

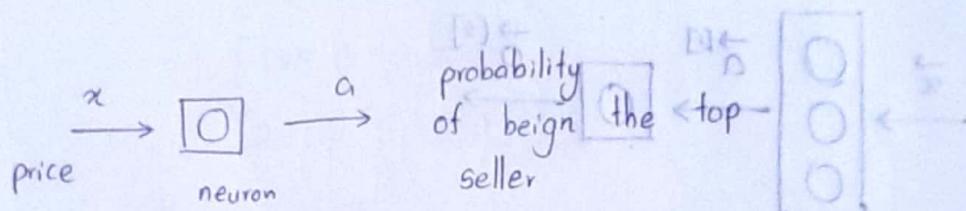
Advanced Learning Algorithms

Demand Prediction

x = price

input \leftarrow time

$$a \underset{\uparrow \text{activation}}{=} f(x) = \frac{1}{1 + e^{-(wx+b)}} \quad \text{output}$$



Neural network layer

$$\vec{x} \begin{bmatrix} 197 \\ 184 \\ 136 \\ 214 \end{bmatrix} \rightarrow \begin{array}{c} \text{neuron } 1 \\ w_1^{[1]}, b_1^{[1]} \\ \text{neuron } 2 \\ w_2^{[1]}, b_2^{[1]} \\ \text{neuron } 3 \\ w_3^{[1]}, b_3^{[1]} \end{array}$$

$a_1 = g(\vec{w}_1 \cdot \vec{x} + b_1) = 0.3$

$a_2 = g(\vec{w}_2 \cdot \vec{x} + b_2) = 0.7$

$a_3 = g(\vec{w}_3 \cdot \vec{x} + b_3) = 0.2$

$\vec{a}^{[1]} = \begin{bmatrix} 0.3 \\ 0.7 \\ 0.2 \end{bmatrix}$

$$a \underset{\uparrow [1]}{=} \begin{array}{c} \vec{w}_1^{[1]} \\ b_1^{[1]} \end{array} \quad a_1 = g(\vec{w}_1 \cdot \vec{a}^{[1]} + b_1)$$

$a^{[2]} = 0.84$

if $[a^{[2]}] \geq 0.5$

yes $\hat{y} = 1$

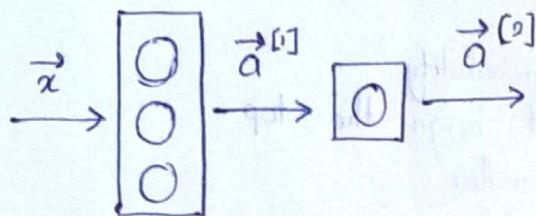
no $\hat{y} = 0$

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

input $\vec{x} = a^{[0]}$

layer $\rightarrow l$
unit $\rightarrow j$

Inference in Code



$x = \text{np.array}([[200.0, 17.0]])$:

$\text{layer_1} = \text{Dense}(\text{units} = 3, \text{activation} = \text{'sigmoid'})$

$a_1 = \text{layer_1}(x)$

$\text{layer_2} = \text{Dense}(\text{units} = 1, \text{activation} = \text{'sigmoid'})$

$a_2 = \text{layer_2}(a_1)$

if $a_2 \geq 0.5$:

$y_{\hat{h}} = 1$

else $y_{\hat{h}} = 0$

$x = \text{np.array}([[200, 17]]) \rightarrow [200, 17] \leftarrow 2D \text{ array}$

$x = \text{np.array}([200, 17]) \rightarrow \begin{matrix} \text{no row} \\ \text{or column} \end{matrix} \leftarrow 1D \text{ vector array}$

Building a neural network architecture

```
layer_1 = Dense (units = 3, activation = "sigmoid")  
layer_2 = Dense (units = 1, activation = "sigmoid")  
model = Sequential ([layer_1, layer_2])
```

```
x = np.array ([[200.0, 17.0],  
[120.0, [5.0],  
[425.0, 20.0],  
[212.0, 18.0]]])
```

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

```
model.compile (...)
```

```
model.fit (x, y)
```

```
model.predict (x-new)
```

Vectorization

loops vs vectorization

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

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

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

units = w.shape[1]

```
def dense (a-in, w, b):
```

```
a-out = np.zeros (units)
```

```
for j in range (units)
```

```
w = w[:, j]
```

```
z = np.dot (w, x) + b[j]
```

```
a[j] = g(z)
```

```
return a.
```

```
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.T, W) + B
```

```
A-out = g(z)
```

```
return A-out
```

