

Pandas:

- will complete this by 260g
- stay on fast track - index
- Pandas is a library for data analysis
- Extremely powerful table (DataFrame) system built on Numpy.
- Fantastic documentation

What can we do with pandas:

- Tools for recoding and working with data between many formats
- Intelligently grab data based on indexing logic, subsetting, and more
- Handle missing data
- Adjust and reshape data

Section - Overview

- Series and DataFrame
- Condition filtering and use full methods.

- Missing data
- Combining DataFrames
- Text method and time method
- Input and Output

Series - Part-one

- A Series is a data structure in pandas that hold an array of information along with a named index
- The named object differentiates this from a simple Numpy array
- Formal definition : One-dimensional n-darray with axis label
- NumPy array has index

Index	Data
0	1776
1	224
2	1460

Label	Index	Data
USA	0	1776
CANADA	1	1867
MEXICO	2	1821

Data is still numerically organized

Code #1

```
import numpy as np  
import pandas as pd  
  
# help(pd.Series)
```

```
pd.Series(data=None,  
          index=None,  
          dtype=None,  
          name=None,  
          copy=False,  
          fastpath=False,)
```

```
myindex = ['un', "canada", 'mexico']  
mydata = [1776, 1867, 1821]  
myseries = pd.Series(data=mydata,  
                      index=myindex)
```

```
myseries[0]  
out: 1776
```

```
myseries['USA']  
out: 1776
```

Series from Dictionary.

```
age = {'suny': 5, 'frank': 10,  
       'spike': 7}  
pd.Series(age)
```

```
out: sum      5  
Frank    10  
Spilke     7  
dtype: int64
```

Pandas - Series part 2

```
# Imagine Sales for 1st and
```

```
2nd Quarters for Col
```

```
q1 = { "Japan": 80, "China": 450, "Brazil": 200, "Russia": 250 }  
q2 = { "Brazil": 100, "China": 500, "India": 110, "Russia": 260 }
```

```
Sales_q1 = pd.Series(q1)
```

```
Sales_q2 = pd.Series(q2)
```

```
Sales = q1 + q2 # Addition
```

```
Sales_q2
```

```
Sales_q1.keys()
```

```
out: Index(['Japan', 'China', 'Brazil'],  
          dtype='object')
```

Broadcasting Occurs in
pandas Series:-

```
[1, 2] * 2 => [1, 2, 1, 2]
```

But in np broadcasting occurs

```
np.array([1, 2]) * 2  
=> [2, 4]
```

then :-

Sales_q1 + q2

Japan 160

China 900

India 600

USA 500

Sales_q1 / 100

Sales_q1 + Sales_q2

Brazil NaN

China 950.0

India 110.0

Japan NaN

USA 510.0

Sales_q1.add(Sales_q2, fill_value=0)

Convert to numerical to
float in numerical
operations

Pandas - Data Frame

" DataFrame is a table of columns and rows in pandas that we can easily restructure and filter.

Formal Definition -

A group of pandas series object that share the same index.

Index	Year	Ref
USA	12	33
CANADA	14	34
MEXICO	17	72

This Series covers first the "basics"

- Create a data frame
- Grab a column or multiple columns
- Grab a row or multiple rows
- Insert a new column or new rows

Code #

```
import numpy as np  
import pandas as pd
```

pt.

```
# np.random.seed(101)
mydata = np.random.randint(0,101,(4,3))

mydata
array([[ 1, 11,  8],
       [ 22, 10, 15],
       [ 95, 9, 77],
       [ 10, 4, 63]]]

myindex = ['Jan', 'NY', 'A2', 'TX'],
mycolours = ['Jan', feb', 'Mar']

df = pd.DataFrame(data=mydata,  
                   index = myindex,  
                   columns = mycolours)

df.info();

# Where my python code available  
pwd → full path  
ls → Report all the files in  
Current directory  
df = pd.read_csv('tips.csv')
```

```
pd.read_csv("C:\pdsca\1\dataset1.csv")
```

