

# Introduction to Data Science

## Lecture 8

### Data Wrangling, preprocessing, and Visualization

#### Part I



### Data Preprocessing (cleaning)

1. Reading the file in a structural format (excel type format)
2. Dealing with missing values
3. Finding outliers and dealing with them
4. normalization
5. possible visualization

### Reading and Handling Data files using Pandas

Function	Description
<code>read_csv</code>	Load delimited data from a file, URL, or file-like object; use comma as default delimiter
<code>read_table</code>	Load delimited data from a file, URL, or file-like object; use tab (' <code>\t</code> ') as default delimiter
<code>read_fwf</code>	Read data in fixed-width column format (i.e., no delimiters)
<code>read_clipboard</code>	Version of <code>read_table</code> that reads data from the clipboard; useful for converting tables from web pages
<code>read_excel</code>	Read tabular data from an Excel XLS or XLSX file
<code>read_hdf</code>	Read HDF5 files written by pandas
<code>read_html</code>	Read all tables found in the given HTML document
<code>read_json</code>	Read data from a JSON (JavaScript Object Notation) string representation
<code>read_msgpack</code>	Read pandas data encoded using the MessagePack binary format
<code>read_pickle</code>	Read an arbitrary object stored in Python pickle format
<code>read_sas</code>	Read a SAS dataset stored in one of the SAS system's custom storage formats
<code>read_sql</code>	Read the results of a SQL query (using SQLAlchemy) as a pandas DataFrame
<code>read_stata</code>	Read a dataset from Stata file format
<code>read_feather</code>	Read the Feather binary file format

## Working on AirPassengers data set:

- The classic Box & Jenkins airline data. Monthly totals of international airline passengers, 1949 to 1960.
- source AARSHAY JAIN

In [2]:

```
1 import pandas as pd
2 import numpy as np
3 import datetime as dt
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 from matplotlib.pyplot import rcParams
7 rcParams['figure.figsize'] = 15, 6
```

Now, we can load the data set and look at some initial rows and data types of the columns:

```
In [16]: 1 filepath = '/Users/martin/Documents/MyLecturesSBU/Fall2018/CSE391/data/2
2 data = pd.read_csv(filepath)
3 data.head(10)
4 #print('\n Data Types:')
5 #print(data.dtypes)
```

```
Out[16]:
```

	Month	#Passengers
0	1949-01	112
1	1949-02	118
2	1949-03	132
3	1949-04	129
4	1949-05	121
5	1949-06	135
6	1949-07	148
7	1949-08	148
8	1949-09	136
9	1949-10	119

```
In [17]: 1 ?pd.read_csv
```

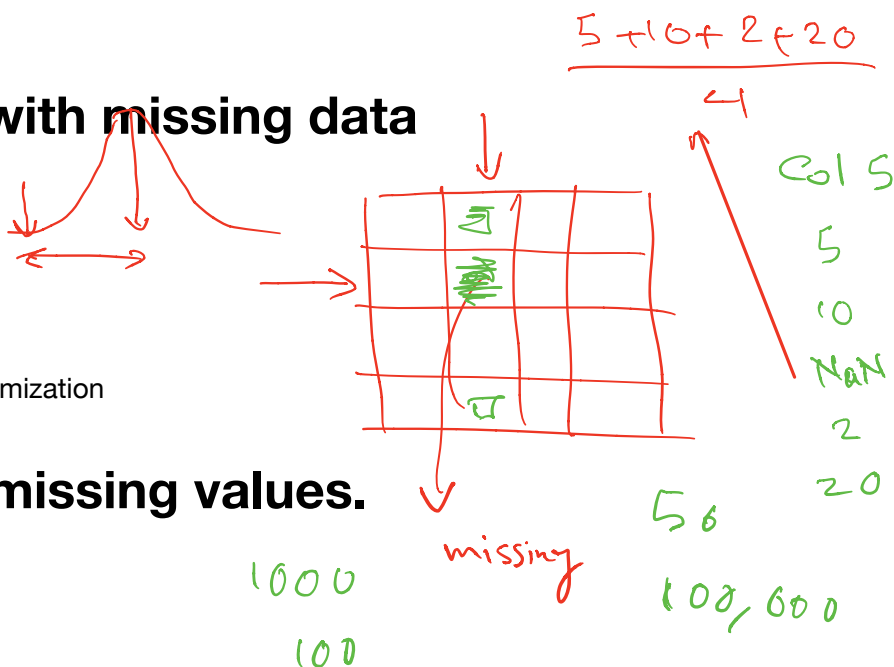
## Dealing with missing values

### Two type of missing data

- missing completely at random
- missing not at random

### How to deal with missing data

- drop out
- imputation
  - expert gussing
  - averaging
  - regression
  - Expectation maximization



