

U-4 → Predictive Big Data Analytics in Python

- Python

- high-level scripting language.
- true object-oriented language
- Guido van Rossum (developed by)
- 2008
- no need to end with special characters.

`print("Hello World")`
Hello World

- features -

- high level programming
- interactive
- object-oriented scripting
- simple & easy
- portable
- free & open source
- perform complex tasks
- run equally on different platforms
- vast range of libraries

- Adv → . ease

- min time
- modular & object-oriented
- large community of users
- large standard & user-contributed libraries

Disadv → . interpreted & therefore slower than compiled language.
. decentralized with packages.

* Python Libraries -

1) Tensorflow -

- high performance numerical computations. with around 95,000 comments
- across various scientific fields
- defining & running computation that involves tensors, ~~possibly~~ produce a value

Features → better computational graph visualization -

- reduce error by 50 to 60% in neural ML.
- libraries backed by Google.
- quicker updates & frequent new releases.
- appln → speech & image recognition.
- text-based applications.
- time-series analysis
- video detection.

2) Scipy -

- free & open-source Python library for data science -
- high level computations.
- 19,000 comments → Github.
- scientific & technical computations.
- extends NumPy & provides many user-friendly & efficient for scientific calculations.

Features → collection of algorithms

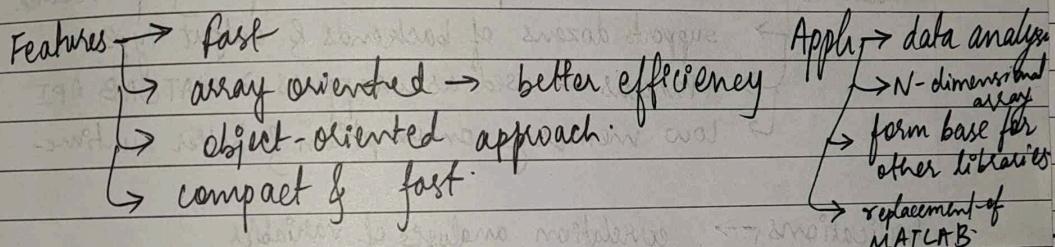
- High level commands
- Data manipulation & visualization
- multidimensional image processing
- Inbuilt-functions → solving differential equations

Applications → multiple dimensional image operations

- diff equations & Fourier
- optimization alg
- Linear algebra

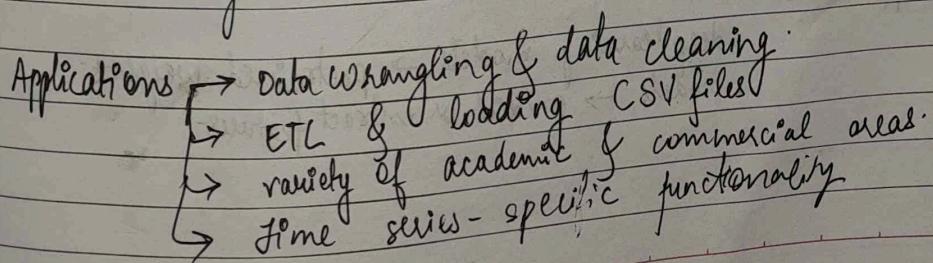
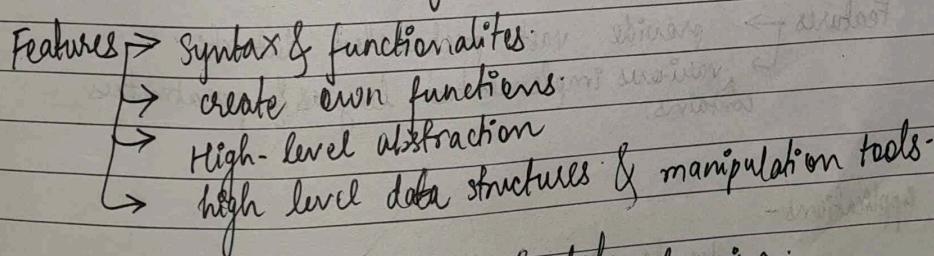
3) Numpy-

- fundamental package for numerical computation.
- n-dimensional array object (powerful)
- 18,000 comments on Github.
- general-purpose array processing package.
- address the slowness problem partly by processing multidimensional arrays.



4) Data Pandas

- python data analysis
- python lib → data science Numpy → matplotlib.
- 17,000 comments → Github.
- fast, flexible, data structures & data frames/EDS.
- structured data easily.



Page No.

Date

5) Matplotlib-

- powerful & beautiful visualizations.
- github community → 26,000
- graphs & plots → extensively used for data visualizations.
- provides an object-oriented API.

Features →

- MATLAB replacement → free & open source.
- supports dozens of backends & output types.
- Pandas → used as wrappers → MATLAB API
- low memory consumption & better runtime.

Applications → correlation analysis of variables

→ 95% confidence intervals of models.

→ outlier detection.

→ visualize distribution of data.

6) Keras.

- used for deep learning & neural network modules.
- supports Tensorflow & Theano backends.

