

Machine Learning Roadmap.

Python language (most preferred by companies)

↓
exploratory data analysis (explore the data gathered, like the nature/features in data)

↓
feature engineering (→ (i) exploratory Data analysis

↓ (ii) handling missing value

(iii) Handling outliers

(iv) Categorical encoding

(v) Normalizing & Standardizing

↓
feature selection (selecting the appropriate feature according to model)

↓
ML Algorithms - Regression & Classification, clustering

↓
(i) Correlation

(ii) Forward Elimination

(iii) Backward Elimination

ML

(iv) Uni variate Selection

(v) Random Forest importance

(vi) Feature selection with Decision Tree

Regression → if continuous nature of data
classification → if discrete - 1/-1
clustering → if nature is not known

↓
Machine learning model is
chosen like Linear, Logistic Regression, Decision Tree, Random Forest, Kmeans,

↓
Hyper parameter tuning

Improving the model, by Improving accuracy.

↓
To improve model still more we use Grid search, Randomised search, hyperopt, Genetic Algorithms.

↓
Docker & Kubernetes to store model

↓
Deploy the model

↓
End to end ML projects.

classification of machine learning

- (i) Supervised Learning → future prediction
- (ii) unsupervised learning → classification task.
- (iii) Reinforcement learning → used in games, when played against computer.

Advantages

- (i) to identify trends & patterns
- (ii) human intervention not needed.
- (iii) Handling multi-dimensional & multi-variety data.

Disadvantages

- (i) Large amount of data & time & Resources reqd.
- (ii) high error susceptibility.

Areas of usage

Automatic language transition

Self driving cars

Product Recommendations

Traffic predictions

Speech Recognition

Image Recognition

stock market trading

Online fraud Detection

Virtual personal Assistant

Email spam & malware filtering

Medical Diagnosis

Libraries in python to learn

Numpy → To solve numerical problems.

matplotlib / seaborn → To analyse the data & have it a graphical structure, to see nature of data.

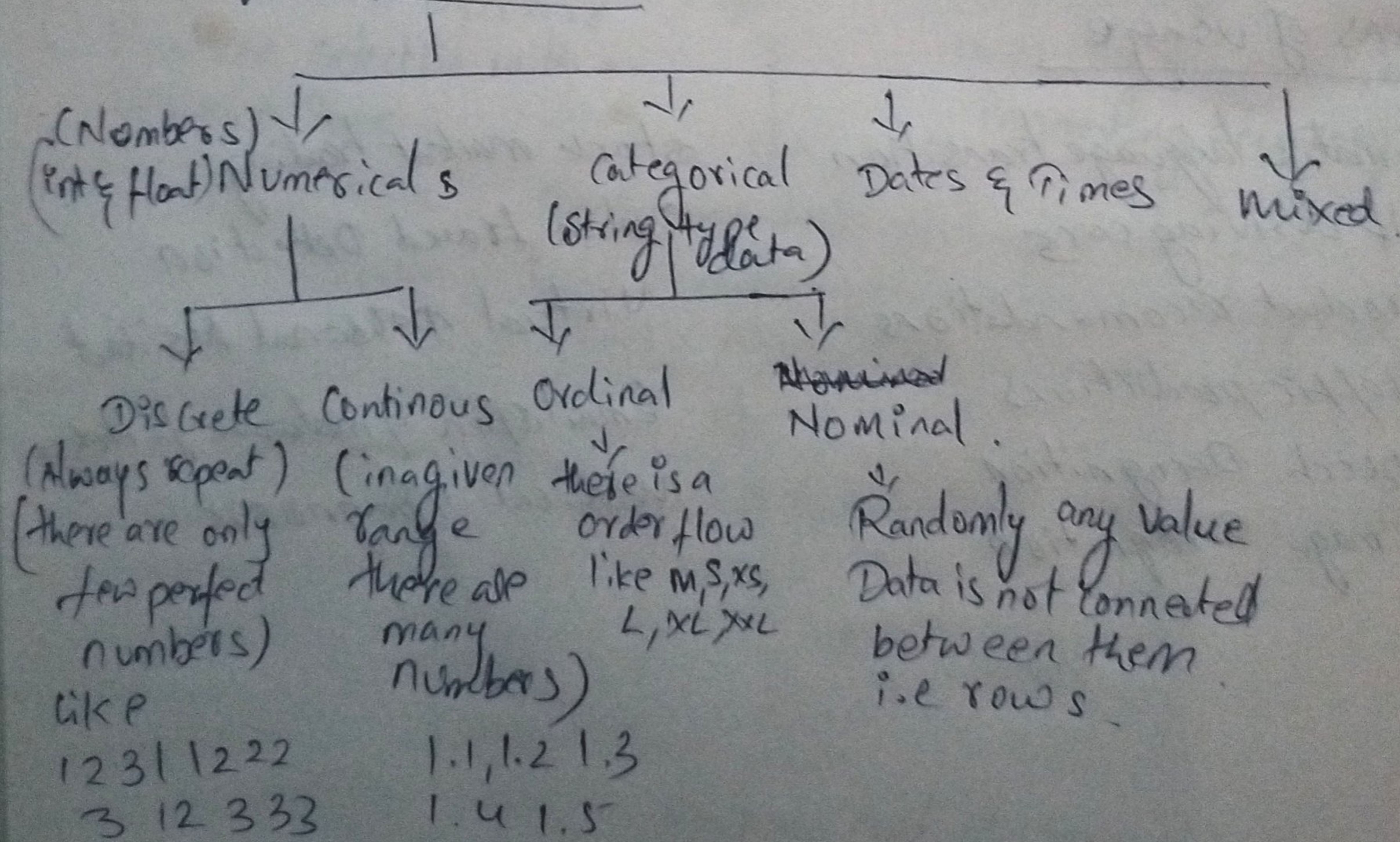
scikit-learn → To build machine learning model

scipy → it provides constants.

TensorFlow → To perform Deep learning

NLTK / opencv → if data is in image or textual format

Data types in Machine Learning types of variables.



Categorical is a object data

Mixed variable is a object data.

Date time variables

like Date of joining or birth, contains date & time.

Mixed variables

it is a mixed mixture of both Numerical & categorical.

Data cleaning

It is the process of preparing data for analysis / ML / DL by removing or modifying data that is incorrect, incomplete, irrelevant, duplicated, or improperly formatted.

Steps in Data cleaning

- (i) Handling Missing Data
- (ii) Outlier Detection & Handling
- (iii) Data Scaling & Transformation
- (iv) Encoding categorical variables
- (v) Handling Duplicates
- (vi) Dealing with inconsistent Data.

What is missing value

- having blank value in the Data set,
- Algorithms don't work on missing value. As it works on math problem. for what $-x2 = ?$, we won't come to know.
- if we want to remove that data or fill any value based on our purpose
- ~~how~~ first identify missing value.
- Then using pandas we would remove.