### 1.1 Drop the missing values.

```

In [7]: 1 filepath = '/Users/martin/Documents/MyLecturesSBU/Fall2018/CSE391/data/
2         titanic_train.csv'
3         passengers = pd.read_csv(filepath)
4         passengers.dropna(axis=0, how='any', thresh=None, subset=None,
5                           inplace=False)
6         passengers.describe()
7         passengers.head()

```

100,000  
Kaggle

```

Out[7]:

```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN

- This will drop any rows with missing values. Clearly this isn't a good idea
- What instead if we wanted to remove any columns with missing values? How?

```

In [19]: 1 passengers.dropna(axis=1).head(5)

```

```

Out[19]:

```

	PassengerId	Survived	Pclass	Name	Sex	SibSp	Parch	Ticket	Fare
0	1	0	3	Braund, Mr. Owen Harris	male	1	0	A/5 21171	7.2500
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	1	0	PC 17599	71.2833
2	3	1	3	Heikkinen, Miss. Laina	female	0	0	STON/O2. 3101282	7.9250
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	1	0	113803	53.1000
4	5	0	3	Allen, Mr. William Henry	male	0	0	373450	8.0500

```

In [ ]: 1

```

```
In [ ]: 1 # Keep only the rows with at least 11 non-na values:
        2 passengers.dropna(thresh = 11)
```

## 1.2 imputation ( a way to replace the missing values)

- Filling the missing values

- mean

- your suggestion?

median

median

```
In [7]: 1 passengers["Age"].fillna(value=passengers["Age"].mean()).head()
```

```
Out[7]: 0    22.0
        1    38.0
        2    26.0
        3    35.0
        4    35.0
        Name: Age, dtype: float64
```

## Summarizing and Computing Descriptive Statistics

Pandas objects are equipped with a set of common mathematical and statistical methods. Most of these fall into the category of reductions or summary statistics, methods that extract a single value (like the sum or mean) from a Series or a Series of values from the rows or columns of a DataFrame.

```
In [9]: 1 df = pd.DataFrame([[1.4, np.nan], [7.1, -4.5], [np.nan, np.nan], [0.75,
        2 .....: columns=['one', 'two']])
        3 print(df)
```

```
   one  two
a  1.40 NaN
b  7.10 -4.5
c   NaN NaN
d  0.75 -1.3
```

- Calling DataFrame's sum method returns a Series containing column sums

```
In [11]: 1 df.sum()
```

```
Out[11]: one    9.25
         two   -5.80
         dtype: float64
```

- computing mean

```
In [ ]: 1 df.mean(axis='columns', skipna=False)
```

m

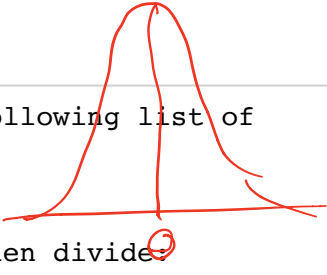
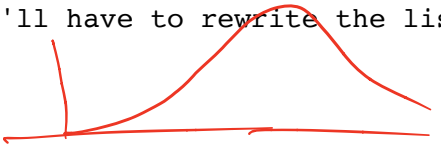
Method	Description
<u>count</u>	Number of non-NA values
<u>describe</u>	Compute set of summary statistics for Series or each DataFrame column
<u>min, max</u>	Compute minimum and maximum values
<u>argmin, argmax</u>	Compute index locations (integers) at which minimum or maximum value obtained, respectively
<u>idxmin, idxmax</u>	Compute index labels at which minimum or maximum value obtained, respectively
<u>quantile</u>	Compute sample quantile ranging from 0 to 1
<u>sum</u>	Sum of values
<u>mean</u>	Mean of values
<u>median</u>	Arithmetic median (50% quantile) of values
<u>mad</u>	Mean absolute deviation from mean value
<u>prod</u>	Product of all values
<u>var</u>	Sample variance of values
<u>std</u>	Sample standard deviation of values
<u>skew</u>	Sample skewness (third moment) of values
<u>kurt</u>	Sample kurtosis (fourth moment) of values
<u>cumsum</u>	Cumulative sum of values
<u>cummin, cummax</u>	Cumulative minimum or maximum of values, respectively
<u>cumprod</u>	Cumulative product of values
<u>diff</u>	Compute first arithmetic difference (useful for time series)
<u>pct_change</u>	Compute percent changes

