Lecture objectives

Main concepts covered:

• Learn Pandas basics

Data import/export

Data cleansing

Data plotting

Data replacing

Basic features

In []: import pandas as pd

Importing/exporting data

Importing data from a file

In []: # import raw data from excel file

df_copy = titanic.copy()

df = titanic.copy()

In []: | # display the first 5 samples

Variables meaning

1st = Upper,

2nd = Middle,

■ 3rd = Lower,

• ticket: Ticket number

• fare: Passenger fare

• cabin: Cabin number

C = Cherbourg,

Q = Queenstown,

■ S = Southampton

Exporting data in a file

In []: # export dataframe to a csv file

Useful functions

In []: # useful info of the dataset

df.info()

df.head()

df.shape

df_copy

In []: | # result as expected

df.groupby(['sex'])

df.groupby(['sex']).mean()

In []: # for a given parameter (column)

Data plotting

df.plot()

df.plot.bar()

• df.hist(bins=...)

2- call the plot you need

In []: # same for the 'survived' column

In []: # customize the number of 'bins' df['age'].hist(bins=20);

> # - the diagonal > histograms # - elsewhere -> scatter plots pd.plotting.scatter_matrix(df);

plot on the same figure

Series and DataFrames

df['column'] = serie

just the first 3 samples

df_copy.set_index('name')

In []: # let's retrieve only the passengers ages

Navigating through the data

• df['age'] = serie:ndarray

• df['age'][0:18]:indexing

• df['age'] < 18: mask

In []: # group function after boolean indexing

Localisation by index

df.iloc[5:10,0:4]

• Localisation (by column)

df.loc[5:10,'column']

In []: # direct nagivation on indexes not allowed

each column is a serie -> 1D array

In []: ser1 = pd.Series([1,2,3],index = ['USA', 'Germany','Japan'])

In []: ser2 = pd.Series([10,20,30],index = ['USA', 'Germany','Japan'])

In []: | df_new = pd.DataFrame(data=[[200, 'cat1'], [250, 'cat2'], [300, 'cat1']],

• df[df['age'] < 18] = dataFrame: boolean indexing

df[df['age'] < 18].groupby(['sex','pclass']).mean()</pre>

Localization methods fo navigating through data

columns=['COL1','COL2'])

df_copy.head(3)

df['age']

ser1

ser2

ser3

In []: df_new

In []: # indexing

In []: | # masking

df['age'][0:18]

df['age'] < 18

df[df['age'] < 18]

In []: # boolean indexing

df[5:10,0:4]

In []: # iloc method allow it

df.iloc[5:10]

df.iloc[5:10,:]

df.loc[5:10]

Navigation by column

In []: # more explicit navigation

df.loc[5:10,:]

df.loc[5:10,0:3]

age <= 20 -> 0

age > 40 -> 3

df_copy

In []: df_new

In []:

In []:

In []:

In []:

In []:

In []: **from** IPython.core.display **import** HTML

return HTML(styles)

In []: $df_{copy.loc}[df_{copy}['age'] \le 20, 'age'] = 0$

df_copy.loc[df_copy['age'] > 40, 'age'] = 3

styles = open("../styles/custom.css", "r").read()

In []: $df_{copy}.loc[(df_{copy}['age'] > 20) & (df_{copy}['age'] <= 30), 'age'] = 1$

plt.hist(titanic['fare'][titanic['survived'] == 0], bins=30,

plt.hist(titanic['fare'][titanic['survived'] == 1], bins=30,

plt.hist(df['age'][df['survived'] == 0], bins=30,

plt.hist(df['age'][df['survived'] == 1], bins=30,

styles = open("../styles/custom.css", "r").read()

alpha=0.5, edgecolor='black', label='Did not survive', color='r')

alpha=0.5, edgecolor='black', label='Did not survive', color='r')

alpha=0.5, edgecolor='black', label='Survived', color='g')

alpha=0.5, edgecolor='black', label='Survived', color='g')

 $df_{copy}.loc[(df_{copy}['age'] > 20) & (df_{copy}['age'] <= 30), 'age'] = 1$ $df_{copy}[(df_{copy}['age'] > 30) & (df_{copy}['age'] <= 40), 'age'] = 2$

def css_styling():

In []: df_copy['age'].value_counts()

plt.ylabel('Fare')

In []: import matplotlib.pyplot as plt

plt.figure(figsize=(10,5))

plt.ylabel('Passengers ages')

In [1]: from IPython.core.display import HTML

return HTML(styles)

def css_styling():

css_styling()

plt.legend();

plt.legend();

In []:

In []:

Out[1]:

In []:

In []: import matplotlib.pyplot as plt

plt.figure(figsize=(10,5))

css_styling()

df_copy

20 < age <= 30 -> 1 # 30 < age <= 40 -> 2

complete code below

Exercise

In []:

In []: # explicit indexing on one axis

Navigation by index

df.iloc[5:10,0:4]

In []: # iloc navigation on a specific axis

with all data on the other axis

In []: # let's make make it more explicit -> recommended

implicitly -> all data on the other axis

df.loc[5:10,['pclass', 'survived', 'sex', 'age']]

For the following exercises, don't delete the cells with the expected results

Create your own in order the chech the matching with your result

In []: # create 4 categories of age using boolean indexing an loc() function:

In []: # let Pandas count the values of each of the created categories of ages

In []: # noy allowed -> only for index localisation

use df_copy (2nd copy of the raw data)

display the result as a dataframe

display the resulting dataframe

Replacing data with new values

• df['sex'].map({'male':0, 'female':1})

• df['sex'].astype('category').cat.codes

In []: df_new['COL2'] = df_new['COL2'].astype('category').cat.codes

In []: df_new['COL2'] = df_new['COL2'].replace([0,1],['cat1','cat2'])

In []: df_new['COL2'] = df_new['COL2'].map({'cat1':0, 'cat2':1})

In []: df_new['COL2'] = df_new['COL2'].map(lambda x:x+1)

Congratulations!