$\text{AT} = \underline{\text{A.T}}$

Transpose

$z = \text{np. matmul}(\text{AT}, w)$

$(z = \text{AT} @ w)$

$\text{AT} = \text{np. array}([[200, 17], [0.01, 0.001]])$

$w = \text{np. array}([[1, [-3, 5], [0.01], [-2, 4, -6]]])$

$b = \text{np. array}([[-1, 1, 2]])$

$z = \text{np. matmul}(\text{AT}, w) + b$

$a_{\text{out}} = g(z)$

$([[1, 200]])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1^2, 1, 1])$

$\therefore (d, W, n, A) \text{ versch. hab}$

$\therefore (W, n, A) \text{ konstan. } q_n = 1$

$(c) e = \text{two. A}$

$\text{two. A} \text{ matrix}$

notwendig zu $n = 1$

$(1, 1, 1)$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

$([1, 1, 1])$

Train a NN in Tensor Flow

```
import tensorflow as tf  
from tensorflow.keras import Sequential  
from tensorflow.keras.layers import Dense  
  
model = Sequential ([  
    Dense (units = 25, 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)
↑
number of steps
in gradient descent

logistic loss → Binary Crossentropy

Regression (Predicting numbers)

Loss → MeanSquareError ()

Choosing Activation Function

Output layer

Binary classification - Sigmoid

Regression - Linear activation function
(+ and -)

activation = 'linear'

Only + - ReLU

activation = 'relu'

Hidden layer

Common - ReLU

Softmax

Softmax regression (N possible outputs)

$$y = 1, 2, \dots, N$$

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j = 1, \dots, N$$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j | \vec{x})$$

$$a_1 + \dots + a_N = 1$$

for 2 outputs

$$a_1 = g(z)$$

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

$$= P(y=1 | \vec{x})$$

$$a_1 + a_2 = 1$$

Cost

Logistic Regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1 + e^{-z}} = P(y=1 | \vec{x})$$

$$a_0 = 1 - a_1 = P(y=0 | \vec{x})$$

$$\text{loss} = -y \log a_1 - (1-y) \log \underbrace{(1-a_1)}_{a_0}$$

$$J(\vec{w}, b) = \text{average loss}$$

Neural Network with softmax

