

Introduction to MACHINE LEARNING

Linear Regression :- ~~Predictions~~ Prediction
Using the simple yet powerful ~~DLR~~ LR algorithm.

Logistic Regression :- Predicting classes for events
Using Probability of odds.

D.T. & R.F. :- Question based approach to
Prediction.

PCA :- Simplify data to lower dimensions.

KNN :- Classifies a data point based on its
nearest neighbours.

Naïve Bayes :- Classifies data based on
Conditional Probability.

K-means :- Cluster a set of objects based on measure of similarity.

SVM :- Predict or classify using support vectors.

Time Series Analysis :- Analysing time series and forecasting future occurrences.

Linear Regression :-

Logistic Regression :- Used when there is a linear relationship b/w variables

DT & RF :- There is no linear relation & completely non-linear.

SVM :- Good; when no. of features are higher than the no. of obs.

KNN :- Good; when everything is mathematical, because it looks at Euclidean distance / Manhattan distance, need to normalise the data really well.

PCA :- Good, when the features have not a lot of interdependence b/w them & dealing with uncorrelated features.

LDA :- Advanced version of PCA used for classification.

Recommendation Engine

A priori algo :- How to come up with recommendation based on time slices. Particular transaction with particular time what item the customer will buy

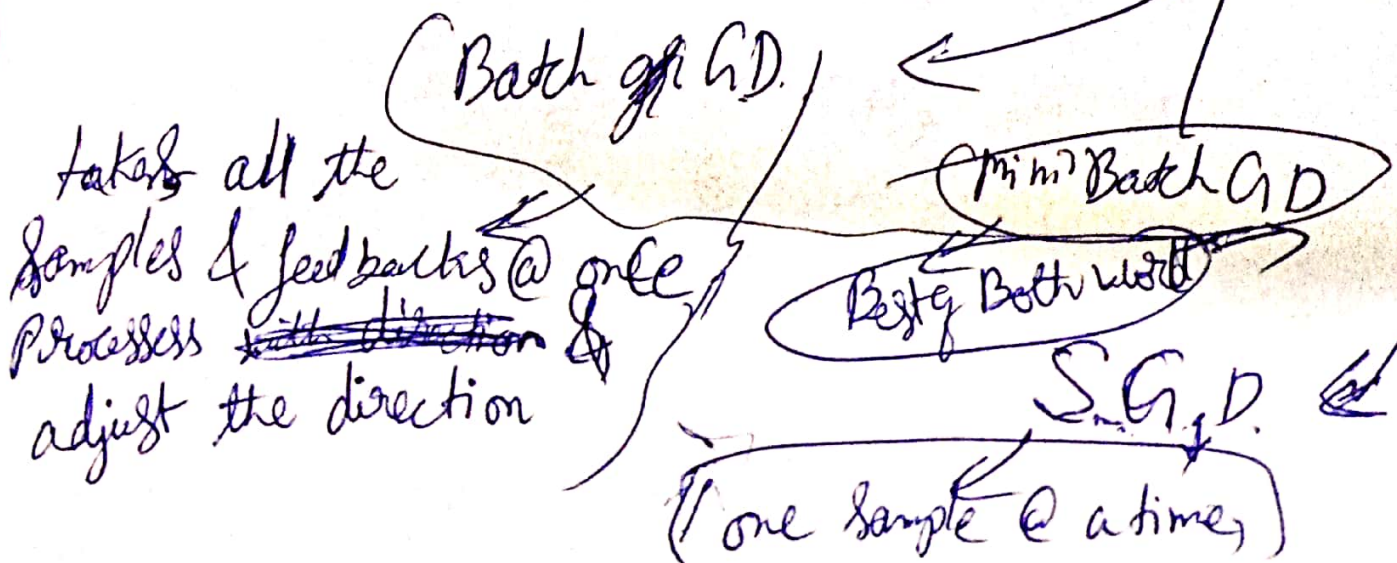
Anomaly detection - detect what is normal & identify what is abnormal. We use (T-digest, windowing algo).

