# 1. Batch Vs Online ML

Wednesday, March 17, 2021     5:30 PM
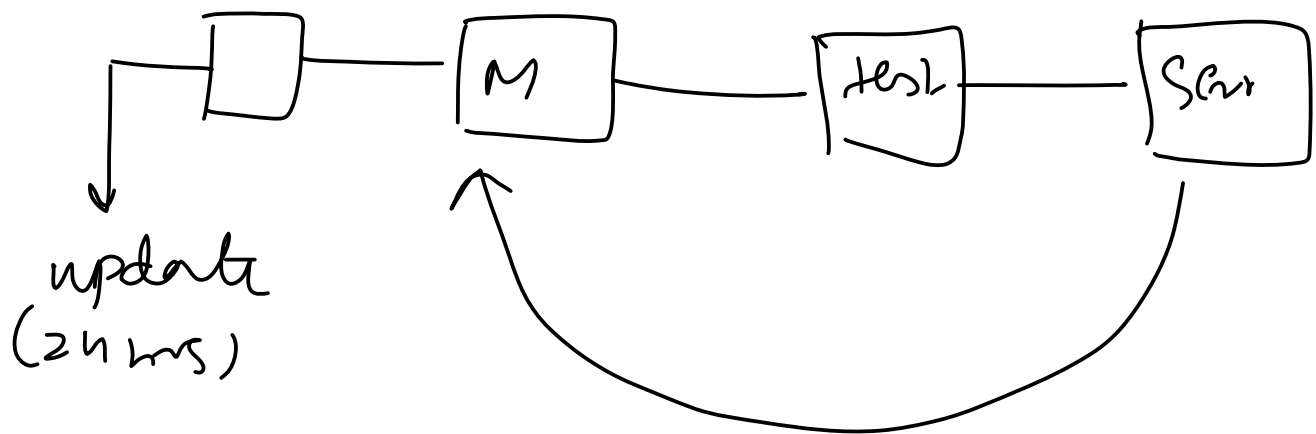
# 2. Batch/Offline ML

