

Axes :- AWS - ML

Date : 02/01/2022.

→ AWS ML X OVERVIEW:

1 → How to Select ML technique to solve given business problems

→ Which AWS Services

→ How to design, Implement, Scale, Secure.

→ cost-effective solutions.

→ 65 questions Scenario → 3 hours

75% pass exam → \$ 300 US dollars.

1. Data Engineering — 20.1.

2. EDA — 24.1.

3. Modeling — 36.1.

4. ML implementation & operation — 100.1.



Date _____

Page _____

DEI.

1. Create Data Repository
2. Data ingestion
3. Data transformation

EDA:-

Prepare data for Modeling

1. Perform feature engineering
2. Analyze & Visualize data from ML

→ Apply Basic ML Security practices to ML Solutions

→ Deploy ML Solutions.

Important Concepts:-

1. Data storage 2. EDA — Ans clue Data pipeline
3. Model 4. MLops. Ans Data migrations

Modeling:

" frame business problems

1. Select appropriate ML model
2. Train
3. Perform hyper parameter
4. Evaluate ML

Kinetics or views. } Amazon Drama features.

Kinetics } Quick right ML solution

EMR

→ Hadoop

Apache spark

features

ML overview

→ Features extracted manually in ML

and in case of DL its done automatically

AE

ML, robotics
and Computer Vision

→ Supervised → Classification
→ Unsupervised → Clustering
→ Reinforcement

ML

→ AWS

→ Cloud platforms
→ 165 fully featured services (40 of them
are not offered anywhere)

Compute
Storage
Database

Data + Model + Compute

Compute → ECR, ECR, ECR

Data Lingo:-

Data
Unlabeled Data ← Labeled Data

ImageNet → 21,841 classes.

Air Ground Paths :-

Data Types

Quantitative
Numerical
Categorical
Qualitative
Ordinal
Mix between
Numerical
and Categorical

↳ discrete
↳ continuous
or continuous

You can multiply them together for example.

Database vs datalake vs data warehouse

structured
Data with defined Schema
to form a Database

to put together

Data lake:-

- Data lake is a centralized repository for Structured and Unstructured data storage.
- Data lakes could be used to store raw data as there is no need to perform any ETL processing. Could be done on export to schema it is defined.

(A)

AWS Key Storage a.a.m.gg1.

storage service allows enterprises / individual extremely durable

1. Amazon S3 → to store and protect any amount of data
2. Aurora RDS → Relational Database, fully managed.
3. Redshift → data warehousing service.
4. DynamoDB.

↳ Redshift uses a unique data warehouse architecture that relies on Massively Parallel Processing (MPP)

↳ Redshift uses Machine learning to optimise performance.

DOMAIN #1

Code:- V1.35%

DATA ENGINEERING :-

1. Amazon S3: Simple Storage Service.
2. Data availability
3. Security
4. Performance.

Data is stored in Buckets. Global Unique Name.
 You can store infinite amount of data
 in bucket in which each object contains upto 5TB.

① Object tags : Confidential / PII / project.

Here an example of an object tag
 1. Project → ML
 Classification → Confidential
 PII → True.

AWS S3 Data Lake:-

x gigabytes → Petabytes

AWS S3 Common formats: CSV, Parquet, ORC, Avro, protobuf, and JSON.

→ Amazon S3 hosts that data the machine learning such as sage maker will use for model training and testing.

training Data → Sage maker.

trained Mode

Bucket Integration with AWS Lambda

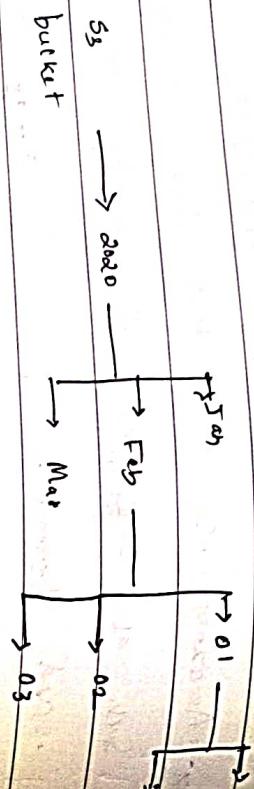
Compute and Data Processing decoupled

→ Centralized Data Architecture.
 → Integration with Cloudless / Serverless AWS services.
 Standardized API's → Hadoop, Spark.

AWS S3 Data Partitioning:-

Date Partitioning is critical when querying data because it could dramatically reduce the cost required for scanning.

General → infrequent → Access → Glacier.



Life Cycles :-

transitions: Policy that governs when an object transitions from one storage class to another.

→ Glacier: Data need to be retained from Compliance purpose.

AWS Storage Classes.

1. S3 Standard → works with general data purpose and frequently used / needed
2. S3 Intelligent : works well with data by varying access.
3. S3 Standard - infrequent access.
4. S3 One Zone - infrequent access
5. S3 Glacier.



