

Supervised

Regression:- Predict a number/infinitely many possible outputs

Classification:- predict categories small number of possible outputs.

Unsupervised

Data only comes with inputs x - but not Output labels y

Algorithm has to find structure in the data

157

Clustering

grp. similar data pts together

Dimensionality reduction

compress data using fewer numbers.

Anomaly detection

Find unusual/unusual data pts

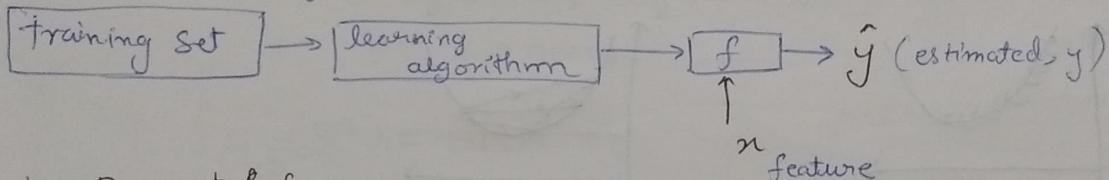
Linear Regression :-

$x \rightarrow$ input variable/feature

$(x, y) \rightarrow$ single training example.

$y \rightarrow$ output variable/target variable, $(x^i, y^i) \rightarrow$ i^{th} training example.

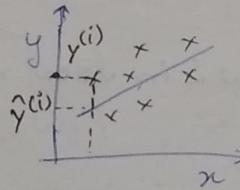
$m \rightarrow$ no. of training examples



How to Represent f

model: $f_{w,b}(x) = w_1x_1 + b$
 $f(x)$

$w, b \rightarrow$ parameters, coefficients



Cost function :-

Squared Error Cost function

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

$m \rightarrow$ no. of training eg.
 $f_{w,b}(x^{(i)})$

goal of Linear Regression \Rightarrow minimise
 $J(w)$
 or $J(w, b)$

Gradient Descent

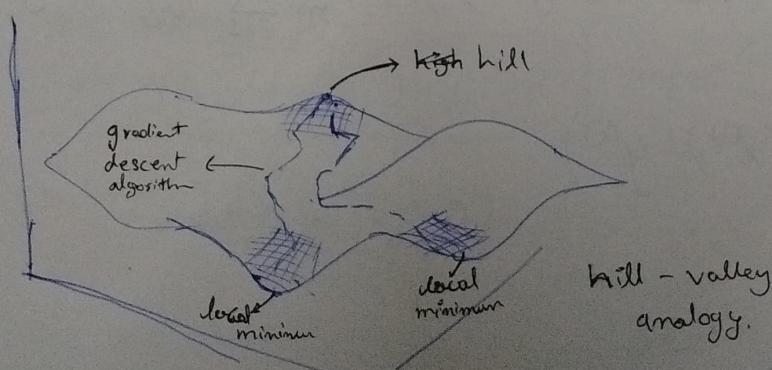
Have $J(w, b)$

or Have $J(w, w_1, w_2, w_3, \dots, b)$

$$\text{Want } \min_{w, b} J(w, b)$$

$$\text{Want } \min_{w, b} J(w, w_1, w_2, \dots, b)$$

Outline \rightarrow \rightarrow Start with some $w, b \rightarrow$ Keep changing w, b to reduce $J(w, b)$



gradient descent algorithm

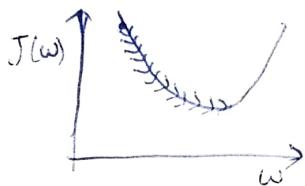
$w = w - \alpha \frac{\partial J(w, b)}{\partial w}$, $b = b - \alpha \frac{\partial J(w, b)}{\partial b}$

assignment operator learning rate repeat until convergence $\left\{ \begin{array}{l} \text{tmp-}w = w - \alpha \frac{\partial J(w, b)}{\partial w} \\ \text{tmp-}b = b - \alpha \frac{\partial J(w, b)}{\partial b} \\ w = \text{tmp-}w \\ b = \text{tmp-}b \end{array} \right\}$

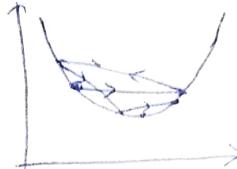
determine ~~too~~ how big step u take down hill

tends to find w, b for which $J(w, b)$ is min

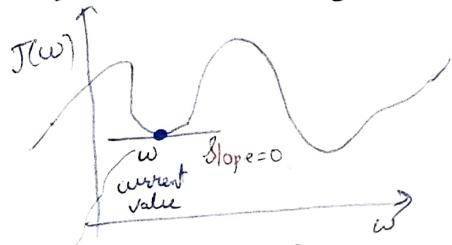
if α is too small



if α is too large



if current value of w is in local minimum



$$w = w - \alpha \frac{\partial J(w)}{\partial w} \Rightarrow \text{update remains unchanged}$$

$$w = w$$

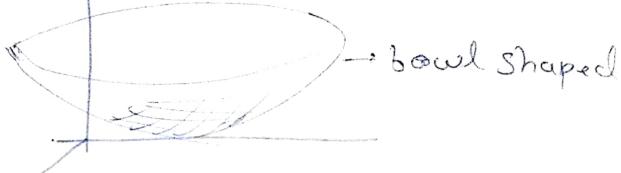
Training Linear Regression for Gradient descent

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

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

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

Squared error cost



feature scaling { Multi-Feature Regression }

$x_j \rightarrow j^{\text{th}}$ feature $n \rightarrow \text{no. of features}$ $\vec{x}^{(i)} = \text{features of } i^{\text{th}}$ training eg.

eg:-

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

✓ w_1 ✓ w_2 ✓ w_3 ✓ w_4 ✓ b
 price #size #bedroom #floor #years old Base price

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

$b \rightarrow$ is a number

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

Vectorization

$$\vec{w} = [w_1 \ w_2 \ w_3], \vec{x} = [x_1 \ x_2 \ x_3], b$$

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

code $\rightarrow f = \text{np.dot}(\vec{w} \cdot \vec{x}) + b$

without vectorizatⁿ:-

without Vectorizatⁿ

for ~~tot~~ j in range(16):

$$w[j] = w[j] - 0.1 * d[j]$$

with Vectorizatⁿ

$$w = w - 0.1 * d$$

Gradient Descent for Multi-Feature Regression

repeat {

one feature $\rightarrow \frac{\partial J(w, b)}{\partial w_i}$

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

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

repeat {

multi-feature $\rightarrow \frac{\partial J(w, b)}{\partial w_i}$

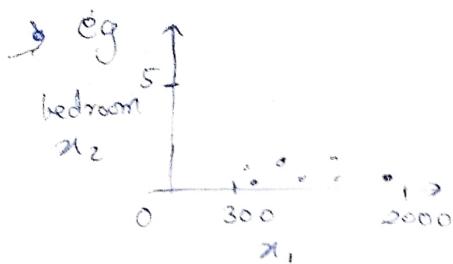
$$w_1 = w_1 - \alpha \left(\frac{1}{m} \sum_{i=1}^m (f_{w,b}(\vec{x}^{(i)}) - y^{(i)}) \vec{x}_1^{(i)} \right)$$

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

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

feature Scaling

(for optimal maths, involved / no abrupt data comes in)