```
import tensorflow as tf
```

```
from tensorflow.keras import Sequential
```

```
from tensorflow.keras.layers import Dense
```

```
model = Sequential([
```

```
Dense(units=25, activation='relu'),
```

```
Dense(units=15, activation='relu'),
```

```
Dense(units=10, activation='sigmoid'))]
```

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

SparseCategoricalCrossentropy

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

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

Softmax regression

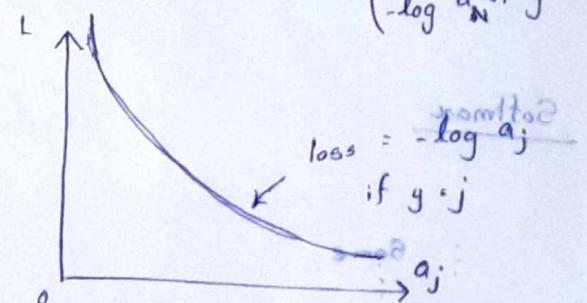
$$a_i = \frac{e^{z_i}}{e^{z_1} + \dots + e^{z_N}} = P(y=i | \vec{x})$$

$$a_N = \frac{e^{z_N}}{e^{z_1} + \dots + e^{z_N}} = P(y=N | \vec{x})$$

Crossentropy loss

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

Crossentropy loss



Tensorflow

Session

sess.run

eval

accuracy

best

loss

min

epoch

batch

step

iteration

gradient

update

weights

parameters

variables

optimizer

learning rate

momentum

decay

epsilon

beta

gamma

clipnorm

clipvalue

epsilon

beta_1

beta_2

epsilon

clipnorm

<p

Improved version of Softmax

Logistic Regression

Numerical Roundoff error

model = Sequential ([

Dense (units = 25, activation = 'relu')

Dense (units = 15, activation = 'relu') 'linear'

Dense (units = 10, activation = 'sigmoid')

model.compile (loss = BinaryCrossEntropy())

model.compile (loss = BinaryCrossEntropy (from_logits = True))

logit : z

Softmax

$\text{softmax} = \frac{e^z}{\sum e^z}$

i.e.:

Same

Dense (units = 10, activation = 'linear')

model.compile (loss = SparseCrossEntropy (from_logits = True))

* get output before softmax

a_1, \dots, a_{10}

z_1, \dots, z_{10}

predict: logits = model(x)

$f_x = \text{tf.nn.softmax}(\text{logits})$

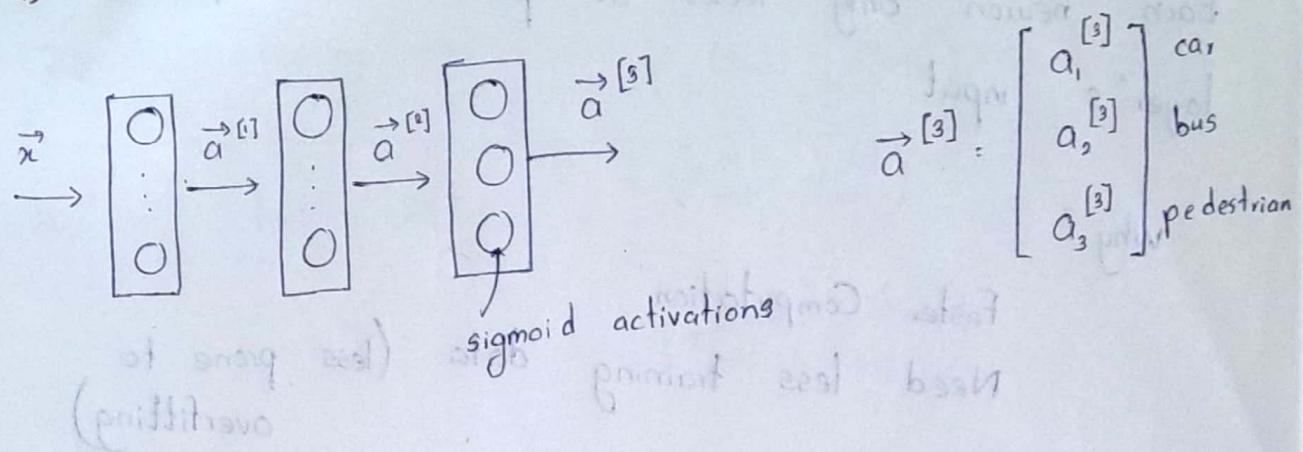
for logistic regression

logit = model(x)

