Academic Year 2023-24 SAP ID: 60004220126



SHRI VILEPARLE KELAVANI MANDAL'S DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING



(Autonomous College Affiliated to the University of Mumbai)
NAAC ACCREDITED with "A" GRADE (CGPA: 3.18)

COURSE NAME: Machine Learning CLASS: Third Year BTech

NAME: Shashwat Shah

BATCH: C22

EXPERIMENT NO. 1

AIM / OBJECTIVE:

To perform data preprocessing in terms of handling, missing data, removing outliers, eliminating duplicate rows and modifying the datatype, etc.

DESCRIPTION OF EXPERIMENT:

Python is an easy-to-learn programming language, which makes it the most preferred choice for beginners in Data Science, Data Analytics, and Machine Learning. It also has a great community of online learners and excellent data-centric libraries. With so much data being generated, it becomes important that the data we use for Data Science applications like Machine Learning and Predictive Modeling is clean. But what do we mean by clean data? And what makes data dirty in the first place? Dirty data simply means data that is erroneous. Duplicacy of records, incomplete or outdated data, and improper parsing can make data dirty. This data needs to be cleaned. Data cleaning (or data cleansing) refers to the process of "cleaning" this dirty data, by identifying errors in the data and then rectifying them. Data cleaning is an important step in and Machine Learning project, and we will cover some basic data cleaning techniques (in Python).

Cleaning Data in Python

We will now separate the numeric columns from the categorical columns.

Missing values

We will start by calculating the percentage of values missing in each column, and then storing this information in a DataFrame.

Drop observations

One way could be to drop those observations that contain any null value in them for any of the columns. This will work when the percentage of missing values in each column is very less.

Remove columns (features)

Another way to tackle missing values in a dataset would be to drop those columns or features that have a significant percentage of values missing.

Impute missing values

Academic Year 2023-24 SAP ID: 60004220126



SHRI VILEPARLE KELAVANI MANDAL'S DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING



(Autonomous College Affiliated to the University of Mumbai)
NAAC ACCREDITED with "A" GRADE (CGPA: 3,18)

There is still missing data left in our dataset. We will now impute the missing values in each numerical column with the median value of that column.

Outliers

An outlier is an unusual observation that lies away from the majority of the data. Outliers can affect the performance of a Machine Learning model significantly.

Duplicate records

Data can sometimes contain duplicate values. It is important to remove duplicate records from your dataset before you proceed with any Machine Learning project. In our data, since the ID column is a unique identifier, we will drop duplicate records by considering all but the ID column.

Fixing data type

Often in the dataset, values are not stored in the correct data type. This can create a problem in later stages, and we may not get the desired output or may get errors while execution.

PROCEDURE:

Describe the procedure that is used to carry out the experiment step-by-step. Describe every line of code with the proper interpretation of the output.

Perform data preprocessing with respect to your case study and discuss results of all the steps.