# 3. The problem with Batch Learning
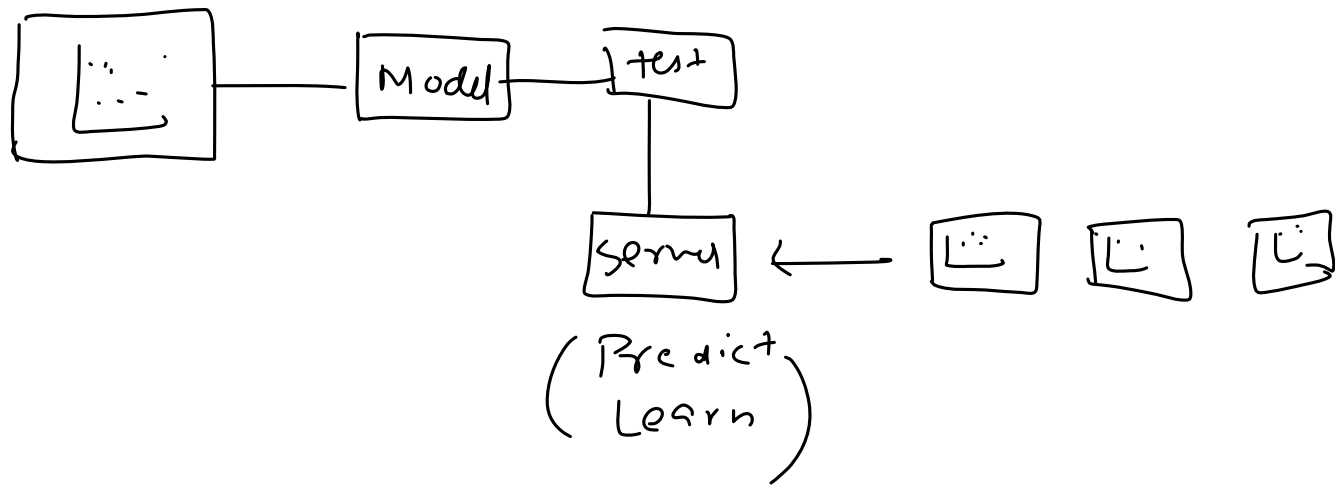
update
(24 hrs)

# 4. Disadvantages of Batch ML

Wednesday, March 17, 2021       5:32 PM

1. Lots of Data
2. Hardware Limitation
3. Availability

# 1. Online Machine Learning

Thursday, March 18, 2021    4:27 PM

# 2. When to use?

1. Where there is a concept drift
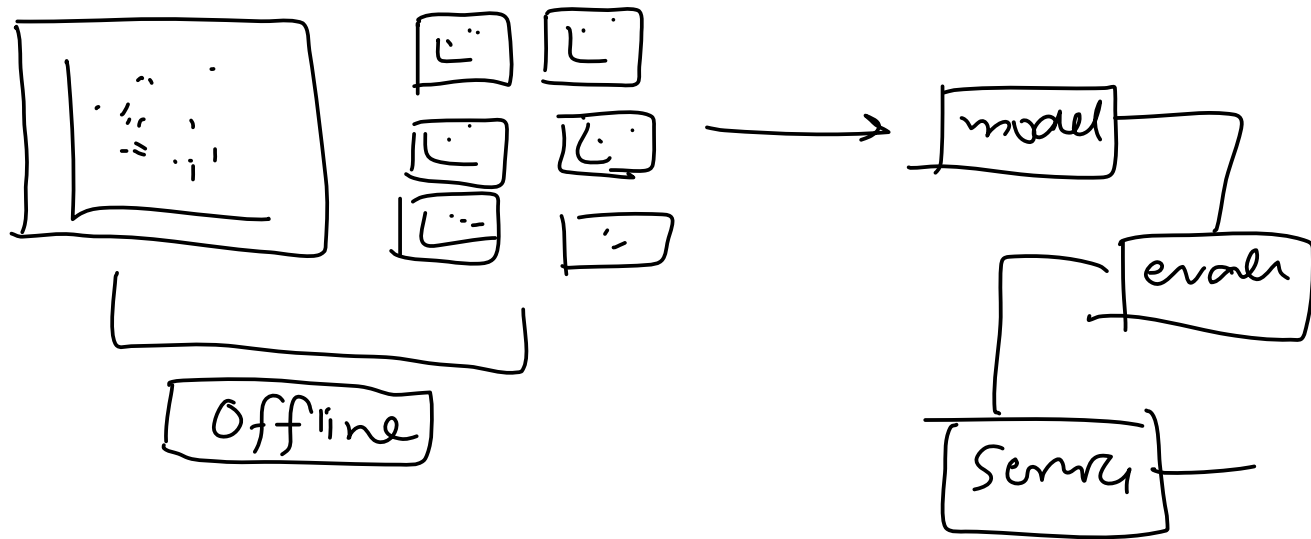2. Cost Effective
3. Faster solution

# 3. How to implement?

Thursday, March 18, 2021     4:28 PM

# 4. Learning Rate

Thursday, March 18, 2021 4:28 PM

# 5. Out of Core Learning

# 6. Disadvantage

Thursday, March 18, 2021     4:29 PM

1. Tricky to use
2. Risky

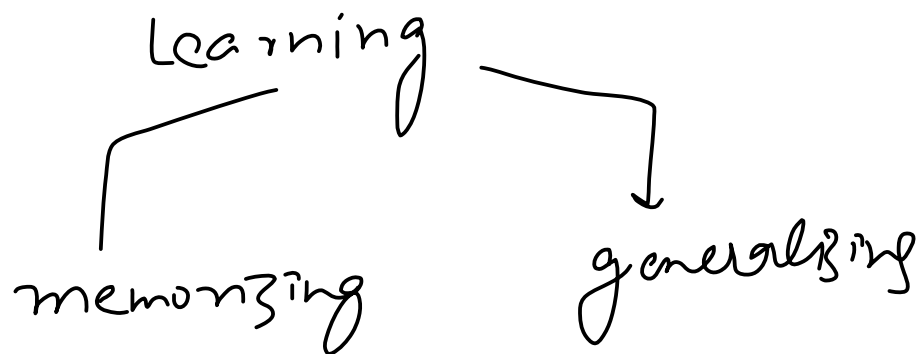# 7. Batch Vs Online Learning

Thursday, March 18, 2021    4:29 PM

| Offline Learning | Features | Online Learning |
|---|---|---|
| Less complex as model is constant | Complexity | Dynamic complexity as the model keeps evolving over time |
