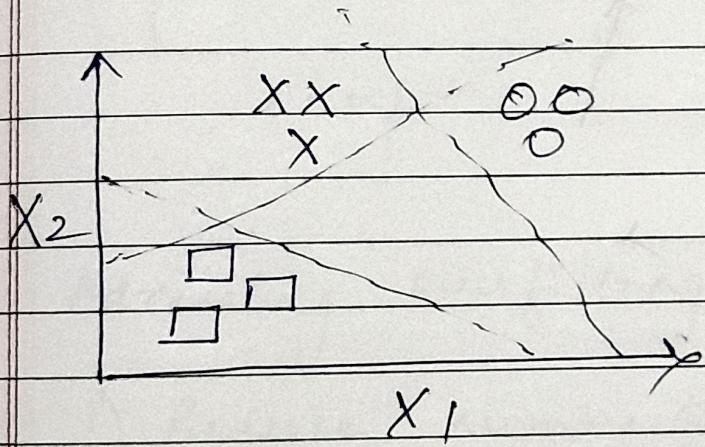
Multi class classification problem -

Email : Work, Friend, Family, Hobby.

$$y = 1 \quad y = 2 \quad y = 3 \quad y = 4$$

Weather \Rightarrow Rainy, Cloudy, Raba, Snow

$$y = 1 \quad y = 2 \quad y = 3 \quad y = 4$$



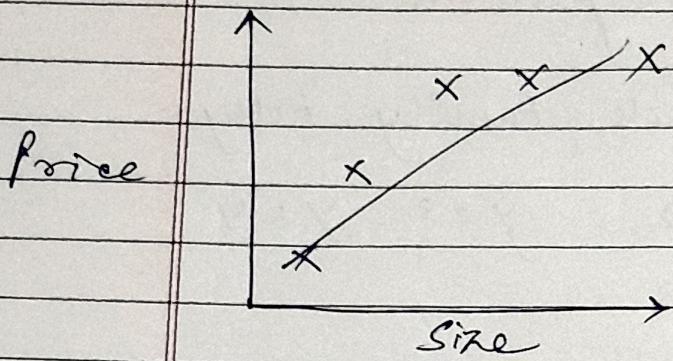
$$\varphi h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta_0 = 6, \theta_1 = -1, \theta_2 = 0$$

$$6 - x_1 + 0 = h_0(x)$$

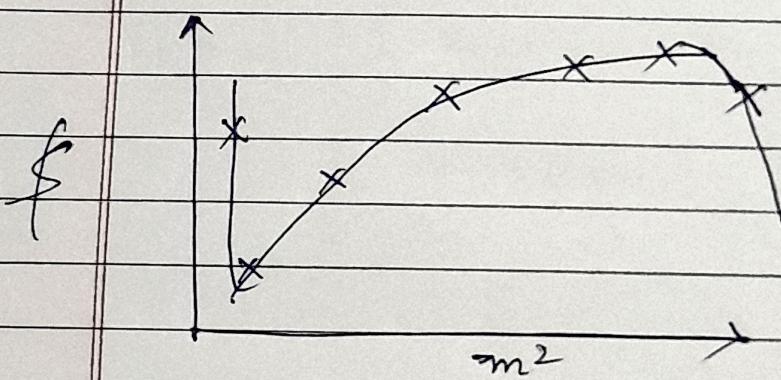
$$6 - x_1 > 0$$

Overfit and Underfit. (Linear Regression)



Underfit

- High bias
- Low variance.

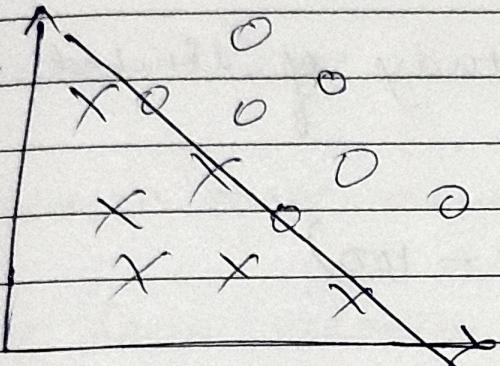


- Overfit
- High Variance
- Low Bias

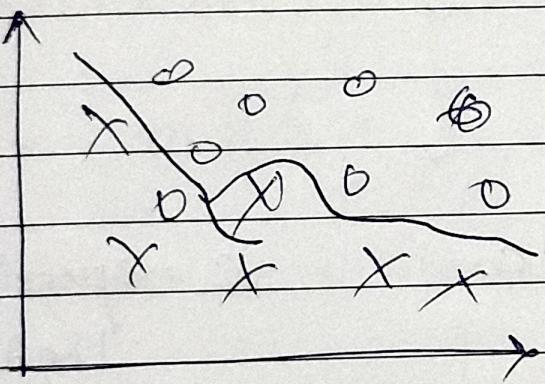
Overfit -

When we have too many features, the learned hypothesis may fit the training set very well but fail to generalize new examples

Logistic Regression.



Underfit



Overfit.

Addressing overfitting -

- 1) Reduce number of features
- 2) Regularization. (Reduce magnitudes/value of parameters).

15/7/22

CLASSMATE
Date _____
Page _____

Regularization

ML Course

$$1 - \frac{\alpha}{m} < 1$$

Normal Equation =

$$X = \begin{bmatrix} x^{(1)\top} \\ x^{(m)\top} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(m)} \end{bmatrix}$$

 $\min(\mathcal{J}(\theta))$

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

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

Logistic regression maybe prone to overfitting if we use a higher order polynomial.