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

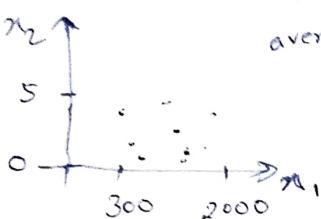
$$x_1, \text{scaled} = \frac{x_1 - 300}{2000}$$

$$x_2, \text{scaled} = \frac{x_2}{5}$$

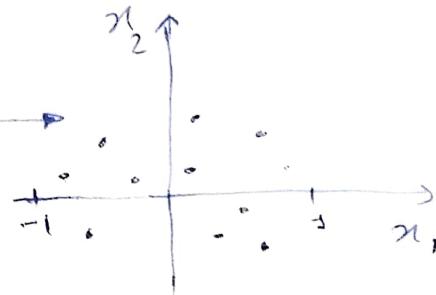
$$0.15 \leq x_1, \text{scaled} \leq 1$$

$$0 \leq x_2 \leq 1$$

mean normalisation

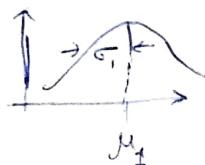
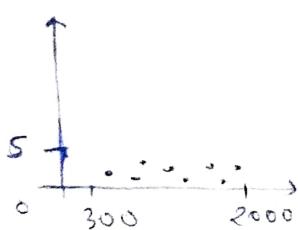


average



$$x_1, \text{scaled} = \frac{x_1 - \mu_1}{\text{max} - \text{min}} \rightarrow -0.18 \leq x_1, \text{scaled} \leq 0.82$$

Z-score normalisation



$$300 \leq x_1 \leq 2000$$

$$x_1, \text{scaled} = \frac{x_1 - \mu_1}{\sigma_1}$$

$$-100 \leq x_1, \text{scaled} \leq 100$$

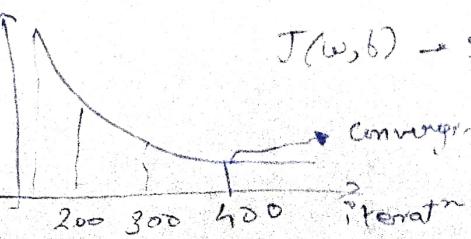
$$-0.001 \leq x_1, \text{scaled} \leq 0.001$$

$$98.6 \leq x_1, \text{scaled} \leq 105$$

make sure gradient descent is working correctly

$J(w, b)$ should decrease after every iteration

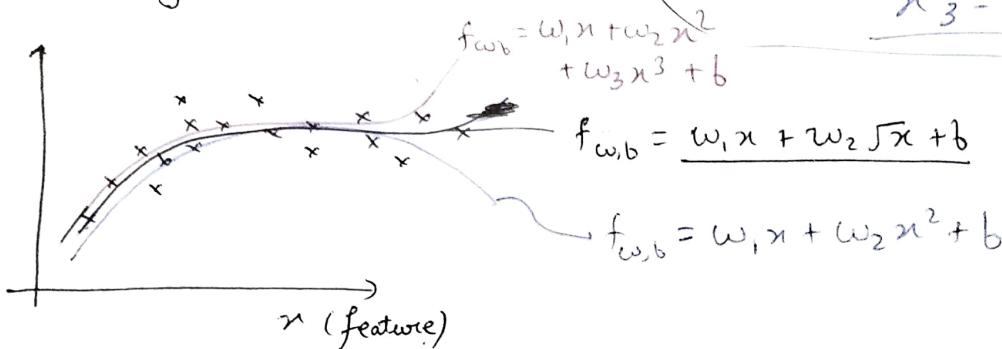
converges after 400 iteration



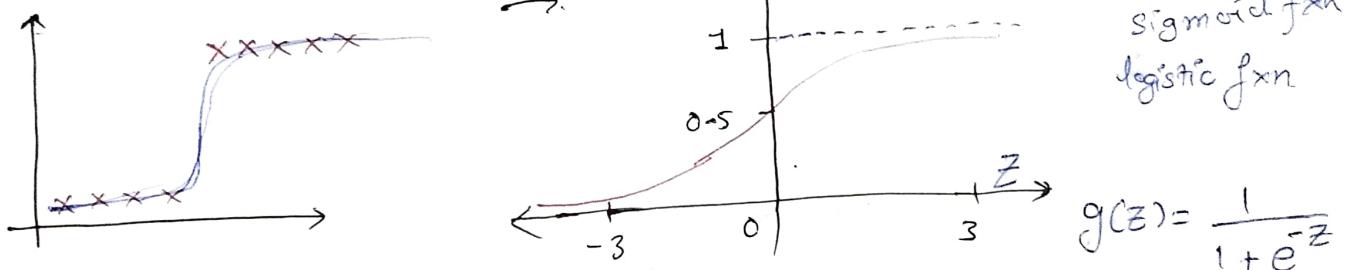
feature Engineering \rightarrow Using intuitions to design transforming or combining original features.

$$f_{w,b} = w_1 x_1 + w_2 x_2 + b \xrightarrow{x_3 \rightarrow \text{new feature}} f_{w,b} = w_1 x_1 + w_2 x_2 + \underline{w_3 x_3 + b}$$

Polynomial Regression



Logistic Regression

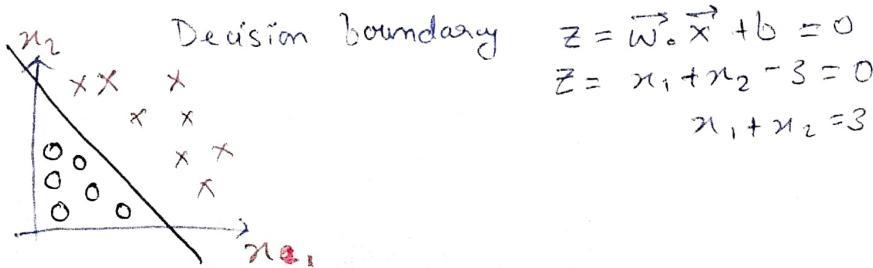


$$f_{\vec{w},b} = \vec{w} \cdot \vec{x} \Rightarrow z = \vec{w} \cdot \vec{x} + b \downarrow z \rightarrow g(z) = \frac{1}{1 + e^{-z}}$$

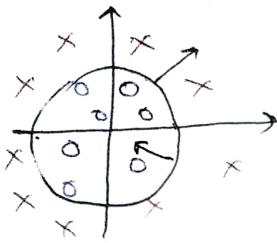
$$f_{\vec{w},b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}} \quad \} \text{logistic Regression Model.}$$

Decision Boundary

$$f_{w,b} = g(z) = g(w_1 x_1 + w_2 x_2 + b)$$



non-linear decision boundaries

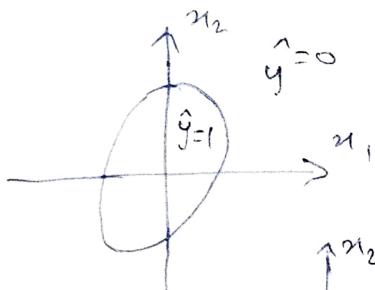


$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1^2 + w_2 x_2^2 + b)$$

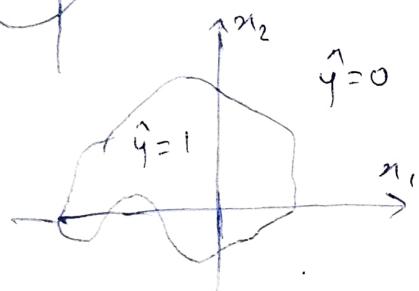
$$\text{decision boundary: } z = w_1^2 + w_2^2 - 1 = 0$$

$$x_1^2 + x_2^2 = 1$$

$$x_1^2 + x_2^2 \geq 1 \quad \hat{y} = 1 \quad \left. \begin{array}{l} x_1^2 + x_2^2 < 1 \\ \hat{y} = 0 \end{array} \right\}$$



