

Lesson-3 (Data Import & Transformation)

Data wrangling:-

↳ process of cleaning and transforming data into much more suitable format for ml algorithms

↳ Iterative process

(Trial & error)

① Explore raw data & check general quality

② Transform raw data

③ Validate & publish the data

manage data:-

+ Data stores (hidden)

+ Datasets (specific data files)

Sets of files that contain test, validation, training data

Data store:- Securely connect to my data

→ abstraction of azure data store

→ Compute location independence.

(Several computers can access data stores)

Dataset:- Explore, transform & manage data

Dataset Versioning:-

- * when new data is available
- * when different feature engineering

Feature Engineering:-

↳ Define new features from existing features

Curse of dimensionality:-

Feature Engineering

* data source (Relational database)

* within python environment

* ing streaming data like spark

* during training the model

Feature Engineering tasks:-

* Aggregation (mean, median)

* part - of (part of data)

* Binning (group of entities into bins.
apply aggregation on bins)

* Flagging (define boolean conditions)

* Frequency-based

* Embedding

* Derived by example

most widely-used data

↳ num (Tabular format -

↳ text

↳ img

Text Embedding:-

Text frequency - inverse document frequency

word embedding ^{TF-IDF}

img × RGB needs to be translated

500 × 400 × color depth ⁽³⁾

Grey scale - ⁽¹⁾

simple vector →

multi dimension matrix

feature engineering naturally happens in hidden layers of neural networks

Convolution neural networks - dedicated to learn complex patterns in image data

Feature selection:-

↳ choosing useful features for a model.

① elimination of irrelevant, redundant, highly correlated data

② Reduction of dimensionality

dimensionality reduction.

Commonly used techniques are.

* PCA (principle component analysis)

* t-SNE (t-Distributed ~~stochastic~~ stochastic Neighboring entities)

* Feature Embedding.

→ linear technique (statistical)

→ probabilistic approach (2 or 3 dimensions are resulted)

↳ used for visualization of data (scatter plot)

↳ multi-dimension x 2-dimensionality

→ Encode large no of features to few features

Filter Based Feature selection:-

Permutation Feature Importance.

data drift:- causes degradation in model's performance.

↳ sensor info

↳ sensor breaks (data quality)

↳ data changes over time

(customer behavior not guaranteed)

↳ changes in relations b/w features (degree of correlation)

dataset monitors - detects drift in time.

models to maintain the level of performance in time.

data drift algorithm - provides the change in the data.

baseline dataset - ~~training~~ training dataset

target dataset - input data for the model

model training:

predict a value of a feature - model aim

What is the problem?

classification? (categorical)

Regression? (numerical value)

↳ choice of algorithm

↳ various approaches to get to the result.

* scale / encode data

* split data

- Training dataset
- validation dataset
- Test dataset

selecting hyperparameters - train the model - Evaluate the model

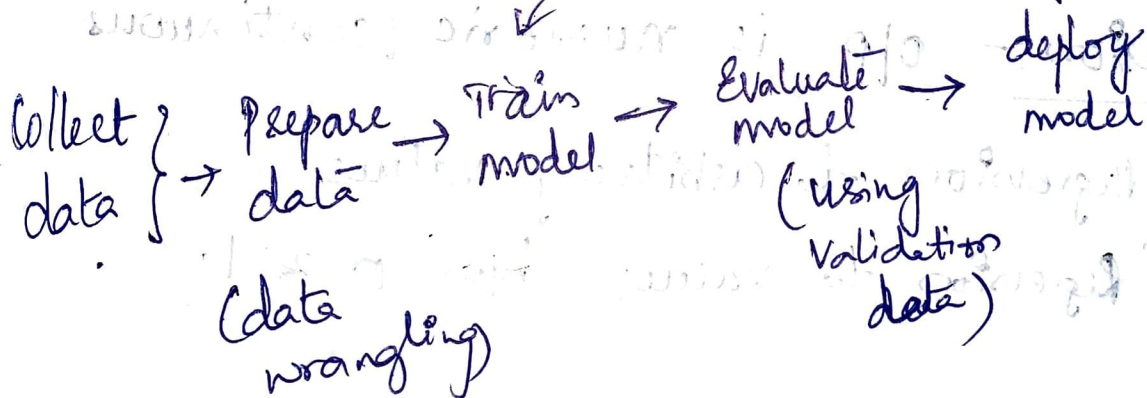
hyperparameters — values are not set before training.

Ex:- no. of layers in deep neural network
no. of clusters in a k-means
learning rate of a model.

splitting data

- Training dataset — used to learn parameters
- Validation dataset — data that checks model's performance
- Test dataset — to finally check the model performance.

→ tune the hyperparameters until it performs better on validation data. ~~test~~ ^{train} model.



Classification:-

Outputs are categorical.

Binary classification

Ex: Fraud detection,
anomaly detection

multi class single label classification

↳ o/p can belong to a single class.

multi class multiple label classification

↳ o/p can belong to multiple classes

Ex:- several tags to texts



Classification algorithms.

↳ logistic regression.

↳ svm (support vector machine)

Regression:- o/p is numeric / continuous

↳ Regression to arbitrary values

↳ Regression to values b/w 0 & 1.