## Note on mean, median, mode and their differences:

```

1 Find the mean, median, mode, and range for the following list of
  values:
2 13, 18, 13, 14, 13, 16, 14, 21, 13
3
4 The mean is the usual average, so I'll add and then divide
5
6  $(13 + 18 + 13 + 14 + 13 + 16 + 14 + 21 + 13) \div 9 = 15$ 
7
8 Note that the mean, in this case, isn't a value from the original
  list. This is a common result. You should not assume that your mean
  will be one of your original numbers.
9
10 The median is the middle value, so first I'll have to rewrite the list
    in numerical order:
11
12 13, 13, 13, 13, 14, 14, 16, 18, 21
13
14 There are nine numbers in the list, so the middle one will be the  $(9 +$ 
     $1) \div 2 = 10 \div 2 = 5$ th number:
15

```

```

16 13, 13, 13, 13, 14, 14, 16, 18, 21
17
18 So the median is 14.
19
20 The mode is the number that is repeated more often than any other, so
    13 is the mode.
21
22 The largest value in the list is 21, and the smallest is 13, so the
    range is  $21 - 13 = 8$ .
23
24 mean: 15
25 median: 14
26 mode: 13
27 range: 8
28
29 Note: The formula for the place to find the median is " $([the\ number\ of\ data\ points] + 1) \div 2$ ", but you don't have to use this formula. You
    can just count in from both ends of the list until you meet in the
    middle, if you prefer, especially if your list is short. Either way
    will work.

```

## Reading Text Files in Pieces

- When processing very large files or figuring out the right set of arguments to correctly process a large file, you may only want to read in a small piece of a file or iterate through smaller chunks of the file.

```

In [ ]: 1 # If you want to only read a small number of rows (avoiding reading the
        2 pd.read_csv('examples/ex6.csv', nrows=5)
        3 # To read a file in pieces, specify a chunksize as a number of rows:
        4 chunker = pd.read_csv('examples/ex6.csv', chunksize=1000)
        5
        6 # The TextParser object returned by read_csv allows you to iterate over
        7 #the value counts in the 'key' column like so:
        8 chunker = pd.read_csv('examples/ex6.csv', chunksize=1000)
        9 tot = pd.Series([])
       10 for piece in chunker:
       11     tot = tot.add(piece['key'].value_counts(), fill_value=0)
       12     tot = tot.sort_values(ascending=False)

```

## Pivot tables in Pandas

### Panada pivot table example

- Adapted from Chris Moffitt (<http://pbpython.com/pandas-pivot-table-explained.html>).

```

In [2]: 1 import pandas as pd

```

In [5]:

```

1 # Read in the data; sale pipeline data
2 df = pd.read_excel("/Users/martin/Documents/MyLecturesSBU/past/
3                      Spring2018/CSE391/data/sales-funnel.xlsx")
4 df

```

3	737550	Fritsch, Russel and Anderson	Craig Booker	Debra Henley	CPU	1	35000	declined
4	146832	Kiehn-Spinka	Daniel Hilton	Debra Henley	CPU	2	65000	won
5	218895	Kulas Inc	Daniel Hilton	Debra Henley	CPU	2	40000	pending
6	218895	Kulas Inc	Daniel Hilton	Debra Henley	Software	1	10000	presented
7	412290	Jerde-Hilpert	John Smith	Debra Henley	Maintenance	2	5000	pending
8	740150	Barton LLC	John Smith	Debra Henley	CPU	1	35000	declined
9	141962	Herman LLC	Cedric Moss	Fred Anderson	CPU	2	65000	won
10	163416	Purdy-Kunde	Cedric Moss	Fred Anderson	CPU	1	30000	presented
11	239344	Stokes LLC	Cedric	Fred	Maintenance	1	5000	pending

In [7]:

```

1 # This isn't strictly required but helps us keep the order we want as
2 #we work through analyzing the data.
3
4 df["Status"] = df["Status"].astype("category")
5 df["Status"].cat.set_categories(["won", "pending", "presented", "declined"],
6                                inplace=True)
7

```



```
In [8]: 1 #The simplest pivot table must have a dataframe and an index .
        2 #In this case, let's use the Name as our index.
        3 pd.pivot_table(df,index=["Name"])
```

