

Unocodile, stones, crabs

## Fraud Analytics

$$\frac{100}{100} = \text{Recall } 100\%$$

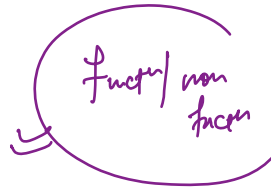
(what you predicted)  
Actually

$$= \frac{\text{Intend to captured (what caught fishes is true)}}{\text{Actual true}}$$

$$\left[ \frac{100}{100 + 50 + 20 + 20} \Rightarrow \frac{100}{190} \Rightarrow \text{Precision } \downarrow \right]$$

AWS Services (End-to-End)

Real time - fraud detection



# System Component

[10,000 Transaction per Seconds]

✓ High throughput.

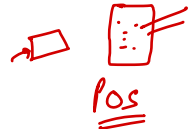
- ① Kinesis Data Streams → Ingesting transaction
- ② Kinesis Data Analytics (Apache Flink)
- ③ Dynamo DB → low latency lookups  
→ Store id 123 ⇒ fetch the details of store (metadata)

## Data

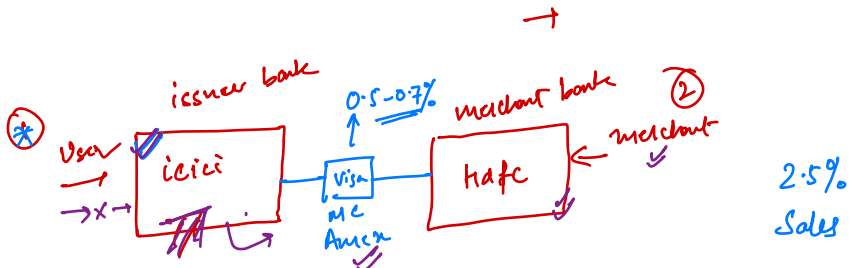
{ "txn id" : 1234  
"User id" : "abcd", "User-1234"  
"amount" : 50,000  
"timestamp" :  
"merchant id" :  
"ip address" :  
"location" : }

⇓

partition  
userid



Issuer bank



Max Share → Issue Bank → Max risk

Configuration of Kinesis  
+ 4 shards (1 shard ⇒ 1 MB/sec input  
1000 records/sec)

+ partition Key = store-id | user-id

② Real time feature engineering (think)

Agg

User-456 → 1 hr rolling spend  
→ Time happened in last 10 minutes

Enrichment

DynamoDB

fetch detailed merchant name,  
account age,  
past fraud record



Credit Score

5 km in 10 min  
150 km/h

```
{
  "transaction_id": "txn_123",
  "features": {
    "user_hourly_spend": 4500.00,
    "merchant_risk_score": 0.85,
    "transaction_frequency_10min": 5,
    "location_velocity": "150 km/h" // Distance from last transaction
  }
}
```

A B  
deli. Noida  
Karolbagh



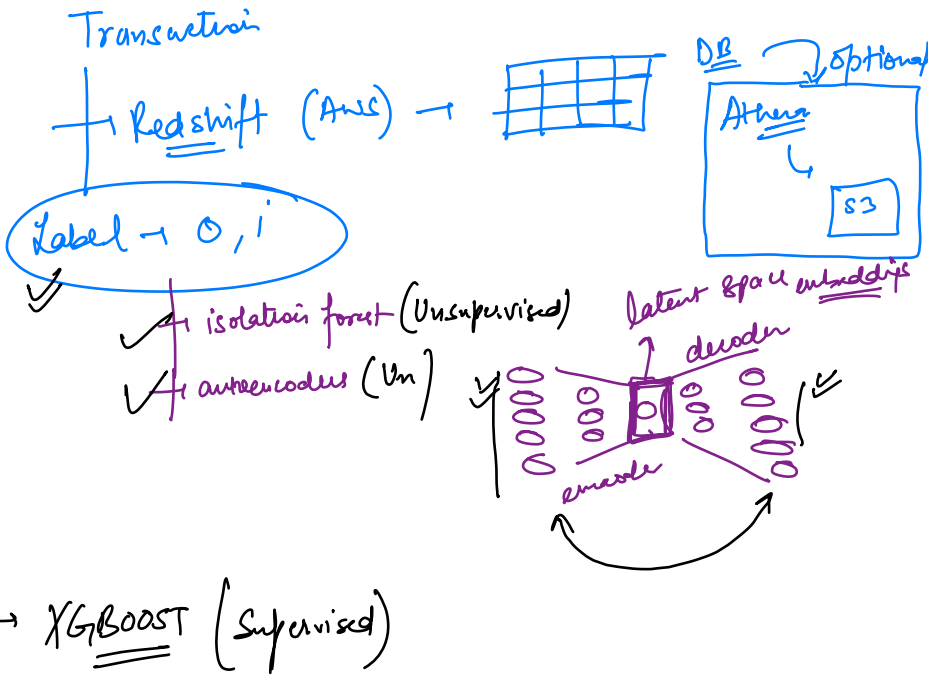
Amazon S3 (training models, artifacts)