Amazon S3 encryption

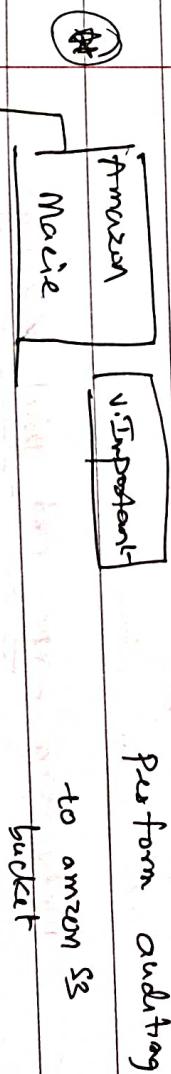
Amazon S3 ensure data protection and durability by

1. Adding redundancy storage over multiple devices across many facilities / regions
2. S3 ensure that any corrupted files are detected and repaired.
3. Using Versioning which allows users to retrieve several versions of the same object.

For ML → below encryption used.

1. SSE-S3 and AWS KMS.

Amazon S3 Security :-



Security of the cloud → AWS

Security in the cloud → user

find sensitive data and alert users

for anomalies.



→ "block public access" → by default, you can't go through the public internet and will stay inside the VPC for maximum security.

1. Cloud watch alarms: certain threshold is exceeded for multiple number of cycles

→ VPC end point will route request to Amazon S3 and back to the VPC.

2. Cloud trail logs: track activity made by users, roles, IP address, record of past requests, timing of the request.

Tagging:-

you can use tagging bucket policies to ensure security.

"Sensitive = TRUE."

3. S3 access logs:-

Networking - VPC and Point gateway:

Amazon Virtual private cloud allows user to create AWS resource inside a virtual network.

Data Warehouse

Data Base

- Used for data analytics
- Used for data translation
- Sources of data?
 - Data gathered from multiple sources.
 - is from one source

- Data writing frequency?
- Bulk write per Many write operations
- Storage optimized for high work optimizability
- frequency Speed query in single-row
- in column format

Elastic search :-

↳ ingestion, enrichment, storage, analysis

and visualization

→ Tools are known as ELK stack

→ fully Managed service that allows that allows for Elastic search deployment easily and securely

→ Click Stream Analytics, data rendering

* AWS - Glue * V. Important Code: 1.22.2.

AWS Glue → ETL

Data generated by AWS Glue could be fed into analytics tools such as Amazon QuickSight.

1. Setup AWS Glue and direct it to your stored data on AWS.

2. AWS Glue extracts metadata from the source and transform it to target. users can then test the code

ETL.

→ AWS Glue → integrated S3, quicksight,

Redshift.

→ AWS Glue → Crawls data storage, extract data schemas and transformations

→ AWS glue generates data transformations codes - on its own.

Step 2: - After identifying the source data. AWS

Glue will crawl the data source (s) and generate the catalogue. AWS Glue works with all kinds of data.

Step 3: - In python / scala prog language, AWS

Glue generates ETL code to extract the data from the source and transform it to

target. users can then test the code in IDE / note book

Step 3: - AWS Glue facilitates running scheduling multiple ETL jobs. AWS Glue efficiently and automatically scales based on demand.



Important :-

AWS DATA PIPELINE:-

- AWS Data Pipeline is an orchestration service that allows AWS user to transfer data easily and securely between various AWS compute and storage services.



Amazon Quick sight

It is used scan the stored data and

infer the schema and partitions.

↓ Once PT Scan it generates the

table in data catalog.

→ You can run the crawlers on demand

→ Machine learning transformation:

• find material: Could be used to find matching data which enables you to find related products, customers

→ Amazon sage maker consumes data, trains the model, update end point for inference

AWS Glue:-



Example:-

→ Data in dynamo DB is copied to Amazon S3(ugly)

continuously as part of Data pipeline

→ JSON file are then converted into CSV to be used by Amazon sage maker.

→ Amazon sage maker consumes data, trains



AWS Data Migrations (DMS)

→ AWS DMS migrates data to AWS or AWS to another service.

→ AWS DMS allows for Database Migration

- 4. Data in CSV format is being queried by Amazon Athena per customer requirements.
- 5. Data pipeline manages the workflow as per customer requirements.

AWS Glue :-

- Glue is an ETL Service that runs on a Serverless Apache Spark Environment.
- Glue contains a data Catalogue for ETL that could interact with Athena and Redshift that could interact with Amazon S3 and Spectrum.
- AWS Glue ETL job uses Scala or Python.

AWS Batch:-

- AWS Batch performs all the scheduling and execution of the batch using Amazon EC2 and Spot instance.

- Cloud Watch Events.

AWS Data Pipeline:-

- AWS Pipeline Managed orchestration service
- AWS Data Pipeline offers more

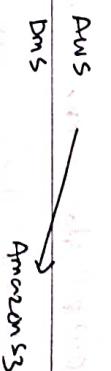


AWS Data Migrations (DMS)

allows to Aurora

allows for Database Migration

4. Data in CSV format is being queried by Amazon Athena
 5. Data pipeline manages the workflow per customer requirements.



