

# ① Data engineering techs

↳ S3 → VPC configuration  
Object encryption

↳ Kinesis → Real-time big data

↳ Kinesis Streams

↳ Low latency streaming

↳ Kinesis Analytics:

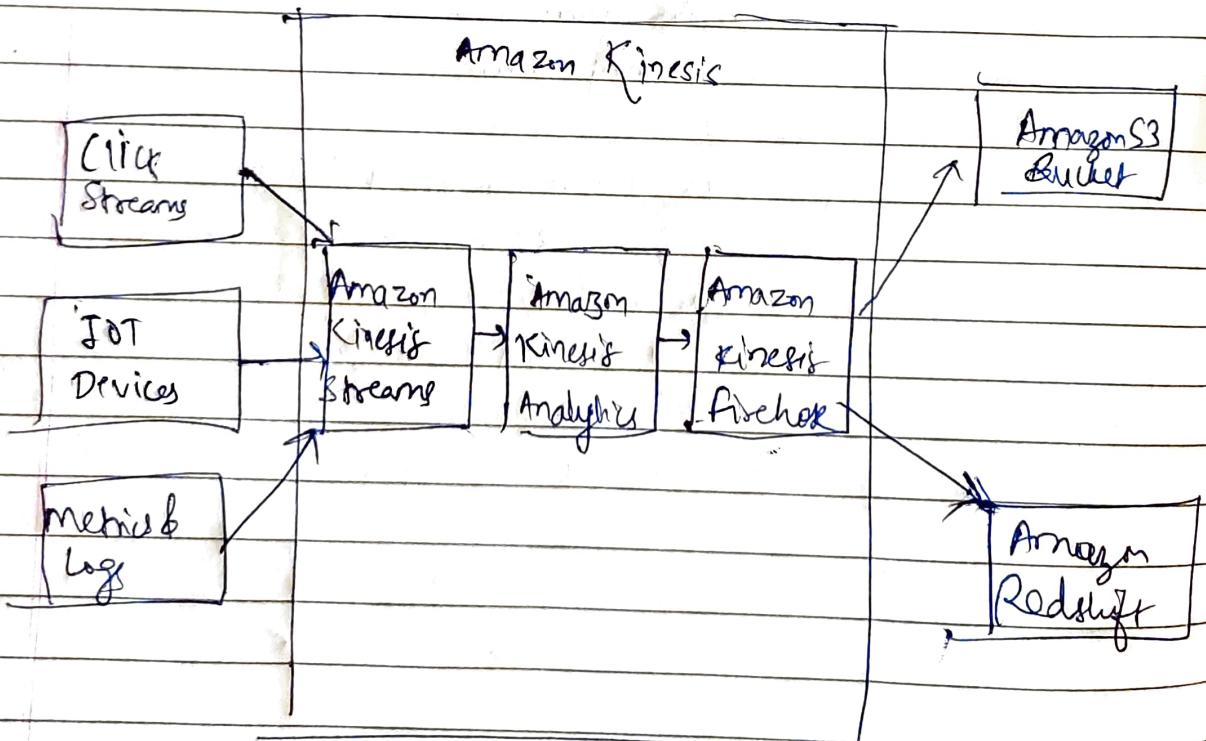
↳ Real-time analytics using SQL

↳ Kinesis Firehose

↳ Load streams to S3, Redshift, etc.

↳ Kinesis Video Streams

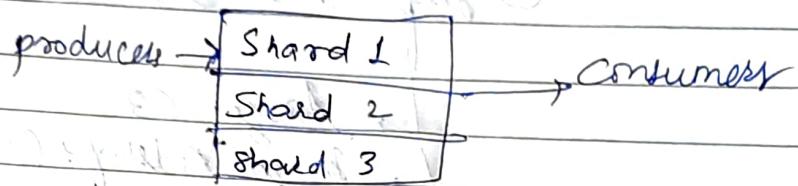
↳ Meant for streaming video in real time



(1)

Kinesis Streams:

↳ Divided in ordered shards/partition.



↳ Retention duration = 24 hrs (default)  
≤ 365 days (max)

- ↳ Ability to reprocess/purge data
- ↳ Multiple applications can consume same stream
- ↳ Once data is inserted in Kinesis, it cannot be deleted.
- ↳ Records can be upto 1MB in size.

### Limits:

Producer:

→ 1MB/s or 1000 msg/sec at write per shard.

↳ "Provisioned Throughput Exception" error if above limit is crossed.

