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) Model Evaluation Model Building **Data Exploration** & Visualization Data Preprocessing **Data Collection** & Assembly 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 conveneince all the data values are numericals. In []: How do we convert a "feeling " into a number? • We can measure a "feeling " into a number through a "scale range" • If the scale is 1 to 4, then we can term -1 → Not Satisfied ■ 2 → Slightly Satisfied 3 → Satisfied 4 → Highly Satisfied **Sports** Scenario 1 Considering the variables that are directty 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. male femal 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 performace (may be or may not be) 1. With this, you 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 peroformance 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. 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. 50% **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. In []: Note - Data Analyst should be wise enough to select the important data variables. This helps to get proper insights pertaining to the problem statement that he/she is assigned to do. In []: Let's make our hands dirty data source → http://wiki.stat.ucla.edu/socr/index.php/SOCR_Data_Dinov_020108_HeightsWeights import pandas as pd In [1]: import numpy as np In [2]: | df = pd.read_csv('students_hw.csv') df.head() Out[2]: Height(Inches) Weight(Pounds) 0 112.99 65.78 71.52 136.49 2 69.40 153.03 68.22 142.34 67.79 144.30 Check the length of the df In [3]: # shape df.shape Out[3]: (200, 2) **Data Preprocessing** Check for NaN In [4]: # isnull().any() df.isnull().any() Out[4]: Height(Inches) True Weight (Pounds) True dtype: bool In [5]: # list of columns df.columns Out[5]: Index(['Height(Inches)', 'Weight(Pounds)'], dtype='object') In [6]: def check_for_nan(dframe): 11 11 11 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 In [7]: check_for_nan(dframe=df) Out[7]: {'Height(Inches)': [10, 32], 'Weight(Pounds)': [19]} • In the column Height (Inches), there are two NaN values at indices 10 and 32. • 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 [8]: # pdf # dropna pdf = df.dropna(axis=0) check the length of pdf In [9]: # shape pdf.shape Out[9]: (197, 2) Since the index of the data frame is not in order, we need to reindex the index values to get the perfect order. In [10]: # head pdf.head(12) Out[10]: Height(Inches) Weight(Pounds) 0 65.78 112.99 1 71.52 136.49 2 69.40 153.03 3 68.22 142.34 67.79 144.30 5 68.70 123.30 6 69.80 141.49 70.01 7 136.46 8 67.90 112.37 120.67 9 66.78 11 67.62 114.14 12 68.30 125.61 Reset the index In [11]: # rdf # reset with drop pdf = pdf.reset_index(drop=True) In [12]: # shape pdf.shape Out[12]: (197, 2) In [13]: # head pdf.head(12) Out[13]: Height(Inches) Weight(Pounds) 0 65.78 112.99 71.52 136.49 1 2 69.40 153.03 3 68.22 142.34 4 67.79 144.30 123.30 5 68.70 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 68.30 125.61 11 Check if Height (Inches) < 40 In [14]: # inch_thresh $inch_thresh = 40$ In [15]: # filter with < pdf[pdf['Height(Inches)'] < 40]</pre> Out[15]: 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 <code>drop()</code> method. In [16]: # drop by index rdf = pdf.drop(index=[68, 93], axis=0) In [17]: # shape rdf.shape Out[17]: (195, 2) Reset the index In [18]: # hdf # drop = True hdf = rdf.reset_index(drop=True) In [19]: # shape hdf.shape Out[19]: (195, 2) Since the index of the data frame is not in order, we need to reindex the index values to get the perfect order. Categorize the data In [20]: # avg_height → mean avg_height = hdf['Height(Inches)'].mean() In [21]: avg_height Out[21]: 67.95584615384617 In [22]: # round() avg_height = round(avg_height) In [23]: avg_height Out[23]: 68 In [24]: # write code # height_cat -> consider round()

height_cat = []

len(height_cat)

print(height cat)