AWS Glue:-

- Glue is an ETL Service that runs on a Serverless Apache Spark Environment.

- Glue contains a data Catalogue for ETL that could work with Athena and Redshift Spectrum.

- AWS Glue ETL job uses Scala or Python.

AWS Batch:-

- Batch Computation Job on AWS.
- AWS Batch performs all the scheduling and execution of the batch using Amazon EC2 and Spot instance.

- Cloud Watch Events.

AWS Data Pipeline:-

- AWS Pipeline Managed Orchestration Service

- AWS Data Pipeline offers more flexibility over EC2 instance.

AWS Step Functions

#4

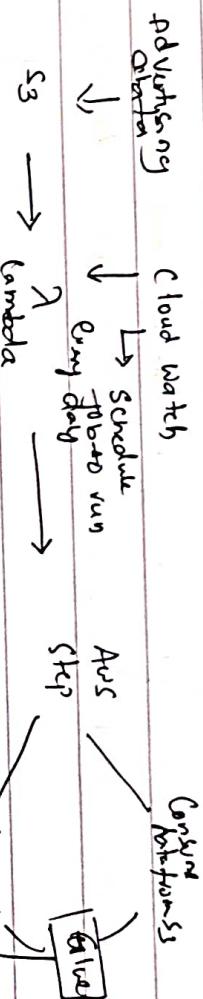
- AWS Step function for Creating Sequential Workflows

→ output from one step is fed as input to the next step

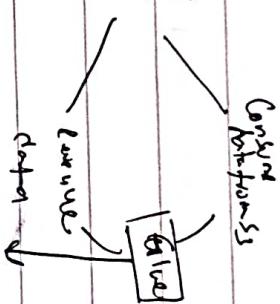
- AWS Step function creates Workflows instead of Machine diagrams that easy to debug and understand.
- AWS Step function allows for performing Resilient Workflow Automation fast with out writing code.

→ it allows for advanced error handling and retrying mechanism.

→ AWS Step Functions



ensure the data is present and available



AWS Step Functions:-

#5

WU
WU

V. Important



DATA STREAMING

B#7

B#7

B#7

Amazon Rekognition

→ Real time

→ Scalable

→ Fully Managed.

→ Kinesis → Apache Kafka

Video stream → Kinesis → monet. tensorflow
Sdk Video Stream File-Based video playback

→ Kinesis Video streams

Data Streams → Real time

Data Firehose → Near real time

Data Analytics → Data Analytics processing Service.

→ Amazon Kinesis to develop smart cities by
reducing crime rates and saving
Smart preventive maintenance

→ B

Amazon
DynamoDB

TOT

AMAZON → AWS Lambda → [AWS Lambda Function]

KINESIS STREAMS

Click streams → Data streams Analytics house

S3 red
diff

Kinesis Data Streams:

Data Producers
Data Consumers
Data Streams.

Shard → base throughput unit

60 seconds

Lambda

KPL
KCL
SDK APP

The KCL allows for ↓
data writing to a KCL allows data

Kinesis data streams
(Kinesis data stream)

→ producer.

producers → shard → consumers

1. Real-time
2. Real-time 60 seconds
3. from 1 to 7 days No data storage Auto scaling
4. Required Shards

Kinesis Firehouse

↳ near-real time analytics
ingested, transformed, encrypted and
loaded into Amazon ss, Redshift, Splunk.

Device

Kinesis Data Stream
Lambda
Splunk
CloudWatch Metrics
CloudWatch Metrics Insights
CloudWatch Metrics Insights Processor
CloudWatch Metrics Insights Processor

Kinesis Data Stream

Firehouse

AWS Kinesis & Cloud Analytics

Real time, automatically scales based on incoming data throughput.

redacting upfront cost, pay per use model

Step 1:-
Setup Data Stream Source
Develop Queries
for processing data.

Kinesis Video → for streaming live video

Kinesis Analytics → To run SQL queries on

Case #1:
Streaming ETL for IoT:

- Amazon Kinesis Data Analytics could be used to transform and filter streaming data from IoT devices such as Sensors then send real time alerts when Variable exceed certain limit.

Kinesis Summary:

Kinesis Data → for real-time streaming of data and gain real-time insights

Kinesis Firehose → for streaming data in near real time and storing them into S3.

Kinesis Video → for streaming live video in real time

Kinesis Analytics → To run SQL queries on streaming data and generate graphs and gain metrics

Domain 2 :- 24/1.

Code : L1X2#

VI
IT

EDA :-

Code :- A#14E

W Date _____
Page _____

→ Athena Athena :-

Code :- A#14E

↳ Cost-effective Query Service.

S3 → SQL Queries

Athena Very fast, Serverless.

It is very cost effective, you only pay per query

you choose to run.