$f_x = \text{tf.nn.sigmoid}(\text{logit})$

Classification with multiple output

(x) Multi-label classification (Self driving car)



Advanced Optimization

Gradient Descent Optimization

Adam : Adaptive Moment estimation

not just one α

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

$$w_{i0} = w_{i0} - \alpha_{i0} \frac{\partial}{\partial w_{i0}} J(\vec{w}, b)$$

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

If w_j (or b) keeps moving
same direction
oscillation

$\alpha_j \uparrow$
 $\alpha_j \downarrow$

code
compile

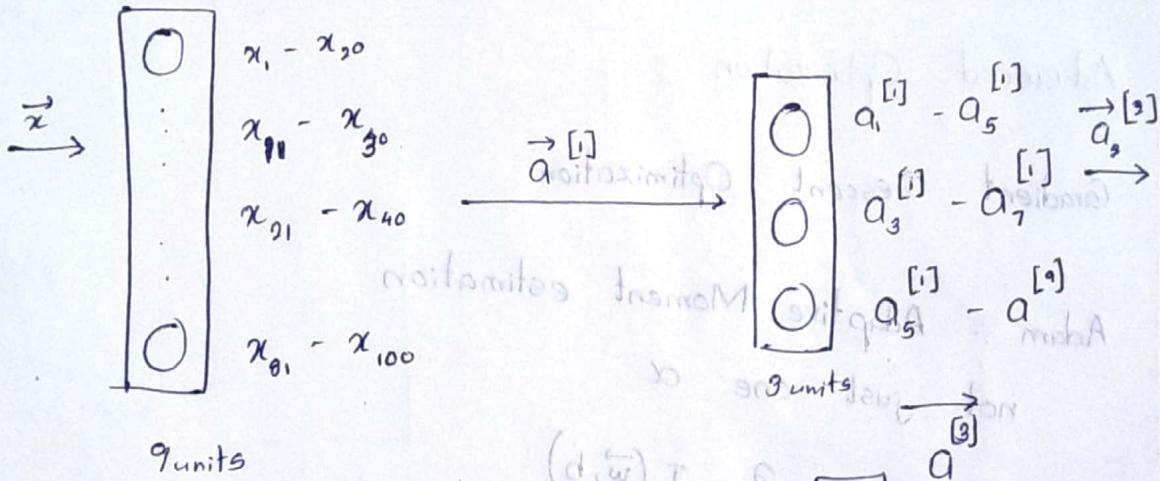
model.compile(optimizer = tf.keras.optimizers.Adam(learning_rate = 1e-3),
loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))

Additional Layer Types
Convolutional Layer
 Each neuron only looks at part of the previous layer's input

Why?

Faster Computation

Need less training data (less prone to overfitting)



$$(d, \vec{w}) \vdash \frac{\partial}{\partial \theta} = d$$

sigmoid

$$(d, \vec{w}) \vdash \frac{\partial}{\partial \theta} = d$$

↑ d Previous epoch $(d, \vec{w}) \vdash$
 ↓ d current epoch $(d, \vec{w}) \vdash$
not taken

stair, piano) make proximity exist. If = proximity) piano, piano
 (8-91)
 (but not piano) piano (but not piano) edge, exist, exist. If = not

Evaluating your model

$$(x^{(i)}, y^{(i)})$$

70% training set $(x^{(m_{train})}, y^{(m_{train})})$

30% test set $(x_{test}^{(i)}, y_{test}^{(i)})$

$$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$$

The test is test vs no broad noiseless soft set

Minimize cost function

for regression

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

Compute test error

$$J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \left[\sum_{i=1}^{m_{test}} (f_{\vec{w}, b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2 \right]$$

