

Data Exploration in Python USING

NumPy stands for Numerical

Python. This library contains basic linear algebra functions Fourier transforms, advanced random number capabilities.

Pandas for structured data operations and manipulations. It is extensively used for data munging and preparation.

Pandas

Matplotlib Python based plotting library offers matplotlib

with a complete 2D support along with limited 3D graphic support. CHEATSHEET -

Contents Data Exploration



3. How to transpose a table?

4. How to sort Data? 5. How to create plots (Histogram, Scatter, Box Plot)?

1. How to load data file(s)? 2. How to convert a variable to different data type?

6. How to generate frequency tables?

7. How to do sampling of Data set? 8. How to remove duplicate values of a variable?

9. How to group variables to calculate count, average, sum'? 10. How to recognize and treat missing values and outliers?

11. How to merge / join data set effectively?

How to load data file(s)?

Here are some common

atemp \

1 9.84 14.395

1 9.02 13.635

1 9.02 13.635

Description Read delimited data from a file. Use Comma as default delimiter read csv read table Read delimited data from a file. Use tab ('\t') as default delimiter read excel Read data from excel file read_fwf Read data in fixed width column format

Function

CODE

Output

CODE

01-01-2011 00:00

1 01-01-2011 01:00

2 01-01-2011 02:00

8.0

8.0



read_clipboard | Read data from clipboard. Useful for converting tables from web pages Loading data from CSV file(s):

loading...

functions used to read data

import pandas as pd
#Import Library Pandas
df = pd.read_csv("E:/train.csv") #I am working in Windows environment
#Reading the dataset in a dataframe using Pandas
print df.head(3) #Print first three observations

datetime season holiday workingday weather temp

32

40

32

0

df=pd.read_excel("E:/EMP.xlsx", "Data") # Load Data sheet of excel file EMP

humidity windspeed casual registered count

Loading data from txt file(s):

df=pd.read_csv("E:/Test.txt",sep='\t')

- Convert character date to Date

from datetime import datetime

Load Data from text file having tab '\t' delimeter print df

- Convert numeric variables to string variables

Loading data from excel file(s):

How to convert a variable to different data type?

print date_obj

- Data set used

ID

1

print df

result

Output

Out[35]:

#Sorting Dataframe

Table A

Product

AAA

Sales

50

#Transposing dataframe by a variable

and vice versa srting_outcome = str(numeric_input) #Converts numeric_input to string_outcome

integer_outcome = int(string_input) #Converts string_input to integer_outcome

float_outcome = float(string_input) #Converts string_input to integer_outcome

How to transpose a Data set?

char_date = 'Apr 1 2015 1:20 PM' #creating example character date

date_obj = datetime.strptime(char_date, '% b % d % Y % I : % M % p')



BBB 45 52 46 52 AAA 46 BBB Code

ID

1

df=pd.read_excel("E:/transpose.xlsx", "Sheet1") # Load Data sheet of excel file EMP

result= df.pivot(index= 'ID', columns='Product', values='Sales')

ID Product

Product AAA BBB

AAA

BBB

AAA

BBB

Table B

AAA

50

BBB

45

Total rows: 4 Total columns: 3

Product Sales

52

AAA

AAA

BBB

BBB

Sorted Table

Sales

123

114

135

139

117

121

133

140

133

133

OutPut

35

Age

35

Age

OutPut

30

40

45

#Add by variable name(s) to sort

Total rows: 4 Total columns: 3

Product

AAA

BBB

AAA

BBB

Orginal Table

ID

df=pd.read_excel("E:/transpose.xlsx", "Sheet1")

print df.sort(['Product','Sales'], ascending=[True, False])

Sales

50

45

52

46

How to create plots (Histogram, Scatter, Box Plot)?

М

М

М

Μ

М

40

37

30

44

36

32

26

32

36

ID

45 1 50 2 52 46 How to sort DataFrame? CODE

Sales

50

45

52

46

EmpID Gender Age E001 М 34

E002

E003

E004

E005

E006

E007

E008

E009

E010

Histogram

#Plot Histogram

fig=plt.figure()

#Variable

#Labels and Tit

plt.xlabel('Age')

plt.show()

Code

fig=plt.figure()

#Variable

#Labels and Tit

plt.xlabel('Age')

plt.show()