Thus we can regularize the cost function of logistic regression \rightarrow

$$\text{by adding, } \frac{1}{2m} \sum_{j=1}^m \theta_j^2$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[-y^i \log(h_\theta(x^i)) - (1-y^i) \log(1-h_\theta(x^i)) \right]$$

$$\text{grad} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^i) \sigma_j$$

Regularized Cost Function \rightarrow

$$\frac{1}{2m} \sum_{j=1}^n \theta_j^2 + J(\theta)$$

1/1/22

classmate

Data
Page

Week - 5

L = Total number of layers in the network.

s_L = no. of units in layer L

Classification

Binary

$$y = 0 \text{ or } 1$$

$$s_L = 1$$

Multi class

$$y \in R^K$$

pedestrian, car, motor

K output units

$$s_L = K$$

Cost function :

Logistic Regression :

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

Back propagation +

$J(\theta)$

$$J(\theta) = \frac{-1}{m} \left[\sum_{l=1}^m \sum_{k=1}^K y_k \log h_\theta(x^{(l)}) + (1 - y_k) \log (1 - h_\theta(x^{(l)})) \right] \\ + \frac{1}{2m} \sum_{l=1}^L \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (o_j^{(l+1)} - t_j^{(l)})^2$$

$\min J(\theta)$

Need to compute :

1) $J(\theta)$

2) $\frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$

Gradient Computation!

$$a^{(1)} = x$$

$$z^2 = \theta^{(1)} \cdot a^{(1)}$$

$$a^2 = g(z^2)$$

$$z^3 = \theta^{(2)} \cdot a^2$$

$$a^3 = g(z^3)$$

$$z^4 = \theta^{(3)} \cdot a^3$$

$$a^4 = g(z^4)$$

$$a^4 = h_\theta(x)$$

Intuition of back propagation algorithm is that for each node we are going to compute the term $\delta_j^{(l)}$

$\delta_j^{(l)}$ = error of node j in layer l

for each output unit $\rightarrow l = 4$

$$\delta_j^4 = a_j^4 - y_j$$

y_j = ground truth

$$\delta^3 = (\theta^3)^T \cdot \delta^4 \cdot g'(z^3)$$

a_j^4 = value of the hypothesis we got.

$$\delta^2 = (\theta^2)^T \cdot \delta^3 \cdot g'(z^2)$$

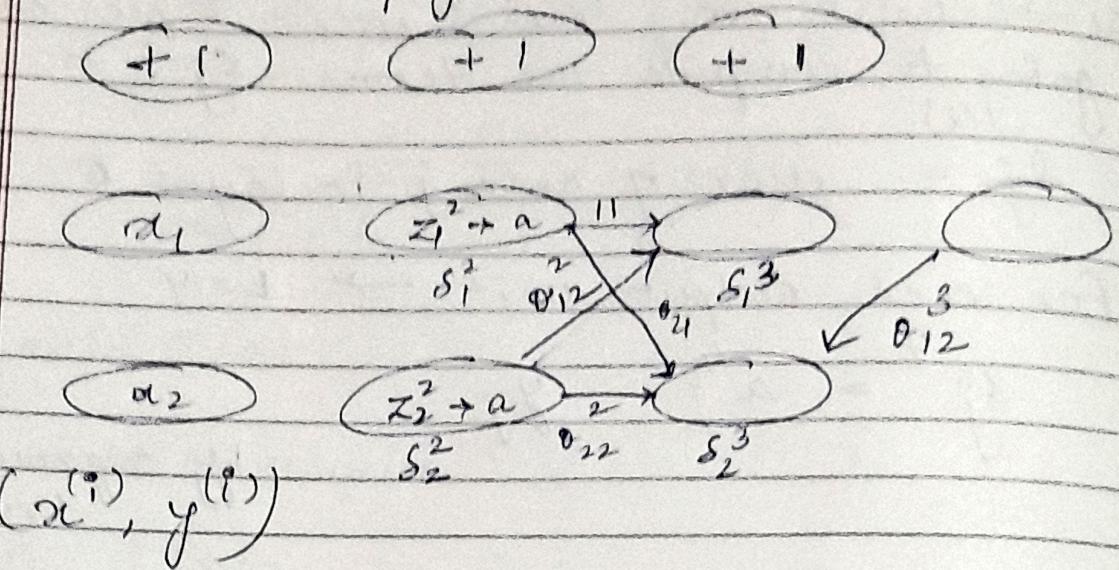
$$\Delta_{i,j} := \Delta_{ij} + a_j^{(l)} \delta_i^{(l+1)}$$

Suppose you have 2 training examples (x^1, y^1) and (x^2, y^2) . What is the correct sequence for computing the gradient.

Ans FP of x^1 , followed by BP of y^1 then.

FP of x^2 , followed by BP of y^2

Forward Propagation.



After we apply activation function to Z , we get the activation values.

Activation function decides whether a neuron should be activated or not by calculating weighted sum and further adding bias with it.

Bias is a sort of a constant which is used to adjust the output along with the weighted sum.

In a NN, we update the weight and biases of the neuron based on the error at the output.

