# Python IV – Lesson 20

Date: Nov 13, 2022

# Agenda

► Google Colab
► Pandas
► Numpy
► Matplotlib

# Proverbs 11:21

- " Be sure of this: The wicked will not go unpunished, but those who are righteous will go free."

# Review + Homework

https://colab.research.google.com/drive/1DF7sFWdp5PfuM_UmxkbHSCw5MPoQx5va#scrollTo=u6OCNcz_gKws

# Pandas Cheatsheet

## Data Wrangling
with pandas Cheat Sheet
http://pandas.pydata.org

Pandas API Reference   Pandas User Guide

## Tidy Data – A foundation for wrangling in pandas

In a tidy data set:

Each **variable** is saved in its own **column**

Each **observation** is saved in its own **row**

Tidy data complements pandas's vectorized operations. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

$M * A$

## Creating DataFrames

```python
df = pd.DataFrame(
        {"a" : [4, 5, 6],
         "b" : [7, 8, 9],
         "c" : [10, 11, 12]},
        index = [1, 2, 3])
```
Specify values for each column.

```python
df = pd.DataFrame(
        [[4, 7, 10],
         [5, 8, 11],
         [6, 9, 12]],
        index=[1, 2, 3],
        columns=['a', 'b', 'c'])
```
Specify values for each row.

```python
df = pd.DataFrame(
        {"a" : [4 ,5, 6],
         "b" : [7, 8, 9],
         "c" : [10, 11, 12]},
index = pd.MultiIndex.from_tuples(
        [('d', 1), ('d', 2),
         ('e', 2)], names=['n', 'v']))
```
Create DataFrame with a MultiIndex

## Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.
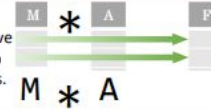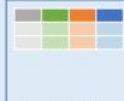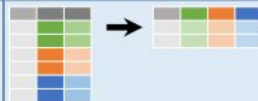```python
df = (pd.melt(df)
        .rename(columns={
                'variable':'var',
                'value':'val'})
        .query('val >= 200')
     )
```

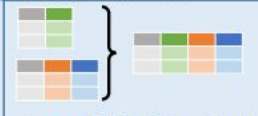## Reshaping Data – Change layout, sorting, reindexing, renaming

pd.melt(df)
Gather columns into rows.

df.pivot(columns='var', values='val')
Spread rows into columns.

pd.concat([df1,df2])
Append rows of DataFrames

pd.concat([df1,df2], axis=1)
Append columns of DataFrames

df.sort_values('mpg')
Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)
Order rows by values of a column (high to low).

df.rename(columns = {'y':'year'})
Rename the columns of a DataFrame

df.sort_index()
Sort the index of a DataFrame

df.reset_index()
Reset index of DataFrame to row numbers, moving index to columns.

df.drop(columns=['Length', 'Height'])
Drop columns from DataFrame

## Subset Observations - rows

df[df.Length > 7]
Extract rows that meet logical criteria.

df.drop_duplicates()
Remove duplicate rows (only considers columns).

df.sample(frac=0.5)
Randomly select fraction of rows.

df.sample(n=10)   Randomly select n rows.

df.nlargest(n, 'value')
Select and order top n entries.

df.nsmallest(n, 'value')
Select and order bottom n entries.

df.head(n)
Select first n rows.

df.tail(n)
Select last n rows.

## Subset Variables - columns

df[['width', 'length', 'species']]
Select multiple columns with specific names.

df['width']   or   df.width
Select single column with specific name.

df.filter(regex='regex')
Select columns whose name matches regular expression regex.

### Using query

query() allows Boolean expressions for filtering rows.

df.query('Length > 7')
df.query('Length > 7 and Width < 8')
df.query('Name.str.startswith("abc")',
          engine="python")

## Subsets - rows and columns

Use df.loc[] and df.iloc[] to select only rows, only columns or both.
Use df.at[] and df.iat[] to access a single value by row and column.
First index selects rows, second index columns.

df.iloc[10:20]
Select rows 10-20.

df.iloc[:, [1, 2, 5]]
Select columns in positions 1, 2 and 5 (first column is 0).

df.loc[:, 'x2':'x4']
Select all columns between x2 and x4 (inclusive).

df.loc[df['a'] > 10, ['a', 'c']]
Select rows meeting logical condition, and only the specific columns.

df.iat[1, 2]   Access single value by index

df.at[4, 'A']   Access single value by label

### Logic in Python (and pandas)

| | | | |
|---|---|---|---|
| < | Less than | != | Not equal to |
| > | Greater than | df.column.isin(values) | Group membership |
| == | Equals | pd.isnull(obj) | Is NaN |
| <= | Less than or equals | pd.notnull(obj) | Is not NaN |
| >= | Greater than or equals | &,\|,~,^,df.any(),df.all() | Logical and, or, not, xor, any, all |

### regex (Regular Expressions) Examples

| | |
|---|---|
| '\.' | Matches strings containing a period '.' |
| 'Length$' | Matches strings ending with word 'Length' |
| '^Sepal' | Matches strings beginning with the word 'Sepal' |
| '^x[1-5]$' | Matches strings beginning with 'x' and ending with 1,2,3,4,5 |
| '^(?!Species$).*' | Matches strings except the string 'Species' |

Cheatsheet for pandas (http://pandas.pydata.org/) originally written by Irv Lustig, Princeton Consultants, inspired by Rstudio Data Wrangling Cheatsheet

# Pandas Cheatsheet

## Summarize Data