part 1S only in jupyter

```
part 2 - pandas DataFrame
```

```
df.columns = attributes
```

```
Index([ "total_bill", ... ])
```

```
df.index RangeIndex(start=0, stop=244, step=1)
```

```
# df.head() - First 5 rows
```

```
# df.head() -> last 5 rows
```

```
# df.info() basic info from info using
```

```
# df.describe() calculate statistical values
```

```
# df.describe().transpose()
```

```
part # 3 - pandas DataFrame
```

```
type(df[ 'total_bill' ])
```

```
pandas.core.series
```

```
mycols = ['total_bills', 'tips']  
df[mycols]
```

```
# New column  
df['tip_percent'] = df['total']/100.
```

$$df['tip_percentage'] = (df['tips']/df['total']) * 100$$

df

```
np.round((df['total_bills']/df['tips']), 2)
```

```
df.head(6)
```

Remove columns

```
df.drop(label=None, axis=1)
```

axis=0 remove rows

axis=1 remove columns

```
df = df.drop(label=None, axis=1)
```

df.shape

(249, 11)
6 1

axis

Range Index

Part 4 Dataframe

```
# df.set_index  
# df.set_index("Payment-ID")  
# Set index to column  
  
df = df.set_index("Payment-ID")  
df.head(6);
```

(f) df.reset_index() permanent change -
df

df.iloc[0] → numeric.

df.loc['Sunday'] → label based

slice for multiple

df.iloc[0:4]

df.loc[['Sunday', 'Sunday260']]

df.drop('Sunday', axis=0)
df =

one-row = df.iloc[0]

df.append(one-row)
df = df.append(one-rows)

Condition Filtering :-

Organizing data

Index	Year	Pop	GDP
USA	1968	32.8	20.5
CANADA	1871	38	1.7
MEXICO	371	126	1.12

df["pop"] > 50

return

Series of boolean values

{ true }
False }

df[df["pop"] > 50]

o Conditional Filtering:

- filter by multiple conditions -
- check against multiple possible values -

Usefull Methods

Dont use - Apply Method

- apply()

apply our custom function to
each row of
the series -

```
def testfour(num):
```

```
    str(1234567890)[4:0]  
    return str(num)[4:0]
```

```
dfflat = df[["CC.Number"]].apply(testfour)
```

```
df["total_bill"].mean()
```

19.785

```
def yelp(price):  
    if price < 10:  
        return '$'  
    elif price >= 10 and price < 30:  
        return '$$'  
    else:  
        return '$$$'
```

```
return '$$$$'
```

```
df['elp'] = df['total_bill'].apply(yelp)
```

Part #1 2

Apply on Multiple Columns:-

Lambda function

```
def simple(num):
```

```
    return num * 2
```

↳ to lambda

```
lambda num: num * 2
```

```
}  
df['total_bill'].apply(lambda num: num * 2)
```

apply with multiple columns:-

```
def quality(total_bill, tip):
```

```
    if tip / total_bill > 0.25:  
        return "generous"
```

```
    else:
```

```
        'others'
```

```
}  
quality(16.99, 1.01)
```

```
df[['total_bill', 'tips']].apply(lambda  
    df: quality(df['total_bill'], df['tips'])  
    , axis=1)
```

Run alot faster:-

```
df[['Quality']] = np.vectorize(quality)  
(df['total_bill'] * df['tips'])
```

Appreciate just a convenience
import tips

Useful methods part #3

Describing and Sorting methods

```
df = pd.read_csv('tips.csv')
```

```
df.describe()
```

```
df.sort_values('tip')
```

a default ascending

```
df.sort_values('size', 'size')
```

first it sort by tip then
if they are same then if
sort by size column.

```
df['total_bill'].max() .min()
```

```
df['total_bill'].idxmax() .idxmin()  
170
```

```
df.iloc[170]
```

df. corr() — correlation
between
numerical column

df[1:sex1], value=counts()

on column

df[1:day1].unique()

↳ unique values

df[1:day1].nunique()

df.head()

df['sex'].replace({'Female': 'F'})

df['sex'].replace({'Female': 'M'}, {'M': 'F'})

mymap = {'Female': 'F', 'Male': 'M'}

df['sex'].map(mymap)

↳ lots of item

df.duplicated() — duplicates()

