

I Data Engineering

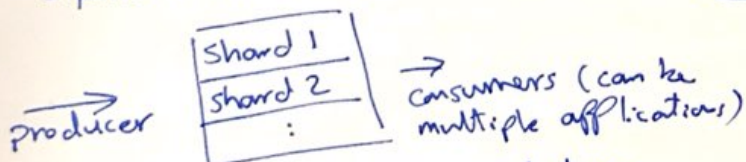
→ S3

- create a "Data Lake"
- You can do data partitioning By date
By product
- data partitioning in S3 can be identified automatically by AWS Glue crawlers
- S3 storage tiers
 - General Standard (freq access)
 - intelligent (infreq access) & One-Zone
 - Glacier → Archiver
- S3 lifecycle rules ⇒ move data between diff tiers
 - General Purpose → infreq Access → Glacier
- S3 Encryption ⇒ SSE-S3
SSE-KMS
- S3 Security ⇒ ~~access~~ user-based ⇒ IAM policies
⇒ resource-based ⇒ Bucket policies
- Use S3 bucket for policy to ⇒ grant public access to the bucket
⇒ force objects to be encrypted at upload
⇒ grant access to another account
- By VPC endpoint gateway → allow traffic to stay within your VPC instead of going through public web

→ Kinesis (managed alternative to Apache Kafka)

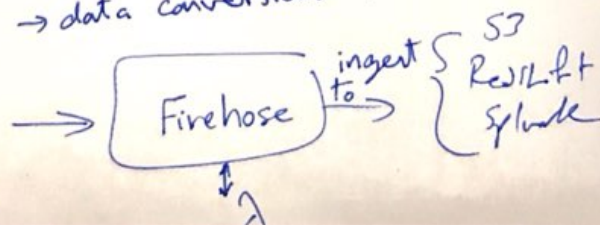
- Kinesis Streams: low latency ingest at scale
- " Analytics: perform real-time analytics on streaming using SQL
- " Firehose: load streams into S3 & Redshift & ElasticSearch & Splunk
- " Video Streams: meant for streaming video in real-time

- Kinesis Streams are divided in ordered shards or partitions



- Kinesis Firehose
 - near real-time
 - auto scaling
 - data conversions from CSV/json → Parquet/ORC

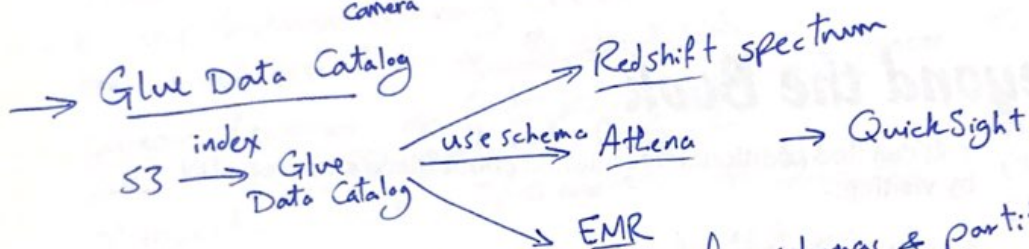
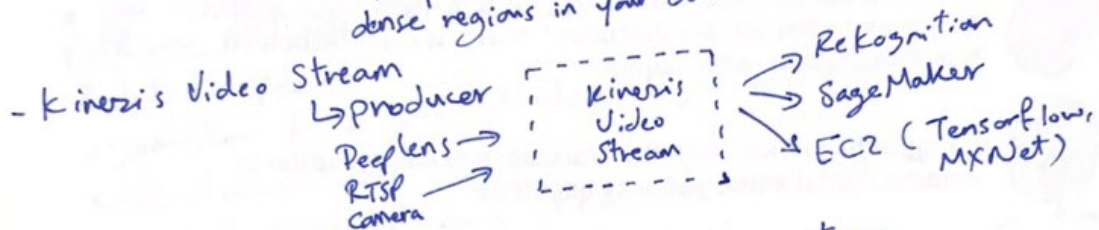
Kinesis Producer Lib (KPL)
" Agent
Kinesis streams
IoT rules or cloudWatch



- Kinesis Analytics
 - use case \rightarrow streaming ETL \Rightarrow select columns, make simple transformations, on streaming data
 - SQL / Flink
 - Serverless

\rightarrow ML on Kinesis Analytics

- * Random cut forest: SQL function used for anomaly detection on numeric columns in a stream
- * Hotspot: locate and return information about relatively dense regions in your data