Athena

(VS) Redshift

Send queries directly to S3

Send queries directly to S3

Designed for easy ad-hoc

Designed for ease of

queries into S3

Redshift

Doesn't require redshift

You will require redshift clusters.

Amazon Athena:-

1. Durability & High availability

2. Very Secure

3. Seamless integration with other AWS services (Glue).

→ AWS QuickSight :- BI tools, visualize

1. Cost effective

2. Highly Scalable, scale from 10 to 10,000 users

3. Allow for easy Data Analytics Embedding

4. End-to-End Business intelligent Solutions.

B

S3

Glue catalog → Quicksight

NL insight

→ Amazon QuickSight: SPICE

→ fast, highly available

→ import the data in SPICE for a great improvement.

Use case:-

→ packaged data products

→ improve apps by offering Analytics

→ integrate data into workflows.

↳ Run dashboards/reports into portals and Company sites.

Embedding Analytics:-

→ provides a complete Analytics experiences

→ Run App

VW
Range

Amazon QuickSight: NL insights:-

→ NL insight feature

1. Find data insights:

→ NL tool can offer an anomaly detection tool.

→ NL tool can perform anomaly detection and notify management.

2. Forecasting:

→ QuickSight NL tool can perform auto forecasting to predict critical business metrics.

3. Generate Auto-Narratives:

→ NL tool powered by NLG, QuickSight NL can offer narrative and tell a story.

→ Quicksight Security:-

→ provides a complete Analytics experiences

→ Run App

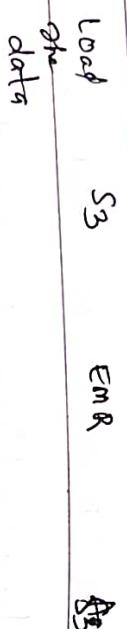
W
Range

Elastic Map Reduce (EMR)

Big Data platform.
Apache Spark, HDFS, Hbase.

→ EMR allows more dynamic Scaling among
Amazon EC2

→ Storage of Amazon S3.



Compute

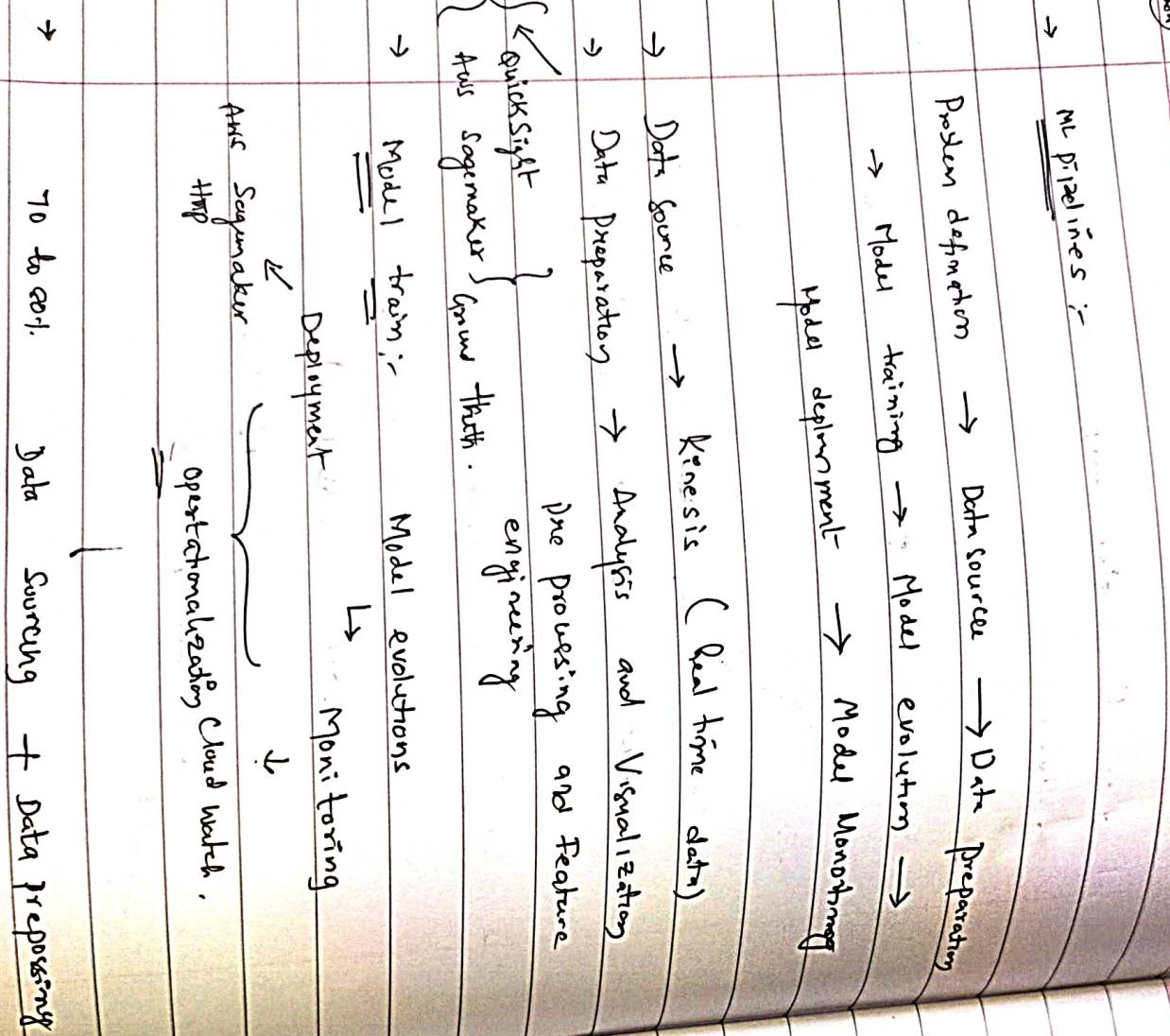
1. Easy to use.
2. Cost Effective
3. Elastic. → Auto Scaling
4. Enhanced Reliability
5. Improved Security
6. Increased flexibility