How to win Kaggle Competition = Ensemble learning

Boosting → Adaboost
Gradient Boosting
XG Boost

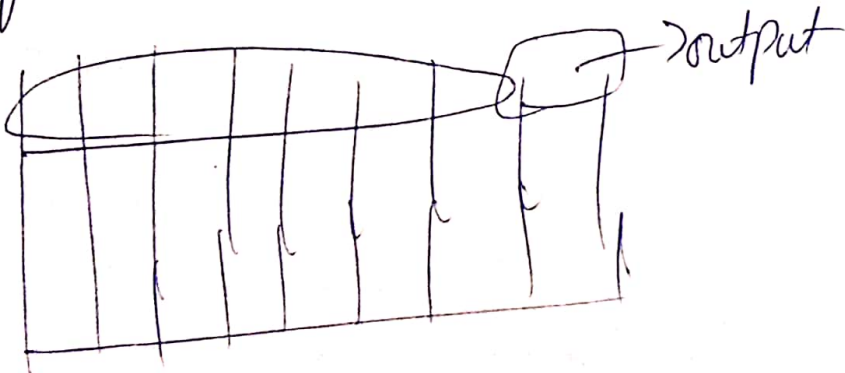
Optimization :- Linear Problem = Use minimize SSE etc

Non-linear problem = we use gradient descent which direction we need to go



Top ML Terminologies

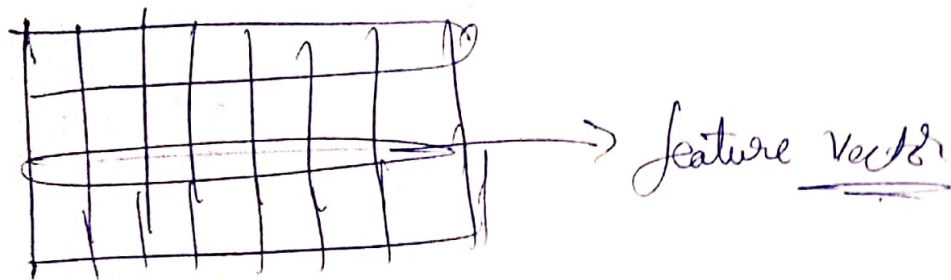
↳ Term - 1 = Features = The No of distinct traits that can be used to describe each item in a quantitative manner.



Term - 2 = A sample is an item to process.
For instance, it could be a document, picture, sound, or a raw database.

Term 3: Feature Vector

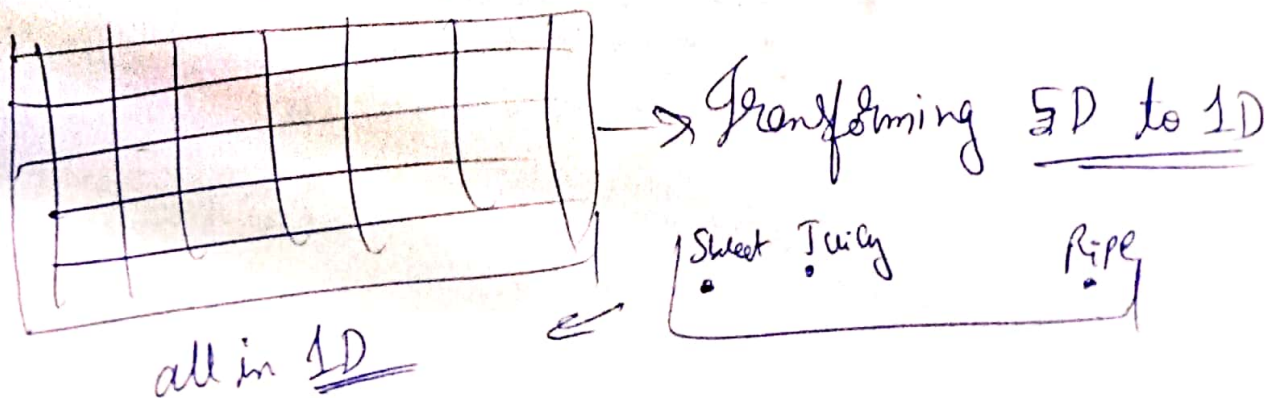
A feature vector is an n -dimensional vector of numerical features that represent some object



Term 4: Feature Extraction.