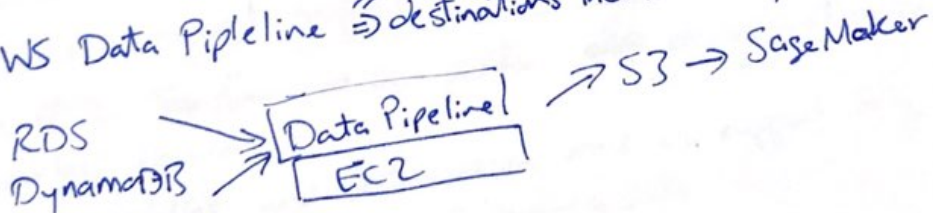
- crawlers go through your data to infer schemas & partitions
- Glue crawler will extract partitions based on how your S3 data is organized
- Glue ETL \Rightarrow Fully managed transformation \Rightarrow Find Matches ML: identify duplicate or matching records in your dataset

\rightarrow Redshift (OLAP)

- \rightarrow use Redshift spectrum to query data directly in S3 (no loading)

\rightarrow AWS Data Pipeline \Rightarrow Orchestration service

\Rightarrow destinations include S3 / RDS / EMR



Exam Review Tips (Data Engineering)

- Content Type: application/x-recordio-protobuf if you want to use data in the protobuf recordIO format for training
- Training Input Mode parameter options while using SageMaker algorithms are File or Pipe
- Kinesis Data Streams PutRecord API uses name of the stream, a partition key and the data blob, whereas Kinesis Data Firehose PutRecord API uses the name of the delivery stream and the data record
- Redshift spectrum
- SageMaker Factorization Machines algorithms supports for inference
- X-recordio-protobuf & json
- Kinesis Producer Library (KPL) provides built-in performance benefits and ease of use on the client side
- Linear learner & factorization machines process the training data in recordIO-protobuf float 32
- IP insights (entity, IPv4 address) & object 2Vec (token, seq) or (token, token) or (seq, seq)
- Using Crawler + Athena is the least effort way to make S3 data available for SQL queries
- AWS Data Pipeline & AWS Glue ETL \Rightarrow from Redshift \rightarrow S3
- Converting the data to recordIO-protobuf file type can significantly improve the training time with marginal increase in cost
- XGBoost only supports text/csv & text/libsvm
- Apache MXNet recordIO is the recommended input format for SageMaker image classification algo.
- AWS Glue does not support Timestream as the source input type
- Glue ETL Transform job can perform ~~data~~ deduplication in a serverless fashion
- Seq2Seq modeling and factorization machines support only the recordIO-protobuf file type for training data

[2] Exploratory Data Analysis (EDA)

- ~~ordinal~~ - ordinal \Rightarrow such as rating
- numerical \Rightarrow such as ~~results~~ age
- categorical \Rightarrow , , gender

- Data Distributions

- \rightarrow normal distribution
- \rightarrow probability mass function
- \rightarrow Poisson distribution
- \rightarrow Binomial
- \rightarrow Bernoulli



- Time Series Analysis

- \rightarrow Trends
- \rightarrow Seasonality
- \rightarrow time series = seasonality + trends + noise

- Amazon Athena

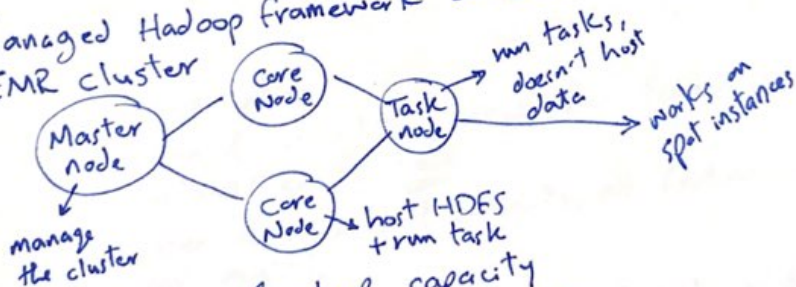
- \rightarrow serverless interactive queries of S3 data (SQL)
- \rightarrow no need to load data
- \rightarrow Athena not for ETL, use Glue instead

- Amazon QuickSight

- \rightarrow Biz analytics & visualizations in the cloud
- \rightarrow Serverless
- \rightarrow data sources: RedShift / RDS / Athena / S3
- \rightarrow data sets are imported to SPICE (Super-fast, Parallel, In-memory, Calculation Engine)
- \rightarrow ML insights: anomaly detection, forecasting, auto-narratives
- \rightarrow QuickSight not for ETL, use Glue instead

- EMR

- \rightarrow managed Hadoop framework on EC2 instances
- \rightarrow EMR cluster



- \rightarrow spot instances reserved for temp capacity
- \rightarrow long-running clusters