Spot instance:-

EMRFS:- EMR:- → allows direct access to
Hadoop environment

allows for flexible, lower-level
access to tools beyond spark.

Dom 21. Exploratory Data analysis



Demo:-

Ans Segmenter :-
Three Dimension → 3,

Statistical:-

{ Demo Statistical Data Analysis using Ans:- }

→ 15% Max. →

* Skewness:- less than 1 → highest skewness

Kurtosis → larger than zero. "2"

→ Ansible and Segmenter:-

Service → Quick sight ←

Manage Data → New data
L → upload file as CSV.

→ Visualize data

→ Sales price →

→ 70 to 80%. Data Sourcing + Data prepossing

→ Data Preparation Using AWS.



imbalanced Data Different scale Incomplete

Difficult Preprocessing Missing Data Outlier

Higher Dimensions Missing Correlated Features

Malformed

Distribution.

→ Possible - Solutions
Scale of features:-

→ Inconsistent scale. "Inconsistent formats."

Curry & us t : Always
solution:- leave distance as

→ Categorical Data

→ Label encoding → unique number, not enough
for larger categories

→ One - hot encoding → sparsity.

→ Missing Data problem → optional fields.

qs. 1. → legit } → [legit]
2.1. → fraudulent } → [legit]

→ legit

Drop them or

Values

Simple Values
can lose critical known words values
information classification trees.

M under security threat.

Impose Missing Values

Mean, mode, Predicting

Outliers:

Nanom Data entry

Data processing

Extremely rare non-representative.

Conditions.

They upset statistical characteristic.

Solution for outliers

Finding outliers → Z-score.

TAR options

Box and Scatterplot → Remove

correction.

→ High Dimensionality :-

cause of dimensionality

→ Very bad

→ challenging to visualize

→ increase risk of over fitting

→ feature engineering.

Feature Selection.

Discriminatory Reduction

Reducing dimension of data

(*) PCA :-

Objective of PCA is to reduce from n-dimension dataset to k-dimension dataset.

by finding K vectors onto which to

project the data so to minimize the projection error.

→ Highly Correlated features:-

Violation of regression.

Find highly correlated features, remove.

Variance Influence factors, VIF

$$x_i = c_1x_1 + c_2x_2 + \dots + c_nx_n$$

$$VIF = \frac{1}{1-R^2}$$

$R^2 = 1$, it will infinity.

VIF = 1 not complicated

VIF = 1 + s = moderately correlated

VIF > 5, highly correlated, remove the feature.

Data ingestion and Transformation

→ Data ingestion

Kinesis Data Stream, Firehose, Data Analytics

↓
Mensive Scale
duration 7 days

In
Manually Sharded

Configurations.

Use case:-

Data streaming
Data storage of streaming
data is needed

Kinesis Video Stream:

need

→ Realtime processing is a

requirement.

→ No manual configuration required
→ Integrate video data with
facial recognition in Amazon Rekognition

Kinesis Data Analytics:
→ Enables queries to run against
the input data

Kinesis fine house:

Fully managed service. doesn't require manual configuration for scale

Data transformation:

- Amazon EMR: Big Data Platforms
 - Spark
 - Hadoop, Hive, Flume, Presto.
- AWS Glue: Fully Managed ETL, Serverless
 - * AWS Glue Crawler → Metadata, Schema
 - * AWS Glue Catalog
- S3, DynamoDB, Redshift
- AWS Data Firehose Data transformer
 - * provides the ability to invoke Lambda function on incoming data
- Doesn't require configuration for scaling
 - * includes a sub-service that supports medical use case.
- > 5 minutes invocation time, payload limit of 6 MB.



AWS AE Services:-

- Language Analysis
 - * Amazon Lex :- Chatbot for customer service.
 - * Amazon Polly :- Speaks → Text to speech.
- Amazon Transcribe :- data from scanned documents
- * Amazon Translate :- NMT → different languages
- * Amazon transcribe :- Speaks to text from audio.

