

SPPU-TE-COMP-CONTENT - KSKA Git

Python is object oriented, high level programming lang. feature.

- integrated.
- interactive
- object oriented.
- free & open source
- scripting lang.
- simple & easy to learn
- portable
- small sol. for complex prob.

Python libraries.

①. A library is a collection of files. that contains func. for other programs.

②. Reusable chunk of code. to include in program.

①. Numpy. (numeric python)

- scientific computing & basic & advanced array operations.
- this lib offers many hand features. to perform operation on n-arrays & matrices in Python.
- makes performing math operations on arr easier.

② Pandas

- It provides support for data structures & data analysis tools.
- This lib is optimized to perform data science tasks. fast & efficiently.
- Best suited for, structured, labelled data i.e. tabular data.
- Pandas core data structure
 - ① Series
 - ② DataFrame

Series - 1D array like structure.

- holds single column

DataFrame - Represents tabular data.

- Organized in columns. (col has single datatype)

SPPU-TE-COMP-CONTENT - KSKA Git

(3) SciPy.

- contains many different packages & modules to assist in Mathematical & scientific computing
- provides packages like
 - (1) Matplotlib } data visualization.
 - (2) Python

(4) SciKit-learn

- contains a lot of efficient tools for ML & statistical modeling.
- used in dimensionality reduction, classification, clustering, & regression.
- provides many supervised learning algo's.
- cross validations - to check the result of supervised learning.
- provides toy dataset ie. iris, Boston housing.
- used to extract features from images & text.

~~Data~~ P eg. input numpy as np.

```
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
# datasets input make classification.
```

~~X, y = make_classification(100, 2, 2, 1.)~~

model = ^{linear} LogisticRegression()

model.fit(X, y) y_pred = model.predict(x)

plt.scatter(X, y, color='blue', label='Actual data')

plt.plot(X, y_pred, label='AB1')

plt.xlabel("...")

plt.title("...")

plt.show()

X = np.array([1], [2], [3], [4], [5])

Y = np.array([2, 4, 5, 4, 5])

Data Preprocessing.

- Data mining technique to transform raw data to understandable format.
- Reduce data size, find relation b/w data & normalize steps.

①. Data cleaning :- fill missing val, smoothing noise, removing inconsistencies.

②. Data integration :- collect from diff source, conflict resolution.

③. Data Transformation :- the data transformed to required structure; normalized, aggregated & generalize

④. Data Reduction :- Normalize, aggregate & generalize

⑤. Data Discretization :- dividing range of val (10-20) in fixed intervals, & assigning each val to one interval

Removing Duplicates.

- in data preprocessing, duplicates are row that appears more than once. causes biased / incorrect results.

- produce false output
- reduces ML model unless / not accurate output
- increases datasets size unnecessarily.

code. import pandas as pd:

```
data = { 'Name': ['A', 'B', 'C'], 'Age': [10, 20, 30] }
```

```
df = pd.DataFrame(data)
```

```
a = df.drop_duplicates()
```

```
print(a)
```

Handling Missing Values.

- Missing val = absence of val. in a dataset.
- NaN is its representation.
- ML cant work with missing vals.

SPPU-TE-COMP-CONTENT - KSKA Git

ways to handle

- ① ignore tuple
- ② use global constant
- ③ use Mean / median interpolation.

Eg.

```
df['A'].fillna(df['sexA'].mean(), inplace=True)
```

Data Transformation.

- process of changing the format, structure, or values of data to make it suitable for analysis or model training.

Types of Transformations

1. Normalization / standardization.
2. log transformation.
3. Encoding
4. Scaling.

Eg. `Scaler = MinMaxScaler()`

```
df['Normal'] = Scaler.fit_transform(df[['Salary']])
```

can be done by ETL (Extract transform & load) tools.
either on premise / cloud based.

Analytics Types.

- ① Predictive analytics
- ② Descriptive analytics
- ③ prescriptive / Diagnostic analytics.

① Predictive analytics.

