

Data Preprocessing - Data Analysis

Data Preprocessing

It is used for maintaining the **Quality** of the data. It includes important factors like -

- Selecting the valid data variables
- Data editing is important in some aspects
- Maintaining uniformity in data values
- Manipulation of the data for achieving the above factors (Data Wrangling)

Credits - Image from Internet

Classroom Management

Imagine you are the new teacher freshly appointed to manage the classroom.

You want to know how many students are good in -

1. sports
2. academics
3. creative work
4. marketing

Based on the students data, you have to conclude **who** can do well in **what**.

Dataset description

You are given the **data of the students** that included the following variables -

- Age
- Gender
- Address
- Father's occupation
- Mother's occupation
- Place of birth
- Height - ft
- Weight - kg
- Prev sports performance
- Prev academics performance
- Voluntary experience
- Extra co-curricular activities
- Arts and Design

Note - For our convenience all the data values are numericals.

Sports

Scenario 1

Considering the variables that are directly related -

- Height
- Weight
- Prev sports performance

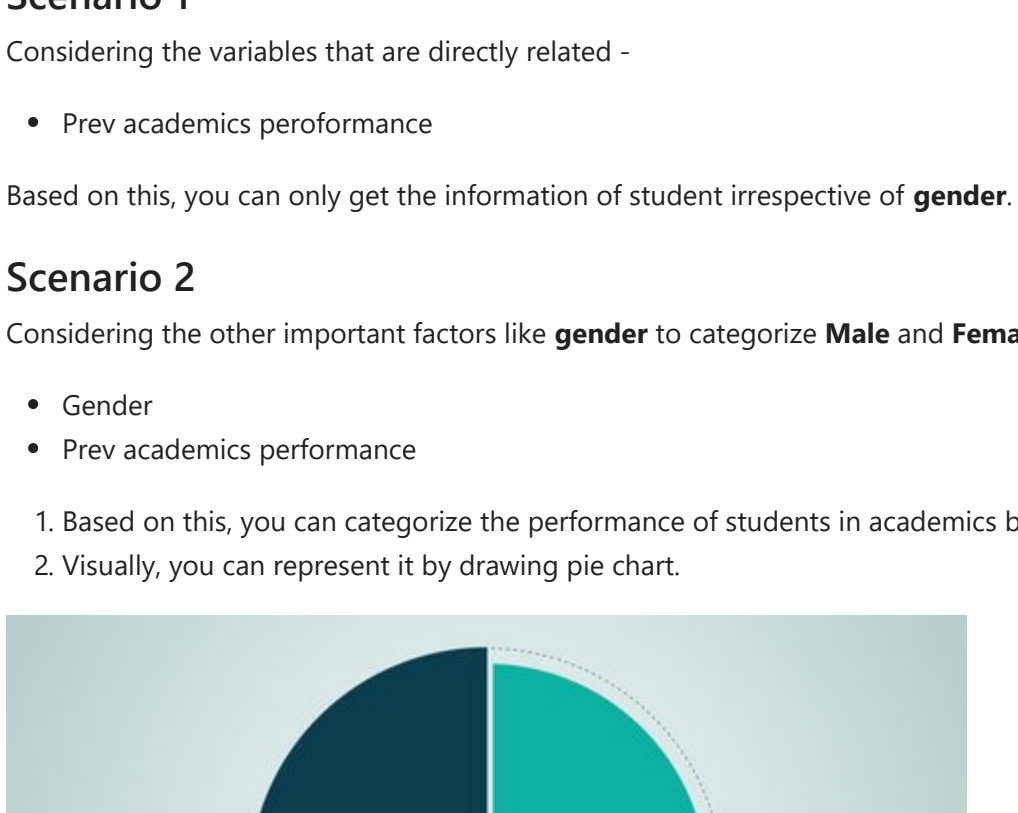
Based on this, you can only get the information of a student irrespective of **gender**.

Scenario 2

Considering the other important factors like **gender** in order to categorize as per **Male related sports** and **Female related sports**.

- Gender
- Height
- Weight
- Prev sports performance

1. Based on this, you can categorize the performance of students in sports by **Male** and **Female**.
2. Visually, you can represent it by drawing pie chart.



Credits - Image from Internet

Scenario 3

If you want to do further research on how good the person is performing in other areas, you can do so by considering -

- Gender
- Height
- Weight
- Prev sports performance
- Voluntary experience
- Extra co-curricular activities
- Prev academics performance (may be or may not be)

1. With this, you can conclude the overall students performance on sports
2. Since you are a kind teacher and well wisher of student, you can give the student a proper career guidance.