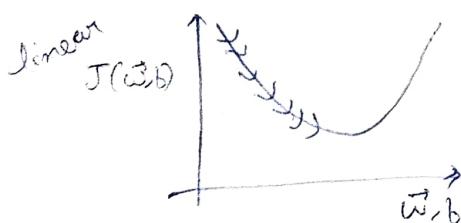
$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + w_3 x_1^2 + \dots + b)$$



Cost fn for Logistic Regression

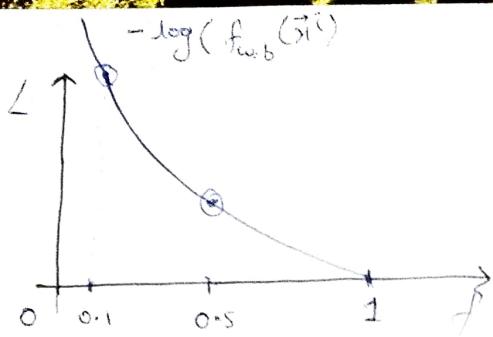
$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

logistic

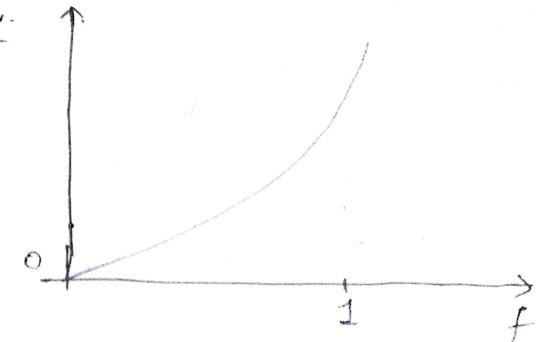


logistic loss function:-

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



$$y \cdot y^{(i)} = 1$$



$$\text{cost} = \frac{1}{m} \sum L(f_{w,b}(x^{(i)}), y^{(i)})$$

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

Training logistic Regression \rightarrow Find \vec{w}, b
gradient descent

repeat $\{$

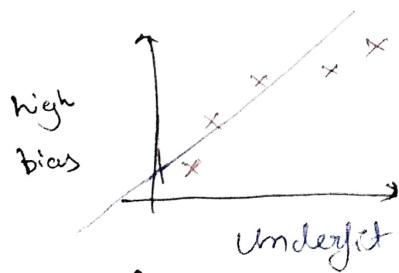
$$w_j = w_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

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

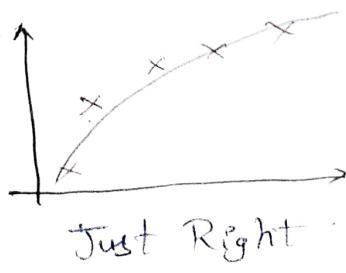
$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

ex training example
ka error har
feature k saath
multiply hogi a

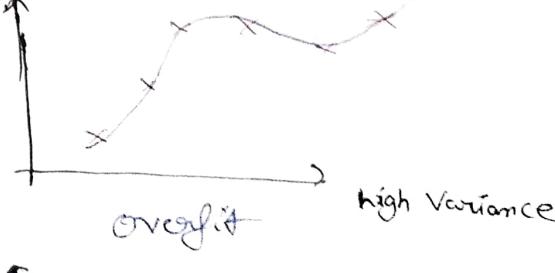
Overfitting



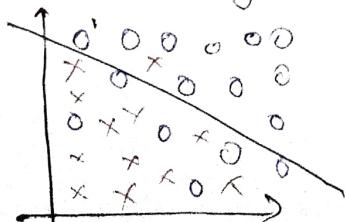
Underfit



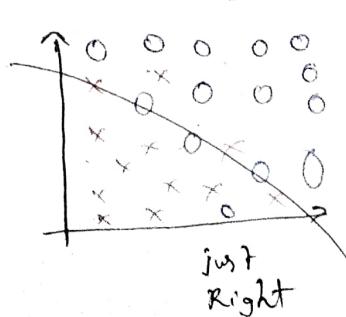
Just Right



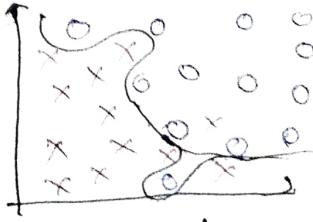
Overfit



underfit



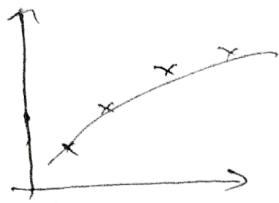
just right



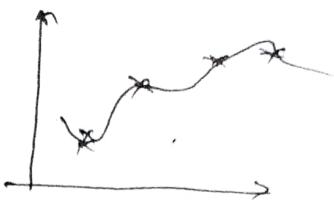
overfit

high variance

Regularization



$$w_1 x + w_2 x^2 + b$$



$$w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

make w, w_4 really small

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

choose $\lambda = 10^{-6}$
then algorithm will
choose the w_j small

fits data \rightarrow Keep w_j small
 λ balances both

gradient descent with regularization

$$f_{\vec{w}, b}(\vec{x}) = w_1 x + w_2 x^2 + w_3 x^3 + \cancel{w_4 x^4} + b$$

$$w_j = w_j - \alpha \left(\underbrace{\frac{1}{m} \sum (f_{\vec{w}, b}(\vec{x}) - y^i) x_j^i}_{\frac{\partial J(\vec{w}, b)}{\partial w_j}} + \frac{\lambda}{m} w_j \right)$$

$$b = b - \alpha \left(\underbrace{\frac{1}{m} \sum (f_{\vec{w}, b}(\vec{x}) - y^i)}_{\frac{\partial J(\vec{w}, b)}{\partial b}} \right)$$