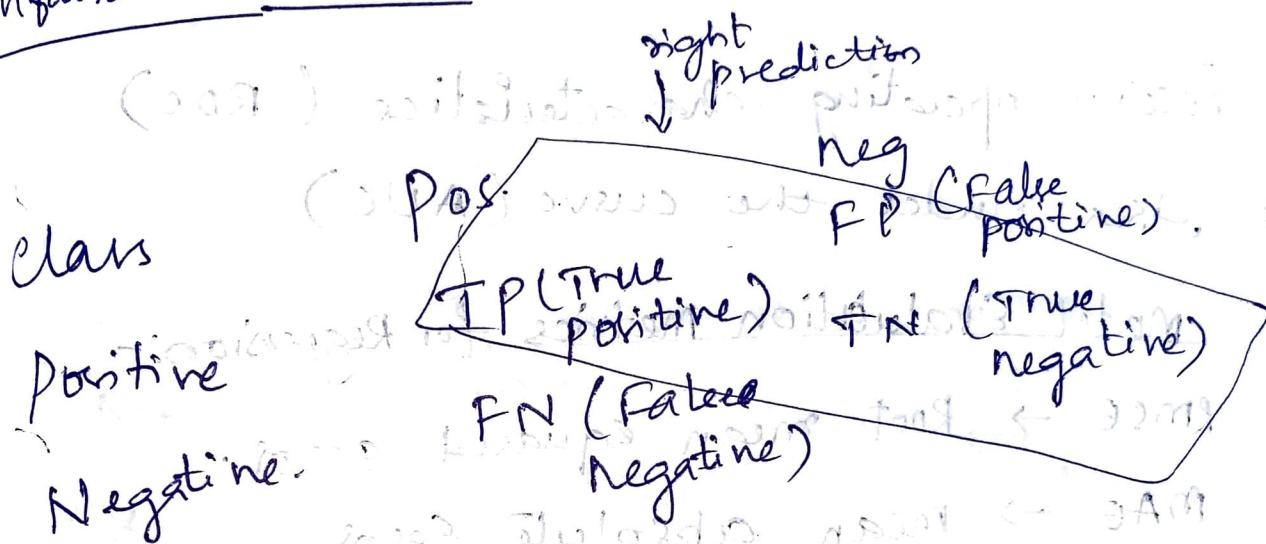
* Linear Regression

* Decision forest Regressor

Evaluating model performance:-

test dataset is a portion of labeled data that is split off and reserve for model evaluation

Confusion matrices:-



$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + FP}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$F_1 \text{ score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

→ Always used in pairs rather than

Receiver operating characteristic (ROC)

Area under the curve (AUC)

Model Evaluation Metrics for Regression:-

RMSE → Root mean squared error.

MAE → mean absolute error.

R-squared → how close the values are to the regression line.

Spearman → strength of the correlation

Histogram of residuals.

less bias.

regressor

↳ decent level of performance



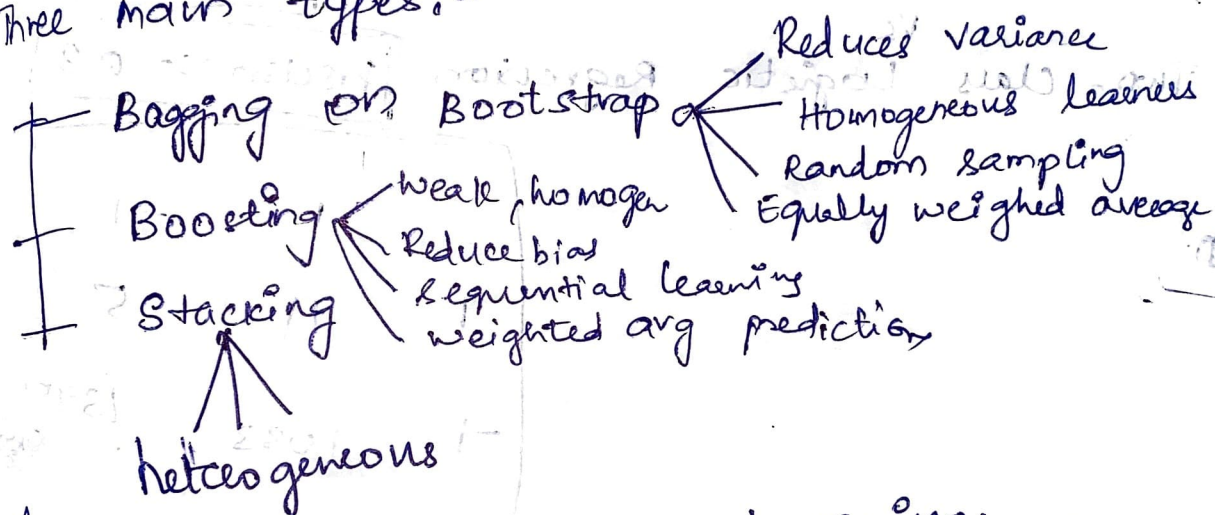
Train and Evaluate a model:-

→ ~~auto~~ Ensemble learning → multiple models
→ Automated machine learning → automate

• train multiple models, and collect the results

Ensemble learning:- combines multiple models to produce one predictive model.

Three main types:-



6-4 Supervised & Unsupervised learning:

Types of classification:-

- + classification on tabular data
- + classification on image (or) sound data
- + classification on text data.