Simple - df.drop_duplicates()

df[['total_bill']].between(10, 20, inclusive='True')

df.nlargest(10, 'Hps')

df.nsmallest(2, 'Hps')

df.sample(5) — return random rows

df.sample(frac=0.1) will return 10% of the sample

Missing Data

- When reading in missing values, pandas will display them as NAN values -
- There are also newer versions specified null pandas value such as pd.NAT to imply the value missing should be a timestamp.
 - Option for missing Data
 - Keep it
 - Remove it
 - Replace it
- Keep the missing data
 - drop
 - Enforce to data
 - Does not manipulate data
 - Change the true data

cons:
many methods does not support
NaN

* Often there are reasonable guesses

o Dropping or Removing the missing data
easy to do
- can be based on Rule

• cons:

- potential to lose a lot of data or useful information -
- limit trained data on Future prediction

• Removing or Dropping missing data -

- dropping a row makes sense when info is missing -

Filling the missing data:-

• pros:

- potential to save a lot of data for use in training a model -
-

• cons:

- Hardest to do and somewhat arbitrary -
- potential to lead to false conclusion

Missing Data

part - two

np.nan — old version
pd.NA

pd.NAT → time stamp

np.nan == np.nan
False

np.nan is np.nan
True

df = pd.read_csv('movies_score.csv')

df.head()

if

df.isnull() → Show null value
df.notnull() → Show not values
df[df['pre-movie score'].notnull()]

if

df['pre-movie score'].isnull() & (df['pre-movie score'].notnull())

df[(df['pre-movie score'].isnull()) & (df['pre-movie score'].notnull())]

```
# keep the data  
# Fill DATA  
# Fill DATA
```

df.dropna(axis=1)
at least not
df.dropna(thresh=1)
value

df.dropna(axis=1)

df.dropna(subset=['last name'])

Fill the Data

df.fillna('New Value! : ')

df['pre-money'].fillna(0)

df

df['pre-money'].mean()

df['pre-money'].fillna(df.mean())

Interpolation to fill missing
values :-

airline_fix = { 'first': 10, 'business': np.nan,
 'economy_plus': 30 }
Sex = pd.Series(airline_fix)
Sex

See

• `Series.interpolate()`

14. Groupby Operations

- A `groupby()` operation allows us to examine data on a per category basis =

category	DataValue
A	10
	5
B	4
	12
C	6

→ use on aggregate function
`sum()`

category	Result
A	15
B	6
C	18

• `mean()`
• `count()`

• Note: in pandas calling `groupby()` by itself creates a "Lazy" groupby object waiting to be evaluated by an aggregate method call -

```
df = pd.read("mpg.csv")
```

```
df['model-year'].value_counts()
```

```
df.groupby('model-year').mean()
```

```
df.groupby(['model-year', 'cylinder']).mean()
```

```
df.groupby(['model-year', 'cylinder']).mean().index
```

```
df.groupby(['model-years']).describe().transpose()
```

→ `year_cyl = df.groupby(['model-year', 'cylinder']).mean()`

```
year_cyl.index.names
```

```
[70, 82]]
```

```
year_cyl.loc[70]
```

```
year_cyl.loc[(70, 4)]
```

Groupby

multilevel index

Cross-section

```
year_cyl_xs(key=70, level='model-year')
```

```
year_cyl_xs([70, 80])
```

```
year_cyl_xs(key=4, levels='cylinders')
```

```
df[df['cylinder'].isin([6, 8])].
```

```
groupby(['model-year', 'cylinder'])  
    .mean()
```

```
year_cyl.swaplevel().swaplevel()
```

```
year_cyl.sort_index(level='model year',  
                     ascending=False)
```

```
df.agg(['std', 'mean']).T[m pg!]
```

```
df.agg({'mpg': ['max', 'mean'], 'weight':  
        ['mean', 'std']})
```

16 - Combining Dataframe - Concatenation

- Often the data needs to exist in two separate sources, fortunately pandas make it easier to ~~concatenate~~ combine together using pd.concat() call.