Features → provide vast prelabelled datasets.

→ various implemented layers & parameters

Applications -

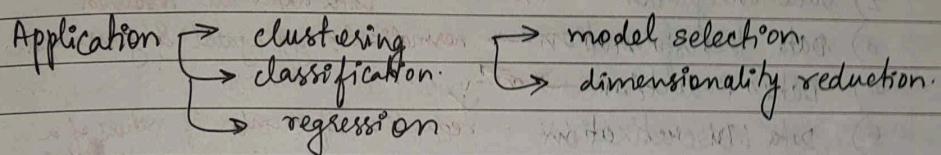
- deep learning models → pretrained weights.
- models → predict or extract features.

Page No.

Date

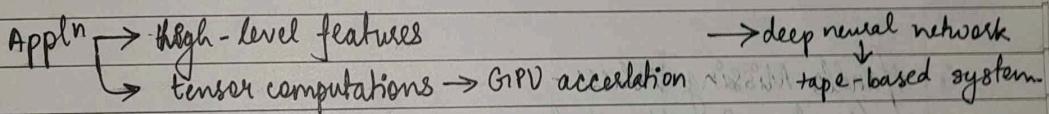
7) Scikit-Learn

- ML algo (provided)
- interpolated into NumPy & SciPy.



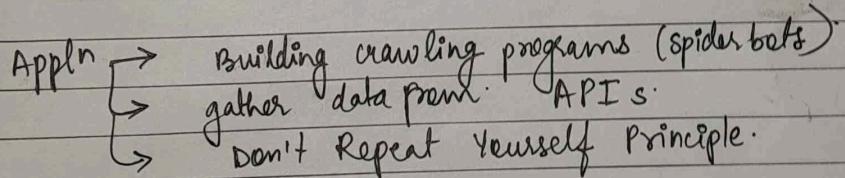
8) PyTorch -

- Python based scientific package → power of graphics.
- deep learning research platform → flexibility & speed.



9) Scrapy -

- Popular, Fast, open-source web crawling framework -
- extract data from web page → X Path -



10) Beautiful Soup -

- Web crawling & data scraping.
- collect data → without proper CSV or API
- help collect & required format (arrange).

* Data Preprocessing -

- 1) Data Cleaning → missing values, noisy data or inconsistency
- 2) Data Integration → diff representation with put together & conflicts are resolved.
- 3) Data Transformation → normalized, aggregated & generalized.
- 4) Data Reduction → reduction of number of values of a continuous attribute.
- 5) Data Discretization → dealing with continuous values - dividing range of attributes.
 - mapping
 - (numerical data)

Handle Missing Values -

- Ignoring the tuple
- Manually filling in missing value.
- Using global constant to fill missing value.

Analytics Type (Data Processing)

Page No.

Date

Predictive

- predict with confidence
 - smarter decisions
 - improve business outcomes
- likelihood different samples.
(find)
- calculates live transactions.
- utilize variety of variable data
- variability → component data
- big basket data
- validate findings.
- monitored → desired results.
- examples → social media analysis, weather, retail, health care & fraud detection.

Descriptive

- gathering, organizing tabulating & depicting data.
- relation b/w product / service.
- model → organize a customer by their personal preferences.
- Business intelligence → sense of communication data → this
- examples → reports that provides historical insights
- data visualization about company
 - ↳ easier communication

Prescriptive

- suggests course of action-
- finding optimal solution to problem.
- what-might-happen Analysis
- determine best course of action
- examples → Traffic applications, Product Optimization & Operational Research

* Association Rule -

- unsupervised learning method.
- descriptive method → discover reln b/w hidden in large dataset.
- used to mining transactions in databases.
- FREQUENTLY BROUGHT TOGETHER.
- goal → discover interesting relationships among items.
- support → popularity of product of all product transactions

$$\text{support}(A) = \frac{\text{No. of transaction in which } A \text{ appears}}{\text{total no. of transactions}}$$

- confidence → likelihood of purchasing both products
A & B.

$$\text{confidence}(A \rightarrow B) = \frac{\text{No. of transaction includes both } A \text{ & } B}{\text{No. of transaction includes only product } A}$$

- Itemset - set of one or more items.

- Frequent Itemset - an itemset whose support is greater than or equal to minsup threshold.



03663 2484523

* Market Basket Analysis -

- determine what products ~~are~~ customers purchase together.
- name → idea of customers throwing all purchases into shopping (market basket) during grocery shopping.
- Association analysis & Frequent itemset mining.
- creates If- Then scenario rules
If item A is purchased then item B is likely to be purchased.
Rule → If $\{A\}$ Then $\{B\}$
- Algorithm → Association Rule & Apriori

Application → • Retail • Telecommunications • Banks
• Insurance • Medical

* Apriori Algorithm -

Improve efficiency → Partitioning
• Reducing Sampling
• Hashing
• Dynamic Counting

- classic ML algorithm.
- designed to work on databases covering transactions.
- aimed to find subsets which are common to at least a min. number.