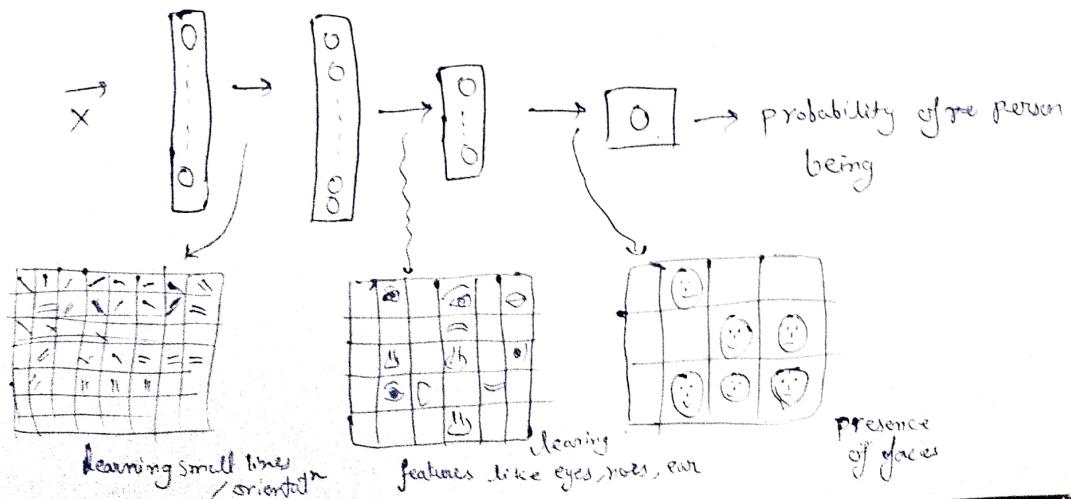
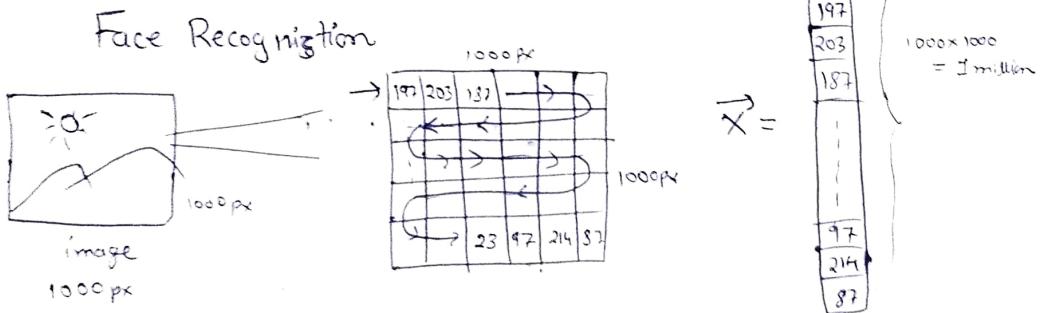
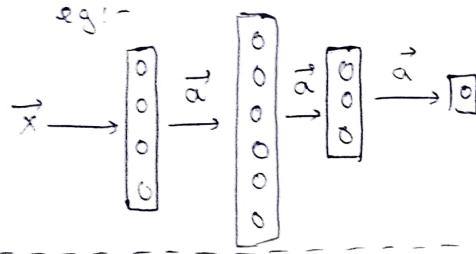
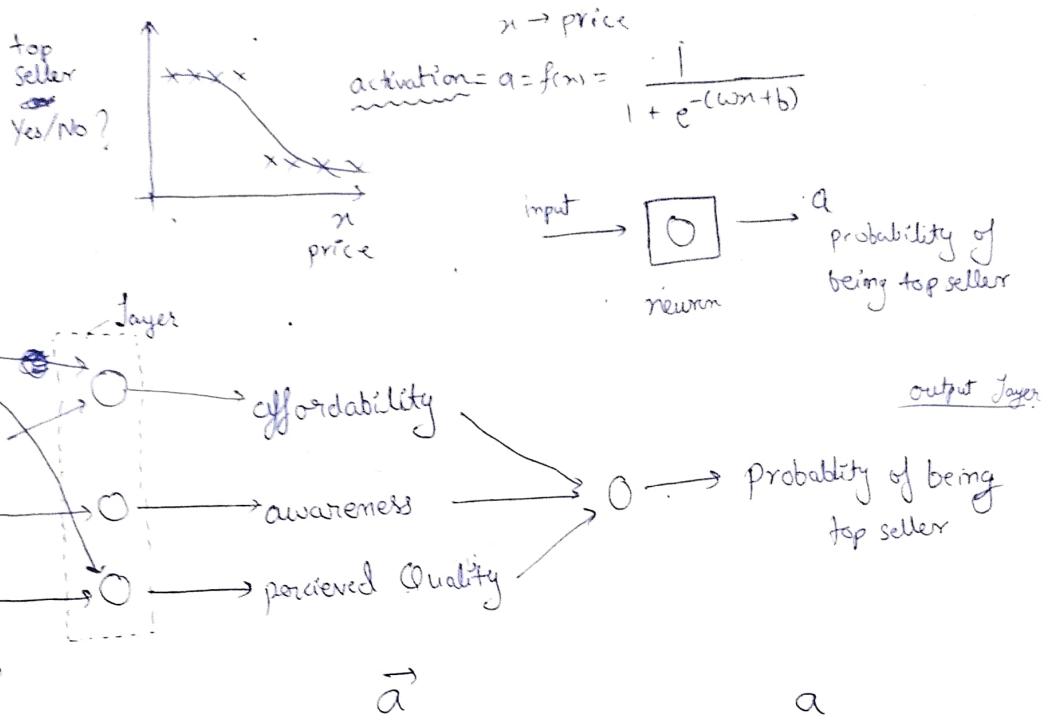
update rule of $w, b \rightarrow$

$$w_j = w_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \underbrace{\frac{1}{m} \sum (f_{\vec{w}, b}(\vec{x}^i) - y^{(i)}) x_j^i}_{\text{usual update}}$$