- Linear
- Sigmoid — Binary
- Tangent Hyperbolic — tanh
- ReLU
- Softmax — Multiclass



$$\delta_1 \stackrel{4}{\Rightarrow} y^{(i)} - a_4$$

$$\delta_2 = \theta_{12}^2 \delta_1 + \theta_{22}^2 \delta_2$$

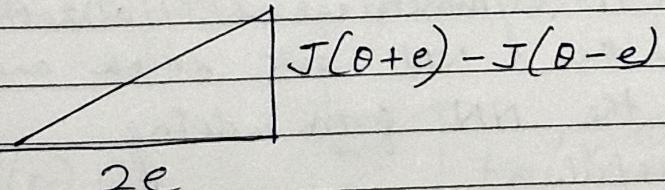
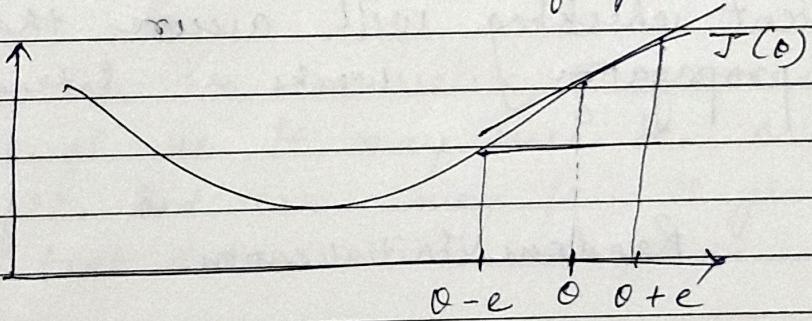
~~$$\delta_2 = \delta_1$$~~

$$\delta_2 = \theta_{12}^3 \times \delta_1$$

Q

~~$$\delta_1 \stackrel{3}{\Rightarrow} \theta_{12}^3 \times \delta_1$$~~

~~18/7/22~~ Numerical estimate of gradients.



$$\frac{dJ(\theta)}{d\theta} \approx \frac{J(\theta + e) - J(\theta - e)}{2e} \quad e = 10e^{-4}$$

$$\text{gradapprox} \doteq \frac{J(\theta + e) - J(\theta - e)}{2 * e}$$

\varnothing $J(\theta) = \theta^3$. $\theta = 1, e = 0.01$

$$\frac{J(\theta + e) - J(\theta - e)}{2e} = \frac{d J(\theta)}{d \theta}$$

$$\frac{(1+0.01)^3 - (1-0.01)^3}{2 \times 0.01} \approx 3$$

$$\doteq 3.0001 \approx 3$$

Gradient checking will assure that back propagation works as intended.

Random Initialization

Zero initialization or symmetrical initialization isn't helpful as the NN gets stuck and it prevents the NN from doing something useful.

So we go for Random Initialization. —

We initialize (θ_i, j) to a random value in $[-e, e]$

E, this is not same as the one from gradient checking.

Training a neural network -

- No of input units : Dimension of features
- No of output units : Number of classes.
- Reasonable Default \Rightarrow 1 hidden unit.

If $f > 1$, then have the same no of hidden units in every layer (more the better).

$$J(\theta) = 3\theta^4 + 4 \quad \theta = 1, \epsilon = 0.01$$

$$\frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

- ~~Just~~ Just a low value of error is not enough as it may even be a case of overfit. and thus may fail to generalize the test data.

Training/Testing procedure for linear regression.

$$J(\theta) \Rightarrow \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \left[h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)} \right]^2$$

Avg Squared error.

set

Why we need a val ~~test~~ →

A validation set is a set of data used to train AI with the goal of finding and optimizing the best model to solve a given ~~prob~~ problem.

It acts as an intermediate phase used for choosing the best model and optimizing it.

Overfitting is checked and avoided due to this.

$J_{CV}(\theta)$ is lower than $J_{test}(\theta)$ because an extra parameter (d) degree of polynomial had been fit to the cross validation set.

$$y = 4 - 3x$$

$$y = e^u$$

$$u = 4 - 3x$$

$$y = e^u$$

$$\frac{dy}{du} = -3$$

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx}$$

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx} = \frac{dy}{du}$$

$$R = 35\Omega, V = 3V, \frac{dV}{dt} = 1.8V/\text{min}$$

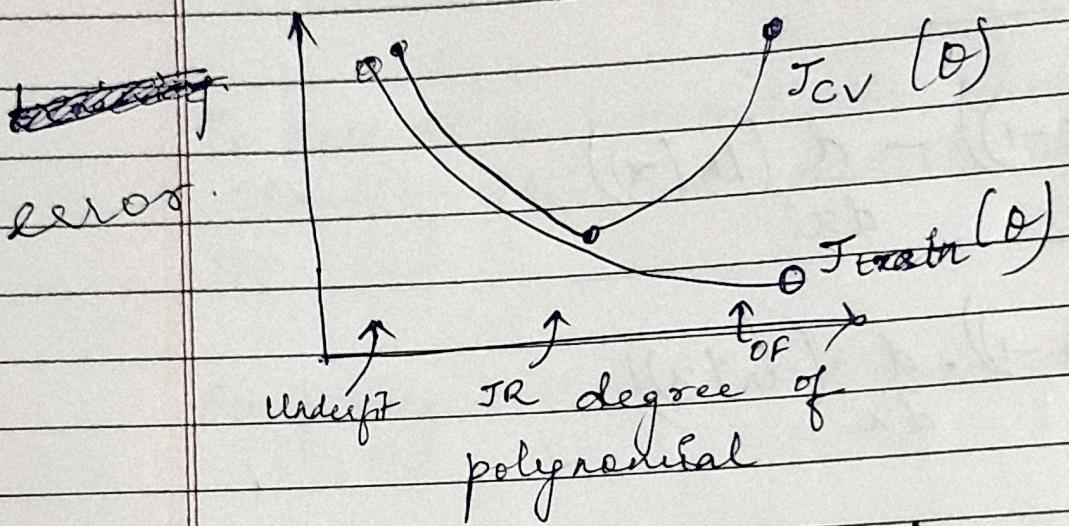
$$\frac{dR}{dt} = 1.7\Omega/\text{min}$$

$$V = IR$$

$$\frac{dV}{dt} = \frac{dI}{dt} \cdot R + I \cdot \frac{dR}{dt}$$

High Bias - Underfit

High Variance - Overfit.