- What could happen in future.
- helps organizations in future to predict future.
- calculates live transactions multiple times to help eval the benefit of a customer transaction

- makes use of variety of variable data to make predictions.
- the variability of the component data will have a relationship with what it is likely to predict
- Predictive analytics is covered in following steps.
 - ①. Project definition → what shall be the outcome of the project
 - ② Data collection → various customer interaction at single view
 - ③ Analysis → data processed & utility is found
 - ④ Statistics → validate (assumption, true output)
 - ⑤ modelling → predict multiple models for precise eval
 - ⑥ Deployment →
 - ⑦ monitoring

Eg:- weather analysis, Healthcare, fraud detection

- ② Descriptive "What has happened"
- looks at past data to understand what has already occurred.
 - helps in identifying trends patterns & summaries
 - used in first phase of analytics. involving gathering & organizing data.
 - Shows relationship b/w product with data
 - to organize customer by their preferences
 - BI fails to make the data communicate to people & hence. Descriptive analytics is used.

Purpose:- Summarize historical data
Eg. Monthly salary report

- ③ Prescriptive / Diagnostic Analytics:- "Why did it happen"
- investigate cause behind past outcomes
 - helps to identify reason or cause of trend or any problem
 - assist user to find optimal solution to make right choice

SPPU-TE-COMP-CONTENT - KSKA Git

- under utilizes the understanding what, why questions. to help user determine best course of action to take
eg:- Sales drop in specific region

Market basket analysis

technique that allows us to discover the relationship between products.

makes use of

Association Rule.

- ① Apriori algorithm
- ② FP growth algorithm.

⇒ Association Rule

- it is useful for discovering interesting relationship hidden in large datasets.
- the uncovered relationships can be represented in the form of association rules or sets of frequent items.
- Association rules are if then statements that helps uncover relationship between seemingly unrelated data.
- eg. "if a customer buys bread, then he is 80% likely to also purchase milk."
- Association rule mining has 2 steps.
 - ①. find all frequent itemsets: (use min support count)
 - ②. generate strong association rule from those itemsets.by definition these rules should satisfy min support & confidence

Support :- how frequently the total collection of items. {A, B} occur together. as a percentage of all transaction

$$\text{Support} = \frac{A + B}{\text{Total}}$$

Confidence :- The ratio of no. of transaction including all items in B & A to the no. of transaction that include all items in A

$$\Rightarrow \text{confidence} = \frac{A + B}{A}$$

SPPU-TE-COMP-CONTENT - KSKA Git

lift: implies the type of relation those objects have

- ie. no relationship
- +ve relationship
- ve relationship

if lift val = 1 \Rightarrow no relationship

lift val $> 1 \Rightarrow$ +ve relationship

lift val $< 1 \Rightarrow$ -ve relationship

Given a table eg

transaction item id	list of buy
A1	I1, I2, I5
A2	I3, I2

Table 1

row fi of each item.

Apriori Algo.

①. find frequency of occurrence of every item

eg. Item	frequency	Support
I1	6	$\frac{6}{\text{total}} \times 100$
I2	8	$\frac{8}{\text{total}} \times 100$
I3	5	
I4	3	
I5	3	

table 1 has 9 rows in this eg.

②. use minimum support val ie 4 to reduce the itemset.
so. anything less than 4 removed.

Item	fi	Support
I1	6	$\frac{6}{9} \times 100$
I2	8	
I3	5	

③. make combinations. & find their frequency.

frequent itemset create rules

I1 \Rightarrow I2
I1 \Rightarrow I3
I2 \Rightarrow I3

Item	fi	Support
{I1, I2}	5	$\frac{5}{9} \times 100$
{I1, I3}	4	$\frac{4}{9} \times 100$
{I2, I3}	4	

I1 & I2 comes in the. Table 1 5 times (together)
like in A1 {I1, I2, I5}

④. make combinations again

{I1, I2, I3}	3
--------------	---

do this till. combination cant be done or till ϕ (empty set)