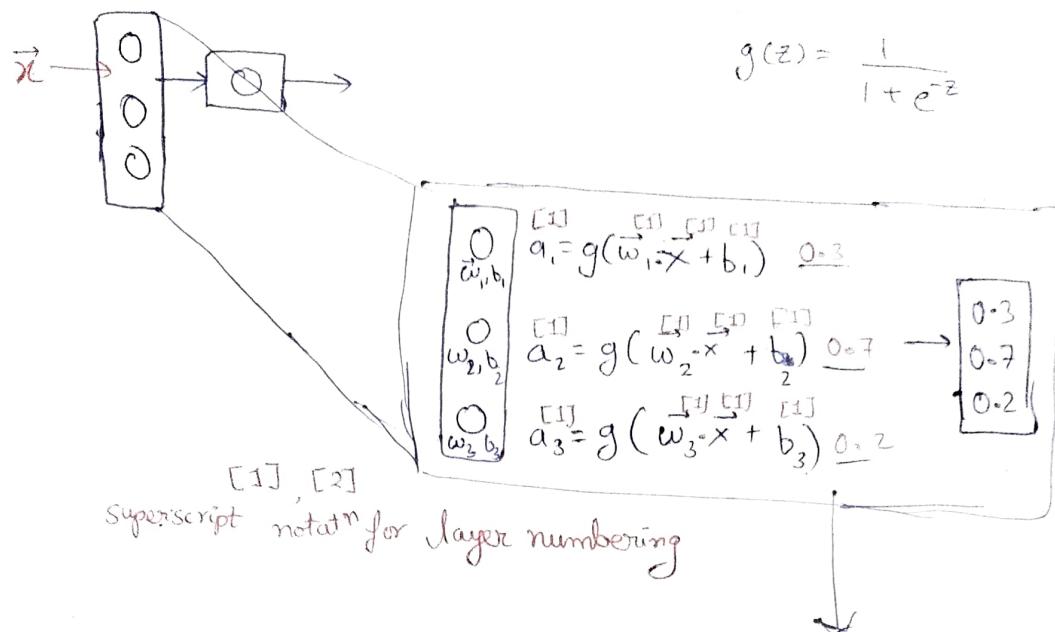
$$b = \text{usual update}$$

Neural Network

→ Demand Prediction



Neural network layer



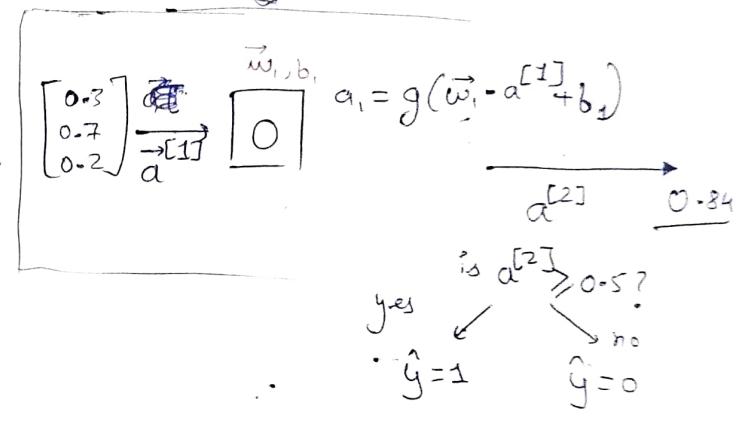
activⁿ value of layer l
unit (neuron) j

$$a_j^{[l]} = g(\vec{w}_j \cdot \vec{a}^{[l-1]} + b_j)$$

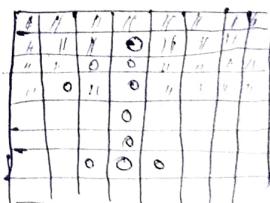
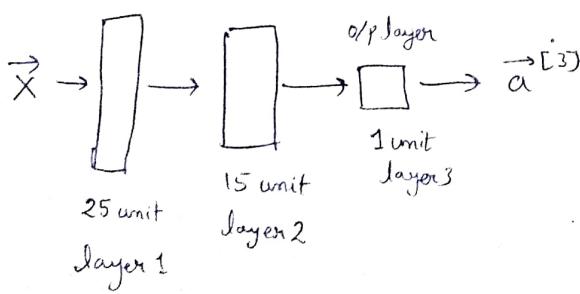
Sigmoid
Activation Function

Parameters w, b of layer l , unit j

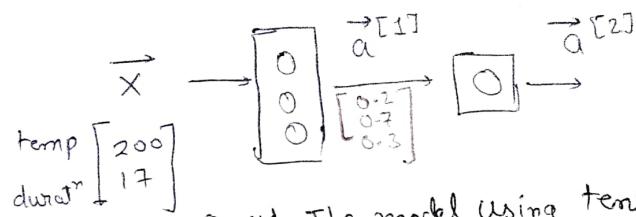
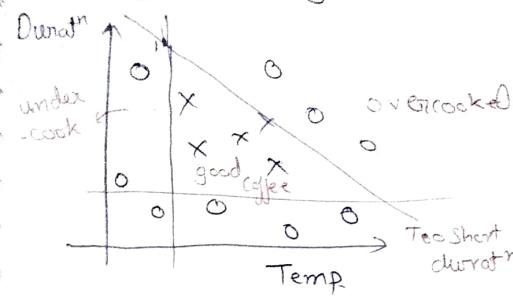
$$g(z) = \frac{1}{1 + e^{-z}}$$



Hand Written Digit Recognition



Coffee Roasting

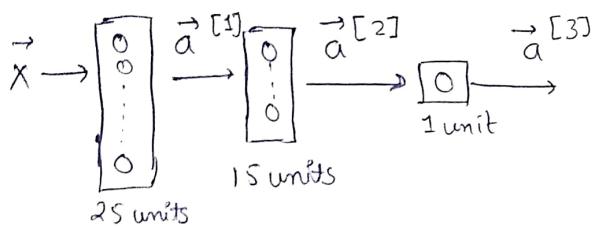


Build The model using tensorflow

```
x = np.array ([[200.0, 17.0]])
```

Tensor \leftarrow layer-1 = Dense (units=3, activation='sigmoid')
 $\left[\begin{bmatrix} 0.2, 0.7, 0.3 \end{bmatrix}\right]$ \leftarrow $a1 = \text{layer-1}(x)$
 $1 \times 3 \text{ matrix}$ if you print out $a1$ $\text{tf.Tensor}(\left[\begin{bmatrix} 0.2, 0.7, 0.3 \end{bmatrix}\right], \text{Shape}=(1, 3), \text{dtype}=\text{float32})$
 $a1$ layer-2 = Dense (units=1, activation='sigmoid')
 $\left[\begin{bmatrix} 0.7, 0.3 \end{bmatrix}\right], \text{dtype}=\text{float32}$ $a2 = \text{layer-2}(a1)$

model for digit classification



```

x = np.array ([[0.0, ..., 255, ..., 240, ..., 0]])
layer_1 = Dense (units=25, activation='sigmoid')
a1 = layer_1(x)
layer_2 = Dense (units=15, activation='sigmoid')
a2 = layer_2(a1)
layer_3 = Dense (units=1, activation='sigmoid')
a3 = layer_3(a2)

```

if $q3 \geq 0.5$
 \downarrow
 $\hat{y} = 1$ else $\hat{y} = 0$

numpy arrays / Matrices

$x = np.array([200, 17])$ $x = np.array([[200, 17]])$ $\rightarrow [200, 17]$ 1×2

$x = np.array([[200], [17]])$ $\rightarrow \begin{bmatrix} 200 \\ 17 \end{bmatrix}$ 2×1

`x = np.array([200, 17])` → 1D vector

Building a neural network architecture

model = sequential ([

```
Dense(units=3, activation="sigmoid"),  
Dense(units=1, activation="sigmoid")])
```

```
x = np.array([ [200, 17],  
               [120, 5],  
               [425, 20],  
               [212, 18] ])
```

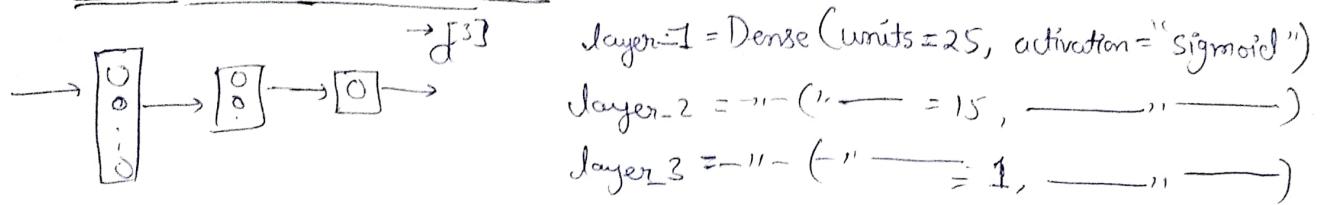
200	17	1
120	5	0
425	20	0
212	18	1

```

y = np.array([1, 0, 0, 1])
model.compile(..., ...)
model.fit(x, y)
model.predict(x_new)

```

→ Digit Classification model :-



model = Sequential([layer1, layer2, layer3])

model.compile(...)

x = np.array([[0, ..., 245, ..., 17, ...],
[0, ..., 200, ..., 184, ...]])

y = np.array([1, 0])

model.fit(x, y)

model.predict(x_new)

Coffee roasting model (using numpy / py)

forward Prop

x = np.array([200, 17])

$a_1^{[1]} = g(\vec{w}_1^{[1]} \cdot \vec{x} + b_1^{[1]})$

$w_{1-1} = np.array([1, 2])$

$b_{1-1} = np.array([-1])$

$z_{1-1} = np.dot(w_{1-1}, x) + b_{1-1}$

$a_{1-1} = \text{Sigmoid}(z_{1-1})$

Similarly

$a_{1-2} \checkmark$

$a_{1-3} \checkmark$

$a_1 = np.array([a_{1-1}, a_{1-2}, a_{1-3}])$

forward Prop in Numpy

$W \rightarrow$ matrix

def dense(a_in, W, b):