High Bias -

- 1) High training error

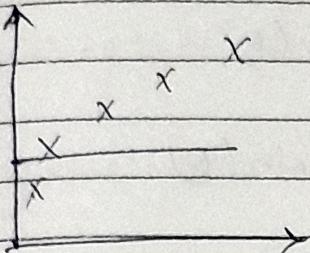
- 2) high cross val error

High Variance

- 1) ~~Low~~ training error
- 2) High cross val error.

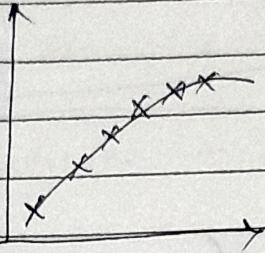
11/7
22

Linear regression with regularization



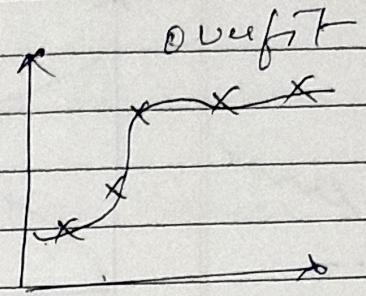
Underfit

$$\lambda = 10000$$



Just Right

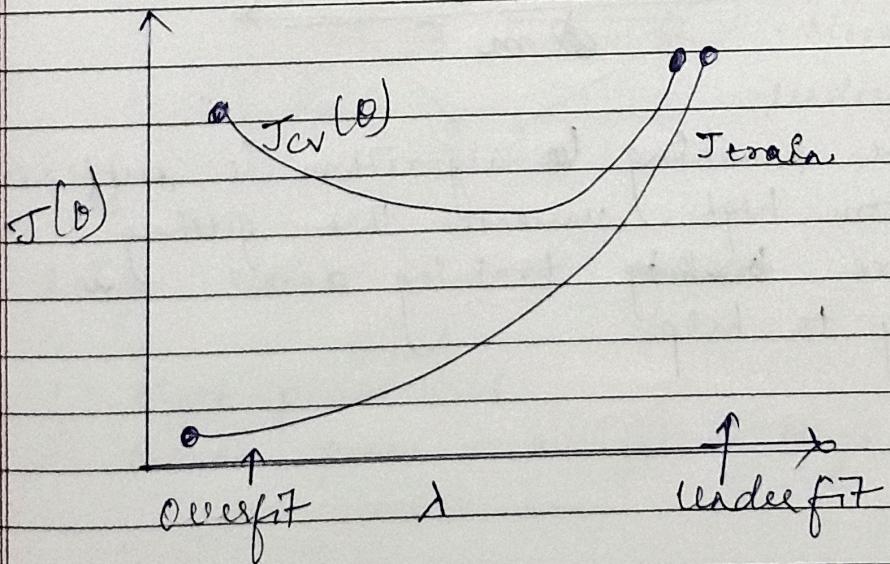
$$\lambda = 500$$



$$\lambda = 0$$

The way we find the value of λ is same as how we find the value of 'd' or degree of polynomial in a NN.

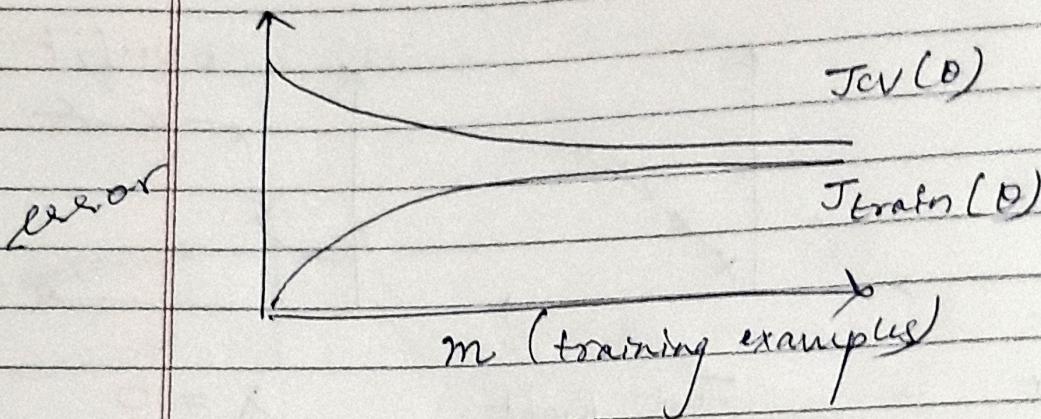
We test how well the model performs across various values of lambda.



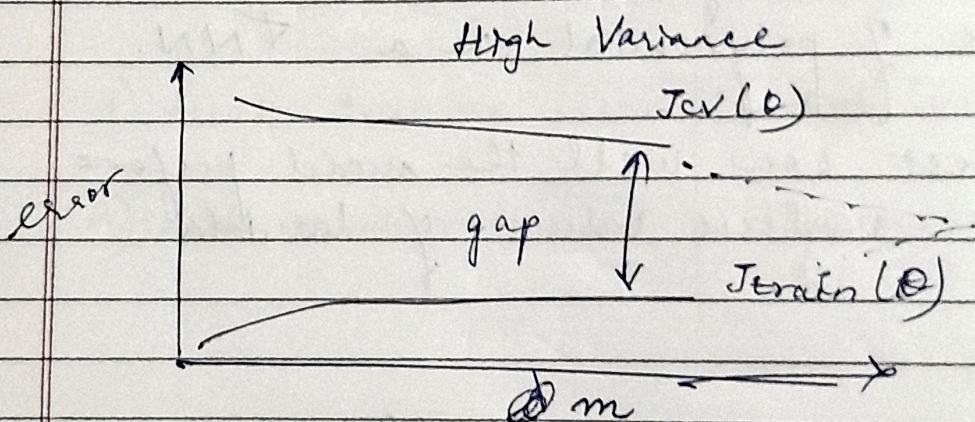
Learning Curve

Case →

High Bias



If a learning algorithm is suffering from high bias, getting more training will not help much.