Note \rightarrow SVD and PCA \Rightarrow to reduce dimensions

\rightarrow Kinesis \rightarrow Spark

- Feature Engineering

- \rightarrow Median ~~is~~ may be better choice than mean when outliers are present
- \rightarrow imputing missing data: ML
 - \rightarrow KNN (numerical data)
 - \rightarrow Hamming distance (categorical data) \downarrow better
 - \rightarrow Regression \rightarrow MICE DL

→ Handling unbalanced data

→ oversampling: duplicate samples from the minority classes

→ undersampling: instead of creating more positive samples, remove negative ones



→ SMOTE: Synthetic Minority Over-Sampling Technique
↳ artificially generate new samples of the minority using KNN

→ adjust thresholds when making predictions about a classification, if you have too many false positive, one way to fix that is to simply increase that threshold

→ Handling Outliers

→ Variance measures how "spread-out" the data is

→ AWS's Random Cut Forest ⇒ made for outlier detection

- Bining: bucket observations together based on ranges of values

⇒ Note

small batch size tend to not get stuck in local minima

- Regularization to prevent overfitting
in DL: - Dropout
- early stopping

- The vanishing Gradient problem
↳ How to fix it:

① Multi-level ~~tree~~ hierarchy

② LSTM

③ Residual N/w (ResNet)

④ Better choice for activation →

ReLU

- L1 & L2 Regularization

L1 → sum of weights, used for feature selection, entire features go to zero

L2 → sum of squared of weights, all features remain considered, just weighted

- Bagging: each resampled model can be trained in parallel

Boosting: Training is sequential, each classifier takes into account the previous one's success

Exam Review Tips (EDA)

- The best way to engineer the cyclic features is to represent them as (x, y) coordinates on a circle using sin & cos functions
- Use exclude pattern `**.*.metadata` in the crawler definition to ignore the metadata. AWS Glue crawler supports exclude pattern
- Removing stop words & tokenizing each word & lowering-case each word are the recommended pre-processing for NLP
- Elastic Search, EMR and EC2 are not "Serverless"
- Rule of thumb: drop a feature that will not help a model to learn. Any feature that has low variance or a lot of missing values or has a low correlation to the target label ought to be dropped

$$IQR = Q3 - Q1$$

$$\begin{aligned} \text{minimum outlier cutoff} &= Q1 - 1.5 * IQR \\ \text{maximum} &= Q3 + 1.5 * IQR \end{aligned}$$

- Histogram is best suited to show distributions of variables ~~with~~ while bar charts are used to compare variables.
Histogram plot quantitative data with ranges of the data grouped into bins or intervals while bar chart plot categorized data

$$\begin{aligned} tf('fox', d1) &= \frac{2}{12} \\ tf('fox', d2) &= \frac{3}{12} \end{aligned}$$

$$c = 2$$

$$idf('fox', D) = \log \frac{2}{2} = 0$$

↳ 'fox' term appears in both documents

$$\begin{aligned} tf-idf('fox', d1, D) &= tf('fox', d1) * idf('fox', D) \\ &= \frac{2}{12} * 0 = 0 \end{aligned}$$

- Box plot ~~& scatter plot~~ ^{histogram} → to spot outliers
- LDA & NTM (Neural Topic Modeling) ⇒ for topic modeling
- PCA → reduce dimensionality
- k-Means → clustering

[3] Modeling

- SageMaker is built to handle the entire ML workflow



- data must come from S3, ideal format varies with algorithm - often it's RecordIO - Protobuf

- SageMaker deploy in 2 ways.

① persistent endpoint for making individual predictions on demand

② Batch Transform to get predictions for an entire dataset

** SageMaker's Built-in Algorithms

① Linear Learner $\begin{cases} \text{Regression} \\ \text{classification} \end{cases}$

- training i/p \Rightarrow RecordIO - wrapped protobuf, float32 only \Rightarrow CSV

- important hyperparameter: - Balance multiclass weights
- learning rate, mini-batch-size
- L1 regularization
- L2 " (Wd \Rightarrow Weigh decay)

② XGBoost $\begin{cases} \text{Reg.} \\ \text{class.} \end{cases}$

- training i/p \Rightarrow CSV \Rightarrow libsvm

\Rightarrow AWS extended it to accept recordIO-protobuf & Parquet

- important hyperparam: - subsample \rightarrow prevents overfitting
- Eta \rightarrow step size shrinkage, prevents overfitting
- Gamma \rightarrow min. loss reduction
- Alpha \rightarrow L1
- lambda \rightarrow L2