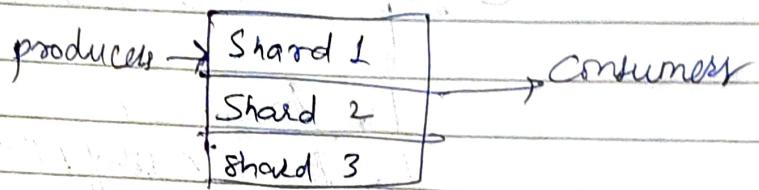
Consumer:

→ 2MB/s of read per shard across all consumers

→ 5 API calls per second per shard across all consumers.

## 1 Kinesis Streams:

↳ Divided in ordered shards/partition.



↳ Retention duration = 24 hours (default)  
≤ 365 days (max)

- ↳ Ability to reproduce/replicate data
- ↳ Multiple applications can consume same stream
- ↳ Once data is inserted in Kinesis, it cannot be deleted.
- ↳ Records can be upto 1MB in size.

### Limits:

#### ↳ Producers:

↳ 1MB/s or 1000 msg/sec at write per shard.

↳ "Provisioned Throughput Exceeded" error if above limit is crossed

#### ↳ Consumers:

↳ 2MB/s of read per shard across all consumers

↳ 5 API call per second per shard across all consumers.

## II Kinesis Data Firehose

↳ AWS Destinations :-

↳ Amazon Kinesis Stream

↳ Amazon S3

↳ Amazon Redshift (Copy through S3)

↳ Amazon Elasticsearch

↳ 3rd party Destinations :-

↳ Data dog

↳ Splunk

↳ MongoDB

↳ New Relic

## III Kinesis Video Stream

• Producer :-

- Security cam, body-wear cam, etc.

- One producer per video stream.

• Video playback capability :-

Consumers

- Build your own (MXNet, TensorFlow, etc.)

- AWS SageMaker

- Amazon Rekognition Video

Keep data for 1hr to 10 years

## \* Glue Data Catalog :-

- Metadata Repository for all your tables :-
- Automated Schema Inference
- Schemas are versioned.
- Integrates with Athena or Redshift Spectrum

↳ Important :-

Glue crawler automatically detects the schema of tables (jsons / csv) in an s3 bucket (specified to it)

It is also able to detect the partitioning scheme of a certain table file.

For eg: If a file is partitioned using date like

s3://bucket-name/file-name/2017/01/01/file-001

## \* Glue ETL :-

Transform data, clean & enrich data  
(before doing analysis)

- → Fully managed, cost effect, pay as you use
- → Jobs are run on a serverless spark platform

## \* Glue Scheduler :- To schedule the job

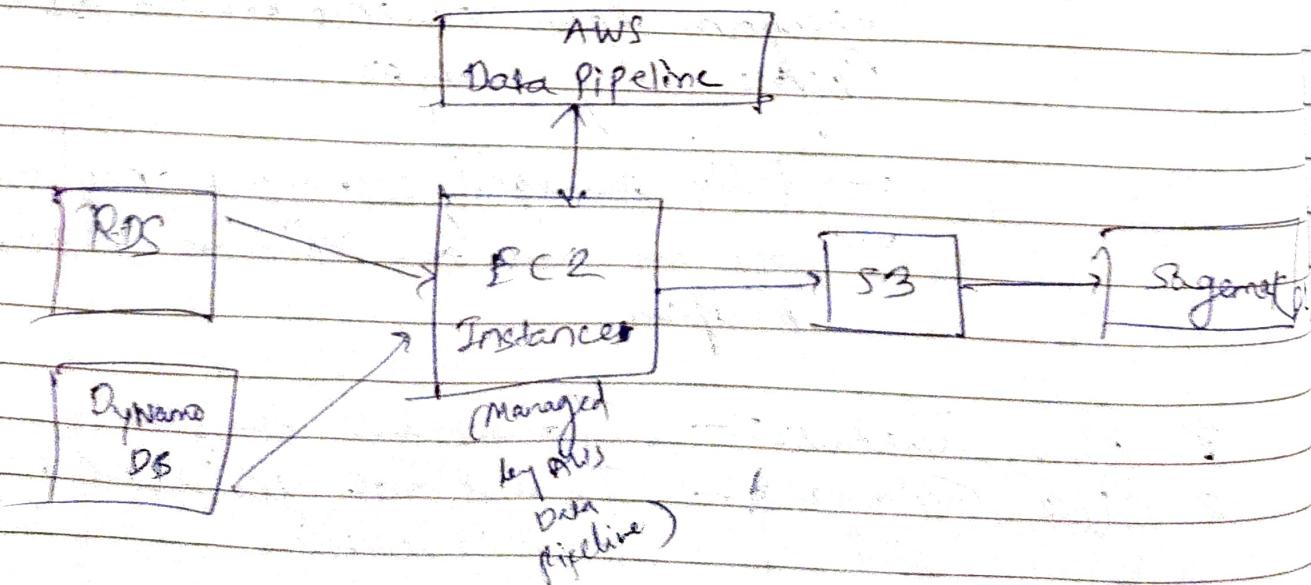
## \* Glue Trigger :- To automate job runs based on "events"

- # Glue ETL - Transformations :-
- Bundled transformations :-
  - DropFields, DropNullFields
  - Filter
  - Join - to enrich data
  - Map - add, delete fields, perform external lookups.
  - ML transformations :-
  - FindMatches ML :-

Identify duplicate or matching records in your dataset.  
 (even when records don't have a common unique identifier & no fields match exactly).