If the learning algorithm is suffering from high variance then getting more ~~less~~ training data is likely to help.

Debugging learning algo -

Getting more training examples - Fix H. Variance

Try smaller set of features - Fix High Variance.

Try getting additional features - Fix High Bias

Try adding polynomial features - Fix High Bias.

Decreasing λ - Fix High Bias

Increasing λ - Fix High Variance

q → High error testing

↳ High variance, overfitting
↳ Increase

Precision / Recall

| | | Actual class | | |
|-----------|---|--------------|----|--|
| | | 1 | 0 | |
| predicted | 1 | TP | FP | |
| | 0 | FN | TN | |

Precision $\Rightarrow \frac{TP}{TP + FP}$ how many had $y = 1$

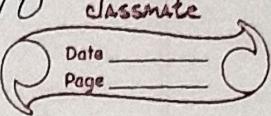
Recall $\Rightarrow \frac{TP}{TP + FN}$ how many did we correctly classify.

| | | Actual class | | |
|-----------|---|--------------|----|--|
| | | 1 | 0 | |
| predicted | 1 | TP | FP | |
| | 0 | FN | TN | |

$$P \Rightarrow \frac{80}{80+20} = 0.8$$

$$R \Rightarrow \frac{80}{80+80} = 0.5$$

| TP | FP | FN | $P = \frac{80}{32} = 93, R = 90$ |
|----|----|----|----------------------------------|
| 30 | 2 | 3 | |
| 30 | 10 | 2 | $P = 75, R = 93$ |



precision = $\frac{\text{true positives}}{\text{no. of predicted positives}}$

recall = $\frac{\text{true positives}}{\text{no. of actual positives}}$.

Trading off precision and recall -

If we wish to increase the ~~threshold~~ threshold of confidence from $0.5 \rightarrow 0.9$.

If we are 90% confident that someone has cancer then we ↑ precision but recall ↓

$$y = 1$$

High Precision, low recall

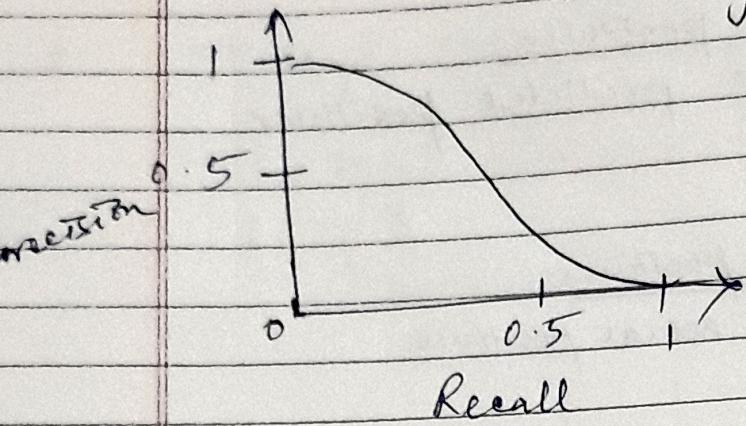
case - 2 To avoid false negatives we may set the threshold to 0.3

Lower Precision, High Recall.

| T | TP | FP | FN | |
|--------|----|----|----|------------------|
| I 0.8 | 30 | 2 | 3 | $P = 93, R = 90$ |
| II 0.3 | 30 | 10 | 2 | $P = 75, R = 93$ |

Date _____
Page _____

$$\text{Accuracy} = \frac{TP + TN}{\text{Total Examples}}$$



F_1 score \rightarrow

| | P | R |
|--------|------|-----|
| Algo 1 | 0.5 | 0.4 |
| Algo 2 | 0.7 | 0.1 |
| Algo 3 | 0.02 | 1.0 |

$$F_1 = \frac{PR \times 2}{P+R}$$

Large data rationale

- Else a learning algorithm with many parameters.
 \rightarrow low bias
- Use a very large training set \rightarrow unlikely to overfit

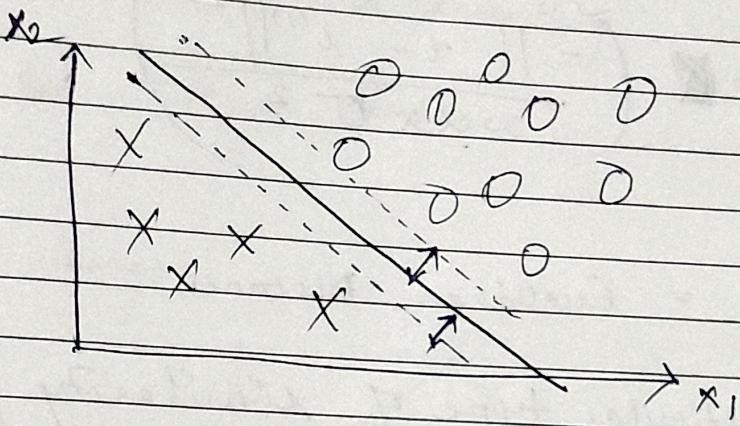
P \Rightarrow

$$\frac{85}{85+890}$$

$$\Rightarrow 0.087$$

22/7/22

SVM

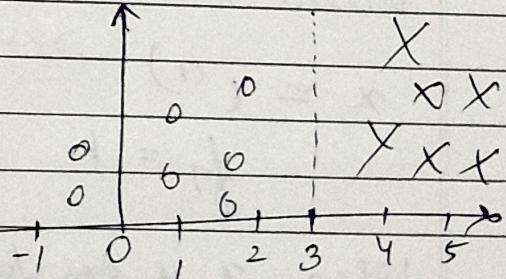


SVM tries to get the largest minimum distance between the two classes.