## Model training & Deployment

Ans

(Training), hosting & monitoring model  
↳ Sagemaker

EMR (Spark) Batch feature engineering for historical data } Optional



## Configuration

+ Deep learning → GPU

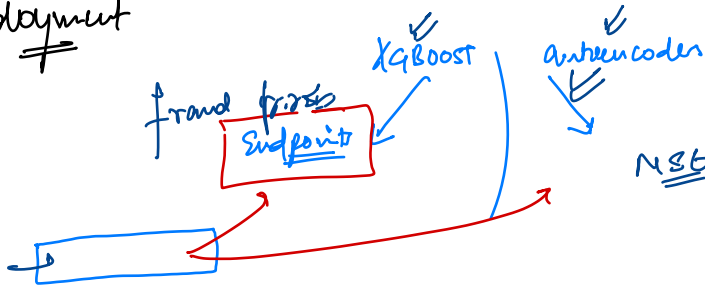
Precision / Recall

ml.p3.2xlarge

+ XGBOOST → CPU

ml.m5.4xlarge } cost efficient

## Deployment



Autoencoders → Computing reconstruction error (MSE)

↳ Anomaly MSE > 0.5

Combined Fraud Score

$$0.7 * \underline{\text{XGBoost-prsb}} + 0.3 \left( \frac{\text{autoencoder-mse}}{\text{max-mse}} \right)$$

# Real time - decision making

## ① Block transactions

AWS  
Lambda

- ✓ ① fraud score  $\rightarrow > 0.9$
- ✓ ② amount  $> \$1000$
- ✓ ③ geo-velocity = "impossible travel" manual review



10 hours

AWS API Gateway → core banking

SMS / email

↳ Amazon SNS

## Monitoring / Retraining phase

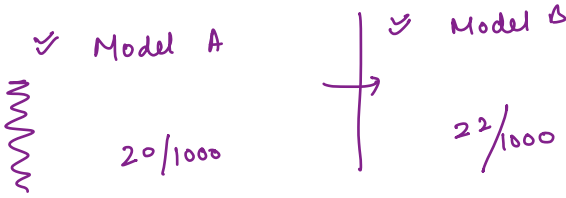
### ① Data Drift

Statistical test

training | production data

- ① KS test
- ② Chi-Sq
- ③ KL divergence

## Experimental



$\rightarrow$  t-test  
 $\rightarrow$  two-proportion Z-test

Q Historically  $\rightarrow$  0.1%

Today  $\rightarrow$  10,000  
Found  $\rightarrow$  100

Observed value  $\rightarrow$  120

$\rightarrow$  t-test  
 $\rightarrow$  Z-test

OK  $\rightarrow$  50 3 20 10 left  
M T W Th Fri  
Today 20 20 20 20 20

HR Head  
Most no. of people take leaves

100

① Real time data processing

$\rightarrow$  Kinesis + Flume

② Hybrid model  $\rightarrow$  Combined Scans  
(Uns + Sup)

③ Automated data

①  $\rightarrow$  Metric Shift (Precision, Recall)

② → Data drift (KL, KS, Chi)

③ Model A → Model B (t-test, z-test, two-prob z test)

④ Cost efficiency

+ GPU  
- CPU

⑤ Alert System → Lambda + API Gateway  
+ SNS