- Format :- CSV, JSON, Avro, Parquet, ORC, XML

## \* AWS Data Pipeline :-



## AWS Batch :-

- Run batch jobs as Docker instances
- Dynamic provisioning

## # Quick Sight

Visualization toolkit for data coming from all kinds of sources —

- On-premises
- AWS S3, Redshift,
- EC2 instances etc.

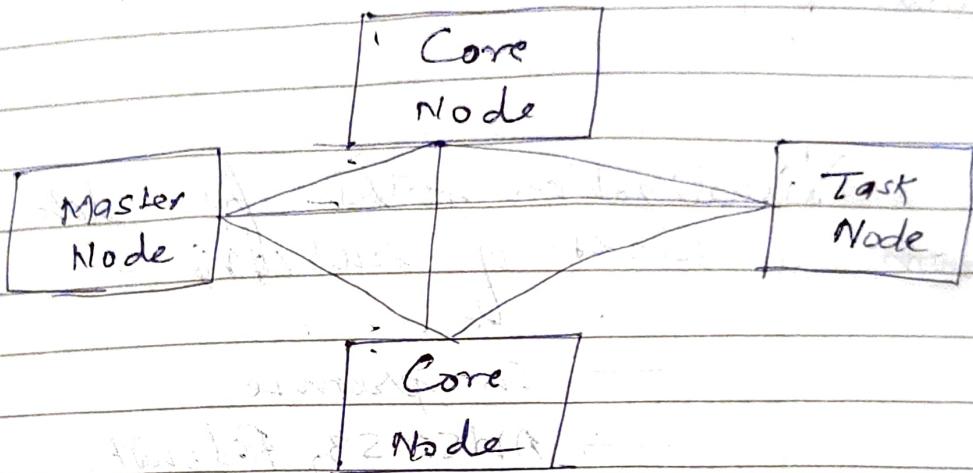
Different types of viz available —

- Bar charts
- Histograms
- Pie charts
- Scatter plots
- Heat maps
- Line plots
- Scatter plots

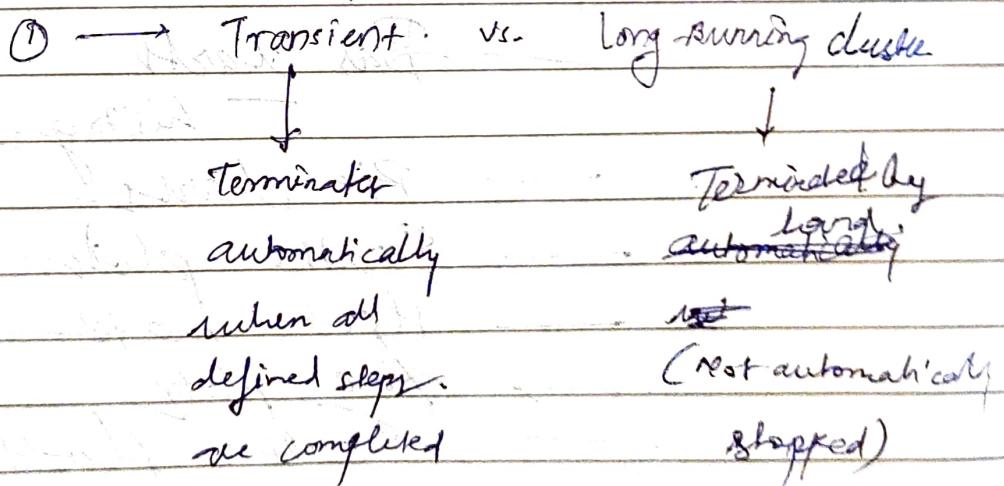
## # EMR

### Elastic Map Reduce

#### EMR Cluster :-



#### EMR Usage :-



#### EMR storage :-

- Temporary storage ←
- HDFS → By default size of a block in HDFS - 128 MB
  - EMRFS — Access S3 as if it were HDFS

## # Feature Engineering :-

→ Applying your knowledge of the data to the model you're using to create better features to train your model with.

→ Which features should I use?

→ Do I need to transform these features?

→ How do I handle missing data?

→ Should I create new features from existing ones?