| Fewer computations, single time batch-based training | Computational Power | Continuous data ingestions result in consequent model refinement computations |
| Easier to implement | Use in Production | Difficult to implement and manage |
| Image Classification or anything related to Machine Learning - where data patterns remains constant without sudden concept drifts | Applications | Used in finance, economics, heath where new data patterns are constantly emerging |
| Industry proven tools. E.g. Sci-kit, TensorFlow, Pytorch, Keras, Spark Mlib | Tools | Active research/New project tools: E.g. MOA, SAMOA, scikit-multiflow, streamDM |

Image courtesy - https://www.iunera.com/kraken/fabric/simple-introduction-to-online-learning-in-machine-learning/
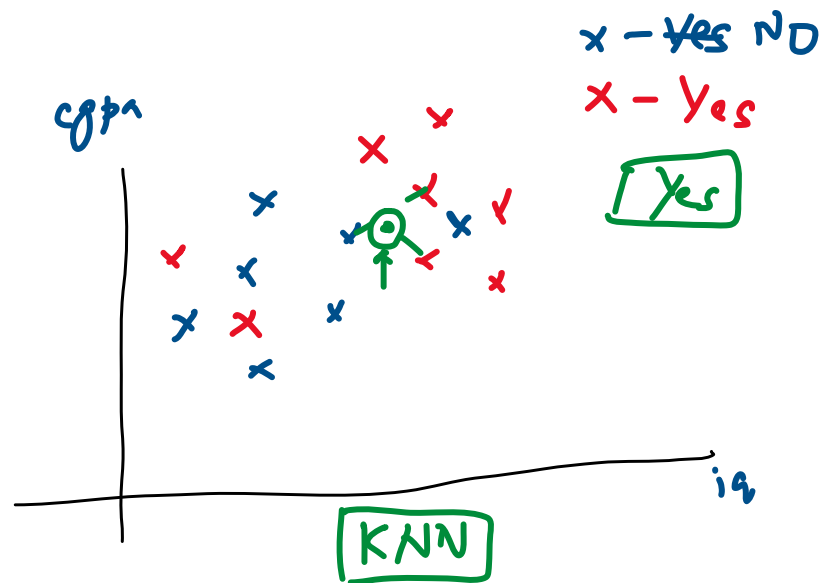
# 1. Instance Vs Model Based Learning

Learning

memorizing        generalizing

# 2. Instance Based

| iq | gpa | placement Y/N |
|----|-----|---------------|
| 80 | 8   | Y             |
| 70 | 7   | N             |

