

Algorithm:

- 1) Import libraries :
import pandas for data handling
import variance & stdev functions from statistic model
- 2) Read data:
Read data from a csv file into a pandas dataframe ('df')
create new dataframe ('data') with selected column
- 3) Calculate mean, median, mode :
calculate mean, median & mode using function .mean()
.median() & .mode() & print them
- 4) Range calculation:
calculate Range of column using formula (max-min)
- 5) Variance & standard deviation:
calculate variance & standard deviation
- 6) Skewness calculation:
calculate skewness using formula $3 * (\text{mean} - \text{median})$
print result.

Algorithm :

- 1) Import libraries :
import pandas for data handling
import the math file for mathematical calculations.
- 2) Read Data :
Read csv file into pandas dataframe.
- 3) calculate means :
calculate mean for 'trip-distance' (\bar{x}) & total distance (\bar{y})
- 4) calculate deviations :
calculate deviation $x - \bar{x}$ & $y - \bar{y}$.
- 5) calculate square & product :
calculate product & square of deviations $(x - \bar{x})^2$, $(y - \bar{y})^2$, $(x - \bar{x})(y - \bar{y})$.
- 6) calculate correlation coefficient r :
$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$
- 7) print pearsons coefficient of correlation r :
print value of r .

for grouped data

- 1) import libraries
- 2) Create dataframe :
create dataframe with two grouped column "pass-count" & "trip-distance" & frequency.
- 3) Use .corr() function to calculate correlation between data
- 4) print r value.

Algorithm :

1) Import libraries

Import pandas for data handling

Import KNN Imputer from sklearn for K Nearest Neighbors Imputation

2) Read Data

Read uber data csv file into pandas Dataframe ('df')

3) Fill NaN values with constant

create new data frame df2 by filling NaN value with a constant value of 1.

4) Fill NaN values with previous value

fill the null value with previous values using method = pad

5) Fill NaN with Next value:

fill null values with next non-NaN value using method = bfill

6) Fill NaN with Mean, Median & Mode:

fill null value with "passenger-count" column mean, median & mode

7) Fill NaN with Max & ~~mean~~ Min value

fill null values with .max() & .min() function on passenger-count column

8) Fill NaN value with frequency occuring value

fill NaN or null value using mode function & taking max value

9) Fill NaN using KNN Imputation

use KNN imputer with 2 Neighbors to impute NaN values in "passenger-count".



Date: _____

formula for calculating T-score.:

$$T_{\text{score}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

where,

$\bar{x} \Rightarrow$ mean of sample data

$\mu \Rightarrow$ mean of population

$\sigma \Rightarrow$ standard deviation

$n \Rightarrow$ The sample size.

Algorithm:

Steps in conducting a T-Test :-

1) Formulate hypothesis :

H_0 : There is no significant difference

H_1 : There is a significant difference.

2) Select significance level (α):

select α (common choices are 0.05, 0.01)

3) calculate T-score :

calculate t score by putting values in formula

4) Determine critical region :

calculate critical value using of α alpha using t-table

5) Make a decision:

If calculated statistics falls in critical region reject the null hypothesis.

6) Draw conclusion.:

Based on result, make inferences about the population



Date: _____

outlier, & more specifically they are considered global outlier.

Primarily KNN rely on distance based outlier detection. In KNN the number of nearest neighbors K is a parameter that need to determine. There are other variants of the KNN algorithm supported by pyOD.

unlike supervised learning approach unsupervised anomaly detection does not require labeled data indicating which instances are anomalies & which are not instead it identify outliers solely based on distribution of data points in the feature space. An appropriate threshold is chosen to differentiate between normal data points & anomalies while those below it are considered normal.

overall unsupervised KNN for anomaly detection provides a flexible & ~~fast~~ interpretable approach for identifying anomalies.

Algorithm:

① Load data & libraries:

- Import numpy pandas & matplotlib.
- Read given csv file into a dataframe
- Extract top-distance as numpy array

② Define function:

- create function to calculate distance to nearest neighbor to each point
- Identify those with mean distance to K nearest neighbor exceeding the threshold as anomalies.



Date: _____

③ Detect anomalies:

call function $\&$ with top_distance , chosen n $\&$ threshold

④ Separate anomalies:

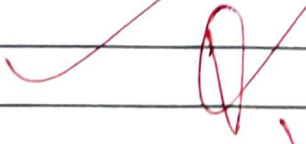
Separate the points in top_distance into anomalies $\&$ normal points.

⑤ Plot Result:

visualize anomalies $\&$ normal points.

Conclusion:

distance base method offer a parameter driven approach for anomaly detection in data. Unsupervised FNN effectively identifies anomalies in distance based data by leveraging the concept of nearest neighbor. Overall FNN provide a valuable unsupervised tool for anomaly detection in various application.


Jul 3/24



SMOTE is careful to create synthetic examples along lines connecting existing minority class instance, ensuring that the synthetic data remain within boundaries of the existing distribution. This prevents introduction of noise & outliers that might adversely affect model performance.

SMOTE is overall valuable tool for addressing class imbalance issue & improving the performance of model on imbalanced dataset.

Algorithm:

Step 1: Import necessary libraries.

Step 2: Read data from 'BSV file' into a dataframe called Data.

Step 3: Create subplot with two panels.

Plot the original data in the first panel.

Step 4: Separate features (X) & target variable (Y) from 'Dataframe'.

Step 5: Initialize SMOTE object.

Step 6: Apply SMOTE oversampling using fit-sample method on X & Y.

Generate synthetic samples for the minority class to balance the dataset.

Step 7: Concatenate resampled features & target variables into a new Dataframe 'resampled_data'.

Step 8: Plot resample data in the second panel.

Step 9: Display plot.



Date: _____

- Integrated (I): Represents the differencing of raw observations to allow the time series to become stationary
- Moving average (MA): Incorporate dependency between an observation & a residual error from a moving average model applied to lagged observations.

ARIMA Parameters:

Each component in ARIMA functions as a parameter with a standard notation. For ARIMA models, a standard notation would be ARIMA with p, d & q where integer values substitute for the parameters to indicate the type of model used. The parameters can be defined as,

- p : The number of lag observations in the model also known as the lag order.
- d : The number of times the raw observations are differenced also known as the differencing
- q : The size of moving average window also known as the order of moving average

Algorithm:

Step 1:

import necessary libraries such as pandas, numpy, matplotlib.pyplot & ARIMA model from statsmodel

Step 2:

Load the time series data from a csv file into pandas dataframe.



Date: _____

step 3 :

visualize the time series data using matplotlib.

step 4:

perform augmented Dickey-Fuller test (ADF) test to check the stationarity of time series data

step 5:

convert the data column to date time format & set it as the index of dataframe

step 6:

fit ARIMA model to the time series data with specific order (p,d,q)

step 7: ~~to~~

generate forecasts for future time periods using the fitted ARIMA model.

Conclusion:

In conclusion, this experiment utilized ARIMA modelling to analyze & forecast time series dataset representing total amount over time. This experiment demonstrate ARIMA modelling for time series analysis & forecasting

10/4/24.