- * support both batch and real.
- Auto machine generating
- includes a sub-service that supports medical use case.

AI Services :-

Computer Vision :- Recognition

Provides auto mark analysis of image

and video data.

1. Amazon forecast :-
→ Deep AR Historical time series data
→ Supply Chain need, demand.

- Predict sales data

2. Amazon fraud detector :-
→ Fraud payment fraudulent payment

fraud transactions.

- Identifying fraudulent account creations

→ Identifying fraudulent actions taken

- 6 legitimate account.

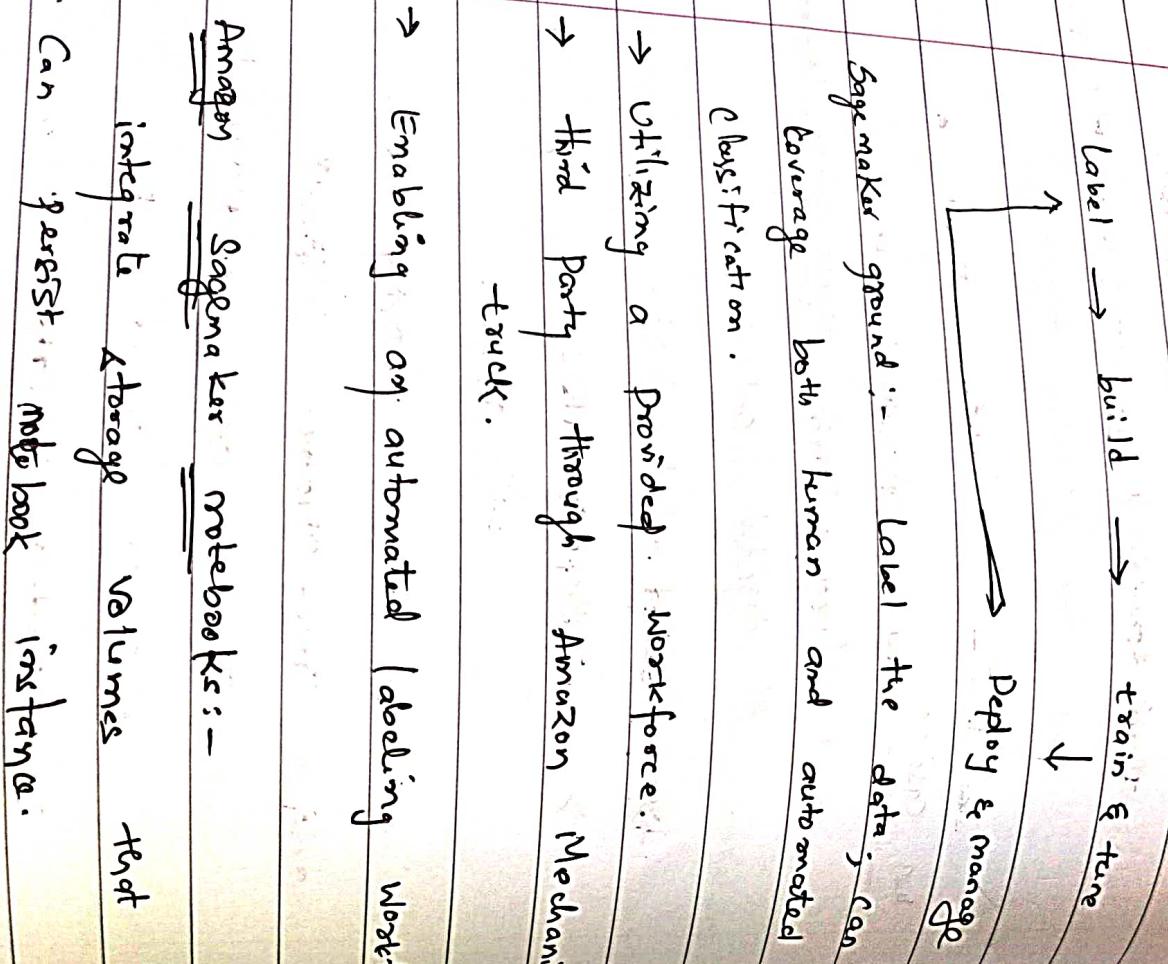
3. Amazon Personalize :-

→ Recommendation Engine.

→ Content or product recommendations.

→ Enable personalized promotions.

1. Sage maker process.



Sage maker ground :- Label the data, can leverage both human and automated classification.

→ Utilizing a provider workforce.

→ Third party through Amazon Mechanical Turk.

→ Enabling of automated labeling workflow

→ Amazon SageMaker notebooks:-

integrate storage volumes that can persist notebook instance.

→ Augmentor studio :- integrated development for ML.

Provides support for building Jupyterlab, deeply integrated to Sage maker

Auto pilot :- building, training and tuning of machine learning model.

- Classification and Regression.