7.5, 103

x — Yes NO

x — Yes

Yes

cgpa

KNN

iq

# 3. Model Based

iq | cgpa | placm



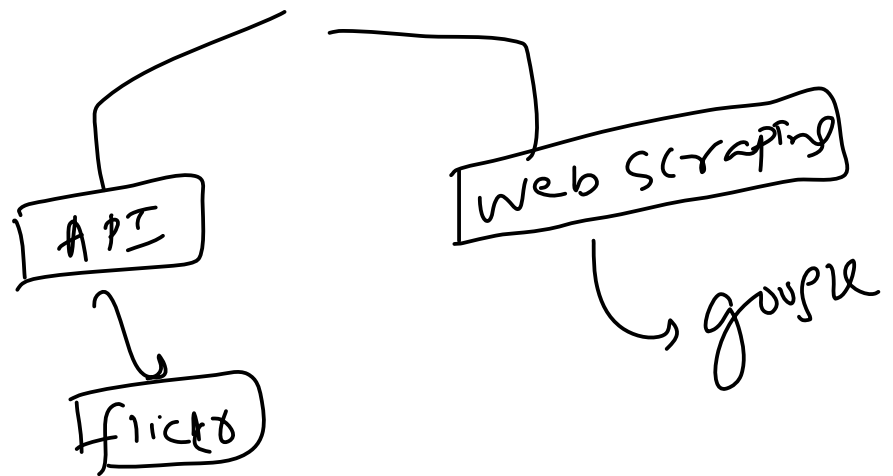cgpa

x        x

iq

yos

la hof | pr

# 4. Differences

Friday, March 19, 2021     4:06 PM

| Usual/Conventional Machine Learning | Instance Based Learning |
|---|---|
| Prepare the data for model training | Prepare the data for model training. No difference here |
| Train model from training data to estimate model parameters i.e. discover patterns | Do not train model. Pattern discovery postponed until scoring query received |
| Store the model in suitable form | There is no model to store |
| Generalize the rules in form of model, even before scoring instance is seen | No generalization before scoring. Only generalize for each scoring instance individually as and when seen |
| Predict for unseen scoring instance using model | Predict for unseen scoring instance using training data directly |
| Can throw away input/training data after model training | Input/training data must be kept since each query uses part or full set of training observations |
| Requires a known model form | May not have explicit model form |
| Storing models generally requires less storage | Storing training data generally requires more storage |

# 1. Data Collection

API

Web Scraping

flickr

google

# 2. Insufficient Data/Labelled Data

$\widehat{NLP}$

A

( 100 )

M1

B

( $10^6$ )

M2 ✓

$|t^o, t^o,$

# 3. Non Representative Data

Sampling Noise

Sampling bias

# 4. Poor Quality Data

60 %

# 5. Irrelevant Features

Garbage In
Garbage Out

age | wt | ht || $\boxed{\text{location}}$ $\times$

bmi $\sim$ feature engineering

# 6. Overfitting

# 7. Underfitting

# 8. Software Integration

Saturday, March 20, 2021　　6:01 PM

# 9. Offline Learning/ Deployment

# 10. Cost Involved

MLops

Devops

# 1. Retail - Amazon/Big Bazaar

Monday, March 22, 2021     6:07 PM

# 2. Banking and Finance

Monday, March 22, 2021        6:07 PM

# 3. Transport - OLA

Monday, March 22, 2021        6:07 PM

# 4. Manufacturing - Tesla

Monday, March 22, 2021 6:08 PM

# 5. Consumer Internet - Twitter

Monday, March 22, 2021          6:08 PM

# Machine Learning Development Life Cycle(MLDLC/MLDC)

SDLC

ML DLC

# 1. Frame the Problem

Tuesday, March 23, 2021     12:10 PM

# 2. Gathering Data

Data

- CSV
- (API)
- Web scraping
- Database
  - ETL → (DW)
- Spark clusters

# 3. Data Preprocessing

Tuesday, March 23, 2021     12:11 PM

→ Remove duplicates

→ Remove missing val

→ Outliers

→ Scale

# 4. Exploratory Data Analysis

Vizs
Univariate/Biavariate
Outlier detection
Imbalance →

# 5. Feature Engineering and Selection

Tuesday, March 23, 2021        12:12 PM

# 6. Model Training,Evalation and Selection

Tuesday, March 23, 2021     12:12 PM

Ensemble
learning

# 7. Model Deployment

pred

API

JSON

binary fila
(pickle)

python

{ Heroku
  AWS
  GCP }

# 8. Testing

A /B testing

# 9. Optimize

Tuesday, March 23, 2021     12:15 PM

Retrain

Rotting

→ Backup

→ Data

→ load balancing

# 1. Various Data Based Job Roles

Wednesday, March 24, 2021    1:25 PM

| DATA ANALYST | DATA ENGINEER |
|---|---|
| DATA SCIENTIST | ML ENGINEER |

Plan → Data → Proces → EDA → mods → eval → deploy → optimze

# 1. Data Engineer

## Job Roles

- Scrape Data from the given sources.

- Move/Store the data in optimal servers/warehouses.

- Build data pipelines/APIs for easy access to the data.

- Handle databases/data warehouses.

OLTP → OLAP

Db          Dw    pipelines

## Skills Required

- Strong grasp of algorithms and data structures

- Programming Languages (Java/R/Python/Scala) and script writing

- Advanced DBMS's

- BIG DATA Tools (Apache Spark, Hadoop, Apache Kafka, Apache Hive)

- Cloud Platforms (Amazon Web Services, Google Cloud Platform)

- Distributed Systems

- Data Pipelines

# 2. Data Analyst

**Responsibilities of a Data Analyst**

- *Cleaning and organizing Raw data.*
- *Analyzing data to derive insights.*
- *Creating data visualizations.*
- *Producing and maintaining reports.*
- *Collaborating with teams/colleagues based on the insight gained.*
- *Optimizing data collection procedures*

**Skills**

- *Statistical Programming*
- *Programming Languages (R/SAS/Python)*
- *Creative and Analytical Thinking*
- *Business Acumen — Medium to High preferred*
- *Strong Communication Skills.*
- *Data Mining, Cleaning, and Munging*
- *Data Visualization*
- *Data Story Telling*
- *SQL*
- *Advanced Microsoft Excel*

# 3. Data Scientist

Wednesday, March 24, 2021     1:26 PM

"A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician".