SPPU-TE-COMP-CONTENT - KSKA Git

⑤. calculate support confidence lift

$$\text{support} = \frac{f_i(A+B)}{\text{total}} \quad \text{confidence} = \frac{f_i(A+B)}{f_i(A)} \quad \text{lift} = \frac{\text{support}(X+Y)}{\text{support}(X) * \text{support}(Y)}$$

$$\text{Support}(J_1, J_2) = \frac{f_i(J_1, J_2)}{\text{total}} \Rightarrow \frac{5}{9} \leftarrow \text{total no. of rows in table}$$

$$= 0.55$$

do same for (J_1, J_3) & (J_2, J_3)

confidence $(J_1, J_2) = \frac{f_i(J_1, J_2)}{f_i(A)} \Rightarrow \frac{5}{6} \leftarrow J_1$'s f_i from step 2

OR $= \frac{\text{Support}(J_1, J_2)}{\text{Support}(J_1)} = 0.83 \leftarrow [\times 100 \text{ for } \%]$

$$\text{lift}(J_1, J_2) = \frac{0.55}{6/9 * 8/9} \leftarrow \text{support}(x+y)$$

$$\frac{f_i(x)}{\text{total}} = 0.94 \Rightarrow \text{-ve relationship}$$

consider 50% confidence if not given.

~~As~~ The apriori algo is fundamental algo used for association rule mining in data mining. It discovers frequent itemset. & derives association rules.

Steps of algo. ①. Create candidate items.

②. Prune infrequent itemsets.

③. Repeat until no new frequent itemsets.

④. Generate association rule.

FP growth - Improvement to apriori

- use min support to find frequent data.

- used to find frequent itemsets in a transaction database without candidate generation.

- in tree format.

SPPU-TE-COMP-CONTENT - KSKA Git

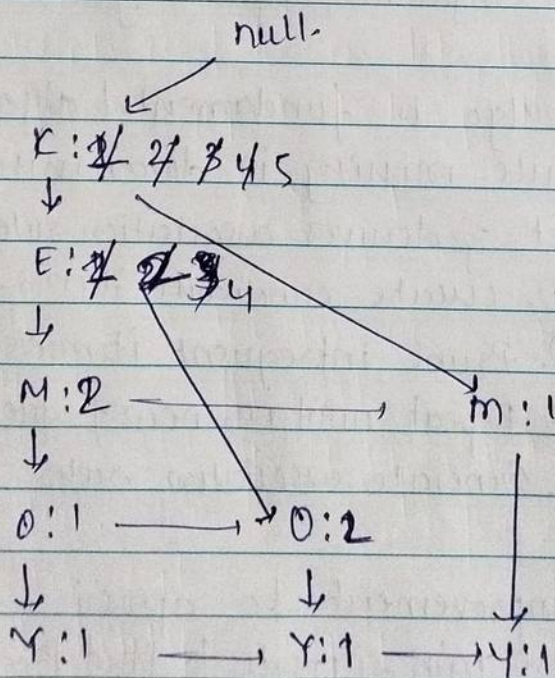
eg.

TID	Item	A	1	X
T ₁	{E K M N O Y}	C	2	X
T ₂	D E I C N O Y	D	1	X
T ₃	A E K M	E	4	
T ₄	C K M O Y	K	5	
T ₅	C E I K O O	M	3	
		N	2	X
		O	4	
		Y	1	X
		Y	3	

min support = 3

arrange in decreasing order of f_i
 K:5 E:4 O:4 M:3 Y:3

TID	Items.	ordered itemset.
T ₁	E K M N O Y	K E M O Y
T ₂	D E I C N O Y	K E O Y
T ₃	A E K M	K E M O
T ₄	C K M O Y	K M Y
T ₅	C E I K O O	K E O O



to reach path.

val	1	1	1
Y	K E M O, K K O, K M		= 3
O	K E M I, K E 2		= 3
M	K E = 2, K = 1		= 3
E	K = 4		4
K			

SPPU-TE-COMP-CONTENT - KSKA Git

Regression.

- data mining func that predicts a number.
- supervised learning technique used to predict continuous numerical vals.
- to find relationship b/w independent var(x) & dependent variables(y).

types.

- ① Linear regression
- ② Logistic regression.

① Linear Regression

- a supervised learning algo used to predict continuous vals.
- models relationship b/w independent & dependent var.