Model training :-

Sage maker experiments:-
↳ autopilot

(Art) → Sage maker Neo :- Provides model optimization and cross-platform compilation

→ Enables model execution in the cloud and on the edge.

→ Hardware intel, Nvidia, ARM

~~Amazon Augmented AI (A2I) :-~~

→ Amazon integrates human review into inference from ECR.

Process.

→ provides a configurable workflow.

Situation when accuracy is critical.

→ Directly integrate into NLP, Text NLP, Recognition.

~~Amazon EMR:-~~

Jupyter notebook, multiple frame nodes.

Tensorflow, PyTorch.

→ On-premise servers,
You can build pre-built SageMaker

Docker image.

Use cases:-

→ Training a model on your hardware
and uploading to use in SageMaker

→ Container Services: Training or inference calls leveraged on any of the AWS Containers

Services.

AWS Deep Learning Containers are

Pre-configured for specific frameworks.

Deployment approaches:-

1. Hosting services:-

Deploy an inference end point to integrate inference into your workflow or application.

2. Batch transform:- Perform an inference job on an AWS Lambda function stored in Amazon S3.

→ post-training Validation

→ fit of model

→ choosing multiple models

→ prediction will be carried outside of sage maker

→

Save API and point

→ Will be integrate inference into an application or workflow.

→ Don't have all of the data you

will need to perform inference in the future.

Batch transform:-

→ Don't need an exposed endpoint for inference.

→ Already have the entire dataset on which you want to perform inference.

→ Need to process an incoming dataset to remove records prior to inference.

→ post-training Validation

→ fit of model

→ choosing multiple models

→ prediction will be carried outside of sage maker

→

Save API and point

→ Will be integrate inference into an application or workflow.

a) Train and Store mode

b) Load test data set → inference into Amazon S3.

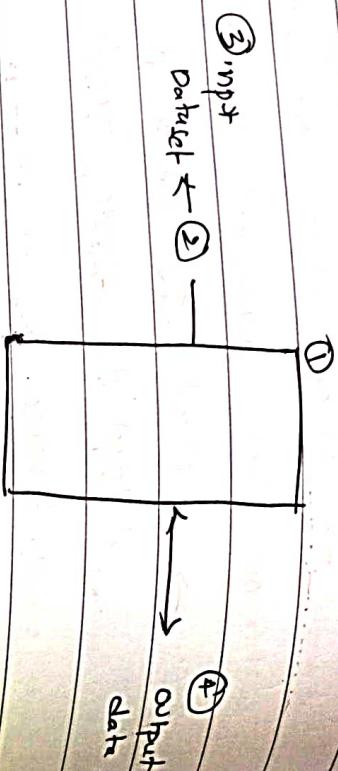
Amazon S3.

c) Specify Configuration for Data transformation Job.

d) Execute job on dataset.

→ Results are stored in S3.

Batch Transform Configuration:



- Infrastructure can be configured with number of instances and container
- their strategy dictates how data is pulled from the input source.
- ⇒ Transform data source input defines the source
- The assemble with parameters dictates how the data will be output

Infrastructure Approaches

1. Managed: SageMaker will launch infrastructure using SageMaker Container using the specified number of instances.

Custom Container:-

2. Custom Container: You can specify a specific container image to use for batch transform inference process.

Amazon Sagemaker Hosting Service :-

Client → API → Lambda → Sagemaker function

application

Load balancer

Many different servers

Trained Model.

(*) Deploy :-

- * Create Model → Create an Endpoint in Sagemaker

↓
Create an Endpoint from Configuration

Model Production Variants

Client

→

Sagemaker hosting services.

→ Endpoints are secured with HTTPS.

→ End point can auto scaled based on the defined Minimum and Maximum

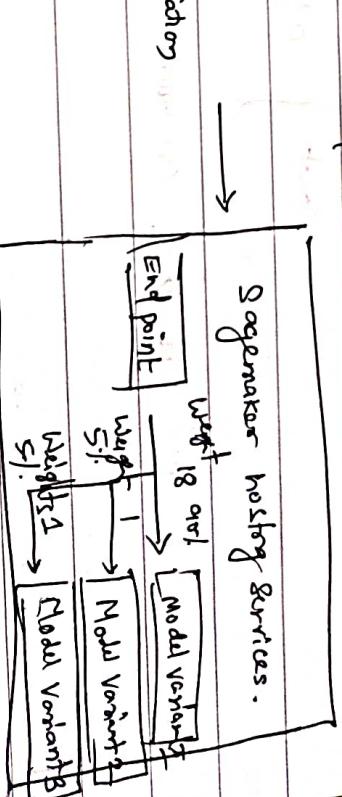
→ End-point can be updated with out downtime.

→ Supports A/B testing through multiple Production Variants.

Production Variants are defined in the endpoint configuration.

→ The initial Variant weight parameter determines the traffic to the variant.

→ Weight per variants can be updated by updating the endpoint configuration.



gather the data, which model is used

V. Important

Deploying a Model using Hosting Service

Steps in Amazon.