3 units = $W.shape[1]$

$a_out = np.zeros(\text{units})$ [0, 0, 0]

for j in range(units): $\rightarrow [0, 1, 2]$ pulling out 1st column from 2D array

~~W~~ $w = W[:, j]$

$z = np.dot(W, a_in) + b[j]$

$a_out[j] = g(z)$

return a_out

eg:-
 $W_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, W_2 = \begin{bmatrix} -3 \\ 4 \end{bmatrix}, W_3 = \begin{bmatrix} 5 \\ -6 \end{bmatrix}$

$W = np.array([1, -3, 5, 2, 4, -6])$

$b_1 = -1, b_2 = 1, b_3 = 2$

$b = np.array([-1, 1, 2])$

$\vec{a} = \vec{x}$

$a_in = np.array([-2, 4])$

def sequential(x):

$a_1 = \text{dense}(x, W_1, b_1)$

$a_2 = \text{dense}(a_1, W_2, b_2)$

$a_3 = \text{dense}(a_2, W_3, b_3)$

$a_4 = \text{dense}(a_3, W_4, b_4)$

$f_x = a_4$

return f_x

Dot prod

$\begin{bmatrix} 1 \\ a \end{bmatrix} \cdot \begin{bmatrix} 1 \\ w \end{bmatrix}$

$z = \vec{a} \cdot \vec{w}$

$\begin{bmatrix} \vec{a} \end{bmatrix} \cdot \begin{bmatrix} \vec{w} \end{bmatrix}$

$z = \vec{a}^T \vec{w}$
 vector
 vector
 multiplication

vectorizatⁿ of above code

$x = np.array([[200, 17]])$

$W = np.array([1, -3, 5, -2, 4, -6])$

$B = np.array([-1, 1, 2])$

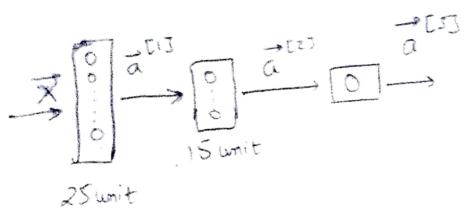
def dense(A_in, W, B):

$Z = np.matmul(A_in, W) + B$

$A_out = g(Z)$

return A_out

Training NN im Tensorflow



```

import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=15, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid')
])

```

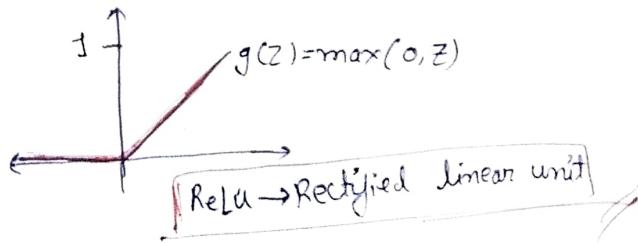
from tensorflow.keras.losses import

BinaryCrossentropy

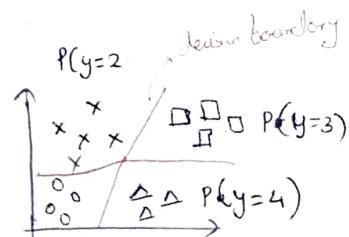
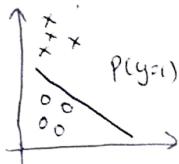
```
model.compile(loss=BinaryCrossentropy())
```

```
model.fit(x, y, epochs=100)
```

→ alternatives to sigmoid fn in Activation



Multiclass classification Problem:-



{Softmax} \hookrightarrow Binary \rightarrow Multi class classification problem.

Regression

$$\times \quad z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$\circ \quad z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$\square \quad z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$\triangle \quad z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1|\vec{x})$$

$$a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}}$$

$$= P(y=2|\vec{x})$$

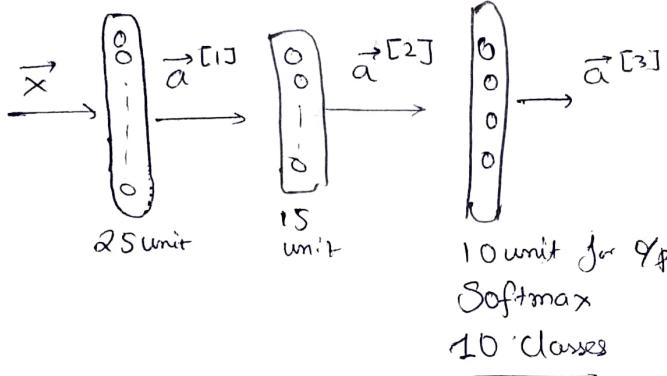
Softmax Regression

Crossentropy loss

$$\text{loss}(a_1, \dots, a_N, y) = \begin{cases} -\log a_1 & \text{if } y=1 \\ -\log a_2 & \text{if } y=2 \\ \vdots & \vdots \\ -\log a_N & \text{if } y=N \end{cases}$$

Neural Network with Softmax Output

Digit Recognition Model.



$$z_1 = \vec{w}_1^{[3]} \cdot \vec{a}^{[2]} + b_1$$

$$\vec{a}_1^{[3]} = \frac{e^{z_1^{[3]}}}{e^{z_1^{[3]}} + \dots + e^{z_{10}^{[3]}}} = P(y=1 | \vec{x})$$

better (recommended) version later.

MNIST with Softmax

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras import layers
```

Don't use the code version here!

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')])
```

```
from tensorflow.keras.losses import
SparseCategoricalCrossentropy
```

```
model.compile(loss=SparseCategoricalCrossentropy(from_logits=True))
model.fit(X, Y, epochs=100)
```

for more
numerically accurate
implementation

Additional Concepts

- ~~ADA~~ Adam Optimisation
- Convolutional layer

Debugging Regression model:-

- get more training examples, → Try smaller sets of features
- try getting additional feature, → Try adding polynomial features
(x_1^2, x_2^2, \dots)
- try decreasing, increasing λ

Evaluating and Choosing model:-

• model selection

$$d=1$$

$$f_{w_b} = w_0 n_0 + b$$

$$\rightarrow J_{\text{test}}^{(1)}$$

$$d=2$$

$$f_{w_b} = w_0 n_0^2 + w_1 n_1 + b$$

$$\rightarrow J_{\text{test}}^{(2)}$$

$$d=3$$

$$f_{w_b} = w_0 n_0^3 + w_1 n_1^2 + w_2 n_2 + b$$

$$\rightarrow J_{\text{test}}^{(3)}$$

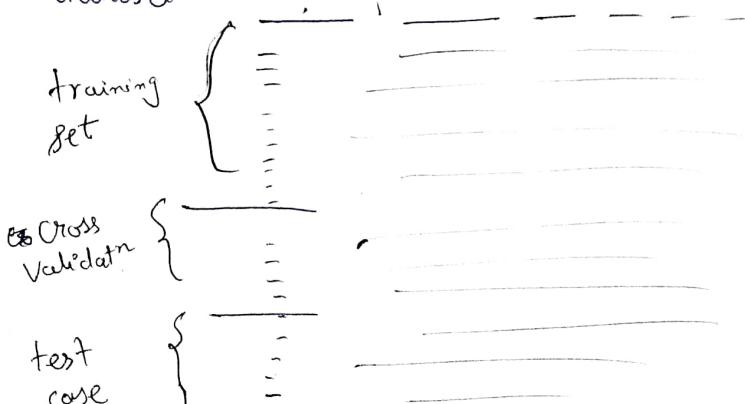
$$d=10$$

$$f_{w_b} = w_0 n_0^{10} + \dots + w_9 n_9 + b \rightarrow J_{\text{test}}^{(10)}$$