Concatination is simply "past[ing]" the two Dataframe together, by columns:

	Year	Pop	GDP	perc
USA	1776	326	USA	2.65
CAN	1867	38	CAN	1.7
MEX	1821	726	MEXICO	1.22

	Year	Pop	GDP	Percent
USA	1776	326	USA	25%
CANADA	1867	38	CANADA	NAN
MEXICO	1821	726	MEXICO	25%

Also can be performed on rows =

MEX	1.7
USA	2.65
CANADA	1.22

Code #

```
dataF1 =  
dataF2 =
```

	A	B	C	D
0	A0	B0	C0	D0
1	A1	B1	C1	D1
2	A2	B2	C2	D2
3	A3	B3	C3	D4

pd.concat([one,two],axis=1)

pd.concat([one,two],axis=0)

↓
One way to resolve it was to
rename

two.columns = one.columns

↓
pd.concat([one,two],axis=0)

④ mydf = pd.concat
mydf.index = range(len(df))

Combining Data Frame

Inner Merge - 17

- merge()

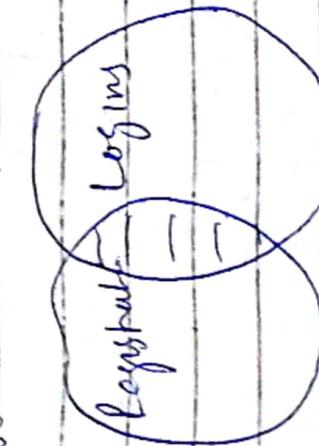
The merge method takes in a key argument labeled 'how' -

- o There are 3 main ways of merging tables together using how parameter,
 - o Inner
 - o Outer
 - o Left or Right

Registrations		Logins	
reg_id	name	log_id	name
1	André	2	André
2	Bob	3	Yolanda
3	Charlie	4	Bob
4	Don	5	

on = "name".

how = "inner"



pd.merge(registrations, logins, how = 'inner',
on = 'name')

Result

reg_id	name	log_id
1	Andrew	2
2	Bob	4

```
registrations = pd.DataFrame([{"reg_id": 1, "name": "Bob"}, {"reg_id": 2, "name": "Andrew"}])
```

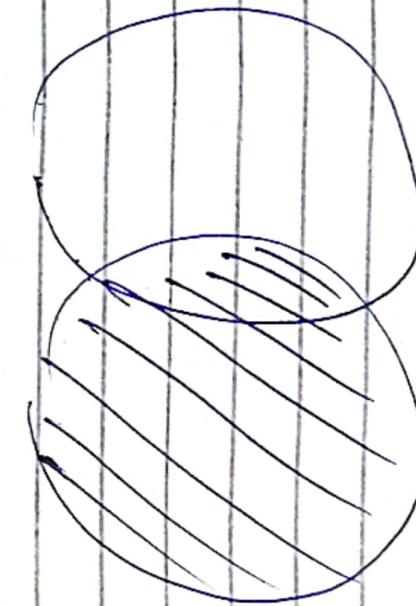
```
pd.merge(registrations, logins, how = "inner", on = "name")
```

reg_id	name	log_id
0	Bob	4
1	Andrew	2

Left - Right (18) merge

- Order matters in different pd.merge()

left merge



everything in left here

column 'name' shows

pd.merge(left, logins, how = "left", on = "name")

Lesson 1

reg_id name / log_id

1	Brienne	2
2	Bob	4
3	Charlie	NaN
4	Daniel	NaN

Right -

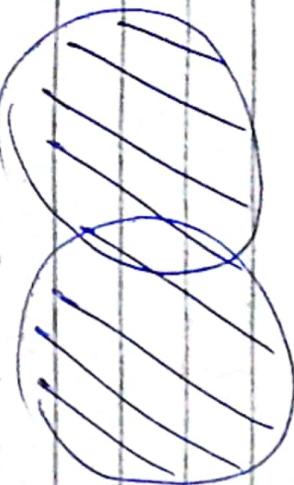
pd.merge(left, right, how='right',
on='name')

Result

	name	2
1	Andrew	4
2	Bob	1
3	Xavier	
4	NaN	3
5	Yorick	

Outer-Merge (1q)

Setting `how = "outer"` allow as
to include everything present in
both table



A: All the
values

order in outer will not matter -

pd.merge(lefts, registrations, how='outer',
on='name')

Registrations: registrations.set_index('name')
registrations

pd.merge(lefts, registrations,

left_index=True,
right_on='name',
how='inner')

registrations.columns = ['reg-name', 'reg-id']

for different column names

pd.merge(lefts, registrations, how='inner',
right_on='name',
left_on='reg-name',
result_drop=True, axis=1)

for same columns

pd.merge(lefts, how='inner',
on='name')

Suffixes = ('_reg', '_log')

20. Text Methods for string Data

- Often text data needs to be cleaned or manipulated for processing.
 - While we can always use a custom apply() function for these tasks, pandas comes with many built-in string methods calls.
 - Let's learn how to use them!