Code and output:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

df = pd.read_csv('/content/dirtydata.csv')
print(df.head(10))
```

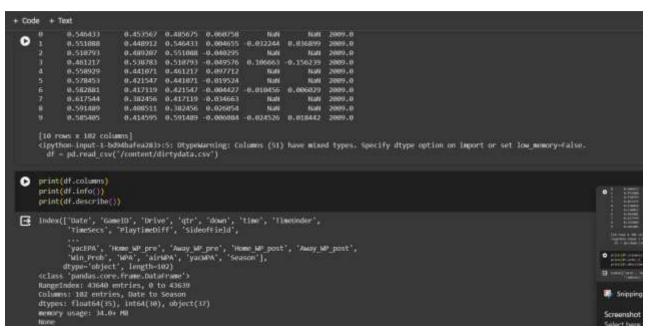


SHRI VILEPARLE KELAVANI MANDAL'S DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING



(Autonomous College Affiliated to the University of Mumbai)
NAAC ACCREDITED with "A" GRADE (CGPA: 3.18)

0	Date	GameID		qtr				er TimeSed	
0		2009091000	1					15 3600.	0
1	2009-09-10	2009091000	1			14:5 3		15 3593.	0
2	2009-09-10	2009091000	1		2.0	14:16		15 3556.	0
3	2009-09-10	2009091000	1		3.0	13:35		14 3515.	
4	2009-09-10	2009091000	1		4.0	13:27		14 3507.	
5	2009-09-10	2009091000	2		1.0	13:16		14 3496.	0
6	2009-09-10	2009091000	2	1	2.0	12:40	,	13 3460.	0
7	2009-09-10	2009091000	2		3.0	12:11		13 3431.	0
8	2009-09-10	2009091000	2	1	4.0	11:34		12 3394 .	0
9	2009-09-10	2009091000	3	1	1.0	11:24		12 33 84 .	0
	PlayTimeDif ⁻	f SideofFiel	d					ay_WP_pre	\
0	0.0				NaN	0.4	85675	0.514325	
1	7.0	ð PI	т	1.146	5076	0.5	46433	0.453567	
2	37.0	9 PI	т		NaN	0.5	51088	0.448912	
3	41.0	ð PI	т	-5.031	L425	0.5	10793	0.489207	
4	8.0	9 PI	т		NaN	0.4	61217	0.538783	
5	11.0	ð TE	Ν		NaN	0.5	58929	0.441071	
6	36.0	ð TE	Ν	0.16	3935	0.5	78453	0.421547	
7	29.0	ð TE	Ν		NaN	0.5	82881	0.417119	
8	37.0) TE	Ν		NaN	0.6	17544	0.382456	
9	10.0	ð TE	Ν	0.541	L602	0.5	91489	0.408511	
		t Away_WP_p		_					
0	0. 54643	3 0.4 53		.485675			NaN		2009.0
1	0.55108	8 0.448	912 0	.546433	0.6	004655	-0.032244	0.036899	2009.0
2	0.51079	3 0.489	207 0	.551088	-0.6	940295	NaN	NaN	2009.0
3	0.46121	7 0.538	783 0	.510793	-0.6	949576	0.106663	-0.156239	2009.0
4	0.558929	9 0.441	071 0	.461217	7 0.6	97712	NaN	NaN	2009.0
5	0.57845	0.421	547 0	.441071	L -0.6	19524	NaN	NaN	2009.0
6	0.58288	1 0.417	119 0	.421547	7 -0.6	004427	-0.010456	0.006029	2009.0
· · · · · · · · · · · · · · · · · · ·		· · · · · · · · · · · · · · · · · · ·	· · · · · · · · · · · · · · · · · · ·	· ·		· · · · · · · · · · · · · · · · · · ·	· · · · · · · · · · · · · · · · · · ·	· · · · · · · · · · · · · · · · · · ·	





SHRI VILEPARLE KELAVANI MANDAL'S DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING



(Autonomous College Affiliated to the University of Mumbai)

NAAC ACCREDITED with "A" GRADE (CGPA: 3.18)

```
43640.000000
      4.364000e+04
                                                  36950.00
       2.009144e+09
                        12.335060
                                       2.560151
                                                      2.015264
                                                                    7.334166
mean
std
      1.7781140+05
                         7.106666
                                       1.126115
                                                      1.009782
                                                                    4.659195
                         1.000000
      2.009091e+09
                                       1.000000
                                                      1.000000
                                                                    0.000000
min
25%
      2.0091010+09
                        6.000000
                                       2.000000
                                                      1.000000
                                                                    3,000000
58%
       2.009111e+09
                        12.000000
                                       2,000000
                                                      2.888888
                                                                    7.000000
75%
       2.009121e+09
                        18.000000
                                       4.000000
                                                      3.000000
                                                                   11.000000
max
       2.010010e+09
                        32.000000
                                        5.000000
                                                      4.000000
                                                                   15.000000
           TimeSecs PlayTimeDiff
                                          yrdln
                                                    yrdline100
                                                                     ydstogo
count 43607,000000 43574,000000 43551,000000 43551,000000 43640,000000
                        20.799904
                                      28.381323
                                                     47,692292
                                                                    7.196471
mean
                                     13.129457
       1057.388873
                        16.910683
                                                    25.187992
                                                                    4.796411
        -893,000000
                         0.000000
                                       1.000000
                                                      1.000000
                                                                    0.000000
min
        803.000000
                         5.000000
                                     20.000000
                                                    30.000000
                                                                    3.000000
50%
       1800,000000
                        17,000000
                                                    49,000000
                                      30.000000
                                                                    9.000000
75%
       2597.000000
                                      39,000000
                                                     69.000000
                                                                   18,86666
                        38,000000
                      234.000000
max
       3600,000000
                                      50.000000
                                                    99.000000
                                                                   36,000000
           yacEPA Home_WP_pre Away_WP_pre Home_WP_post \
16595.000000 40762.000000 40762.000000 40584.000000
count
              -0.400900
                             0.534956
                                            0.465486
                                                           0.535202
mean
                2.008798
                              0.289938
                                             0.289991
                                                           0.292223
             -14.000000
                              0.000000
                                             0.000000
                                                           0.000000
             -0.957404
                              0.327666
                                            0.221636
                                                           0.323830
               0.000000
                              0.530724
                                             0.469555
                                                           0.533272
               0.479230
                              0.778942
                                            0.672906
                                                           0.783487
                              1.000000
                                            1.000000
                                                           1.000000
               9.059733
max
                                                                      yacWPA \
                         Win_Prob
       Away_WP_post
                                            WPA.
                                                       airWPA
count 40584.000000
                    40755.000000
                                   4.296800e+04 16597.000000 16569.000000
           0.465212
                         0.505817 1.841291e-03
                                                      0.014292
                                                                    -0.010349
```