## # Handling (Imputing) missing data.

↳ Replacing by mean

↳ Replacing by mode

↳ Replacing by median

→ Using MICE → Multiple Imputation by Chained Equations

↳ Replacing values using ML

↳ By K-Means for ~~cat~~  
missing col (Categorical)

↳ By Random-cut forest for  
missing column (numerical)

## ↓ # Dealing with unbalanced data

↳ # Replicate minority class label entries

↳ a) Delete majority class label entries  
some

↳ (This might lead to loss  
of data).

↳ m.

# Binning, Transforming, Encoding,  
Scaling & Shuffling.

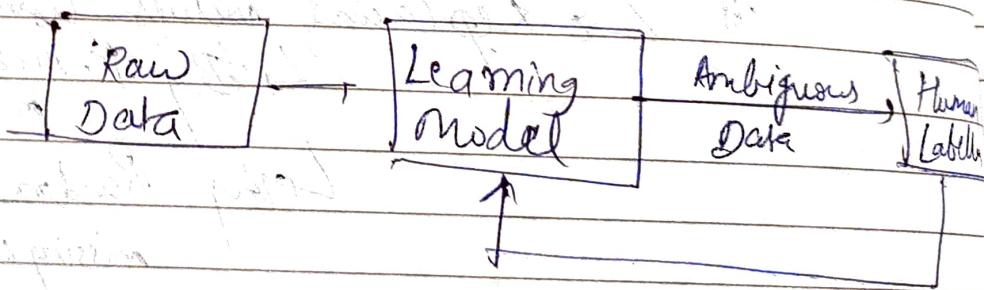
## A. WAYS TO GENERATE TRAINING LABELS

# Sagemaker Ground Truth

↳ AWS Service to use humans  
to label the missing data.

# Ground truth creates its own model  
labelled by people

If this model learns only images  
it is not sure about the sent to  
human labellers.



# Other ways to generate training labels

↳ AWS Rekognition

↳ AWS Comprehend

↳ Any pre-trained model or unsupervised  
technique that may be helpful

## # TF-IDF

- ↳ Term Frequency - Inverse Document Frequency
- ↳ Tools for data search - figures out what terms are most relevant for a doc.
- ↳ TF measures how often a word occurs in a document.
- ↳ DF measures how often a word occurs in a whole set of documents.

. Assumptions of TF-IDF :-

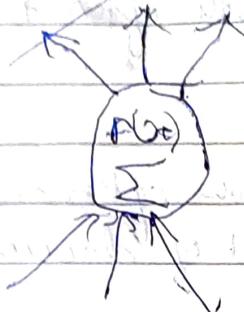
- ↳ A document is just a bag of words
- # Unigrams — Individual words
- # Bi-grams — Set of two words taken at a time
- # n-grams — Set of  $n$  words taken at a time

## # Preparing data for TF-IDF

## II DEEP LEARNING

activation function

→ Defines the output of a node /  
neuron given its input signals



### I Linear Activation Function

- 1. Doesn't really do anything
- 2. Doesn't allow back propagation

### II Binary Step Function

↳ It's on or off.

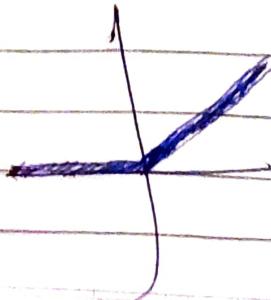
↳ Can handle multiple

↳ Vertical slopes don't work well  
with calc calculus

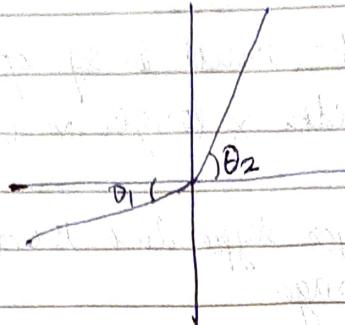
### III Sigmoid / Logistic / TanH

### IV Rectifier Linear Unit (ReLU)

↳ easily  
calculable



## Leaky ReLU



$$\boxed{\theta_1 > \theta_2}$$

## Parametric ReLU (PReLU)

↳ the  $\theta_i$  is decided by RNN

↳ Complicated & computationally expensive

## # Softmax Activation Function

↳ Used on the final output layer of a multiple classification problem.

↳ Basically converts off to probabilities of each classification.

## # CHOOSING AN ACTIVATION FUNCTIONS

① For multiple classifications, use softmax on the output layer.

② RNN's do well with tanh

③ For every thing else -

    ① → start with ReLU

    ② If you need better try Leaky ReLU

    ③ Last resort = PReLU, Maxout

④ SWISH for really deep networks

## # CNNs with Keras / TensorFlow

① source data must be of appropriate dimension  
ie width x length x color channels

② Conv2D layer type does the actual convolution  
on a 2D image.