Academics

Scenario 1

Considering the variables that are directly related -

- Prev academics performance

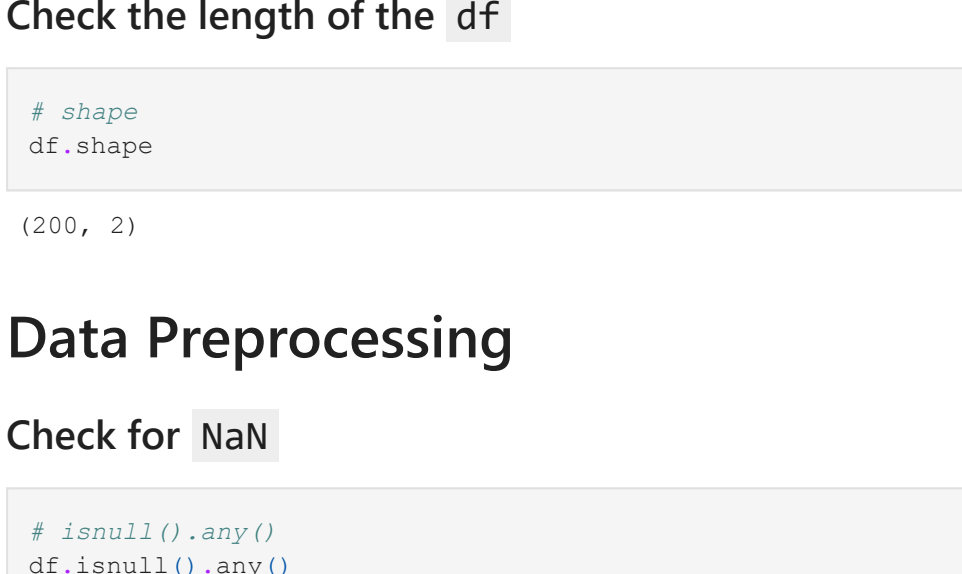
Based on this, you can only get the information of student irrespective of **gender**.

Scenario 2

Considering the other important factors like **gender** to categorize **Male** and **Female** separately.

- Gender
- Prev academics performance

1. Based on this, you can categorize the performance of students in academics by **Male** and **Female**.
2. Visually, you can represent it by drawing pie chart.



Credits - Image from Internet

Scenario 3

If you want to further research on why a particular student is **lagging behind** or **excelling ahead**, you can do so by considering -

- Address
- Father's occupation
- Mother's occupation
- Prev academics performance
- Gender (for categorizing in terms of gender)

and later on, you can decide whether to change your teaching methodology or not.