# 4. ML Engineer

Wednesday, March 24, 2021    1:26 PM

## Responsibilities

- Deploying machine learning models to production ready environment
- Scaling and optimizing the model for production
- Monitoring and maintenance of deployed models

## Skills

- Mathematics
- Programming Languages (R/Python/Java/Scala mainly)
- Distributed Systems
- Data model and evaluation
- Machine Learning models
- Software Engineering & Systems design

# 5. Comparison

Wednesday, March 24, 2021    1:26 PM

|  | ANALYTICAL SKILLS | BUSINESS ACUMEN | DATA STORYTELLING | SOFT SKILLS | SOFTWARE SKILLS |
|---|---|---|---|---|---|
| **DATA ANALYST** | HIGH | MEDIUM TO HIGH | HIGH | MEDIUM TO HIGH | MEDIUM |
| **DATA ENGINEER** | MEDIUM | LOW | LOW | MEDIUM | HIGH |
| **DATA SCIENTIST** | HIGH | HIGH | HIGH | HIGH | MEDIUM |
| **ML ENGINEER** | MEDIUM TO HIGH | MEDIUM | LOW | HIGH | HIGH |

# 1. What are Tensors

Thursday, March 25, 2021    4:44 PM

# 2. 0D Tensor/Scalar

$$(2) \quad (3)$$

Scalars $\longrightarrow$ $\boxed{\text{0D Tensor}}$   $0$

# 3. 1D Tensor/Vector

Thursday, March 25, 2021          4:45 PM

$[1, 2, 3, 4]$ → 1D Tensor

Vector

→ 1D array / array

1D Tensor / Vector

4D

$[0, 1, 2, 3]$

↳ scalars → Vector

nDim → 1

Axis

2 Dim

↳

No. of axes = rank
= dim

$[1, 2]$ → Vector (2)

↳ 1D Tensor

# 4. 2D Tensor/Matrices

$[1, 2, 3]$  $[4, 5, 6]$  $[7, 8, 9]$

$$\begin{bmatrix} [1, 2, 3] \\ [4, 5, 1] \\ [7, 8, 9] \end{bmatrix} \rightarrow 2D$$

Rank = 2 = ndim

# 5. ND Tensors

4D

$4 \times 3 \times 3$

3D Tensor

depth

4D tensor

3D

0 - 5D

1D - 5D

# 6. Rank, Axes and Shape

No. of axis = Rank = No. of dim

vector $\longrightarrow$

[1,2,3]

(3,)

3

Shape

row

2,

(3,3)

a

ndim rank

(2)

Shape (2,3)

6

Size of tensor

Size = 1

(4,2)

8

# 7. Example of 1D Tensors

Students ⟨ $\underline{10000}$ ⟩

WB = 0
KR = 1

1D Tensor | Vector

| cgpa | iq | state | placement |
|------|-----|-------|-----------|

$[ 8.1, 91, 0 ]$ → (1D)

→ Vector

→ 3D

8.1   91   (WB)
            0

$[ 7.2 \quad 102 \quad KR ]$

$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ , \\ 0 \end{bmatrix}$

(iq)

(State)

cgpa

1D Tensor

$[ 1 \; 0 \; 1 \; 01 \; 1 \; 01 \dots ]$

↳ 10000 no.