Conv1D & Conv2D layers can be used - doesn't  
have to be image data

(V1.3.1) (V1.3.1 Software)

③ Max Pooling 2D layers can be used to reduce  
2D layers down by taking the maximum  
value in a given block.

④ Flatten layers will convert 2D layers  
to a 1D layer for passing into a flat  
hidden layer of neurons

⑤ Typical stages

Conv2D → Max Pooling 2D → Dropout →  
Flatten → Dense → Dropout → Softmax.

# CNNs are very resource-intensive (CPU, GPU, RAM, ...)

# Lots of hyperparameters:  
→ Kernel sizes.

→ Many layers with diff no. of units

→ Amount of pooling

→ No. of layers

→ choice of optimizer.

# Getting the training data is the hardest part

# Storing & Accessing it.

CNNs

→

problems

#

## # Specialized CNN Architectures

- Define specific arrangements for layers, padding & other hyperparameters

### I LeNet-5

- Good for handwriting recognition.

### II AlexNet

- Image classification, deeper than LeNet

### III GoogleLeNet

- Even deeper, but with better performance
- Introduces Inception module

~~ResNet~~

### IV ResNet (Residual Network)

- Even deeper
- Maintains performance via skip connections

## # Recurrent Neural Networks

What are they for

→ Time Series Data

→ Data that consists of sequences of arbitrary length.



## # RNN Topologies

→ Sequences to sequence

— predict stock prices based on series of historical data

→ Sequence to Vector

— words in sentence to sentiment

→ Vector to sequence

— create captions from image

→ Encoder → Decoder

— Sequence → Vector → sequence

— machine translation

## II. GPU for Instance types on AWS for Deep Learning

- P3: 8 Tesla V100 GPUs
- P2: 16 K80 GPUs
- P3: M60 GPUs (All NVIDIA chips)

## III. Tuning Neural Networks:

### Learning Rate

- NNs are trained by gradient descent.

### Batch Size

IMPORTANT  
① Large batch sizes tend to get stuck inside a "local minima" as random instead of a correct solution

Counter-intuitive points:

- ② Smaller batch sizes can make it easier to get out of the "local minima" more easily
- ③ Batch sizes that are too large can end up getting stuck in the wrong solution
- ④ Random shuffling at each epoch can make this look very inconsistent results from run to run.

### Important Recap points:

- (good) ① Small batch sizes tend not to get stuck in local minima
- (bad) ② Large batch sizes can converge on the wrong solution at random
- (bad) ③ Large Learning rates can overshoot the correct solution.
- (bad) ④ Small learning rates increase training times



Drop out layers force the network to spread out its learning throughout the network. It can prevent overfitting from learning concentrating on one spot.  
 ↳ Early stopping also does the same.

## # Regularization :-

- Prevents overfitting
- A regularization term is added as weights are learned.

### # L1 & L2 Regularization:

↳ L1 term is the sum of weights,

$$\text{Final } \lambda \sum_{i=1}^k |w_i|$$

↳ L2 term is the sum of squares of weights

$$\text{Final } \lambda \sum_{i=1}^k w_i^2$$

↳ same idea is applied to loss function

### Differences b/w L1 & L2 Regularization

L1 = Sum of weights

① performs feature selection - entire feature is good or bad

② computationally inefficient

③ sparse output

L2 = squared sum of weights

① All features are considered, just weighted

② computationally efficient

③ dense output

H. In sklearn, Lasso Regressor  $\Rightarrow$  L1 regularized model

In sklearn, Ridge Regressor  $\Rightarrow$  L2 regularized model

# ~~RNN~~ ~~LSTM~~ is a ~~recurrent~~ type of RNN that solves the problem of items losing their weight over time. In NLP, words in a sentence may be significant, regardless of their position.

Why would I want L1?

- Feature selection can reduce dimensionality.
- But if you think all features are important, L2 is probably a better choice.

# The Vanishing Gradient Problem :-

- When the slope of the learning curve approaches zero, things can get stuck.
- We end up working with very small numbers that slow down the training, or even introduce numerical errors.
- This becomes a problem with deeper networks & RNN's as these "vanishing gradients" propagate to deeper layers.
- Opposite of this problem is called "Exploding gradients".