<b>Barton LLC</b>	740150.0	35000.0	1.000000
<b>Fritsch, Russel and Anderson</b>	737550.0	35000.0	1.000000
<b>Herman LLC</b>	141962.0	65000.0	2.000000
<b>Jerde-Hilpert</b>	412290.0	5000.0	2.000000
<b>Kassulke, Ondricka and Metz</b>	307599.0	7000.0	3.000000
<b>Keeling LLC</b>	688981.0	100000.0	5.000000
<b>Kiehn-Spinka</b>	146832.0	65000.0	2.000000
<b>Koepp Ltd</b>	729833.0	35000.0	2.000000
<b>Kulas Inc</b>	218895.0	25000.0	1.500000
<b>Purdy-Kunde</b>	163416.0	30000.0	1.000000
<b>Stokes LLC</b>	239344.0	7500.0	1.000000
<b>Trantow-Barrows</b>	714466.0	15000.0	1.333333

```
In [9]: 1 # You can have multiple indexes as well
        2 pd.pivot_table(df,index=["Name", "Rep", "Manager"])
```

```
Out[9]:
```

			Account	Price	Quantity
	Name	Rep	Manager		
	<b>Barton LLC</b>	<b>John Smith</b>	<b>Debra Henley</b>	740150.0	35000.0 1.000000
	<b>Fritsch, Russel and Anderson</b>	<b>Craig Booker</b>	<b>Debra Henley</b>	737550.0	35000.0 1.000000
	<b>Herman LLC</b>	<b>Cedric Moss</b>	<b>Fred Anderson</b>	141962.0	65000.0 2.000000
	<b>Jerde-Hilpert</b>	<b>John Smith</b>	<b>Debra Henley</b>	412290.0	5000.0 2.000000
	<b>Kassulke, Ondricka and Metz</b>	<b>Wendy Yule</b>	<b>Fred Anderson</b>	307599.0	7000.0 3.000000
	<b>Keeling LLC</b>	<b>Wendy Yule</b>	<b>Fred Anderson</b>	688981.0	100000.0 5.000000
	<b>Kiehn-Spinka</b>	<b>Daniel Hilton</b>	<b>Debra Henley</b>	146832.0	65000.0 2.000000
	<b>Koepp Ltd</b>	<b>Wendy Yule</b>	<b>Fred Anderson</b>	729833.0	35000.0 2.000000
	<b>Kulas Inc</b>	<b>Daniel Hilton</b>	<b>Debra Henley</b>	218895.0	25000.0 1.500000
	<b>Purdy-Kunde</b>	<b>Cedric Moss</b>	<b>Fred Anderson</b>	163416.0	30000.0 1.000000
	<b>Stokes LLC</b>	<b>Cedric Moss</b>	<b>Fred Anderson</b>	239344.0	7500.0 1.000000
	<b>Trantow-Barrows</b>	<b>Craig Booker</b>	<b>Debra Henley</b>	714466.0	15000.0 1.333333

```
In [10]: 1 pd.pivot_table(df,index=["Manager","Rep"])
```

```
Out[10]:
```

		Account	Price	Quantity
Manager	Rep			
	<b>Craig Booker</b>	720237.0	20000.000000	1.250000
<b>Debra Henley</b>	<b>Daniel Hilton</b>	194874.0	38333.333333	1.666667
	<b>John Smith</b>	576220.0	20000.000000	1.500000
	<b>Cedric Moss</b>	196016.5	27500.000000	1.250000
<b>Fred Anderson</b>	<b>Wendy Yule</b>	614061.5	44250.000000	3.000000

```
In [11]: 1 # the Account and Quantity columns aren't really useful.
2 #Let's remove it by explicitly defining the columns
3 #we care about using the values field.
4
5 pd.pivot_table(df,index=["Manager","Rep"],values=["Price"])
```

```
Out[11]:
```

		Price
Manager	Rep	
	<b>Craig Booker</b>	20000.000000
<b>Debra Henley</b>	<b>Daniel Hilton</b>	38333.333333
	<b>John Smith</b>	20000.000000
	<b>Cedric Moss</b>	27500.000000
<b>Fred Anderson</b>	<b>Wendy Yule</b>	44250.000000

```
In [13]: 1 import numpy as np
2 # The price column automatically averages the data
3 # but we can do a count or a sum. Adding them is simple
4 #using aggfunc and np.sum .
5 pd.pivot_table(df,index=["Manager","Rep"],
6                 values=["Price"],aggfunc=np.sum)
```

```
Out[13]:
```

		Price
Manager	Rep	
	<b>Craig Booker</b>	80000
<b>Debra Henley</b>	<b>Daniel Hilton</b>	115000
	<b>John Smith</b>	40000
	<b>Cedric Moss</b>	110000
<b>Fred Anderson</b>	<b>Wendy Yule</b>	177000