- Steps →
- 1) computing support for each individual item.
 - 2) Deciding on support threshold
 - 3) selecting frequent items
 - 4) Finding support of frequent itemsets
 - 5) Repeats for larger sets
 - 6) Generation Association Rules & compute confidence
 - 7) compute lift

$$\text{lift} = \frac{P(X \cap Y)}{P(X) * P(Y)}$$

Adv → easy
Join & Prune
(easy on large itemsets)

Disadv → requires high computation
→ entire database has to be scanned

Appln → education • Medical Field • Forestry

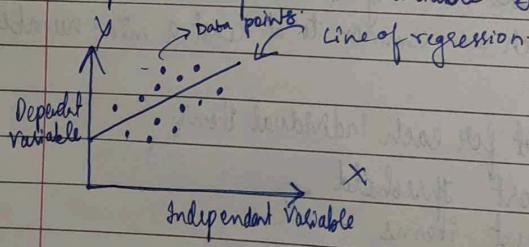
* Regression

- data mining function → predicts a number.
- profit, sale, house values, temperature or distance can be predicted.
- dataset in which target values are known.
- good choice → all predictor variables are continuous values.
- for input x → output continuous → regression problem.

- Dependant Variable
- Independant Variable
- Outliers
- Multicollinearity
- Underfitting & Overfitting

1) Linear Regression -

- ML algorithm → supervised learning.
- predict dependent variable based on independant



Disadv → outliers affect
• over simplifies
real-world
problems.

- regression line is best fit line for model.
- $X \rightarrow$ independent variable
- $Y \rightarrow$ output.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_n X_n + \epsilon$$

2) Simple Linear Regression (SLR) -

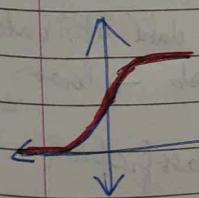
- one independent (or input) variable -
- no. of litres of petrol & kilometers driven. (Univariate Regression).

3) Multiple Linear Regression (MLR) -

- more than one independent variables.
- no. of litres of petrol, age of vehicle, speed & kilometers driven.
- Multivariate Regression.
- multiple inputs & multiple possible outputs

3) Logistic Regression -

- used to model & estimate probability of an event.
- estimate probability of an event occurring \rightarrow value
independent variables
- Is a curve.
- sigmoid curve or S-curve in short
- outcome data points \rightarrow 0 or 1
- 0 \rightarrow totally uncertain & 1 \rightarrow certain
- solve classification issues.
- predictive analytics.
- predict likelihood of categorical dependent variable



- 1) Binomial
- 2) Multinomial
- 3) ordinal

* Naive Bayes.

- uses relation between probabilities of events for classification
- supervised learning algorithm
- solves classification problems
- Based on Bayes Theorem.
- used for Text Classification - (High dimensional training datasets)
- simple, effective, quick predictions.
- probabilistic classifier
- spam filtering, sentiment analysis of classifying articles

* Bayes Theorem -

determine probability of hypothesis with prior knowledge
(conditional probability)

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

likelihood prob
 hypothesis true
 (prob)

Posterior prob
 prob of hypothesis A
 on event B

Prior Probability
 before evidence

Marginal
 Probability → prob. of
 evidence.

→ Adv-

- handle missing data attribute values
- Handle irrelevant input variables
- simple
- fast
- better → categorical variables.

DPsadv -

- assume → data are independent
- not reliable for prob. estimate
- only with categorical variables
- lot of data attributes
prob → less or new

Appln →

- Credit Scoring
- Medical data classification
- sentiment analysis
- spam filtering

Decision Trees:-

- uses concept of trees to structure the given info in sequence of decisions & sequences
- ML algorithms → use a hierarchical structure of decisions to model reln input features & target variables
- tree where node → feature (attribute)
each link (branch) → decision (rule) & leaf represent an outcome (categorical or continuous)
- consists of
 - Nodes → Test for value of certain attribute
 - Edges → corresponds to outcome of test.
 - Leaves → predict outcomes

Steps →

- splitting
- pruning
- Tree Selection

Adv →

- simple & easy
- both nominal & numerical attributes
- handling → errors
- handling → missing values
- non parametric method
- self-explanatory

Disadv →

- target attribute → discrete values
- difficult → solve XOR
- less appropriate → estimation
- prone to overfitting
- classification

Page No.

Date

* More algo's -

- 1) ID3 Algo → Iterative Dichotomiser 3
 - divide into 2 parts, classes or groups
 - entropy → measure of impurity → decision at each
 - selects features with high info → gain to split data & grows tree reaches a stopping criteria.
 - no backtracking
- 2) C4.5 → gain ratio to address bias towards attribute → large no. of val.
 - discrete & continuous attributes & support missing values.
- 3) CART → (Classification & Regression Trees)
 - splitting data based on GIN impurity or min squared error.
- 4) Random Forest → ensemble method → combines multiple trees
 - random subsets of data & of features
 - prediction of individual trees.
- 5) XGBoost → finds optimal splits at each other.
- 6) Gradient boosting → another ensemble learning method that combines multiple decision trees.