Training error

$$J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \left[\sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}_{train}^{(i)}) - y_{train}^{(i)})^2 \right]$$

for Classification

minimize

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

Test error - only testing set

Training error - only training set

Model Selecting & Training

Training 60%

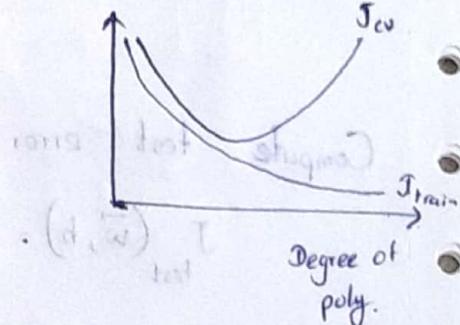
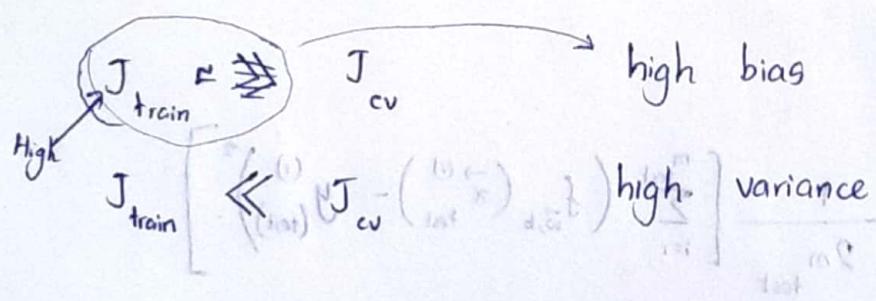
Cross Validation 20% ← Dev set (Development)

Testing 20%

Select the model which has least CV error.

(*) Take the decision based on CV test or test set.

Bias and Variance



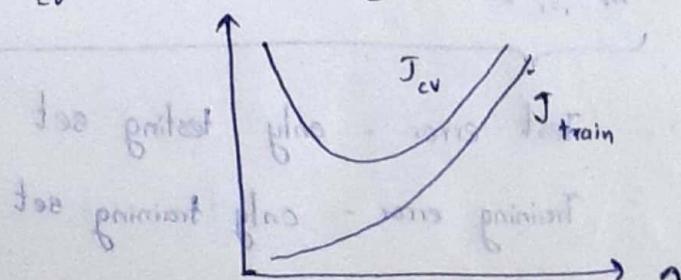
J_{train} high High bias & High Variance

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

Regularization and bias / variance

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

$\lambda = 0 \rightarrow 10$

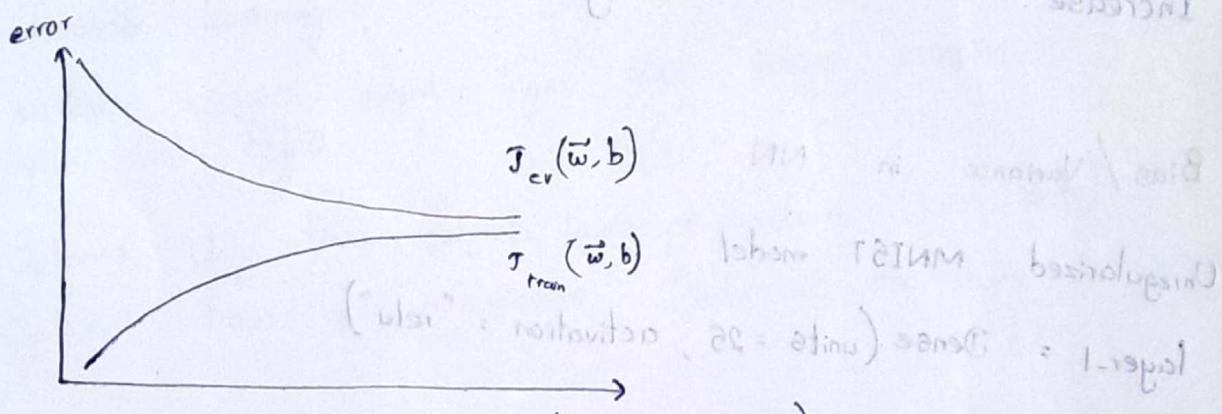


Establishing a baseline level of performance

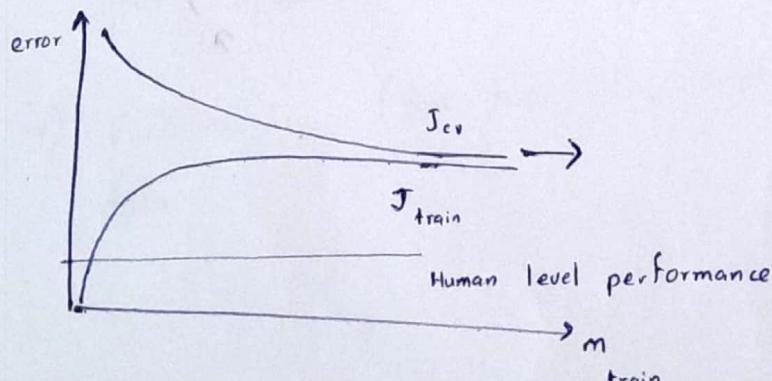
What is the level of error you can reasonably hope to get to?

- Human level performance
- Competing algorithms performance
- Guess based on experience.

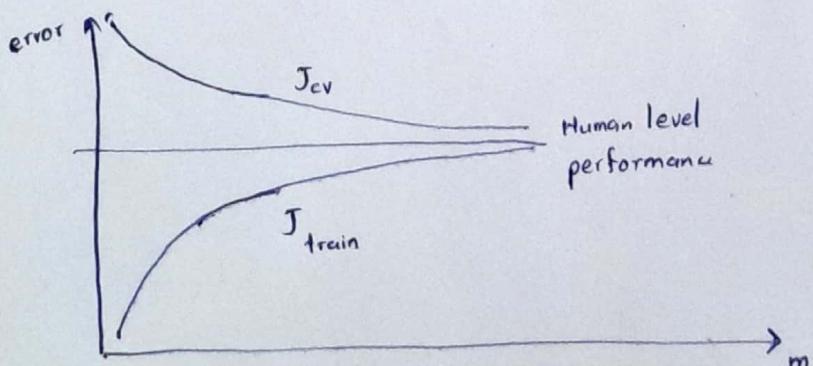
Learning curve



High bias



High Variance



High variance നീ
Data ഒരു നബ്ദാ.

High bias നീ
Data ഒരു നബ്ദാ.
ഒരു വ്യ.

Debugging a learning Algorithm
and ~~and~~ no map wants to load soft or hard
large error ~~(high bias)~~

- Get more training example - fixes high variance
Try smaller sets of features - fixes high variance
Try getting additional features - fixes high bias
Adding polynomial features - fixes high bias
Decrease λ - fix high bias
Increase λ - fix high variance

Bias / Variance in NN $(d, \omega)_\lambda$

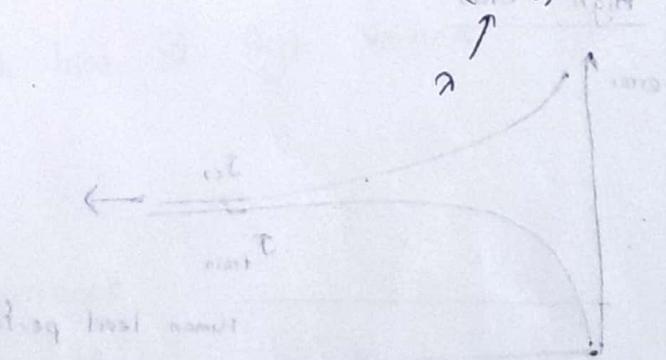
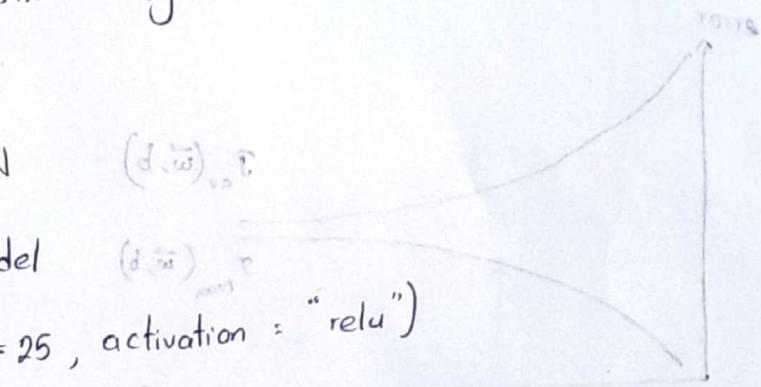