• df['sex'].replace(['male', 'female'], [0,1])

df_copy['age'].value_counts()

In []: # navigating through many columns at once

In []: ser3 = ser1 + ser2

In []: # let's check the type type(df['age'])

In []: # let's recall the data for the illustration

In []: # let's change the default index -> use the name

df['age'].hist();

In []: # scatter matrix plot

Exercises

In []: # complete code here

In []: # complete code here

df['pclass'].value_counts()

Pandas is built on top of Matplotlib!

• df.scatter(x=..., y=...)

1- count the number per category

uggly message above the figure

df['pclass'].value_counts().plot.bar()

df['survived'].value counts().plot.bar();

In []: # display the histogram of the passengers ages

pd.plotting.scatter_matrix(df)

In []: # display a barplot of numbers of passengers per 'pclass'

NB: add the ';' add the end of the call to avoid

In []: """ dont'remove this cell -> create your own to check the result """

(1) 1st plot -> the histogram passengers ages who survived # (2) 2nd plot -> the histogram passengers ages who did not # is there any correlation between 'age' and 'survival'?

In []: """ dont'remove this cell -> create your own to check the result """

is there any correlation between passengers 'fare' and their 'survival'?

• A Pandas **serie** is a NumPy 1D array which has an index axis (which is independent from Numpy).

use the original (raw) dataframe -> 'titanic'

same exercise with passengers 'fare'

• A Pandas **dataframe** is a dictionary of series.

Note: Pandas is composed of **series** and **dataframes**.

df copy

df.describe()

df_copy.describe()

In []: df.info()

In []:

df.describe()

Data cleansing

In []: # let's remove unnecessary columns

Parent = mother, father,

• embarked: Port of Embarkation

Export to an excel file: df.to_excel()

Export to a csv file: df.to_csv()

• Export to a json file: df.to_json()

Export to an html file: df.to_html()

Export to a laxtex file: df.to_latex()

df.to_csv('../DATA/titanic_export.csv')

number of non null values, type of each variable ...

• Remove columns: df.drop(['col1','col2',...]), which may cause loss of information

• Fill missing with a default value: df.fillna(df['age'].mean()), which may introduce biases

• Delete underlying rows: df.dropna(axis=0), which may cause loss of information

'embarked', 'boat', 'body', 'home.dest'], axis=1)

'embarked', 'boat', 'body', 'home.dest'], axis=1, inplace=True)

df = df.drop(['name', 'sibsp', 'parch', 'ticket', 'fare', 'cabin',

Same result with the following instruction (without assignment):

df.drop(['name', 'sibsp', 'parch', 'ticket', 'fare', 'cabin',

In []: # let's check the result by viewing the first 5 samples

In []: # remove of all samples with missing values

df = df.dropna(axis=0)

In []: # new infos on the modifed dataset

In []: # display the data to see the changes

In []: # setting 'inplace' parameter to True

Group_by() and value_counts()

df.groupby(['sex','pclass']).mean()

In []: # group of data (1st copy the raw data) by 'sex'

In []: # with the grouped data (by sex), compute the mean()

count the number of values per category

In []: # count the number of values per category of 'pcalss'

df[df['age'] < 18]['pclass'].value_counts()</pre>

and display the shape to see the change

In []: # statistics on the second copy of the raw data

In []: # let's remove few columns, keeping the 'name' column

df_copy.drop(['sibsp', 'parch', 'ticket', 'fare', 'cabin',

nothing happened -> need to specify 'inplace' to True

df_copy.drop(['sibsp', 'parch', 'ticket', 'fare', 'cabin',

the result is a dataFrame (we will see the meaning later)

In []: # same with nested group -> by 'sex' -> for each 'sex' -> by 'pclass'

only for passengers aging less than 18 (boolean indexing)

'embarked', 'boat', 'body', 'home.dest'], axis=1)

'embarked', 'boat', 'body', 'home.dest'], axis=1, inplace=True)

In []: # basic statistics : total, mean ... of each variable

Export to an sql file: df.to_sql()

In []: | # display the data shape

df.shape

df.head()

Import an excel file: pd.read_excel()

Import a csv file: pd_read_csv()

Import a json file: pd. read_json()

Import an sql file: pd. read_sql()

Import an html file: pd.read_html()

titanic = pd.read excel('../DATA/titanic3.xls')

- to see the last 5 samples -> df.tail()

pclass: a proxy for socio-economic status (SES)

• **sibsp:** the dataset defines family relations in this way...

Sibling = brother, sister, stepbrother, stepsister,

• parch: the dataset defines family relations in this way...

Child = daughter, son, stepdaughter, stepson,

make 2 copies of the raw data for further purposes

- to change the default value i.e. 5, enter another value

• age: age is fractional if less than 1. If the age is estimated, is it in the form of xx.5

■ Some children travelled only with a nanny, therefore parch=0 for them.

Spouse = husband, wife (mistresses and fiancés were ignored)

API reference: https://pandas.pydata.org/docs/reference/index.html

Navigating