In [25]: # check len

In [26]: # display

In [27]: hdf.head()

1

2

3

In [28]: | # new column

0

2

3

check head

hdf.head()

Height(Inches) Weight(Pounds)

112.99

136.49

153.03

142.34

144.30

112.99

136.49

153.03

142.34

144.30

Short

Tall

Tall

Average

Average

Take value_counts() Of height_category variable

df pie = hdf['height cat'].value counts().to frame()

65.78

71.52

69.40

68.22

67.79

65.78

71.52

69.40

68.22

67.79

how many are shorthow many are tall

Plotting the pie chart to show

hdf['height_cat'].value_counts()

Name: height cat, dtype: int64

74

42

height_cat

74

42

In [33]: # plot pie of df_pie with size (width=10, height=6)

Out[33]: array([<AxesSubplot:ylabel='height_cat'>], dtype=object)

df_pie.plot(kind='pie', figsize=(10, 6), subplots=True)

Short

Short

Average

Tall Average

hdf['height_cat'] = height_cat

Height(Inches) Weight(Pounds) height_cat

Out[27]:

In [29]:

Out[29]:

In [30]:

In [32]:

Out[32]:

In []:

What did we learn?

· Data Preprocessing

sportsacademicscreative workmarketing

• Real life scenario Example

· Classroom Management problem

· Getting hands dirty by writing code

· Plotting the pie chart for showing categorization

Out[30]: Short

Tall Average

In [31]: # df pie \rightarrow to frame

df pie

Short Tall

Average

display df_pie

Out[25]: 195

for i in hdf['Height(Inches)'].to_list():

height_cat.append('Short')

height cat.append('Tall')

height cat.append('Average')

t', 'Tall', 'Short', 'Average', 'Average', 'Tall']

Make a new column height category in the dataframe - hdf

['Short', 'Tall', 'Tall', 'Average', 'Average', 'Tall', 'Tall', 'Tall', 'Average', 'Short', 'Average', 'Average', 'Short', 'Average', 'Tall', 'Short', 'Tall', 'Tall', 'Average', 'Tall', 'Short', 'Average', 'Average', 'Short', 'Tall', 'Short', 'Short', 'Tall', 'Short', 'Tall', 'Tall', 'Short', 'Short', 'Tall', 'Average', 'Short', 'Short', 'Short', 'Tall', 'Short', 'Tall', 'Average', 'Short', 'Tall', 'Average', 'Tall', 'Average', 'Tall', 'Short', 'Tall', 'Tall', 'Short', 'Tall', 'Average', 'Tall', 'Tall', 'Tall', 'Tall', 'Tall', 'Tall', 'Short', 'Average', 'Short', 'Average', 'Short', 'Average', 'Tall', 'Tall', 'Tall', 'Tall', 'Tall', 'Tall', 'Short', 'Tall', 'Tall', 'Tall', 'Short', 'Tall', 'Tall', 'Short', 'Tall', 'Tall', 'Short', 'Short', 'Tall', 'Tall', 'Short', 'Short', 'Short', 'Tall', 'Tall', 'Tall', 'Short', 'Short', 'Short', 'Tall', 'Tall', 'Tall', 'Tall', 'Average', 'Short', 'Short', 'Short', 'Tall', 'Tall', 'Tall', 'Tall', 'Average', 'Short', 'Short', 'Tall', 'Tall', 'Tall', 'Short', 'Short', 'Short', 'Tall', 'Tall', 'Short', 'Short', 'Short', 'Tall', 'Tall', 'Short', 'Short', 'Short', 'Tall', 'Short', 'Tall', 'Short', 'Short', 'Tall', 'Short', 'Tall',

'Tall', 'Tall', 'Short', 'Short', 'Short', 'Tall', 'Short', 'Tall', 'Average', 'Tall', 'Short', 'Short

if round(i) < avg_height:</pre>

elif round(i) == avg_height:

Data Preprocessing - Data Analysis