- if we want to see sales broken down by the products, the columns variable allows us to define one or more columns.

```
In [14]: 1 table = pd.pivot_table(df,index=["Manager","Rep"],values=["Price"],
2           columns=["Product"],aggfunc=[np.sum])
3 table
```

```
Out[14]:
```

			sum		
				Price	
	Product	CPU	Maintenance	Monitor	Software
Manager	Rep				
	Craig Booker	65000.0	5000.0	NaN	10000.0
Debra Henley	Daniel Hilton	105000.0	NaN	NaN	10000.0
	John Smith	35000.0	5000.0	NaN	NaN
	Cedric Moss	95000.0	5000.0	NaN	10000.0
Fred Anderson	Wendy Yule	165000.0	7000.0	5000.0	NaN

- Remember, columns are optional - they provide an additional way to segment the actual values you care about. The aggregation functions are applied to the values you list.

## Join and Merge Pandas Dataframe

- adapted from [Chris Albon](https://chrisalbon.com/python/data_wrangling/pandas_join_merge_dataframe/)  
([https://chrisalbon.com/python/data\\_wrangling/pandas\\_join\\_merge\\_dataframe/](https://chrisalbon.com/python/data_wrangling/pandas_join_merge_dataframe/))



```
In [ ]: 1 #Create a dataframe
2 raw_data = {
3     'subject_id': ['1', '2', '3', '4', '5'],
4     'first_name': ['Alex', 'Amy', 'Allen', 'Alice', 'Ayoung'],
5     'last_name': ['Anderson', 'Ackerman', 'Ali', 'Aoni', 'Atiches']]
6 df_a = pd.DataFrame(raw_data, columns = ['subject_id', 'first_name', 'last_name'])
7 df_a
```

```
In [ ]: 1 # Create a second dataframe
2 raw_data = {
3     'subject_id': ['4', '5', '6', '7', '8'],
4     'first_name': ['Billy', 'Brian', 'Bran', 'Bryce', 'Betty'],
5     'last_name': ['Bonder', 'Black', 'Balwner', 'Brice', 'Btisan']]
6 df_b = pd.DataFrame(raw_data, columns = ['subject_id', 'first_name', 'last_name'])
7 df_b
```

```
In [ ]: 1 # create a third dataframe
2 raw_data = {
3     'subject_id': ['1', '2', '3', '4', '5', '7', '8', '9', '10', '11'],
4     'test_id': [51, 15, 15, 61, 16, 14, 15, 1, 61, 16]}
5 df_n = pd.DataFrame(raw_data, columns = ['subject_id', 'test_id'])
6 df_n
```

- Join the two dataframes along rows

```
In [ ]: 1 df_new = pd.concat([df_a, df_b])
        2
        3 df_new
```

- Join the two dataframes along columns

```
In [ ]: 1 pd.concat([df_a, df_b], axis=1)
```

## Reading big files using Pandas

```
In [ ]: 1 import pandas as pd
        2 import numpy as np
        3 #SQLAlchemy is the Python SQL toolkit and Object Relational Mapper
        4 # that gives application developers the full power and
        5 #flexibility of SQL.
        6 from sqlalchemy import create_engine
```

```
In [ ]: 1 file = '/Users/martin/Documents/TCGA/tcga_mRNA/COAD_mRNA/COAD.uncv2.mRNA
        2
        3
        4 #This command uses pandas' "read_csv" command to read
        5 #in only 5 rows (nrows=5)
        6 # and then print those rows to the screen.
        7 pd.read_csv(file, sep='\t', nrows=1)
```

- Before we can actually work with the data, we need to do something with it so we can begin to filter it to work with subsets of the data. This is usually what I would use pandas' dataframe for but with large data files, we need to store the data somewhere else. In this case, we'll set up a local sqllite database, read the csv file in chunks and then write those chunks to sqllite.
- To do this, we'll first need to create the sqllite database using the following command.

```
In [ ]: 1 csv_database = create_engine('sqlite:///csv_database.db')
```

```
In [ ]: 1 #Next, we need to iterate through the CSV file in chunks and store the c
        2
        3 chunksize = 1000
        4 i = 0
        5 j = 1
        6 for df in pd.read_csv(file, chunksize=chunksize, iterator=True):
        7     df = df.rename(columns={c: c.replace(' ', '') for c in df.columns})
        8     df.index += j
        9     i+=1
       10     df.to_sql('table', csv_database, if_exists='append')
       11     j = df.index[-1] + 1
```

```
In [ ]: 1 # To access the data now, you can run commands like the following:
        2
        3 df = pd.read_sql_query('SELECT * FROM "table"', csv_database)
        4
```

```
In [ ]: 1
```