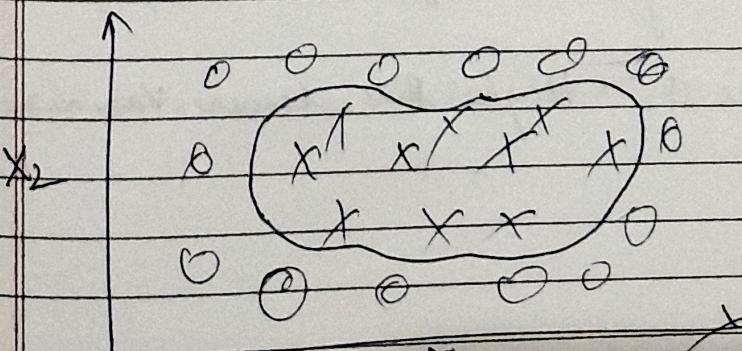
Q

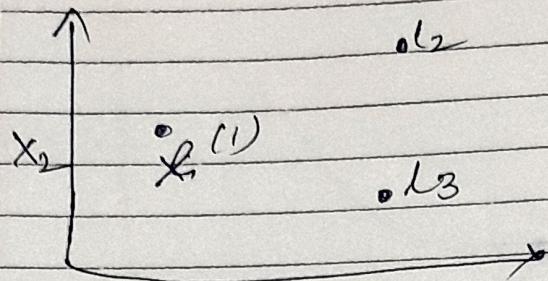
$$\theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$\theta = -3, \theta_1 = 1, \theta_2 = 0$$



Kernels are used to adapt SVM in order to develop complex non linear classifiers.





$f_1 = \text{similarity}(x, l^{(1)}) \rightarrow \text{Kernel}$

$$\exp\left(-\frac{\|x - l^{(1)}\|^2}{2 \times \sigma^2}\right)$$

$x - l^{(1)}$ = Euclidean Distance.

This particular type of similarity function
is known as Gaussian kernel.

When \rightarrow

$$x \approx l^{(1)}$$

$$f_1 = 1$$

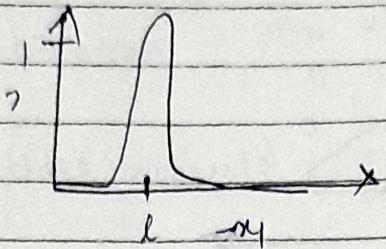
When x is far from l

$$f_1 = 0$$

$C = \frac{1}{\lambda}$ Larger C : High Variance, Low Bias

Smaller C : High Bias, Lower Variance.

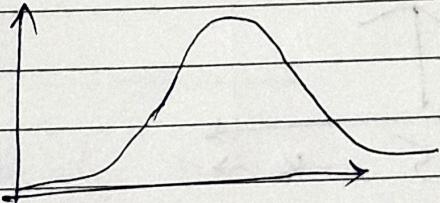
1) $\sigma^2 = \text{smaller}$



Features vary less smoothly.

Low Bias, High Variance

2) $\sigma^2 = \text{Larger}$



Features vary smoothly,

High Bias, Low Variance.

Week - 8

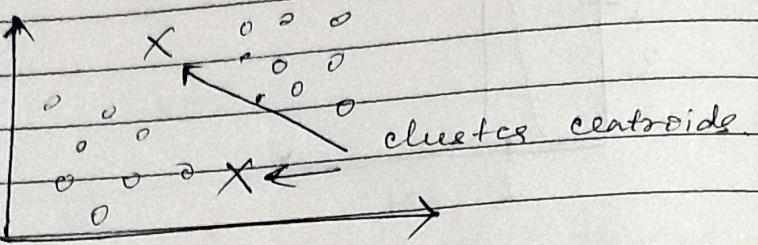
Applications of clustering \rightarrow

a) Market Segmentation

b) Social network analysis.

c) Astronomical data analysis.

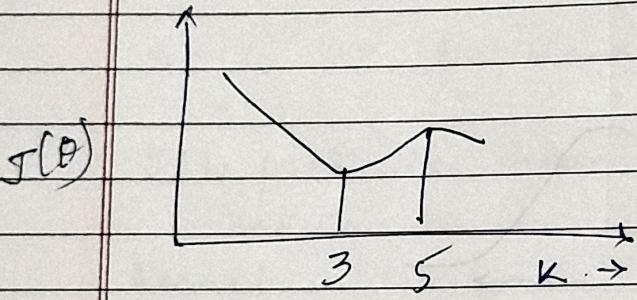
k-means is the most popular clustering algorithm.



Steps →

a) Cluster Assignment

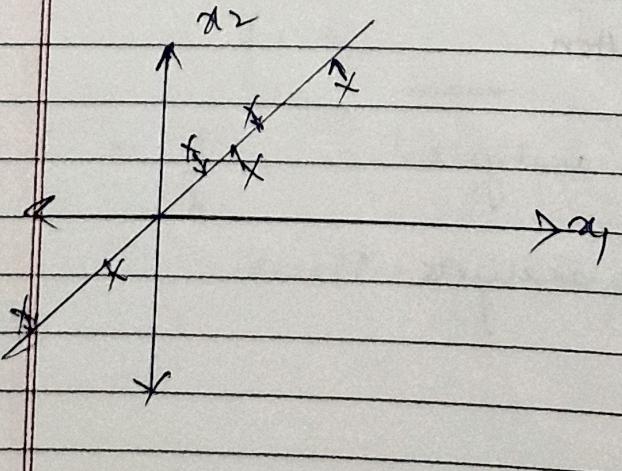
b) Move centroid step. ↘ Loop



PCA

Principal Component Analysis.