Unregularized MNIST model
layer-1 = Dense (units = 25, activation = "relu")

Regularized MNIST model $(d, \omega)_\lambda$
layer-1 = Dense (units = 25, activation = "relu", kernel_regularizer = L2(0.01))

See solid dip

See flat spot

at $m=1000$

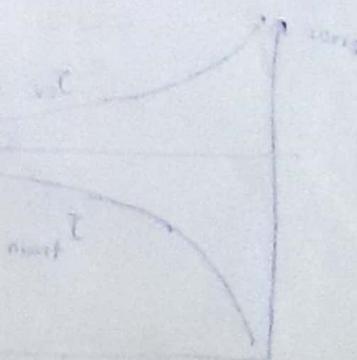


denominator last moment

m^2
moment

denominator dip

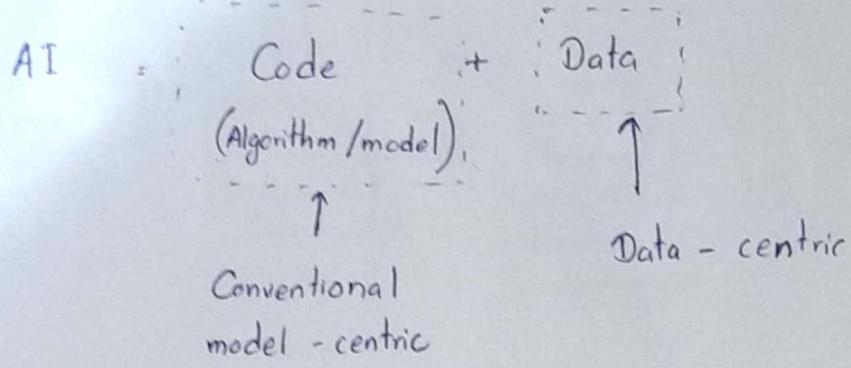
last moment
numerating



See another dip

where See spot

Engineering the data used by your system



Transfer Learning

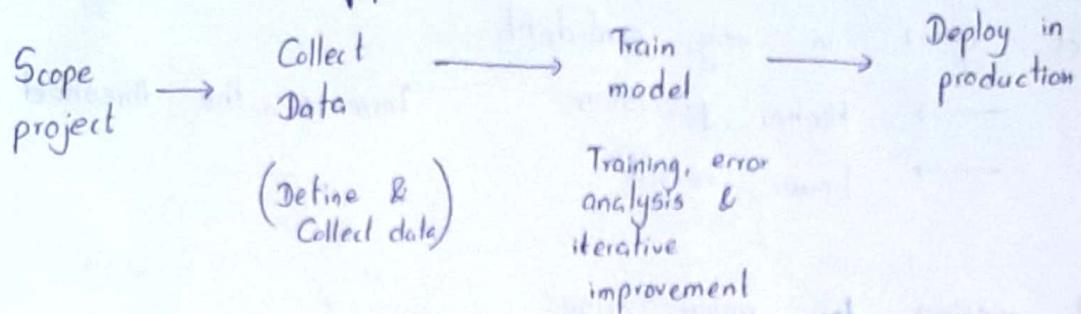
கேள்வி கொடுத்து output layer கிட வேண்டும் என்றால்
ஏதும் விட விரும்பும்

Option 1: Only train output layer parameters. (Less data)

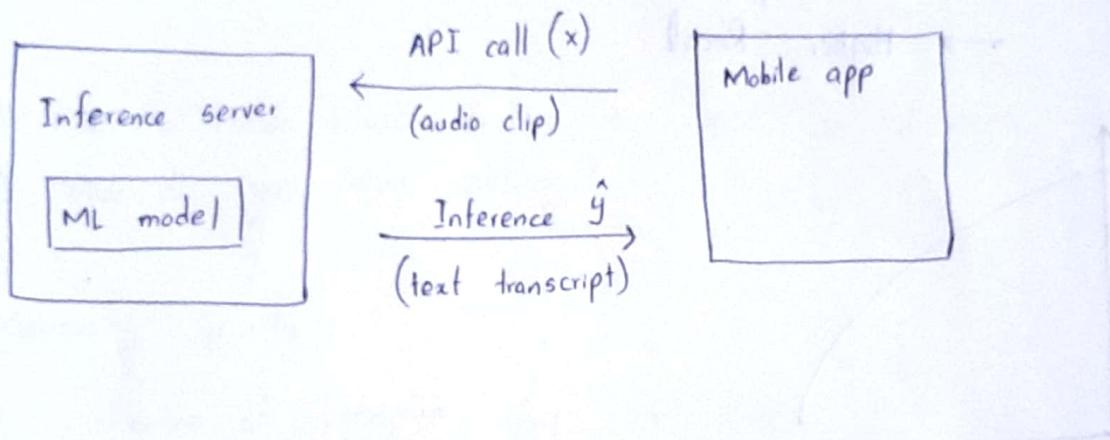
Option 2: Train all parameters (Data set சிறப்பாக அல்லது அதிக)

- 01) Download NN parameters pretrained on a large dataset with same input type (eg: images, audio, text) as your application (or train your own)
- 02) Further train (fine tune) the network on your own data.

Full cycle of ML project



Deployment



Fairness, bias & Ethics

Follow the guidelines.

Precision / Recall

		Actual class	
		1	0
Predicted class	1	True Positive	False Positive
	0	False Negative	True Negative

$$\begin{aligned} \text{Precision} &= \frac{\text{TP}}{\text{Predicted Positive}} \\ &= \frac{\text{TP}}{\text{TP} + \text{FP}} \end{aligned}$$

$$\begin{aligned} \text{Recall} &= \frac{\text{TP}}{\text{Actual Positive}} \\ &= \frac{\text{TP}}{\text{TP} + \text{FN}} \end{aligned}$$

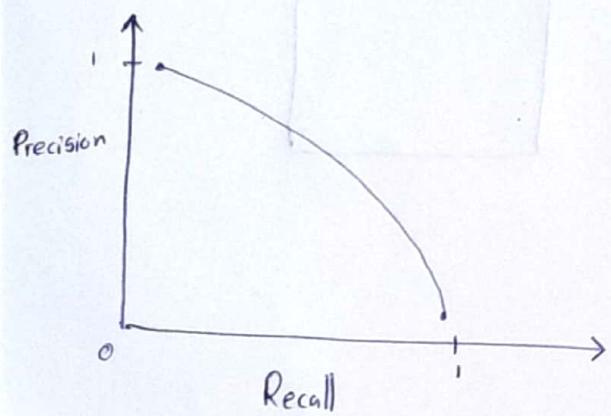
Trading off precision & Recall

predict $y = 1$ in very confident

- Higher precision Increase the threshold
- Lower recall

Avoid missing too many cases

- Lower precision Lower the threshold
- Higher Recall



F1 Score

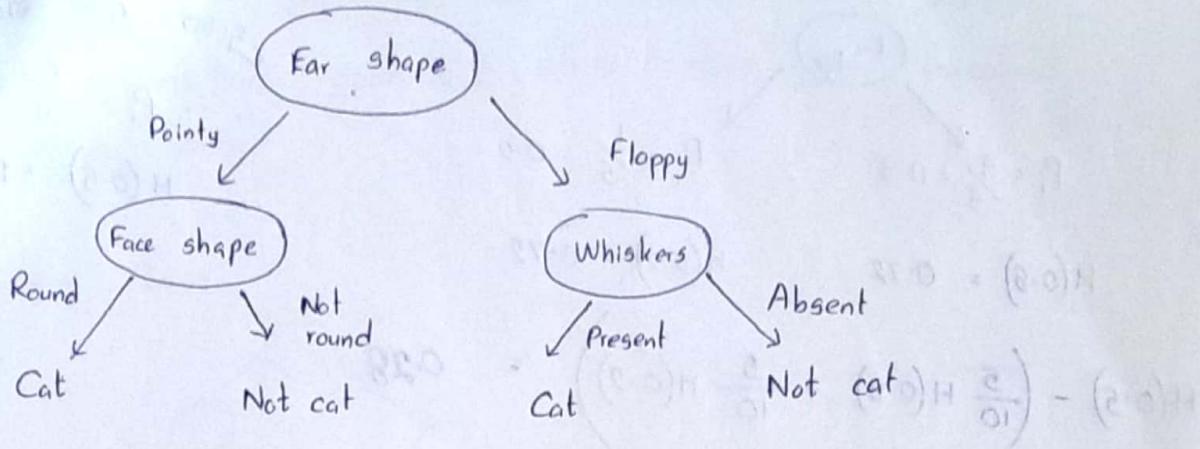