```
df['w'].value_counts()
  Count number of rows with each unique value of variable
len(df)
  # of rows in DataFrame.
df.shape
  Tuple of # of rows, # of columns in DataFrame.
df['w'].nunique()
  # of distinct values in a column.
df.describe()
  Basic descriptive and statistics for each column (or GroupBy).
```

pandas provides a large set of summary functions that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

```
sum()                          min()
  Sum values of each object.     Minimum value in each object.
count()                        max()
  Count non-NA/null values of    Maximum value in each object.
  each object.                 mean()
median()                         Mean value of each object.
  Median value of each object. var()
quantile([0.25,0.75])            Variance of each object.
  Quantiles of each object.    std()
apply(function)                  Standard deviation of each
  Apply function to each object.  object.
```

## Group Data

```
df.groupby(by="col")
  Return a GroupBy object, grouped
  by values in column named "col".

df.groupby(level="ind")
  Return a GroupBy object, grouped
  by values in index level named
  "ind".
```

All of the summary functions listed above can be applied to a group.
Additional GroupBy functions:

```
size()                         agg(function)
  Size of each group.            Aggregate group using function.
```

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

```
shift(1)                       shift(-1)
  Copy with values shifted by 1.  Copy with values lagged by 1.
rank(method='dense')           cumsum()
  Ranks with no gaps.            Cumulative sum.
rank(method='min')             cummax()
  Ranks. Ties get min rank.      Cumulative max.
rank(pct=True)                 cummin()
  Ranks rescaled to interval [0, 1]. Cumulative min.
rank(method='first')           cumprod()
  Ranks. Ties go to first value.  Cumulative product.
```

## Windows

```
df.expanding()
  Return an Expanding object allowing summary functions to be
  applied cumulatively.
df.rolling(n)
  Return a Rolling object allowing summary functions to be
  applied to windows of length n.
```
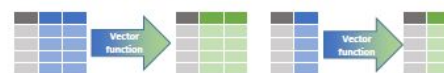
## Handling Missing Data

```
df.dropna()
  Drop rows with any column having NA/null data.
df.fillna(value)
  Replace all NA/null data with value.
```

## Make New Columns

```
df.assign(Area=lambda df: df.Length*df.Height)
  Compute and append one or more new columns.
df['Volume'] = df.Length*df.Height*df.Depth
  Add single column.
pd.qcut(df.col, n, labels=False)
  Bin column into n buckets.
```

pandas provides a large set of vector functions that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

```
max(axis=1)                    min(axis=1)
  Element-wise max.              Element-wise min.
clip(lower=-10,upper=10)       abs()
  Trim values at input thresholds  Absolute value.
```

## Plotting

```
df.plot.hist()                 df.plot.scatter(x='w',y='h')
  Histogram for each column      Scatter chart using pairs of points
```

## Combine Data Sets

adf + bdf =

**Standard Joins**

```
pd.merge(adf, bdf,
          how='left', on='x1')
  Join matching rows from bdf to adf.

pd.merge(adf, bdf,
          how='right', on='x1')
  Join matching rows from adf to bdf.

pd.merge(adf, bdf,
          how='inner', on='x1')
  Join data. Retain only rows in both sets.

pd.merge(adf, bdf,
          how='outer', on='x1')
  Join data. Retain all values, all rows.
```

**Filtering Joins**

```
adf[adf.x1.isin(bdf.x1)]
  All rows in adf that have a match in bdf.

adf[~adf.x1.isin(bdf.x1)]
  All rows in adf that do not have a match in bdf.
```

ydf + zdf =

**Set-like Operations**

```
pd.merge(ydf, zdf)
  Rows that appear in both ydf and zdf
  (Intersection).

pd.merge(ydf, zdf, how='outer')
  Rows that appear in either or both ydf and zdf
  (Union).

pd.merge(ydf, zdf, how='outer',
          indicator=True)
  .query('_merge == "left_only"')
  .drop(columns=['_merge'])
  Rows that appear in ydf but not zdf (Setdiff).
```

Cheatsheet for pandas (http://pandas.pydata.org/) originally written by Irv Lustig, Princeton Consultants, inspired by Rstudio Data Wrangling Cheatsheet

# Numpy and Matplotlib

➢ NumPy (Numerical Python) is a powerful, and extensively used, library for storage and calculations. It is designed for dealing with numerical data. It allows data storage and calculations by providing data structures, algorithms, and other useful utilities. For example, this library contains basic linear algebra functions, Fourier transforms, and advanced random number capabilities. It can also be used to load data to Python and export from it.

➢ Matplotlib is widely used for data visualization like for plotting histograms, line plots, and heat plots.

Let's explore more in Colab!