# Intro. Comp. for Data Science (FMI08)

Dr. Nono Saha

May 31, 2023

Max Planck Institute for Mathematics in the Sciences
University of Leipzig/ScaDS.AI

Spring 2023

## more pandas: exercice

Create a DataFrame from the data available at `../data/rent.csv` using `pd.read_csv()`. These data come from the 2017 American Community Survey and reflect the following values:

- `name` - name of state
- `variable` - Variable name: income = median yearly income, rent = median monthly rent
- `estimate` - Estimated value
- `moe` - 90% margin of error

Using these data, find the state(s) with the lowest income-to-rent ratio.

**pandas**: split-apply-combine

By "group by", we are referring to a process involving one or more of the following steps:

- **Splitting** the data into groups based on some criteria.
- **Applying** a function to each group independently.
- **Combining** the results into a data structure.

Groups can be created within a DataFrame via `groupby()` - these groups are then used by the standard summary methods (e.g. `sum()`, `mean()`, `std()`, etc...).

```
cereal = pd.read_csv("../data/cereal.csv")

cereal.groupby("type")
## <pandas.core.groupby.generic.DataFrameGroupBy object at 0
    x143e2a460>

cereal.groupby("type").mean()
cereal.groupby("mfr").size()

```

Groups can be accessed via `get_group( )` or the DataFrameGroupBy can be iterated over,

```
cereal.groupby("type").get_group("Hot")

cereal.groupby("mfr").get_group("Post")

for name, group in cereal.groupby("type"):
  print(name)
  print(group)
  print("")
```

# Group by: named aggregation

It is also possible to use special syntax to aggregate specific columns into a named output column,

```
cereal.groupby("mfr", as_index=False).agg(
min_cal = ("calories", "min"),
max_cal = ("calories", "max"),
med_sugar = ("sugars", "median"),
avg_rating = ("rating", "mean"))

##                 mfr  min_cal  max_cal  med_sugar  avg_rating
## 0  General Mills      100      140        8.5     34.485852
## 1 Kellogg's 50        160      7.0    44.038462
## 2           Maltex    100      100        3.0     54.850917
## 3          Nabisco     70      100        0.0     67.968567
## 4             Post     90      120       10.0     41.705744
## 5      Quaker Oats     50      120        6.0     42.915990
## 6   Ralston Purina     90      150        5.5     41.542997
```

Tuples can also be passed using `pd.NamedAgg()` but this offers no additional functionality.

The `transform()` method returns a DataFrame with the aggregated result matching the size (or length 1) of the input group(s),

```
1 cereal.groupby("mfr").transform(np.mean)
2 # For the new version of pandas, this may not work
3
4 cereal.groupby("type").transform("mean")
5 # For the new version of pandas, this may not work
6
7 <string>:1: FutureWarning: Dropping invalid columns in
      DataFrameG...
8
```

Note that we have lost the non-numeric columns, in case it works. And there will be a warning message.

`transform()` will generally be most useful via a user-defined function. The lambda argument is for each column of each group.

```
1 ( cereal.groupby("mfr").transform(
2 lambda x: (x - np.mean(x))/np.std(x)) )
3
```

Above, we are standardizing each numerical column of each manufacturer

# Group by: filtering groups

filter() also respects groups and allows for the inclusion/exclusion of groups based on user-specified criteria,

```
cereal.groupby("mfr").size()

## mfr
## General Mills     22
## Kellogg's 23
## Maltex           1
## Nabisco          6
## Post             9
## Quaker Oats      8
## Ralston Purina   8
## dtype: int64

cereal.groupby("mfr").filter(lambda x: len(x) > 10)

(cereal.groupby("mfr").filter(lambda x: len(x) > 10)
.groupby("mfr").size())
```

matplotlib

> `matplotlib` is a comprehensive library for creating static, animated, and interactive visualizations in Python.

```
1 import matplotlib as mpl
2 import matplotlib.pyplot as plt
3
```
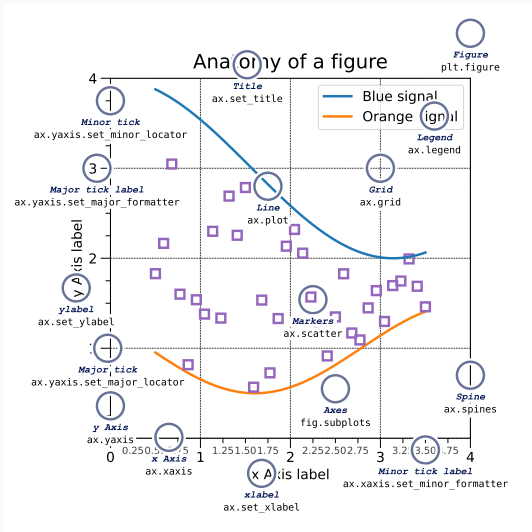
Why do we usually import only `pyplot` then?