- Training of model
- Save train Model
- Creating End-point Configuration
- Deploying End-point
- Implementing Multiple production Variants :-
- im an end point.
- 1. Create a Model → Registering.
- 2. Create End-point Configuration
- Instance type }
instantiation count } all traffic.
- Model Name

→ Implementing Multiple production Variants:-
 A/B testing :-
 Variant 1 }
 Variant 2 } → initial weights
 1 → 50%
 2 → 50%

Production ready codes.
 Model 1 working very high accuracy

Amazon
 Elastic Inference } Multi model inference
 End point } pipelines.

GPU-(EEI) } End point can

Serve multiple } support two
 distinct models } for fine control.
 → Support multiple } HTTP
 Production Variants → requests
 multiple models.

Does not support general → E2 E3
 Inference } Complete

Can reduce } workload



Securing Sagemaker implementation

→ Security and Compliance is Shared Responsibility between Customer and the

→ Data like nobody is watching. Encrypt like everyone is.

→ Reviewing the shared Responsibility model in terms of Sagemakers

→ Learning High-level Security and compliance considerations

→ Exploring VPC configurations for Sagemaker

→ Controlling access to resources and data

using IAM.

→ Reviewing how to secure Sagemaker notebooks

* Data at Rest & Data In transit.
SG, Key use KMS key for encryptions.

→ Sage maker verified that its data outputs are encrypted at rest.

→ Specific Nitro-based instance

are encrypted but not with KMS.

AWS Responsibility Customer Responsibility

→ AWS is training for AWS Employee

→ Global data centres and underlying network.

→

1. Notebook instance
2. Processing Jobs
3. Training Jobs
4. Data transformations
5. Hosting

* Data in Transf:

Within Sagemaker inter-node communication is encrypted in transit with TLS 1.

* [TLS 1.2]

→ Some areas are not encrypted by default.

- Communication between Control Plane and

Training Job instance

- Node communication in distributed

processing jobs.

→ Node communication in distributed training jobs.

* Data Security :-

Amazon Macie :- Fully Managed Service

focused on data security

→ Utilizes Machine learning to

discover and report on data

stored in Amazon S3.

* Configuring for VPC for ML:-

→ To properly secure the infrastructure as well as data the infrastructure can access, infrastructure should be launched inside of a VPC. Efforts should be taken to minimize or eliminate data travelling over the public internet.

Alerts organization on several conditions

1. usual data access patterns
2. Configuration error for sensitive data.

Enables automated Compliance Checking for data storage.

→ Using VPC End Point for AWS S3.



→ Controlling Access with IAM.

They say that are by validating the submitted credentials.

Authorization :- Controls what a specific validated user can do on the platform based on their configured permissions.

→ Multi-factor Authentication

Provides additional security in the authentication process by leveraging another

a physical or virtual device that

generates token for login to AWS.

Services

Inference End points leverage AWS private

→ IAM identities

links.

↳ User

VPC's with no public internet need

to do not require a NAT or internet gateway

→ Roles : assumes permission for task

Groups → Allows to manage permissions for a group of IAM users

Soc 1,2,3 → Sagemaker.



- * Policies in AWS IAM. It defines what actions can be taken by who on which service.
- A JSON document that defines permission for AWS IAM identity for AWS services that identity defines both the AWS services that identity can access and what actions can be taken on that service.
- Can be either customer managed or managed by AWS.
- You should define least privilege access for each user using IAM policies.
- Securing Sagemaker notebooks:-
 - Sagemaker notebooks are internet-enabled by default.
 - Internet access can be eliminated by launching notebook in a VPC.
 - If not launched in VPC instance is launched in Sagemaker-managed VPC.

Implementing Highly availability and fault tolerance: -

Reliability :-

Being able to support the failure of components within the your architecture.

High availability :-

keeping your entire solution running in the expected manner despite issues that may occur.

SageMaker End points :-



Scaling SageMaker End points :-

→

Vertical Scaling:-

You scale up your instance type to a larger instance type with additional resources.

→

Horizontal Scaling:-

You scale out on additional instances to handle the demand of your applications.

→

CPU Utilization.

→

Additional Metrics.

Each instance invoked after a variant.

→

An initial Metric.

To utilize for scaling the SageMaker variant invocations per instance.

→

The average times per minute that

the average times per minute that each instance invoked after a variant.

→

Compute Optimized Instances.

* Standard instances:— t₂, t₃, m₄, m₅, m₆. (t₄, (s, c₆))

→

Configuring autoscaling for SageMaker End Point:-

→

Memory Optimized Instances.

* Accelerated Computing (P₂, P₃, g₄-dn, mlts) * Deployment Methodologies:—

→

Deployment Methodologies:-

→

Elastic Inference:-

Continuous — a development practice that integrates requires developers integrate code into shared repository for several times in day.

→

SageMaker Auto Scaling:-

Can be configured, Cloud Watch metrics to define scaling policy