Code #

```
import pandas as pd  
import numpy as np
```

```
name = 'jose'
```

```
email = 'jose@mail.com'  
email.split('@')
```

```
name = pd.Series(['andrew', 'bobo', 'claire',  
                  'david', 'ss', 'zz'])
```

```
names.str.upper()
```

```
email.isdigit()
```

```
names.str.isdigit()
```

```
tech-finance = ['GOOG', 'APPL', 'AMZN',  
                'FB', 'BAC', 'GS']
```

```
tech-finance  
len(tech-finance)  
tickers = pd.Series(tech_finance)
```

```
tickers.str.split(',')  
str[0]
```

```
tickers.str.split(',', expand=True)
```

```
messy_names = pd.Series(['Andrew',  
                        'Bob', 'Bob', 'Clair', 'Clair'])
```

```
messy_names.str.replace(' ', '_').str.strip()  
str.capitalize()
```

21 - Time Methods in pandas

- Basic Python have a date-time object.
- Containing date and time information
- pandas allow us to easily extract information from a date-time object to use feature engineering

- dt methods

Extract information from the timestamp, such as:

- Day of the week
- Weekend vs weekday
-

from datetime import datetime

myyear = 2015

mymonth = 1

myday = 1

myhour = 2

mymin = 30

mysec = 15

mydate = datetime(myyear, mymonth, myday)

mydate

myser = pd.Series(['Nov 3, 1990',
'2000-01-01', None])

[]

`pd.to_datetime(myse0)`

0 1990-11-03
1 2000-01-01
2 NaT

`time series of year`

`euro_date = '31-12-2000'`

`pd.to_datetime(euro_date)`

`pd.to_datetime(euro_date,format=True)`

`style_date = '12--Dec--2000'`

`pd.to_datetime(style_date,format`

`= '%d--%b--%Y')`

`custom_date = "12th of Dec 2000"`

`pd.to_datetime(custom_date)`

`sales = pd.read_csv('sales.csv')`

22 - Pandas Input and Output

- Note #
 - You need to know exactly directory location and correct files name -
 - You may need passwords (or permission for certain data inputs (e.g. a SQL database password)).

Video lecture

- CSV files
- HTML Tables
- Excel files
- SQL Database

Code # :

```
import pandas  
pd, ls
```

```
import os  
os.getcwd()
```

df = pd.read_csv('example.csv')

df.to_csv('C:\Users\Maravalli\newfiles.csv')

```
df.to_csv('newfiles.csv', index=False)
```

23 - Pandas HTML Table

Website displays tabular information through the use of HTML tables.

- <table>
 - pandas has the ability to automatically convert these HTML tables into a DataFrame.
 - Note every table in a website is dividable through HTML tables

```
import pandas as pd  
#and install lxml  
pip install lxml
```

```
url = "https://...  
pd.read_html(url)
```

```
table = pd.read_html(url)  
len(tables)
```

```
tables[1]
```

```
worldTables[0]
```

```
worldTables[1]
```

```
worldTables
```

```
worldTables = worldTables[1]  
axis=1
```

world_topten = pd.read_csv('2000_1_2015_2030_Census_top10.csv')

```
world_topten
```

world_topten.to.html('Samptable.html',
index=False)

24 - Excel files

• pandas can read and write excel files -

Important Note!

- pandas can read and write in raw data . it is not able to read in across , visualization , or formulas : created inside of spreadsheet.