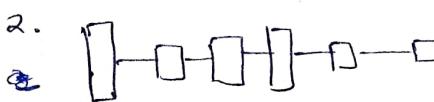
✓ (works well).

$$\xleftarrow{\text{lowest } J_{\text{test}}^{(5)}} w_0 n_0 + w_1 n_1^2 + w_2 n_2^3 + \dots + w_5 n_5^5 + b$$

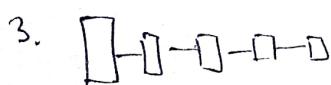
• dataset



$$J_{\text{cv}}(w^{(1)}, b^{(1)})$$



$$J_{\text{cv}}(w^{(2)}, b^{(2)})$$



$$J_{\text{cv}}(w^{(3)}, b^{(3)})$$

Train \curvearrowleft

low error \rightarrow
model selected

Pick $w^{(2)}, b^{(2)}$

↳ generalisatⁿ error using test set

$$J_{\text{test}}(w^{(2)}, b^{(2)})$$

→ regularization and bias/variance

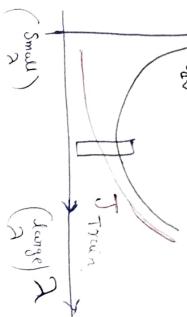
$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

minimizing parameter

Try $\lambda = 0$
Try $\lambda = 0.01$

$\rightarrow \lambda = 0.08$ (goes well)

Try $\lambda \approx 10$



→ Establishing baseline level of Performance.

- Human level performance
- Competing algorithms —
- Curves based on experience

→ Scanning Curves

$$f_{w, b} = w_1 x_1 + w_2 x_2 + b$$

$$J_{\text{cv}}(\vec{w}, b)$$

$$J_{\text{train}}(\vec{w}, b)$$

m_{train}
(train set size)

Debugging algo - revisited.

- get more training set → fines high variance
- try smaller set of feature → fines high variance
- adding additional features → fines high bias.

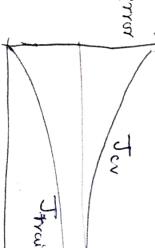
- try adding poly. features → fines high bias → helps to do better on the training set
- decreasing λ → fines high bias
- increasing λ → fines high variance

$$J_{\text{cv}}$$

$$J_{\text{train}}$$

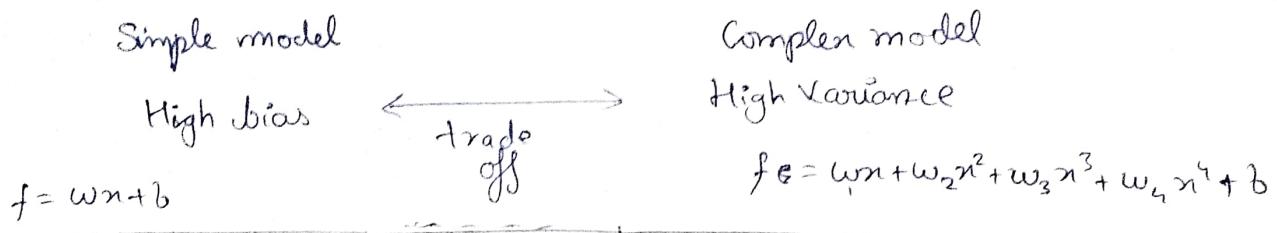
$$m_{\text{train}} (\text{train set size})$$

with high variance

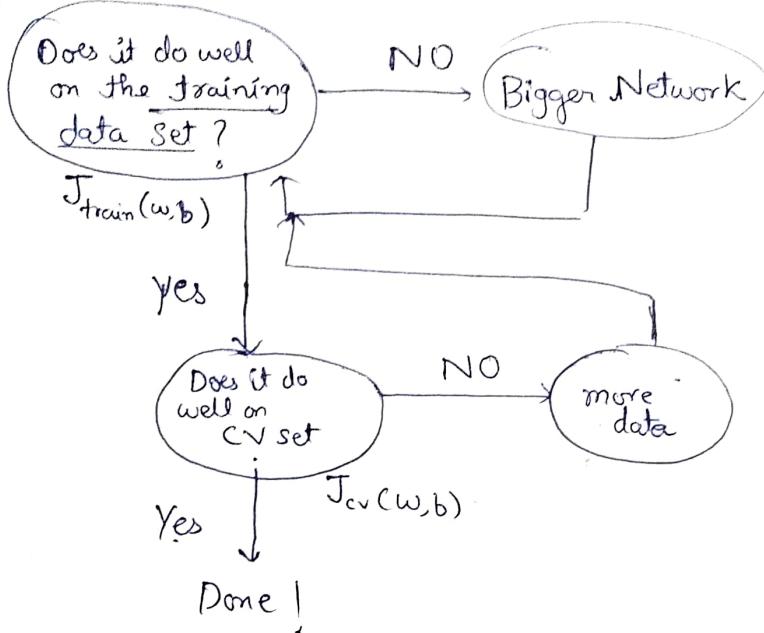


m_{train} (train set size)

with high variance



→ Large N^2 are ~~are~~ slow bias bias machine.



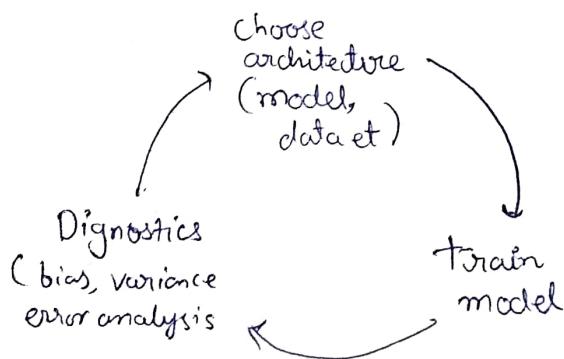
N. N. regularization

$$J(W, B) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum (w^2)$$

all weight w

ML Development process :-

• Iterative loop of ML dev.