↳ Preparation of feature vector.

↳ Transforms the data from high dimensional space to a space of fewer dimensions.



Term 5: Training Set

Set of data to discover potentially predictive relationships.

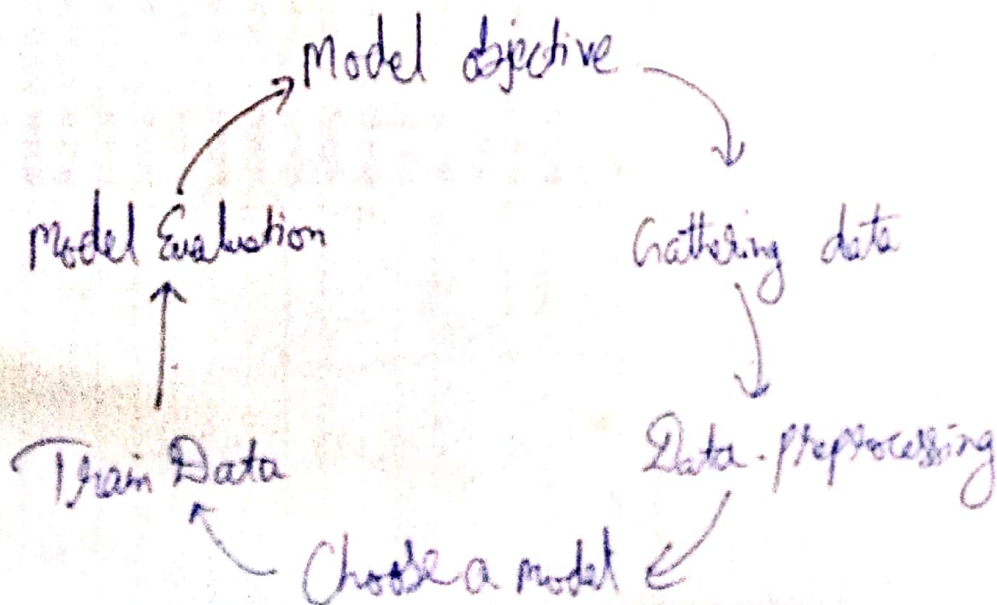
Term 6: Test Set

Set of data to validate the predictive relationship.