# 8. Example of 2D Tensors

2D

$\pm 0000$

$$\overline{cgpa \mid iq \mid state \mid place}$$

$$\left[\begin{array}{c} [\text{—} \quad \text{—} \quad \text{—}] \\ [\text{—} \quad \text{—} \quad \text{—}] \end{array}\right\} \rightarrow vector$$

$\searrow$ $\boxed{b}$

2D Tensor

$$\left[\begin{array}{c} [\text{—} \quad \text{—} \quad \text{—}] \\ [\text{—} \quad \text{—} \quad \text{—} \quad \text{—}] \\ [\text{—} \quad \text{—} \quad \text{—} \quad \text{—}] \\ [\text{—} \quad \text{—} \quad \text{—}] \end{array}\right]$$ $$\left[\begin{array}{c} \phantom{\text{—}} \end{array}\right]$$

1D Tensor

2D Tensor

NLP

Hi Nitish
Hi Rahul
Hi अंकित

| Hi | Nitish | Rahul | अंकित |
|----|--------|-------|-------|
| 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 |

[ [ [1,0,0,0], [0,1,0,0] ]      2u      (3, 2, 4)

  [ [1,0,0,0], [0,0,1,0] ]    → 3D Tensor

  [ [1,0,0,0], [0,0,0,1] ] ]

Timeseries Data

(2)

Highest | Lowest

→ 10 years

(365, 2)

| | | |
|---|---|---|
| Day 1 | — | — |
| Day 2 | — | — |
| Day 365 | — | — |

→ 2D → 10 →    time axis

365

(10, 365, 2)

↳ 3D Tensor

# 10. Example of 4D Tensors

images $\rightarrow$ [ C V ]

$(3, 1200, 800)$

$\downarrow$

3D Tensor

$(1200, 800)$

R      G      B

$\rightarrow$ 50 color images

$(50, 3, 1200, 800)$

$\downarrow$ 4D Tensor

## 12. Example of 5D Tensors

Thursday, March 25, 2021    4:47 PM

Videos
↳ frames

60 sec ──→ 4 videos
  ↳ 30 fps
    ↳ 480p → 480 × 720 (3 channels)

4 (60 sec)
videos
  ↳ 27 GB
      ↳ video
         ↳ mkv → mpeg mp4

(4, (1800, 480, 720, 3))
              ↑
          5D tensor        ↳ float 32

# 1. Installing Anaconda

Friday, March 26, 2021      5:40 PM

# 2. Jupyter Notebook Intro

Friday, March 26, 2021     5:40 PM

# 3. Virtual Env

Friday, March 26, 2021     5:40 PM

# 4. Using Kaggle

Friday, March 26, 2021      5:41 PM
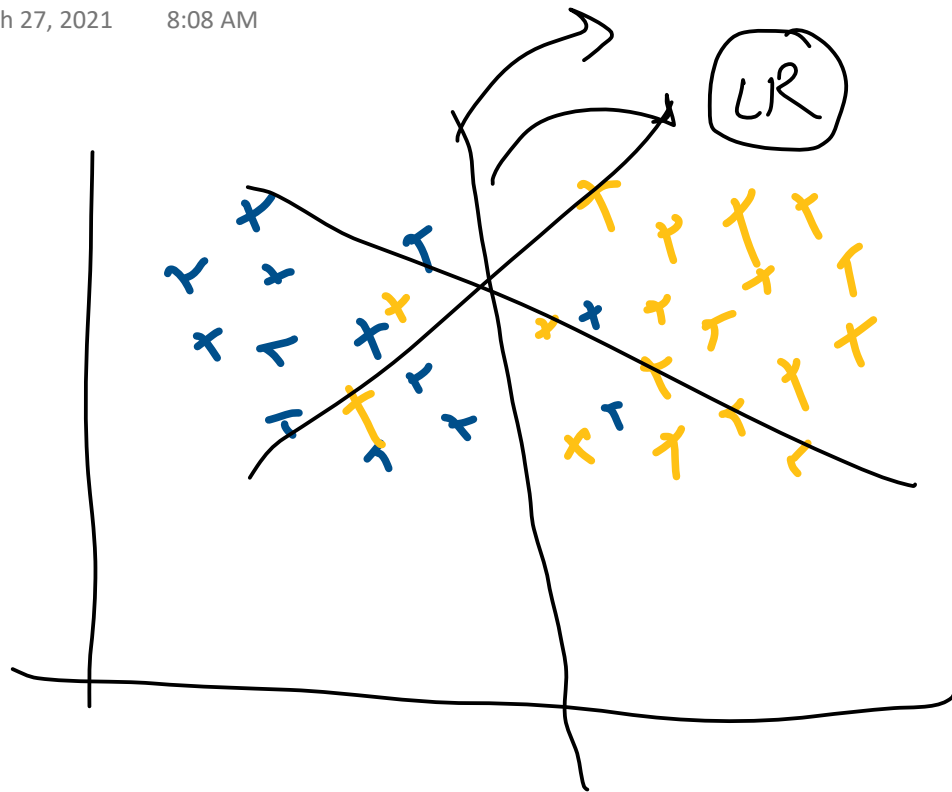
# 5. Using Google Colab