Reduce from 2D to 1D : Find a direction onto which to project the data so as to minimize the projection error.



$1, 2, 3, 4, 5 \Rightarrow \text{Mean} \Rightarrow 3.$

$$\text{Mean} \Rightarrow \frac{1}{m} \sum x_j - M_j$$

Normalization

$$\Rightarrow \frac{1}{5} [(-2) + (-1) + 0 + 1 + 2] = 0$$

$\therefore \text{New mean} = 0.$

Feature scaling \Rightarrow

$$x_j^* \Rightarrow \frac{x_j - M_j}{S_i} \quad S_i = \text{range}$$

$M_j = \text{mean}$

Anomaly Detection

1) Credit Card Fraud Detection.

2) Manufacturing. $P_n < \epsilon$

3) Data centres. Every time P_n goes $< \epsilon$, we flag it as an anomaly.

On a cross val / test set \rightarrow

$$y = \begin{cases} 1, & p(x) < \epsilon \text{ (anomaly)} \\ 0, & p(x) \geq \epsilon \text{ (normal)} \end{cases}$$

Use $\rightarrow F_1, P/R$, as metrics.

Anomaly vs Supervised.

Very small number
of positive
examples.
(0-20)

Large number of
negative examples

We don't know what
they will be like
to tomorrow.

I.e. for an aircraft
engine a lot of
things may
go wrong.

Large no of positive
and negative examples.

Here we know what
the positive examples
may be like tomorrow.

Q

User = j | rated movie = i

$$\therefore g(2,1) = y(i,j)$$

$i = 2, j = 1$

r_{21} User 2, Movie 1

$r(i,j)$

Q

$r(i,j) \Rightarrow i = \text{rated movie}$
 $j = \text{User}$

24/7/22

Stochastic gradient descent will allow us to scale GD to a much larger training set.

SGD

- Randomly shuffle the dataset
- Repeat for $i = 1 \dots m$.

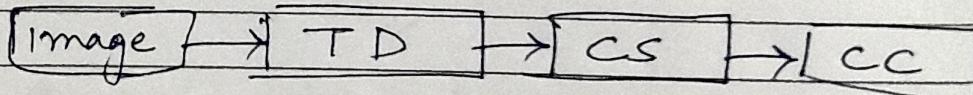
$$\text{cost}(\theta; x^i, y^i) = \frac{1}{2} (h_\theta(x^i) - y^i)^2$$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta; (x^i, y^i))$$

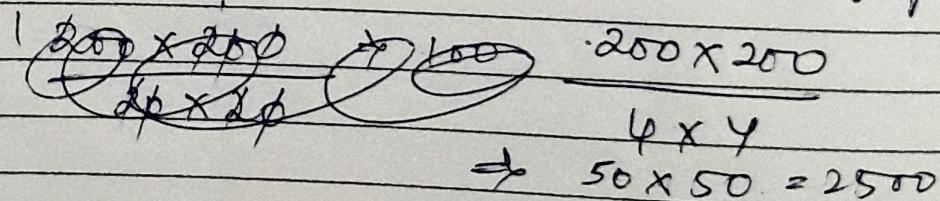
Mini Batch Gradient Descent - Take a batch size of $b = 10$, and perform in iteration.

Photo OCR Pipeline

1. Text Detection
2. Character Segmentation.
3. Character Classification.



Expander operator: Brighten nearby pixels.



10 sec to get 1 ex

$$1 \text{ TE} = \frac{10}{60 \times 60}$$

$$= 2.77 \times 10^{-3} \text{ hrs to get 1}$$

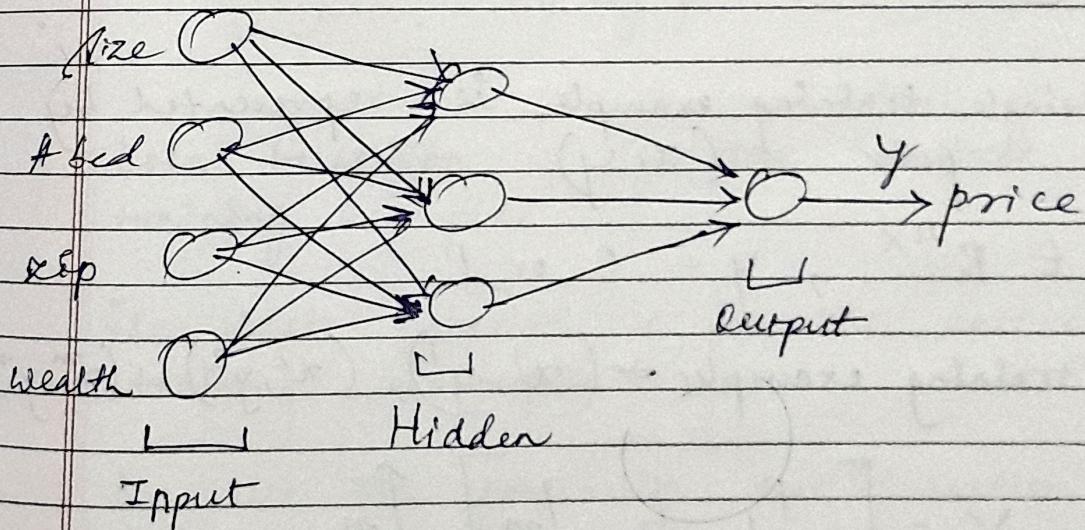
∅ \$10 / hour.

$$4 \text{ ex/min} \Rightarrow 240 \text{ ex/hour} = \$10$$

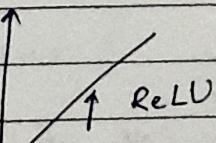
$$\frac{1000 \times 1000}{2 \times 2}$$

$$\$ 500 \times 500$$

Neural Network and DL



ReLU



The ReLU activation function allows better computation by allowing gradient descent to be equal to 1 and converge faster when gradients are used to descend.