Given by.

$$y = B_0 + B_1 x \quad \text{--- (1)}$$

where $B_1 \rightarrow$ slope of line

$B_0 \rightarrow y$ intercept.

$$B_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{(x - \bar{x})^2}$$

$$B_0 = \bar{y} - B_1 \bar{x}$$

put val of x in eq 1. to get new val of y

table to solve/get B_1 & B_0 . eg.

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(x - \bar{x})(y - \bar{y})$
4	3	1	1	1	1
0	4	-3	2	9	-6
3	0	0	-2	0	0
5	1	2	-1	4	-2
<u>12</u>	<u>10</u>				<u>-7</u>
$\bar{x} = \frac{12}{4} = 3$	$\bar{y} = \frac{10}{4} = 2.5$				

SPPU-TE-COMP-CONTENT - KSKA Git

- In linear reg. predicted val. are the mean of target var. at the given val. of the input var.
- used to solve regression problems.
 - graph is a straight line.

② Logistic Regression.

- ① - form of regression analysis in which the outcome is binary
- ② - Supervised learning algo used for binary classification
- ③ - it predicts the possibility of class using a sigmoid function.

~~③~~

$$P(Y=1|X) = \frac{1}{1 + e^{-(mx+c)}}$$

eg. spam detection.

- ④ - The core component of logistic regression is logistic func. aka Sigmoid func. which is used to model the probability that a given input belongs to a particular class.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

where z is the linear combination of input features & their coefficients.

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots \beta_n x_n$$

this func. takes a val for z & gives out as 0/1

Why to use it

① Probability Output.

- logistic func. maps real val input to 0 or 1 which can be interpreted as probability

SPPU-TE-COMP-CONTENT - KSKA Git

②. Non linear transformation.

logistic regression uses logistic func. to squash value to a valid probability range.

③. Decision Boundary

By applying the threshold (commonly 0.5) logistic regression classifies the output in 0 & 1 classes.

④. Interpretability

- helps in understanding confidence level of classification

⑤ Mathematical convenience

The logistic func is smooth & differentiable

Role of sigmoid

- converts linear output (z) into a probability.
- if $O_p > 0.5 \Rightarrow$ class 1 (+ve class)
- if $O_p < 0.5 \Rightarrow$ class 0 (-ve class).
- helps in binary classification using threshold based decision.

Linear. Reg.

Logistic. Reg.

Output

Continuous value.

Binary output.

func used

$$y = mx + c$$

$$\sigma = \frac{1}{1 + e^{-z}}$$

problem type

Regression type

Classification type.

O/p range

$-\infty$ to ∞

0 to 1

curve

straight line

S shaped curve / non linear.

#

Classification.

it predicts categorical labels (Classes),
prediction models continuous valued func.
= considered to be supervised learning

2 types

(1) Naive Bayes

(2). Decision tree

(1) Naive Bayes.

- Supervised learning algo based on bayes theorem.
- with a strong assumption of feature independence
- used to classify tasks

$$P(A | x_1 \cap x_2 \cap x_3 \cap x_4) =$$

$$\frac{P(x_1|A) * P(x_2|A) * P(x_3|A) * P(x_4|A) * P(A)}{P(x_1) * P(x_2) * P(x_3) * P(x_4)}$$

$$\text{Bayes theorem} = P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

$P(A|B) \rightarrow$ posterior probability

$P(B|A) \rightarrow$ likelihood.

$P(A) \rightarrow$ prior. prob of A

$P(B) \rightarrow$ ——— B

working.

1. calculate prior probability of each class.
2. calculate likelihood of each feature
3. Apply Bayes theorem. to complete posterior for each class
4. assign the class with higher. posterior probability.

SPPU-TE-COMP-CONTENT - KSKA Git

Adv.

Simple & fast.

works well with high dimensional data.

Requires less training data.

Application.

Spam detection.

Text classification

Sentiment analysis

Eg.