③ Seq2Seq $\begin{cases} \text{ML Translation} \\ \text{Speech to text} \end{cases}$

- training i/p \Rightarrow RecordIO / Protobuf (tokens must be integer)
 \Rightarrow must provide training data, validation data, and vocab files

- important hyper: - batch size
- optimize type (adam, sgd, rmsprop)
- learning rate
- # layer encoder
- # " decoder
- can optimize on $\begin{cases} \text{accuracy} \\ \text{BLEU score} \\ \text{perplexity (cross-entropy)} \end{cases}$

④ DeepAR → forecasting 1D time series data

- training i/p ⇒ json ~~(with timestamps)~~
⇒ files can be in gzip or Parquet
- important hyper:
 - context length
 - epoch
 - mini batch size
 - learning rate
 - num cell

⑤ Blazing Text → Text classification (web searches, info retrieval, ranking)^{Supervised}

- Word2Vec (unsupervised learning), word embeddings
- only works on individual words, not sentences or documents
- training i/p ⇒ for text classification (supervised mode)
 - ↳ one sentence / line
 - ↳ --label- (first "word" in sentence)
- ⇒ for Word2Vec (unsupervised)
 - ↳ skip-gram
 - ↳ CBow
 - ↳ Batch-skip-gram
 - ↳ text file with one training sentence per line
- important hyper ⇒ Word2Vec:
 - mode (batch-skipgram, skipgram, CBow)
 - learning_rate
 - window size
 - vector dim
 - negative samples
- ⇒ Text classification:
 - epoch
 - learning rate
 - word ngrams
 - vector dim

⑥ Object2Vec

- Note → currently support learning embeddings of pairs of tokens, pairs of sequences, and pairs of token and sequence
- use cases → collaborative recommendation system
- multi-label document classification
- sentence ~~embed~~ embedding
- training i/p ⇒ JSON lines

⑦ Object Detection

- training i/p ⇒ RecordIO or image format (jpg or png)
- uses a CNN with SSD (Single Shot multiple Detector) algorithm
- important hyper:
 - mini batch size
 - learning rate
 - optimizer (adam, SGD, RMSPROP, ADADelta)

⑧ Image Classification

- input training \Rightarrow RecordIO or image files (jpg or png)

⑨ Semantic Segmentation \rightarrow ~~document~~ Pixel-level

- training i/p \Rightarrow image files JPG, PNC for annotation

- choice of 3 algorithms \rightarrow FCN
 \rightarrow PSP
 \rightarrow DeeplabV3

- choice of backbones \rightarrow ResNet50
 \rightarrow ResNet101
 \rightarrow B.12 trained on ImageNet

⑩ Random Cut Forest \rightarrow anomaly detection (unsupervised)

- training i/p \Rightarrow RecordIO - Protobuf
 \Rightarrow CSV

- hyper: - num trees
- num samples per tree

⑪ Neural Topic Modeling \rightarrow organize documents into topics (unsupervised)

- training i/p \Rightarrow RecordIO - protobuf
 \Rightarrow CSV

- hyperparameter: - mini batch size & learning rate
- num topic

⑫ LDA \rightarrow topic modeling algorithm (unsupervised)

\rightarrow can be used also in cluster customers based on purchases & harmonic analysis in music

- training i/p \Rightarrow RecordIO - protobuf
 \Rightarrow CSV

- hyperparameter: - num topics
- alpha0

⑬ KNN $\begin{cases} \text{reg.} \\ \text{class.} \end{cases}$

input training \Rightarrow ~~Record~~ RecordIO - protobuf
 \Rightarrow CSV

hyperparameter: - k
- sample size

⑭ K-Means \rightarrow unsupervised clustering

- i/p training \Rightarrow RecordIO - protobuf
 \Rightarrow CSV

⑮ PCA \rightarrow dimensionality reduction (unsupervised)

\rightarrow the reduced dimensions are called components

- i/p training \Rightarrow RecordIO - Protobuf
 \Rightarrow CSV

→ How PCA is used:

- Covariance matrix is created, then Singular Value Decomposition (SVD)
- Two modes → Regular = sparse data
→ Randomized = large "

→ hyperparameters: - Algorithm mode
- Subtract mean

⑥ Factorization Machines

- good choice for tasks dealing with high dimensional sparse datasets, such as click prediction and item recommendation

→ training i/p ⇒ recordId - protobuf Float32

⑦ IP insights

→ unsupervised learning for IP address usage patterns
→ identifies suspicious behavior from IP address

i/p training ⇒ (entity, IPv4 address) CSV only
↳ hashed and embedded

- Hyperparameter Tunning

↳ Automatic Model Tunning ⇒ Don't run too many training jobs concurrently, because it learns as it goes