Date _____
Page _____

Logistic regression is an algo for binary classification

RGB \rightarrow

say we have an image of 64×64 .

$$\therefore R = 64 \times 64$$

$$B = 64 \times 64$$

$$G = 64 \times 64$$

$$X = \begin{bmatrix} 255 & X \neq 3 \times 64 \times 64 \\ 231 & n_x \neq 12288 \\ \vdots & \\ 255 & \text{dimension of this input} \\ \vdots & \text{feature vector} \\ 255 & \\ 231 & \end{bmatrix}$$

A single training example is represented by
as pair (x, y)

$$x \in \mathbb{R}^{n_x}, y = 0 \text{ or } 1$$

m training examples $\Rightarrow (x^1, y^1), (x^2, y^2) \dots (x^m, y^m)$

$$X = \begin{bmatrix} | & | & | \\ x^1 & x^2 & \dots & x^m \\ | & | & | & \end{bmatrix} \begin{matrix} n_x \\ \downarrow \\ m \end{matrix}$$

$$Y = \begin{bmatrix} y^{(1)}, y^{(2)} \dots y^{(m)} \end{bmatrix} Y = (1, m)$$

$$\text{Output } \hat{y} = \sigma(w^T x + b)$$

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

w is an $n \times 1$ dimensional vector

b is a real number.

Loss Function \rightarrow

$$L(\hat{y}, y) = \frac{1}{2} [\hat{y} - y]^2$$

Squared error loss
Function.

Loss Func of Log Reg \rightarrow

$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

Individual Loss Function is for ~~the~~ entire set of training examples.

The cost function is for the entire set of training examples

We use the particular cost function for logistic regression since we want a convex function.

Repeat {

λ = learning rate.

$$w := w - \frac{\partial J(w)}{\partial w}$$

Computing Derivatives

Chain Rule \rightarrow

$$\frac{dJ}{da} = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

 $a \rightarrow v \rightarrow J$

Change in J = Change in J * Change in v
 due to change due to change due to change
 in a in v in a .

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db}$$

$$b=3 \quad u=bc$$

$$c=2 \quad u=6$$

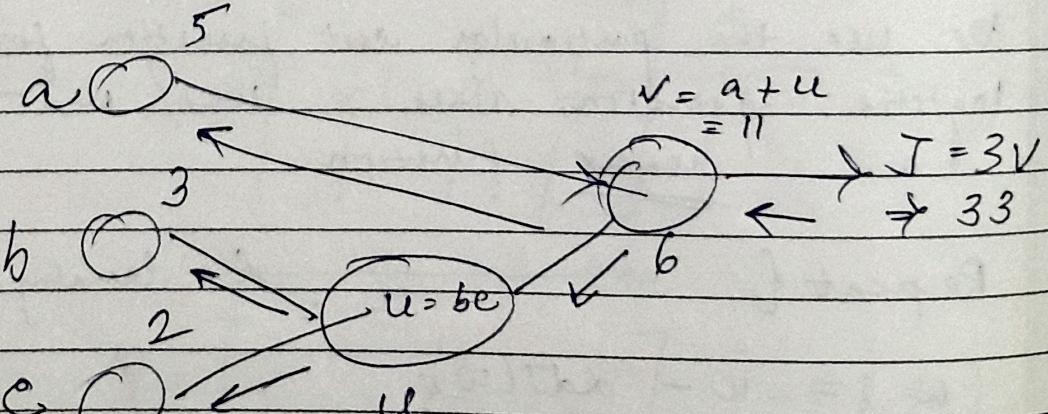
$$b=3.007 \quad u=6.002$$

$$C=2$$

$$\Delta u = 0.002$$

$$\frac{du}{db} = 2$$

$$\frac{dJ}{du} = 3. \quad \frac{dJ}{db} = 3 \times 2 = 6.$$



x_1 w_1 x_2 w_2 b

$$z = w_1 x_1 + w_2 x_2 + b \rightarrow a = \sigma(z) \rightarrow L(a, y)$$

$$\frac{dz}{dL} = \frac{dL}{dz}$$

$$\frac{da}{dL} = \frac{dL}{da}$$

$$dz = a - y$$

$$\frac{da}{dL} = \frac{-y}{a} + \frac{1-y}{1-a}$$

Vectorizing Logistic Regression

$$\begin{array}{|c|c|c|} \hline z^1 & = w^T x^1 + b & z^2 = w^T x^2 + b \\ \hline a^1 = \sigma(z^1) & a^2 = \sigma(z^2) & a^3 = \sigma(z^3) \\ \hline \end{array}$$

$$X = \begin{bmatrix} | & | & | & | \\ x^1 & x^2 & \dots & x^m \\ | & | & | & | \end{bmatrix} (N \times m)$$

~~$$Q \quad [z^1, z^2, \dots, z^m] = w^T X + [b \ b \ \dots \ b]$$~~

$$[z^1, z^2, \dots, z^m] = w^T X + [b \ b \ \dots \ b]$$

$$[z^1, z^2, \dots, z^m] = [w^T x_1, w^T x_2, \dots, w^T x_m] + [b \ b \ \dots \ b]$$

$$\begin{bmatrix} 1, 2, 3, 4, 5 \\ 1 \times 5 \end{bmatrix} - \begin{bmatrix} 4, 6, 4, 3, 2 \\ 1 \times 5 \end{bmatrix}$$