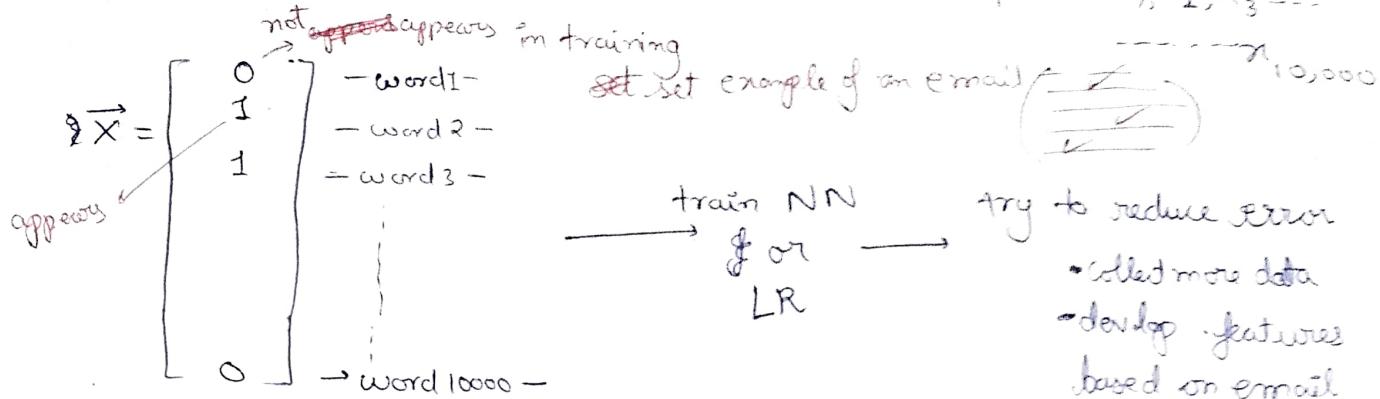


e.g. Building a spam classifier

Super Supervised learning : \vec{x} = features of email

$y = \text{Spam}(1)$ or not $\text{Spam}(0)$

features: list the top 10,000 words to compute. x_1, x_2, x_3, \dots



- Define features from email body
- Design algo. to detect misspellings

eg: W4tches, m0rtgage

→ Error Analysis

$m_{cv} = 500$ examples. in CV set

Algorithm misclassifies 100 of them
manually examine 100 eg. and categorize them on common traits.

- eg of common traits
- {
 - ~ deliberate misspellings (w4tches, etc)
 - ~ Unusual routing
 - ~ Steal passwords
 - ~ Spam message in embedded image

→ B Adding Data:

- Artificial Data Synthesis
- Data augmentation → create new data by modifying ~~the~~ previous data

eg,

$$A \rightarrow A \vee A \wedge A$$

$$B \rightarrow B \otimes B \otimes B$$

→ by introducing distortions.

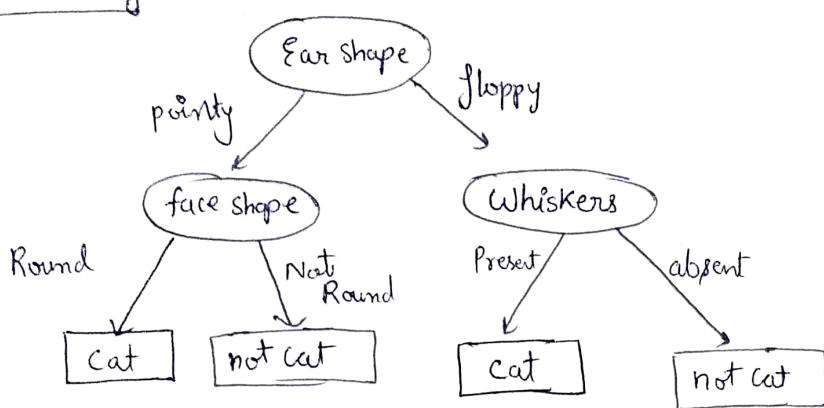
- → used for speech recognition by adding different - different background noise effects.

AI = Code + Data.

Decision Tree Model

eg

cat identifier



- How to choose what feature to split on at each node?

→ ensure max. maximum purity.

- When to stop splitting?

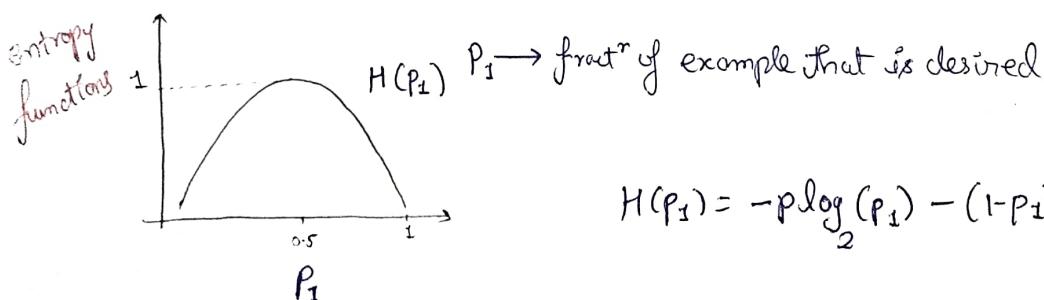
→ when a node is 100% one class

→ When splitting a node will result in the tree exceeding a maximum depth.

→ when no. of egs in a node is below a threshold.

→ when improvements in purity score are below a threshold

Measuring purity



$$H(p_1) = -p_1 \log_2(p_1) - (1-p_1) \log_2(1-p_1)$$

information Gain

$$\text{info. Gain} = H(p_1^{\text{root}}) - \left(w_{\text{left}}^{\text{root}} H(p_1^{\text{left}}) + w_{\text{right}}^{\text{root}} H(p_1^{\text{right}}) \right)$$

p_{true}
Situation
of classifier

Weighted avg
fraction
 $\left(\frac{\text{no. egs in this}}{\text{Total no. of egs}} \right)$

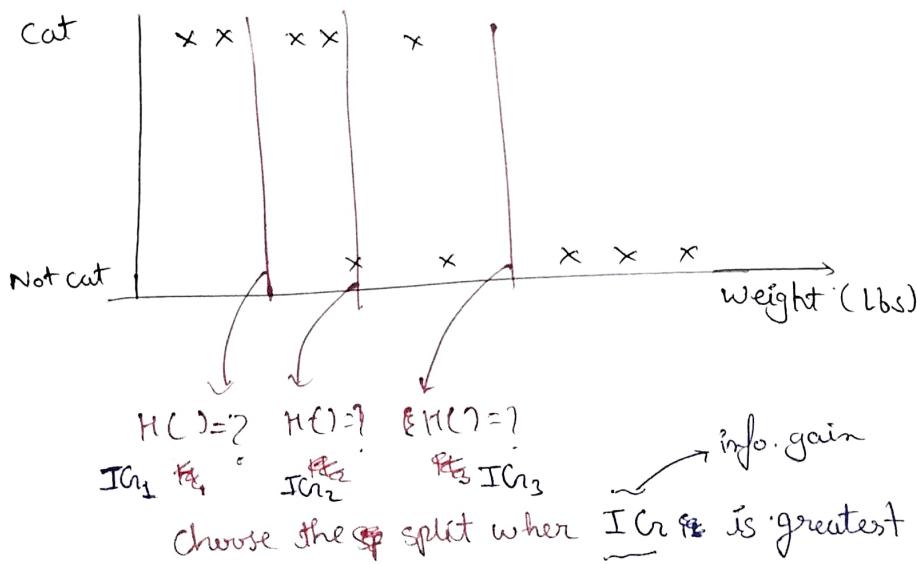
→ Decision tree learning:-

- Start with all examples at the root node.
- Calculate information gain for all possible features, and pick the one with the highest information gain.
- Split dataset according to selected feature, and create left and right branches of tree.
- Keep repeating, splitting process until stopping criteria is met.

One-hot encoding:- if a feature take 'K' values,

→ Create K features with (0 or 1).

Splitting a ~~on~~ continuous Variable:-



XGB Boost

from xgboost import XGBClassifier

model = XGBClassifier()

model.fit(X_train, Y_train)

Y_pred = model.predict(X-test)

Regression

from xgboost import XGBRegressor

model = XGBRegressor()

model.fit(X_train, Y_train)

Y_pred = model.predict(X-test)

Decision trees Vs Neural Network

- works well on tabular Data
- not recommended for unstructured data
 - (image, audio, txt)
- fast
- small decision are human interpretable

- works well on all types of data
- maybe slower than a decision tree.
- works with transfer learning