# Fixing this problem:

- Instead of training all layers together, we can train multi-level layers training.
- Break up levels into their own sub-levels trained individually.
- LSTM, RNNET, Better Activation functions

## CONFUSION MATRIX:-

		Actual YES	Actual NO
Predicted	YES	TRUE POSITIVES	FALSE NEGATIVES
	NO	FALSE POSITIVES	TRUE NEGATIVES

True positive and false negative are called -

True negative and false positive are called -

		Actual YES	Actual NO
Predicted	YES	TRUE POSITIVES	FALSE POSITIVES
	NO	FALSE NEGATIVES	TRUE NEGATIVES

True positive and true negative are called -

False positive and false negative are called -

True positive rate and False positive rate

True positive rate =  $\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

False positive rate =  $\frac{\text{False Positives}}{\text{True Negatives} + \text{False Positives}}$

True negative rate and False negative rate

True negative rate =  $\frac{\text{True Negatives}}{\text{True Negatives} + \text{False Positives}}$

False negative rate =  $\frac{\text{False Negatives}}{\text{True Positives} + \text{False Negatives}}$

- # Precision should be focussed when you care about False positives more.
- # Recall should be focussed when you care about False negatives more.

#### Measuring your models :-

(i)

Recall  $\Rightarrow$

Used called as

$$\frac{TP}{TP + FN}$$

$$\frac{FP}{FP + TN}$$

False Positive Rate

True Positive Rate

Sensitivity

Completeness

(ii)

Precision  $\Rightarrow$

$$\frac{TP}{TP + FP}$$

Correct Positive Rate

Percent of relevant results (relevancy)

(iii)

$$\text{Specificity} = \frac{TN}{TN + FP}$$

True Negative Rate

(iv)

F1 Score :-

$$\frac{2 \cdot TP}{2 \cdot TP + FP + FN} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

= Harmonic mean of precision & sensitivity

(v)

RMSE :-

- Root mean squared error:

- Accuracy measurement

- Only cares about right & wrong answers

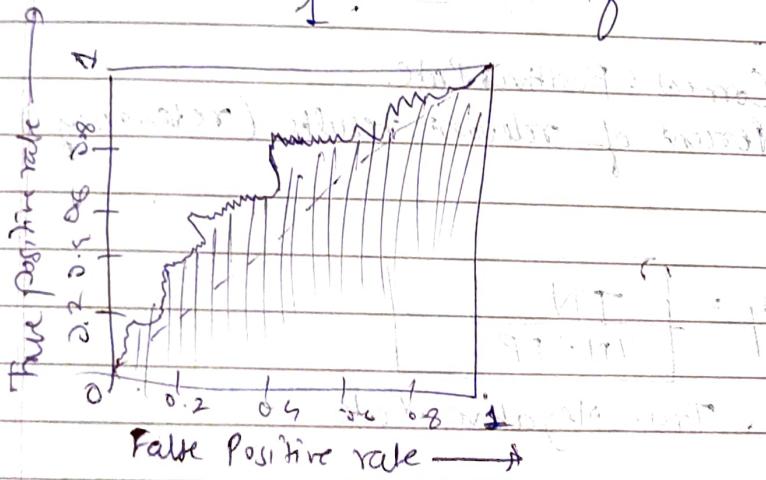
# ROC Curve :- Receiver Operating Characteristic Curve

(1) Plot of TPR vs. FPR at various threshold settings.

(2) Points above the diagonal represent good classification (better than random).

(3) Area under the curve of ROC should be  $> 0.5$ . If it is less than that then it's useless.

Ideal value of ROC curve AUC is 1.



# ENSEMBLE LEARNING:-

Bagging: Generate  $N$  new training sets by random sampling with replacement.

Random sampling with replacement

## # Boosting :-

- ① Observations are weighted
- ② Some will take part in new training sets more often

③

- Training is sequential; each step may a classifier takes into account and learn beyond the previous one's success.

When to choose what?

## Bagging vs. Boosting

① XGBoost has the latest hottest

②

Boosting generally yields better accuracy.

③

Bagging avoids overfitting

④

Bagging is easier to parallelize

⑤

So, depends on your goal

What about which is better  
Bagging or Boosting

## II. Built-in Algorithms in Sagemaker:

### ① Linear Learner Algorithms -

#### ↳ Linear Regression :-

↳ Fit a line to your train data

↳ Predictions based on that line

↳ Can handle both regression & classification predictions

→ for classification a linear threshold is used.

→ Can do binary or multi-class.

I/O Format for this:

→ RecordIO-wrapped protobuf

→ CSV format (1<sup>st</sup> row assumed as header)

→ file or pipe mode both supported.

## # Special things about Linear Learner

### ① Preprocessing :-

- ① Training data must be normalized so that all features are weighted the same.
- ② Linear Learner can do this automatically.
- ③ Input data should be shuffled.

### ② Training:-

- ① Uses SGD (stochastic Gradient Descent)
- ② Choose an optimization algo. (Adam, Adagrad, etc)
- ③ Multiple models are optimized in parallel.
- ④ Tune L1 & L2 regularization.

### ③ Validation:-

- ① Most optimal model is selected.