- pandas treats an Excel Workbook as a dictionary with the key being the sheet name and the value being the Dataframe representing the sheets itself

Note! Using pandas with excel requires additional libraries

- openpyxl → xlst
- xlrd . xls

lecture #1

```
df = pd.read_excel('my.xlsx',  
                   sheet_name='firstsheet')
```

```
nb: pd.ExcelFile('my.xlsx')
```

```
with sheet_names  
    ['firstsheet']
```

```
excel_sheets:dict = pd.read_excel('my.xlsx')
```

```
sheet_name=None
```

if set to

None then

Create a dictionary
and saves everything
to it

```
type(excel_sheet_dict)
```

```
excel_sheets.keys()
```

```
our_df.to_excel('example.xlsx', sheetname  
                ='firstsheet', index=False)
```

25 - Pandas Input and Output SQL - DataBase

- o Pandas can read and write to various SQL engines through use of a drivers and the sqlalchemy Python library -
 - o See how ~~XXXXXXXXXX~~ this works?

Step 1:

- o Figure out what SQL engine you are desiring connecting to, for just a few examples :-
 - o PostgreSQL
 - o MySQL
 - o MS SQL Server

Step 2:

- o Install appropriate Python driver library (most likely requires a Google search);
 - PostgreSQL - psycopg2
 - MySQL - PyMySQL
 - MS SQL Server - pyodbc

Step 3:

Use sqlalchemy library to connect to your SQL database with the drivers.

Step 4:

- Use the SQLAlchemy driver connection with pandas `read_sql` method.
 - pandas can read in entire tables.
- a DataFrame or actual `pyspark` or SQL query through the connection:
 - `SELECT * FROM table;`

Using SQLite we can easily install database



```
pip install SQLAlchemy  
# temp database inside ram
```

```
from sqlalchemy import create_engine  
temp_db = create_engine('sqlite:///memory:')
```

~~Step~~ pd.DataFrame(data=temp_db.read_sql)

```
(low=0, high=100,  
size=(4,4))  
columns=['a','b','c','
```

```
df.toSQL(name='new-table', con=temp_db)
```

```
pd.read_sql('new-table', con=temp_db)
```

New Table A b c d

0	0	0	0	80
1	1	1	1	98
2	2	2	2	63

```
pd.read_sql_query(sql = 'SELECT a,c FROM  
newtable', con = temp_db)  
result =
```

26 - Pivot Tables

- Allows you to reorganize data, refactoring cell based on columns and a new index
- This is best shown visually....

- A DataFrame with repeated value can be pivoted for a reorganization and clarity -
`df.pivot(index='foo', columns='bar', values='baz')`

Code #

```
pd.read_csv('sales-funnel.CRM.csv')
```

`licences = df[['Company', 'product', 'License']]`

```
pd.pivot(index=licences, index='Company',  
columns='products', values='Licence')
```

```
pd.pivot_table(df, index='Company', aggfunc='sum',  
values='Sales', sales_price)  
values = pd['Sales', 'Sales.Price']
```

```
pd.pivot_table(df, index='Company',  
values='Sales', sales_price)  
aggfunc='sum', fill_value=0)  
columns='products', values='Licence')
```

Practice Exercise:-

- 1) How many rows are there?
- 2) Is there missing data? If so which column has the most missing data?
- 3) Drop the "Company" column from the dataset?
- 4) What are the most common types of country codes in dataset?
- 5) What is the name of person who paid the highest total charge (dairy rate)? How much was their AGR?
- 6) The average daily rate for a person is stay at hotel. What is the mean and across all the hotel stays in the dataset?
- 7) What is the average number of (mean) night for a stay across the entire data set, feel free to round this to 3 decimal points.

- 8) What is the average total cost for a stay in dataset? Not average daily cost, but total stay cost. Feel free to round this to 2 decimal points.
- 9) What are the names and emrks of people who made 5 "special requests"
- 10) What percentage of hotel stay were classified as "repeat guest"?
- (Do not sort this on name of person but instead of 15 - repeated guest - column)
- 11) What are the top 5 common last name in the dataset?
- (a) What are the names of the people who have booked the most number of children and babies for their stay? (Don't worry about cancelled just r.)
- (b) Top 5 most common area codes in phone numbers? (First 3 digit)

14) How many arrivals took place between first and the 15th of the month, inclusive of 1 and 15? can you do it in one line of pandas code?

Hard Bonus Task

15) Create a table for counts for each day of week that people arrived.
(e.g.) Good arrival were on a Monday, 3000 were on Tuesday, etc.