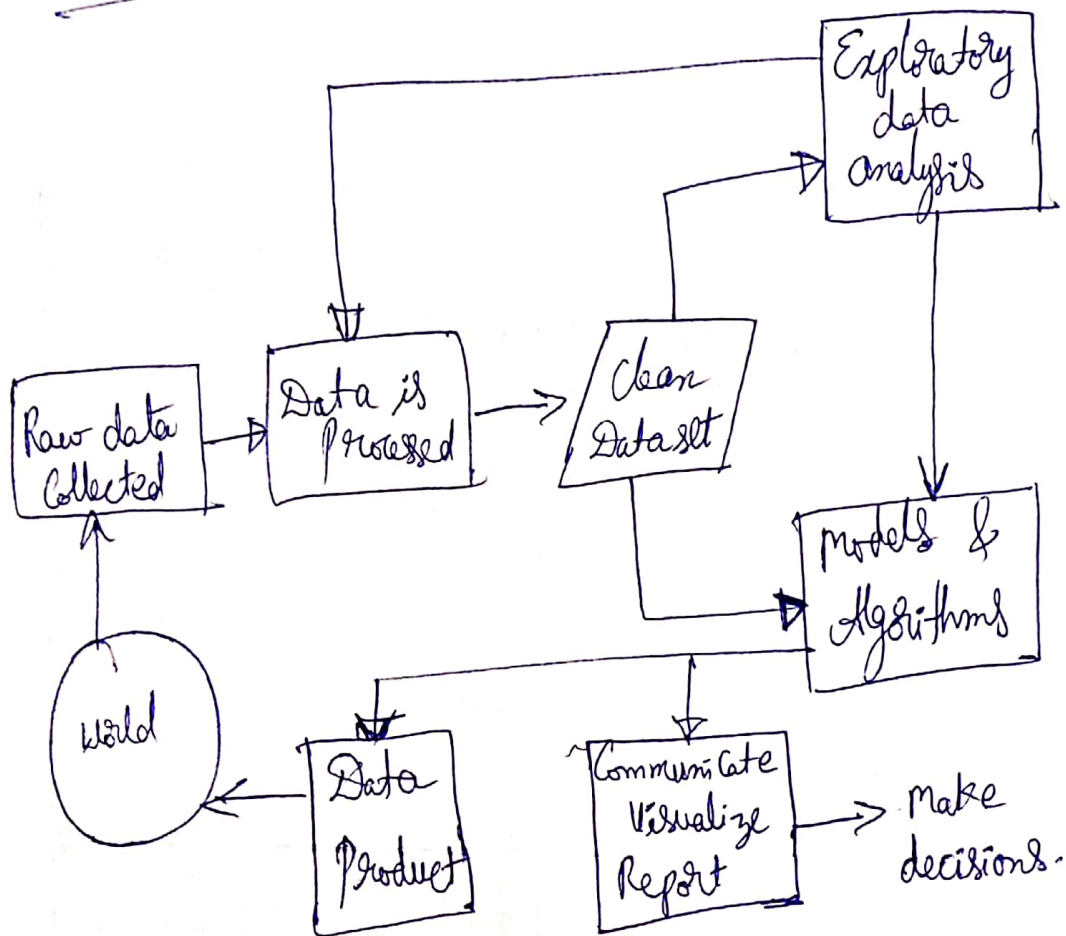
Term 7: Scoring

Evaluating the performance of a machine learning based on any statistical metric.

ML workflow



Putting it all together



Step 1: Collect Raw Data

To solve a given problem, as a data scientist you need data.

Step 2: Store Raw Data

Raw data means data that has not been changed since ~~so~~ acquisition. This raw data is stored in your storage systems.

Step 3: Data Pre-processing.

As a Data Scientist, a lot of your time will go in Data Pre-processing. Also known as Data Cleaning.

This step includes

- ↳ Removing outliers
- ↳ Replacing missing data
- ↳ Malicious Data.
- ↳ Erroneous Data
- ↳ Irrelevant Data
- ↳ Inconsistent Data
- ↳ Formatting.

Once Data is Cleaned, it needs to be processed to make it ready for use.

↳ This stage includes

- Sorting
- Summarization.
- Aggregation
- Validation
- Classification

→ Data Preprocessing (Data Cleaning) is at times considered to be part of Data Processing

Step 4: Exploratory Data Analysis

What are the key concepts about EDA?

2 types of Data Analysis

- ↳ Confirmatory data analysis
- ↳ Exploratory data analysis.

4 objectives of EDA

- ↳ Discover patterns
- ↳ Spot anomalies
- ↳ Frame hypotheses
- ↳ Check assumptions

2 methods for Exploration

- ↳ Univariate analysis
- ↳ Bivariate analysis

Stuffs done during EDA

- Trends
- Distributions
- Mean
- Median

- Outliers
- Spread measurement (SD)
- Correlations
- Hypothesis testing
- Visual exploration

Objectives of EDA

- ↳ Discover patterns
- ↳ Spot anomalies
- ↳ Frame hypotheses
- ↳ Check assumptions