Friday, March 26, 2021      5:41 PM

# 6. Running Kaggle Data on Google Colab

Friday, March 26, 2021        5:41 PM

# End to End Example

Saturday, March 27, 2021        8:08 AM

# 1. Business Problem to ML Problem

Netfix

Churn rate ↓
    ↳ 4·/·
        ↳ (3.75)
            ↳ (3.5)

700 → 2·/·
    ↳ 98 ↑
        ↳ 2/·
            ↳

Increase revenue →
    ↳ 4 /· → (3.75/·)

# 2. Type of Problem

Big picture → end product ⤳ prediction

prediction ↓

1) Superv

↳ Classification

Regression

# 3. Current Solution

Churn rate → 5%  ⟶  ± 10%

6%

↳ factors

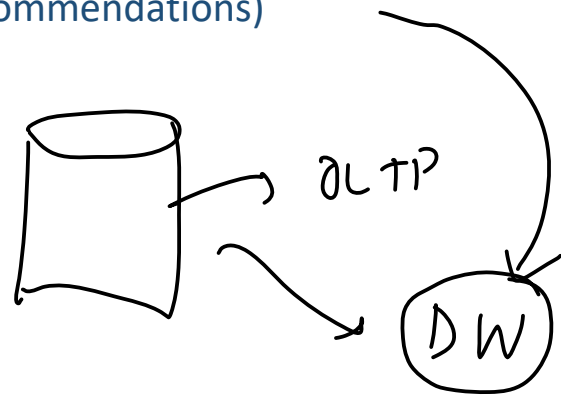# 4. Getting Data

1. Watch time
2. Search but did not find
3. Content left in the middle
4. Clicked on recommendations(order of recommendations)

# 5. Metrics to measure
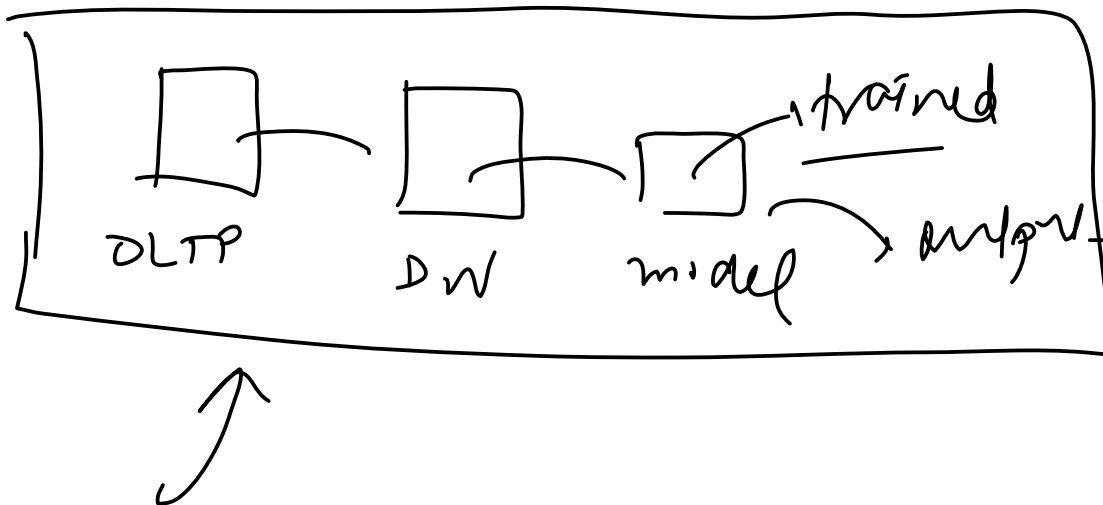
Monday, March 29, 2021     7:30 AM

# 6. Online Vs Batch?

OLTP → DW → model → trained → output

# 7. Check Assumptions