$$\text{F1 Score} = \frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right) = \frac{2PR}{P+R}$$

F1 Score ~~also~~ ~~also~~ ~~also~~ ~~also~~ ~~also~~

When do you stop splitting?

- * When a node is 100% one class
- * When splitting a node will result in the tree exceeding a maximum depth
- * When improvements in purity score are below a threshold
- * When number of examples in a node below a threshold

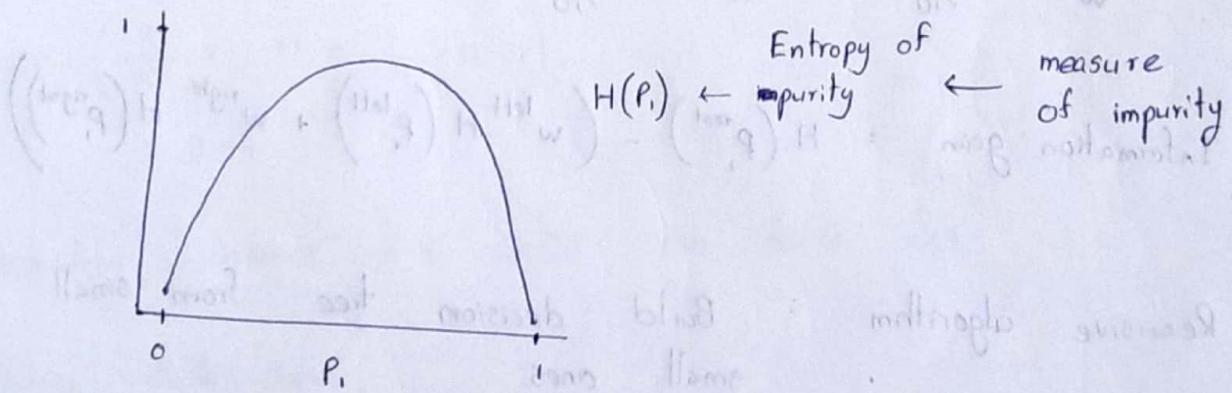
Decision Trees



- o1). How to choose what feature to split on at each node?
- o2). When do you stop splitting?

Measuring purity

P_i = fraction of examples that are cats



$$P_o = 1 - P_i$$

$$\begin{aligned} H(P_i) &= -P_i \log_2 (P_i) - P_o \log_2 (P_o) \\ &= -P_i \log_2 (P_i) - (1 - P_i) \log_2 (1 - P_i) \end{aligned}$$

Note. $0 \log(0) = 0$

Choosing a split : Information Gain

Ear shape

$$P_i = \frac{9}{10} = 0.8$$

$$P_i = \frac{1}{10} = 0.2$$

$$H(0.8) = 0.72$$

$$H(0.2) = 0.72$$

∴ $P_i = \frac{5}{10} = 0.5$

$$H(0.5) = 1$$

$$H(0.5) - \left(\frac{5}{10} H(0.8) + \frac{5}{10} H(0.2) \right) = 0.28$$

∴ calculation എന്തെന്ന് ഉള്ളേം value എന്ന ഗുണം.

Ear shape

$$P_i^{\text{root}} = \frac{5}{10} = 0.5$$

pointy

$$P_i^{\text{left}} = \frac{4}{5}$$

$$w^{\text{left}} = \frac{5}{10}$$

Floppy

$$P_i^{\text{right}} = \frac{1}{5}$$

$$w^{\text{right}} = \frac{5}{10}$$

$$\text{Information gain} = H(P_i^{\text{root}}) - \left(w^{\text{left}} H(P_i^{\text{left}}) + w^{\text{right}} H(P_i^{\text{right}}) \right)$$

Recursive algorithm : Build decision tree from small ones.

One hot encoding

Categorical features എന്ന് എന്ന്.

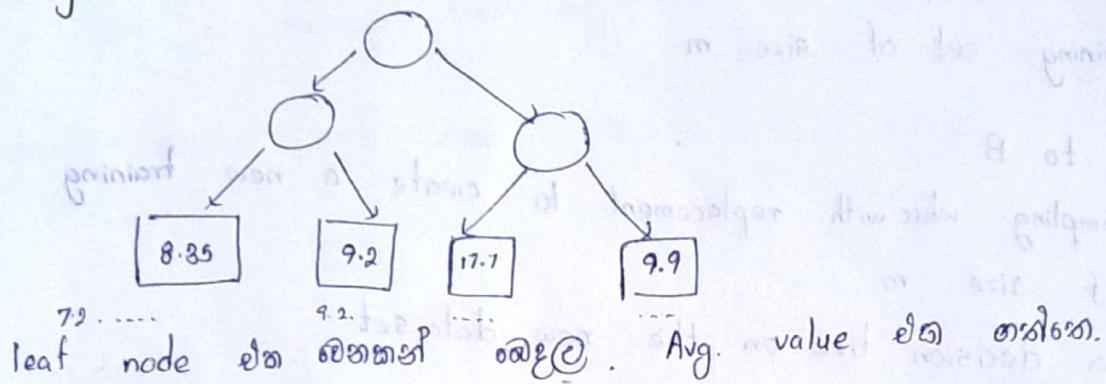
Pointy ears	Floppy ear	Round ear	Face shape	Whiskers present
1	0	0	Round	1

Continuous valued feature

Information Gain එක සොයුනු threshold value ඇතුව.