Step 5: Modelling & Algorithms

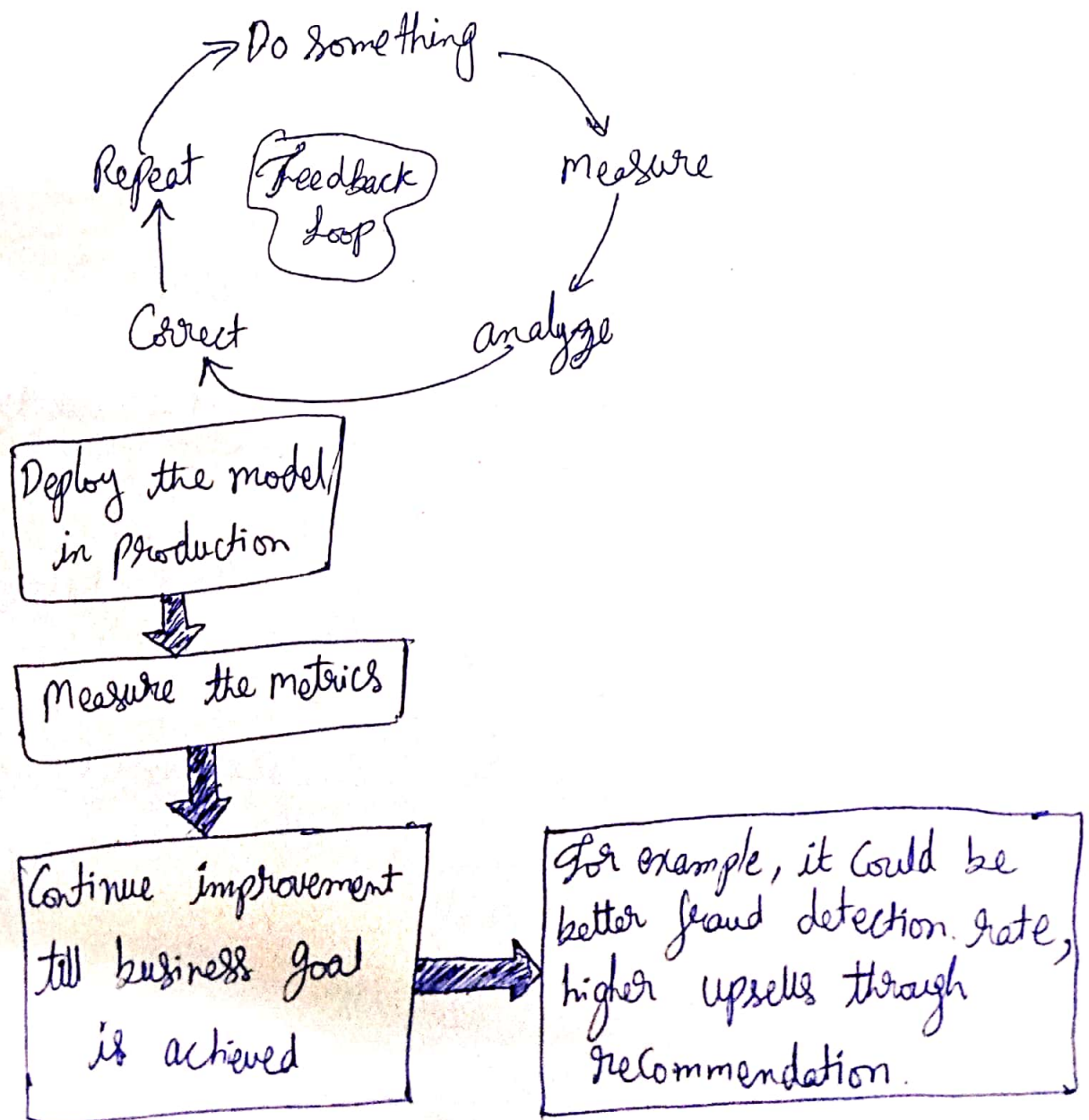
- ① Split data into train and test.
- ② Build multiple models to solve the Business problem.
- ③ Evaluate models using test data and error metrics.
- ④ Choose "Best" model by comparing models to see which one comes closest to answering business objective.

Step 6: Communicate Visualize & Report

- ↳ Brainstorm with management & showcase the benefits, the analysis, and models bring to the plate.

↳ Seek management's Consideration for deploying the solution to real world to help make the business more optimized & beneficial.

Step 7: Take Action & deploy the findings in real world



Types of Machine learning

Supervised Learning enables Computers to learn from labeled data without being explicitly programmed.

Unsupervised Learning draws ~~inferences~~ inferences from data without labeled responses.

<u>Supervised</u>	<u>Unsupervised</u>
<u>Direct feed back</u>	<u>No feedback</u>
<u>Labeled data</u>	<u>No Labeled data</u>
<u>Predictions as output</u>	<u>Find hidden structure in data</u>

Reinforcement Learning :- Learn by interaction with environment.