## # Important Hyperparameters for Linear Learner :-

- i) balanced\_multiclass\_weights
- ii) learning\_rate, mini\_batch\_size
- iii) l1  $\Rightarrow$  regularization
- iv) wd  $\Rightarrow$  weight decay (L2 regularization)

## II XGBoost : Algorithms

Instances for XGBoost  
• uses CPU's only for training  
• memory bound, not compute bound.  
→ So, M5 is a good choice

For XGBoost 1-2,  $\Rightarrow P_3$  (GPU train.)

→ Boosted group of decision trees.

→ Models are serialized / deserialized using pickles.

### Important Hyperparameters:

- i) subsample - prevents overfitting
- ii) eta - step size shrinkage, prevents overfitting
- iii) gamma - min. loss reduction to create a partition

= larger = more conservative

(~~alpha~~, method median (N)) - alpha is same as L1 regularization term.

~~lambda~~ - larger = more conservative

~~lambda~~ - L2 regularization;  
= larger = more conservative

vi) eval\_metric = optimizes AUC, error, rmse, etc.

Ex: if you care more about false positives choose, AUC.

vii) scaled\_pos\_weight =

adjusts balance of true & -veg weights.

scale factor other weight

← helpful for unbalanced classes.

viii) max\_depth =

- max depth of the tree.
- too high = you may overfit.

## Seq2Seq Algorithm.

- \* IP is a sequence of tokens, OP is a seq. of token
- \* Used maybe for:
  - i) M/c translation
  - ii) Text summarization
  - iii) Speech to text
  - iv) Implemented with CNN & RNNs.

IP format for Seq2Seq  $\Rightarrow$  Record IO - Protobut  
 Start with Tokenized text files

### Hyperparameters of Seq2Seq, Important

- i) batch\_size
- ii) optimizer\_type (adam, sgd, rmsprop)
- iii) learning\_rate
- iv) num\_layers\_encoder in AG
- v) num\_layers\_decoder

## DeepTR:

- 1) Forecasting time-series
- 2) Uses RNNs
- 3) Allows you to train the same model over several related time series
- 4) finds series freq & learnability.

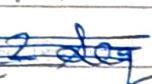
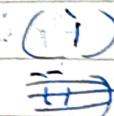
### Hyperparameters:

- i) context\_length
- ii) epochs
- iii) batch\_size
- iv) learning\_rate (v) num\_cells

(VI)

## Blazing Text

+ Word2Vec



(VII)

## ~~Blazing Text~~ → Object2Vec

(VIII)

## Object Detection in Sagemaker

(IX)

## Image Classification in Sagemaker

(X)

## Semantic Segmentation in Sagemaker (Pixel level obj. chg.)

(XI)

## RandomCutForest in Sagemaker - (Anomaly detection) - Unsupervised

(XII)

## Neural Topic Model in Sagemaker - (For organize documents into topics)

(XIII)

## LDA in Sagemaker - (Topic modelling Algo.)

(XIV)

## kNN in Sagemaker -

For simple classification or segregation.



Find the 'K' closest points to a sample point & return the most frequent label.

Find the 'K' closest points to sample point & return the avg. value.

(XV) K-Means  $\rightarrow$  For unsupervised clustering

(XVI) PCA In Sagemaker  $\rightarrow$  For Dimensionality Reduc<sup>n</sup>.

(XVII) Dealing with sparse Data  $\rightarrow$  in factorization machines

Click prediction  
Item recommendations

(XVIII) Reinforcement Learning in Sagemaker:

Q-Learning (Implementation of R.I. Learning)

- i) A set of environment states 's'
- ii) A set of possible actions in those states 'a'
- iii) A value of each state/action  $\delta$   
Start with  $\delta$  values as 0

Good things

$\rightarrow$  after given action  $\Rightarrow Q\uparrow$

Bad things after given action  $\Rightarrow Q\downarrow$

## # Sagemaker autopilot

Automates:

- i) Algorithm selection
- ii) Data preprocessing
- iii) Model tuning
- iv) All infrastructure

- # It does all the trial & error
- # More broadly called as autoML

## # Sagemaker Model Monitor:

L → Get alerts on quality deviations on your deployed models

- L → Detect anomalies & outliers
- L → Detect new features
- L → No code is needed.
- L → Visualize data drift

~~It's also useful for~~

~~→ It's also useful for~~

## # AWS AI/ML SERVICES :-

### 1) Amazon Comprehend :-

- i) NLP + Text Analytics
- ii) Extract Keyphrases, entities, sentiments, syntax, topics, document classifying.
- iii) can train your own data.

### 2) Amazon Translate

### 3) Amazon Transcribe :-

↳ Speech to text.

### 4) Amazon Polly

↳ Text to speech.

### 5) Amazon Rekognition

↳ Computer Vision Tool

↳ Object & Face Detection  
Scene

↳ Facial Analysis :-

(Age, emotion on the face, gender, etc.)