වෙත නිවැරදි Information gain එක හිසේ value එක threshold එක විසින් ගණනා කළ යුතුව.

Regression Trees



වෛශ්‍ය පිළිස්ථාපනය split කළ ඇත් සොයුනු variance එක නිශ්චිත නිවැරදි.

$$20.51 - \left(\frac{5}{10} \times 1.47 + \frac{5}{10} \times 21.87 \right)$$
$$= 8.84$$

Reduction variance එක වෙත නිවැරදි නිවැරදි.

Using multiple trees.

Tree ensemble : Trees ගොඩන් ගොයෙල / යුතු වෙශ්‍යාකෘති තුළ output එක result එක ගුණවා.

Sampling with replacement

Data set එකට භුෂණ ඇල්බිජ් data randomly replace කළ යුතුව

Random Forest Algorithm

- * Sampling & Replacement Data Set
- * Create Decision trees.

60000 trees construct by prediction error minimization

Given training set of size m

For $b=1$ to B

Use sampling with replacement to create a new training set of size m

Train a decision tree on the new data set

At each node, when choosing a feature to use to split, if n features are available, pick a random subset of $k < n$ features and allow the algorithm to only choose from that subset of features.

$$k = \sqrt{n}$$

XGBoost

Given training set of size m

For $b = 1$ to B

Use sampling with replacement to create a new training set of size m

Train a decision tree on the new dataset.

Learned features and Decision tree stages.

Classification

```
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)
```

When to use Decision Trees

1. Decision Trees

Decision Trees and Tree ensembles

- Works well on tabular (structured) data
- Not recommended for unstructured data (images, audio, text)
- Fast development
- Small decision trees may be human interpretable

Neural Networks

- Works well on all types of data, including tabular (structured) and unstructured data.
- May be slower than a decision tree
- Works with transfer learning
- When building a system of multiple models working together, it might be easier to string together multiple NNs.