```
#Dropping duplicates
    df.drop_duplicates(inplace = True)
[ ] # get the number of missing data points per column
    missing values count = df.isnull().sum()
    # look at the # of missing points in the first ten columns
    missing values count[0:10]
    Date
                     0.0
                     0.0
    GameID
    Drive
                    0.0
    qtr
                    0.0
    down
                    0.0
    time
                    0.0
    TimeUnder
                    0.0
    TimeSecs
                     0.0
    PlayTimeDiff
                     0.0
    SideofField
                     0.0
    dtype: float64
    total_cells = np.product(df.shape)
    total missing = missing values count.sum()
    # percent of data that is missing
    (total_missing/total_cells) * 100
    24.974658974497224
```

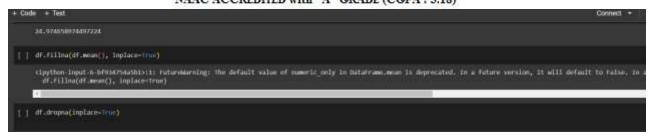
Academic Year 2023-24 SAP ID: 60004220126



SHRI VILEPARLE KELAVANI MANDAL'S DWARKADAS J. SANGHVI COLLEGE OF ENGINEERING



(Autonomous College Affiliated to the University of Mumbai)
NAAC ACCREDITED with "A" GRADE (CGPA: 3.18)



CONCLUSION:

So, we successfully implemented cleaning of data using python

_		Experiment 1 Snashwat Stab
		6000 4210120
		TY Brech Comps &
-		
		Am: To perform data preprocessing in terms of handling
		mission late remarked outliers entrained in
		eliminating duplicate rows & modifying the datatypes etc.
		Theory: Python an easy to clean programming language
	11	which makes it the most popular I prejected carryings
		for beginners in data science, data analytics and machine
		Ceanning, It also has a greate community of online
		earners and excellent data centric libraries. With
		so much double being generated it becomes important
		that the data we use you data science explications
		like machine leasing & predictive modelling is clear.
	1	Data cleaning region to the process of cleaning the
	0	listy date by identifying the errors in data and
	-	sectifying them. Data cleaning is hence a very
	1	mportant step in ML.
		and the same of th
	00	ta cleaning consists of!
١	11	ssing Values .
_	11	will stood by calculating the percentage of values
		ssing in each column and storing information in the
		vlaset.
×		op Observation - One way
-		can decop observators that contain only null values
	10	then for any of the columns, this works when
1		FOR EDUCATIONAL USE

the percontage of mischy volves in each column is less.
* Ramone Columns -
Drop column / features which have significant percentage
of missing values.
Inp missing values -
we can also fill missing value is numerical columns, very
mean, median, mode and other such structures;
Outries
It's an unusual der observation that is random & wrong
They can affect the ML model significantly.
Duplicate records
Data can sometimes contain deplicate values. Its
impostant duplicate records from dataset before me
process to Mc model.
Fising a dalatpe
Often values in deseaset went stored in the c correct
data type.
the second of th
Procedure
Frostly we load the dateset in the datagrams tun
we list columns & the no. of null vernes in that column
Its observed that ge has 177 null values, aun has
687 null values. Next we jill the null -ales. we
also obeen that the cleaned data includes column age
and encoded values of columns embasked & p class
as dependent values & survived column as largest.
Conclusion - Here we conclude that data cleaning is
way important before applying date on ML models.
FOR EDUCATIONAL USE