↳ Celebrity Analysis / Recognition

## 6) Amazon Forecast:

↳ Time Series analysis

↳ automatically chooses

optimal model to train  
data from

↳ ARIMA

↳ DeepAR

↳ ETS

↳ NPTS

↳ Prophet

## 7) Amazon Lex:

↳ Natural Language  
Chatbot Engine

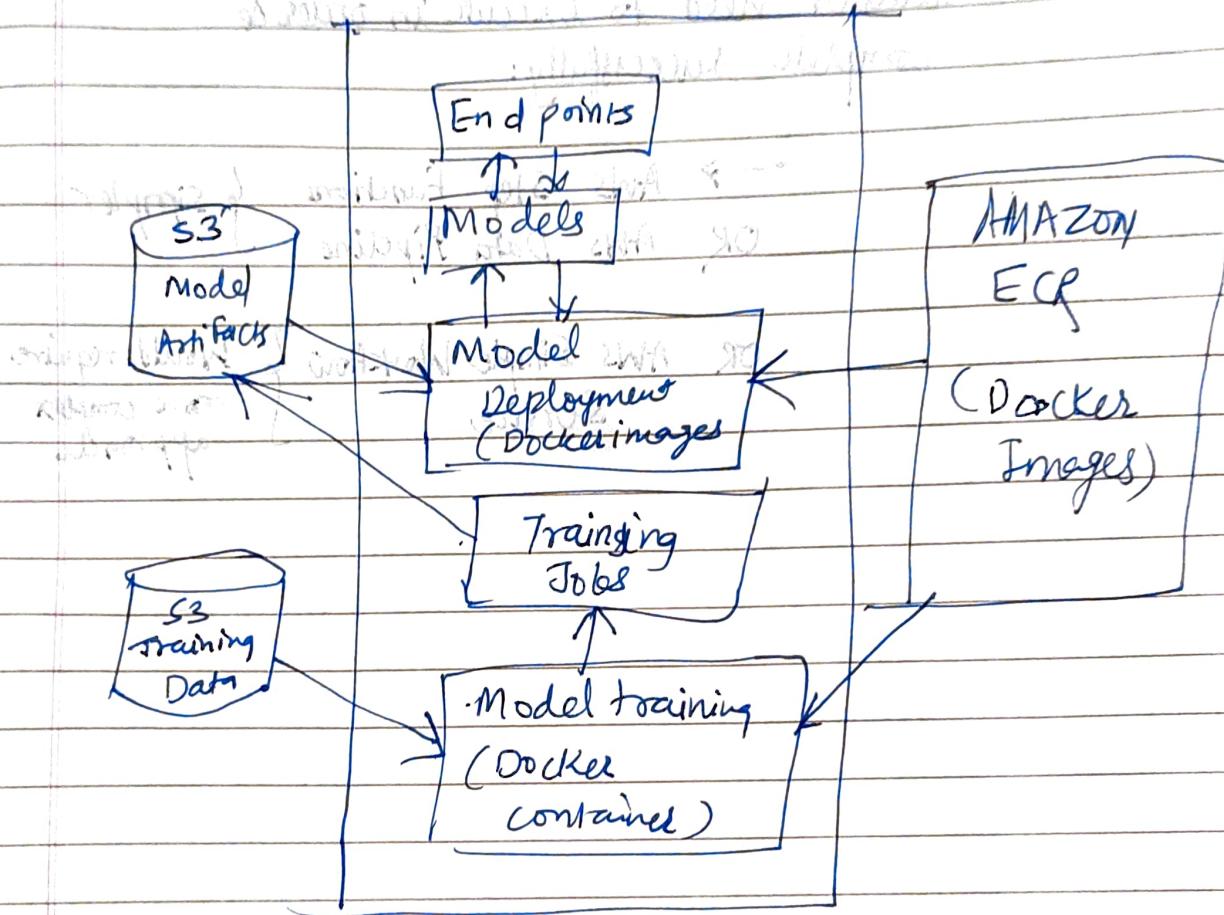
## 8) Amazon Personalize:

↳ Fully managed

Recommendation Engine

## # MLOps: Using Dagerunner with Docker

### Using Dagerunner with Docker



$$2 \cdot \frac{2}{3} \times 0.8 = 2 \cdot \frac{0.72 \times 0.8}{0.72 + 0.8} = \frac{2 \times 0.576}{1.5} = \frac{1.152}{1.5}$$

## # CASES & SERVICES THAT MIGHT BE CHOSEN:

(I)

Multiple long-running ETL jobs which need to execute in order to complete successfully:

→ AWS Step Functions  
OR AWS Data Pipeline

OR AWS Simple Workflow

{ Would require more complex approach }

beginning

step

process [ ]

end

process [ ]

process [ ]

end