Let's make our hands dirty

data source - https://raw.githubusercontent.com/msameeruddin/Data-Analysis-Python/main/3_DA_Preprocessing/students_hw.csv

```
In [1]: import pandas as pd
import numpy as np

In [2]: data_source = 'https://raw.githubusercontent.com/msameeruddin/Data-Analysis-Python/main/3_DA_Preprocessing/s
# df
df = pd.read_csv(data_source)

In [3]: type(data_source)

Out[3]: str

In [4]: type(df)

Out[4]: pandas.core.frame.DataFrame

In [5]: # head
df.head()

Out[5]:
   Height(Inches)  Weight(Pounds)
0             65.78           112.99
1             71.52           136.49
2             69.40           153.03
3             68.22           142.34
4             67.79           144.30

Check the length of the df

In [6]: # shape
df.shape

Out[6]: (200, 2)

Data Preprocessing

Check for NaN

In [7]: # isnull().any()
df.isnull().any()

Out[7]:
Height (Inches)    True
Weight (Pounds)    True
dtype: bool

In [8]: # isnull().sum()
df.isnull().sum()

Out[8]:
Height (Inches)    2
Weight (Pounds)    1
dtype: int64

In [9]: # list of columns
df.columns

Out[9]: Index(['Height (Inches)', 'Weight (Pounds)'], dtype='object')

In [ ]:

Things to read

• What is dictionary in Python?
  └─ Keys and Values pairing. Refer to this link.

• What is a function?
• How to define functions?
• How to call functions?
• Types of functions

In [10]: for col in df.columns:
print(col)

Height (Inches)
Weight (Pounds)

In [11]: def get_nan_indices(dframe):
"""
dframe -> pandas data frame object
returns 'nan_places' a dictionary of column names and the 'nan_indices'
"""
nan_places = {}

for col in dframe.columns:
    indices = list(np.where(dframe[col].isnull())[0])
    nan_places[col] = indices

return nan_places

• Hey Python, take help of numpy to locate the NaN values for each column in dataframe called df and finally save it in a dictionary.

In [12]: # function call
get_nan_indices(dframe=df)

Out[12]: {'Height (Inches)': [10, 32], 'Weight (Pounds)': [19]}

1. In the column Height \(Inches\) , there are two NaN values at indices 10 and 32 .
2. In the column Weight \(Pounds\) , there is one NaN value at index 19 .

What can we do for those?

• Remove the entire row which ever column has a NaN .

For this, we will remove the rows which ever column has NaN . In total, there are 3 rows that need to be removed.

Remove 3 rows

• axis (0) -> row
• axis (1) -> column

In [13]: # dropna
df_1 = df.dropna(axis=0)

In [14]: df_1.shape

Out[14]: (200, 2)

In [15]: df_1.shape

Out[15]: (197, 2)

In [16]: df_2 = df.drop(index=[10, 19, 32], axis=0)

In [17]: df_2.shape

Out[17]: (200, 2)

In [18]: df_2.shape

Out[18]: (197, 2)

In [19]: pdf = df.dropna(axis=0)

check the length of pdf

In [20]: # shape
pdf.shape

Out[20]: (197, 2)

Since the index of the data frame is not in order, we need to reorder the index values to get the perfect order.

In [21]: # head(12)
pdf.head(12)

Out[21]:
   Height(Inches)  Weight(Pounds)
0             65.78           112.99
1             71.52           136.49
2             69.40           153.03
3             68.22           142.34
4             67.79           144.30
5             68.70           123.30
6             69.80           141.49
7             70.01           136.46
8             67.90           112.37
9             66.78           120.67
11            67.62           114.14
12            68.30           125.61

Reset the index

In [22]: # rdf
# reset with drop
rdf = pdf.reset_index(drop=True)

In [23]: # shape
rdf.shape

Out[23]: (197, 2)

In [24]: # head
rdf.head(12)

Out[24]:
   Height(Inches)  Weight(Pounds)
0             65.78           112.99
1             71.52           136.49
2             69.40           153.03
3             68.22           142.34
4             67.79           144.30
5             68.70           123.30
6             69.80           141.49
7             70.01           136.46
8             67.90           112.37
9             66.78           120.67
10            67.62           114.14
11            68.30           125.61

In [ ]:

Check if Height(Inches) < 40

In [25]: # inch_thresh
inch_thresh = 40

In [26]: # filter with <
rdf[rdf['Height (Inches)'] < inch_thresh]

Out[26]:
   Height(Inches)  Weight(Pounds)
68            30.84           134.02
93            36.29           120.03

Remove the rows where Height(Inches) < 40

• In the above case, we can see two values where height is less than 40.
• We remove by specifying the index values in drop\(\) method.

In [27]: # drop by index
rdf = rdf.drop(index=[68, 93], axis=0)

In [28]: # shape
rdf.shape

Out[28]: (195, 2)

Reset the index

In [29]: # hw_df
# drop = True
hw_df = rdf.reset_index(drop=True)

In [30]: # shape
hw_df.shape

Out[30]: (195, 2)

Since the index of the data frame is not in order, we need to reorder the index values to get the perfect order.

In [ ]:

Categorize the data

• Refer to -> https://pandas.pydata.org/docs/reference/api/pandas.cut.html

In [31]: # height_cat
# weight_cat

height_cat = pd.cut(x=rdf['Height (Inches)'], bins=3, labels=['short', 'average', 'tall'])
weight_cat = pd.cut(x=rdf['Weight (Pounds)'], bins=3, labels=['under', 'normal', 'obesity'])

Make a new column height_cat and weight_cat in the dataframe - hw_df

In [32]: hw_df.head()

Out[32]:
   Height(Inches)  Weight(Pounds)
0             65.78           112.99
1             71.52           136.49
2             69.40           153.03
3             68.22           142.34
4             67.79           144.30

In [33]: hw_df['height_cat'] = height_cat
hw_df['weight_cat'] = weight_cat

In [34]: # head
hw_df.head()

Out[34]:
   Height(Inches)  Weight(Pounds)  height_cat  weight_cat
0             65.78           112.99      short      under
1             71.52           136.49       tall      normal
2             69.40           153.03      average      obesity
3             68.22           142.34      average      obesity
4             67.79           144.30      average      obesity

In [ ]:

Plotting the pie chart to show

• how many are short
• how many are tall
• ...

Take value_counts() of height_cat variable

In [35]: hw_df['height_cat'].value_counts()

Out[35]:
average    117
short       60
tall        16
Name: height_cat, dtype: int64

Take value_counts() of weight_cat variable

In [36]: hw_df['weight_cat'].value_counts()

Out[36]:
normal     112
under       44
obesity     37
Name: weight_cat, dtype: int64

In [37]: # hdf_pie -> to_frame
# wdf_pie -> to_frame
hdf_pie = hw_df['height_cat'].value_counts().to_frame()
wdf_pie = hw_df['weight_cat'].value_counts().to_frame()

In [38]: # display hdf_pie
hdf_pie

Out[38]:
   height_cat
average    117
short       60
tall        16

In [39]: # display wdf_pie
wdf_pie

Out[39]:
   weight_cat
normal     112
under       44
obesity     37

In [40]: # plot pie of hdf_pie with size (width=10, height=6)
hdf_pie.plot(kind='pie', figsize=(10, 6), subplot=True)

Out[40]: array([<AxesSubplot:ylabel='height_cat'>], dtype=object)


```

In [41]: # plot pie of wdf_pie with size (width=10, height=6)
wdf_pie.plot(kind='pie', figsize=(10, 6), subplot=True)

Out[41]: array([<AxesSubplot:ylabel='weight_cat'>], dtype=object)

What did we learn?

- Data Preprocessing
- Real life scenario Example
- Classroom Management problem
 - sports
 - academics
 - creative work
 - marketing
- Getting hands dirty by writing code