Email	offer	free	spam	
1	1	0	No	
2	0	1	yes	
3	1	1	yes	o✓ s✓
4	0	1	No	oX sX
5	1	1	yes	o✓ s✓

to find. offer=1 & free=1 is spam or not

$$p(\text{spam is yes}) = 3/5$$

$$p(\text{spam is no}) = 2/5$$

$\frac{2}{3} \rightarrow 0.5$
 $\frac{2}{3} \leftarrow$ spam yes for only 3.

	yes	No.
offer	$\frac{2}{3}$	$\frac{1}{2}$
free	$\frac{3}{3}$	$\frac{1}{2}$

$$p(\text{yes} | \text{offer, free}) = \frac{p(\text{offer, free} | \text{yes}) * p(\text{yes})}{p(\text{offer}) * p(\text{free})}$$

$$p(\text{No} | \text{offer, free}) = \frac{p(\text{offer, free} | \text{no}) * p(\text{no})}{p(\text{offer}) * p(\text{free})}$$

no. need to solve denominators as it is same for both. yes & no.

SPPU-TE-COMP-CONTENT – KSKA Git

$$p(\text{yes} | \text{offer, full}) = \frac{2}{3} \times \frac{2}{3} \times \frac{3}{5} = \frac{2}{5} = 0.4$$

$$p(\text{No offer, full}) = \frac{1}{2} \times \frac{1}{2} \times \frac{2}{5} = \frac{1}{10} = 0.1$$

Since the Denominators are same we can say.

• Hence, as $p(\text{yes} | \text{offer, false}) > p(\text{No} | \text{offer, false})$

\therefore the new msg is spam.

十一

Decision Tree

- used to make decisions.
- build by choosing the best attribute to split the data at each step, based on information gain, which is calculated using entropy.

Steps. yes / No \rightarrow from last col.

①. Calculate the entropy of entire dataset.

values {+ve for yes, -ve for no}.

eg val $\{2+9, -5\} = \frac{-9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$

②. select random attribute, ie whether it has 3 types.

calculate entropy for them. & do the same for other attributes.

Sunny. $\left\{ \begin{matrix} 2 \text{ yes} & 3 \text{ No} \\ +2 & -3 \end{matrix} \right\} \quad \frac{-2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$

cloudy $\{+4, -0\}$ $-\frac{4}{2} \log_2 \frac{0}{2}$ will lead to $0 = 0$

Rainy $\frac{2}{5} + 3, -24$ $-\frac{2}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$

③. Calculate info gain for all attributes. i.e. weather, Temp, humidity & wind.

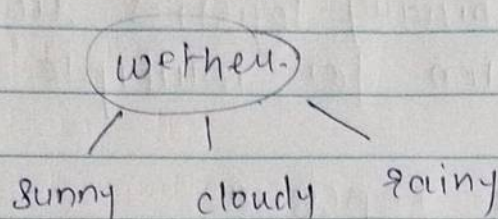
for weather

eq. info gain =

chance is, sunny
cloudy & rainy.

EC whole dataset) = $\frac{\text{no. of Row's (Sunny)} \times E(\text{Sunny})}{\text{total rows in table}} - \frac{\text{no. of Row(Cloud)} \times E(\text{Cloud})}{\text{total rows in table}} \leftarrow \text{Same}$

- ④. select the highest value of info gain
- ⑤. set that node as root node.
- ⑥. select the sub val from the root node i.e.



select sunny.

- eg prepare a table with all sunny values.
- eg find entropy of the entire table
- eg info gain for sunny.

- ⑦ repeat this for all the attributes. manually.

entropy

$$H(S) = - \sum p_i \log_2(p_i)$$

p_i = probability of class i

Info gain

$$IG(S, A) = H(S) - \sum \left(\frac{|S_v|}{|S|} \cdot H(S_v) \right)$$

Types of logistic Regression.

- ①. Binary logistic Regression.
- ②. Multinomial logistic Regression.
- ③. Ordinal logistic Regression.

①. Binary logistic Regression.

used when target variable has 2 categories.

eg. Pass /s fail & yes/no

Output :- probability b/w 0/1 classified using threshold.

②. Multinomial logistic regression
used when the target variable has more than two unordered categories.
eg. fruit (Apple, Banana, orange)
- it extends binary logistic regression using a Softmax function for multiple classes.

③ Ordinal logistic Regression.
- used when the target variable has more than 2 ordered categories.
- eg:- poor, avg, good, excellent -
- it considers the order/rank of categories.