Code

Box-plot:

sns.despine()

import seaborn as sns

sns.boxplot(df['Age'])

plt.ylabel('Sales')

ax = fig.add_subplot(1,1,1)

ax.scatter(df['Age'],df['Sales'])

plt.title('Sales and Age distribution')

create blank figure

ax = fig.add_subplot(1,1,1)

ax.hist(df['Age'],bins = 5)

plt.title('Age distribution')

import matplotlib.pyplot as plt

object, use plt.figure to create new figure

#Create one or more subplots using

add_subplot, because you can't

Code

import pandas as pd Age distribution 3.0 df=pd.read_excel("E:/First.xlsx", "Sheet1") 2.5 #Plots in matplotlib reside within a figure

2.0

1.5

1.0

0.5

0.0 L 25

#Employee

plt.ylabel('#Employee') Scatter plot #Plots in matplotlib reside within a figure object, use plt.figure to create new figure Sales and Age distribution 145 #Create one or more subplots using 140 add_subplot, because you can't create blank figure 135

130

120

115

110

42

Sales

import pandas as pd

df=pd.read_excel("E:/First.xlsx", "Sheet1")

test= df.groupby(['Gender','BMI'])

Code

print df

100%

test.size()

get 5 random rows from df dfr = df.ix[rindex] print dfr

Code

Code test= df.groupby(['Gender']) test.describe()

40 38 36 34 32 28 26 How to generate frequency tables with pandas? **OutPut** Age

E001

E002

E003

E004

E005

E006

E007 E008 E009 E010 Out[84]: Gender Normal

34

Μ

123

114

135

139

121

133

140

133

133

1

EMPID Gender Age Sales BMIUnderweight E005 117 E003 135 Obesity E008 26 140 Normal E009 32 133 Μ Normal E006 36 121 M Normal How to remove duplicate values of a variable? Output

EMPID Gender

F 44

Output

4.000000

36.750000

7.719024

Age

count

Sales

4.000000

126.500000

12.922848

E001

E003

E004

E005

E007

E008

Gender

How to group variables in Python to calculate count, average, sum?

Sales

123

135

133

140

114

BMI

Normal

Obesity

Obesity

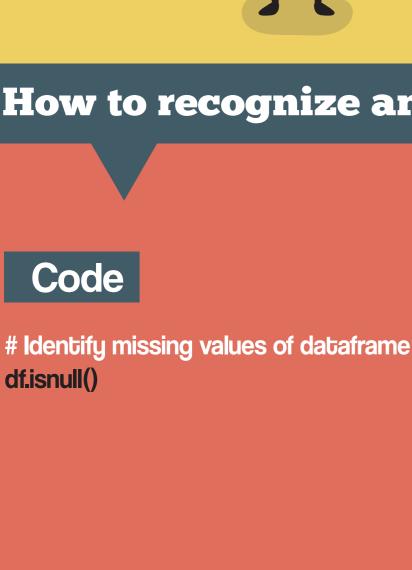
Normal

Overweight

139 Underweight

117 Underweight

#Remove Duplicate Values based on values of variables "Gender" and "BMI" rem_dup=df.drop_duplicates(['Gender', 'BMI']) print rem_dup



Code

Output df.isnull() EMPID Gender Age 0 False False False False 1 False False False False False False False | False 3 False False False False 4 False False False | False

#Create Sample dataframe import numpy as np import pandas as pd from random import sample # create random index rindex = np.array(sample(xrange(len(df)), 5))

import numpy as np meanAge = np.mean(df.Age)

Sales BMI False False

#replacing missing values in the DataFrame df.Age = df.Age.fillna(meanAge)

Code df_new = pd.merge(df1, df2, how = 'inner', left_index = True, right_index = True)

merges df1 and df2 on index # By changing how = 'outer', you can do outer join. # Similarly how = 'left' will do a left join # You can also specify the columns to join instead of indexes, which are used by default. To view the complete guide on Data Exploration in Python visit here - http://bit.ly/1KWhaHH

BMI

Normal

Obesity

Normal

Obesity

Normal

Normal

Overweight

Underweight

Underweight

Underweight

1 Obesity Overweight Underweight Normal Obesity Underweight dtype: int64 How to do sample Data set in Python? OutPut

Code

75% 41.000000 44.000000 max 6.000000 count 33.333333 2.422120 std M

Out[116]: False False False False 6 False False False #Example to impute missing values in Age by the mean False False False 8 False False False | False #Using numpy mean function to calculate the mean value 9 False False False | False | False How to merge / join data sets?

26.000000 114.000000 25% 34.250000 116.250000 38.500000 126.000000 136.250000 140.000000 6.000000 130.333333 6.889606 30.000000 121.000000 32.000000 | 125.500000 33.000000 133.000000 35.500000 133.000000 36.000000 | 139.000000 How to recognize and Treat missing values and outliers? In [116]: # Identify missing values of dataframe