④ ML Ops

→ structure of a training container

/opt/ml

└─ input
 └─ config
 └─ data
└─ model
└─ code
└─ output

→ Structure of your Docker image

.WORKDIR
nginx.conf
predictor.py
server
train/
wsgi.py

→ structure of a Deployment Container

/opt/ml

└─ model

→ Protecting Data in-transit in SageMaker

↳ TLS/SSL

↳ inter-node training communication may be optionally encrypted, may increase time, and cost

→ SageMaker + VPC ⇒ Setup S3
VPC endpoint

→ Notebooks are internet-based by default

↳ if disabled, your VPC needs an interface endpoint
to allow outbound connections, for training and hosting to
work

→ PrivateLink
→ NAT GW

→ Spot instances can be interrupted, use checkpoints to S3 so training can
resume

→ Elastic Inference accelerators may be added alongside a CPU instance

→ inference pipelines:

- linear seq of 2-5 containers
- any combination of pre-trained built-in algorithms or
your own algorithms in docker containers
- Combine (pre-processing → predictions → post-processing)
- Spark ML serialized into MLeap format

01003519014
01050819014

Exam Reviews (Modeling)

- AUC/ROC metric doesn't require you to set a classification threshold. It's also useful when there's high class ~~imb~~ imbalance
- The residuals plot would indicate any trend of underestimation or overestimation. Both Mean Absolute Error and RMSE would only ~~give~~ give the error magnitude. AUC is a metric used for classification models.
- Object2Vec can be used to find semantically similar objects such as ticket.
BlazingText Word2Vec can only find semantically similar words
- PCA & Random Cut Forest are parallelizable ***
- Transfer learning \Rightarrow 3rd party pre-trained model $\xrightarrow{\text{then}}$ remove the original classifier
 $\xleftarrow{\text{then}}$ add the new classifier
- Model is ~~considered~~ the mandatory hyperparameter for Word2Vec & Text classification
- Object2Vec, ~~semantic~~ Semantic segmentation and Seq2Seq modeling are NOT parallelizable
- AWS Glue job can't write output in recordIO-protobuf format
- When a model underfits, it exhibits low accuracy on both training & testing data
- Amazon Personalize is ML service to create individualized recommendations for customers using their app
- \times predict() method for SageMaker Python SDK
- \times create_endpoint_config() & create_endpoint() used in Boto3 library while deploying the model
- \times invoke_endpoint() to call inferences
- \times XGBoost & Random Cut Forest \rightarrow only use CPU
- \times Seq2Seq \rightarrow i/p is a seq of tokens (text, ~~audio~~ ^{audio}), and the o/p is another seq of tokens.
for ex; machine translation, text summarization, speech-to-text
 \rightarrow Seq2Seq, it uses RNN and CNN models with attention as encoder-decoder architectures
- create_transform() \rightarrow to create Batch Transform jobs
- describe_transform() \rightarrow to get a status about the progress of the job

- Think of L1 as reducing the number of features in the model altogether.
- L2 "regulates" the feature weights instead of just dropping them

Exam Review (ML Ops)

- How early stoppy works in hyperparameter tuning job
 - after each epoch of training, get the value of the object metric
 - compute the running avg of the objective metric for all previous training jobs up to the same epoch, and then compute the median of all of the running averages
 - if the value of the objective metric for the current training job is worse than the median value of running avg of the objective metric for previous training jobs. SageMaker stops the current job
- object Detection, Image classification, Semantic segmentation
- Using the training data in protobuff recordIO along with pipe mode can significantly improve the training job performance
- Training image & inference Image Registry paths are region based
- VPC has an interface endpoint (PrivateLink) + NAT GW for outband
- num-round & num-class are mandatory hyper for XGBoost
- SageMaker models are store in model.tar.gz format in S3 bucket
- Within an inference pipeline model, Amazon SageMaker handles invocations as a sequence of HTTP Req.
- N/W isolation is not supported by the following managed Amazon SageMaker containers as they require access to S3 → Pytorch, Reinforcement Learning, Sklearn, Chainer
- People Pathing by AWS Rekognition
- There's no 15 min gap for logs to appear on CloudWatch. The default IAM has the required permissions to writes logs into cloudWatch
- SageMaker supports t, m, c, p and r family instances
- predict() → real-time inference
- deploy(), fit() → for models hosted on SageMaker Hosting Services
- SageMaker ~~flips~~ inference pipelines can be considered for real-time prediction & process batch Transform
- SageMaker Neo supports image classification
- PCA (randomized mode) for larger datasets