> `matplotlib` is the whole package; `matplotlib.pyplot` is a module in `matplotlib`; and pylab is a module that gets installed alongside `matplotlib`.
> `pyplot` provides the state-machine interface to the underlying object-oriented plotting library. The state-machine implicitly and automatically creates figures and axes to achieve the desired plot.
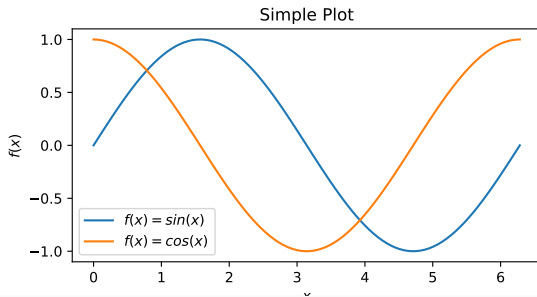
- **Figure** - The entire plot (including subplots)
- **Axes** - Subplot attached to a figure, contains the region for plotting data and axis'
- **Axis** - Set the scale and limits, generate ticks and ticklabels
- **Artist** - Everything visible on a figure: text, lines, axis, axes, etc.
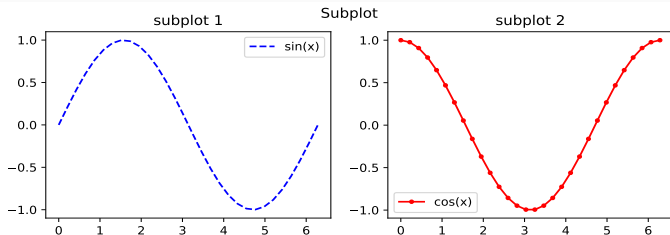
# matplotlib: basic plot

```python
x = np.linspace(0, 2*np.pi, 100)
y1 = np.sin(x)
y2 = np.cos(x)
fig, ax = plt.subplots(figsize=(6, 3))
ax.plot(x, y1, label="sin(x)")
ax.plot(x, y2, label="cos(x)")
ax.set_title("Simple Plot")
ax.legend()
```

## matplotlib: subplot

```python
x = np.linspace(0, 2*np.pi, 30)
y1 = np.sin(x)
y2 = np.cos(x)
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(9, 3))
ax1.plot(x, y1, "--b", label="sin(x)")
ax2.plot(x, y2, ".-r", label="cos(x)")
fig.suptitle("Subplot")
ax1.set_title("subplot 1")
ax2.set_title("subplot 2")
ax1.legend()
ax2.legend()
```

# matplotlib: more subplots and fancy

```python
1 x = np.linspace(-2, 2, 101)
2 fig, axs = plt.subplots(2, 2, figsize=(6, 4))
3 axs[0,0].plot(x, x, "b", label="linear")
4 axs[0,1].plot(x, x**2, "r", label="quadratic")
5 axs[1,0].plot(x, x**3, "g", label="cubic")
6 axs[1,1].plot(x, x**4, "c", label="quartic")
7 fig.legend(loc='upper right', bbox_to_anchor=(1.12, 0.9))
8 fig.suptitle("More subplots")
9
```
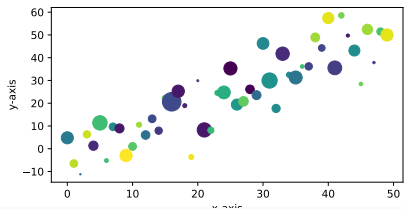
```python
1 x = np.linspace(0, 2*np.pi, 30)
2 y1 = np.sin(x)
3 y2 = np.cos(x)
4 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(9, 3))
5 ax1.plot(x, y1, "--b", label="sin(x)")
6 ax2.plot(x, y2, ".-r", label="cos(x)")
7 fig.suptitle("Subplot")
8 ax1.set_title("subplot 1")
9 ax2.set_title("subplot 2")
10 ax1.legend()
11 ax2.legend()
12
```

# matplotlib: plotting data

Beyond creating plots for arrays (and lists), addressable objects like dicts and DataFrames can be used via data.

```python
np.random.seed(19680801)  # seed the random number generator.
d = {'x': np.arange(50),
  'color': np.random.randint(0, 50, 50),
  'size': np.abs(np.random.randn(50)) * 100}
d['y'] = d['x'] + 10 * np.random.randn(50)
plt.figure(figsize=(6, 3))
plt.scatter('x', 'y', c='color', s='size', data=d)
plt.xlabel("x-axis")
plt.ylabel("y-axis")
```

Data can also come from `DataFrame` objects or series,

```
1 df = pd.DataFrame({"x": np.random.normal(size=10000)
2 }).assign(y = lambda d: np.random.normal(0.75*d.x, np.sqrt
      (1-0.75**2), size=10000))
3 fig, ax = plt.subplots(figsize=(5,5))
4 ax.scatter('x', 'y', c='k', data=df, alpha=0.1, s=0.5)
5 ax.set_xlabel('x')
6 ax.set_ylabel('y')
7 ax.set_title("Bivariate normal ($\\rho=0.75$)")
8
```

Series objects can also be plotted directly. The index is used as the *x* axis values,

```
1 s = pd.Series(np.cumsum( np.random.normal(size=100) ),
2 index = pd.date_range("2022-01-01", periods=100, freq="D"))
3 plt.figure(figsize=(3, 3), layout="constrained")
4 plt.plot(s)
5 plt.show()
6
```

# matplolib: scales

Axis scales can be changed via `plt.xscale()`, `plt.yscale()`, `ax.set_xscale()`, or `ax.set_yscale()`, supported values are "linear", "log", "symlog", and "logit".

```python
y = np.sort( np.random.sample(size=1000) )
x = np.arange(len(y))
plt.figure(layout="constrained")
scales = ['linear', 'log', 'symlog', 'logit']
for i, scale in zip(range(4), scales):
    plt.subplot(221+i)
    plt.plot(x, y)
    plt.grid(True)
    if scale == 'symlog':
        plt.yscale(scale, linthresh=0.01)
    else:
        plt.yscale(scale)
    plt.title(scale)
plt.show()
```

# matplolib: categorical data

```
1 df = pd.DataFrame({"cat": ["A", "B", "C", "D", "E"], "value": np
       .exp(range(5)) })
2
3 plt.figure(figsize=(4, 6), layout="constrained")
4 plt.subplot(321)
5 plt.scatter("cat", "value", data=df)
6 plt.subplot(322)
7 plt.scatter("value", "cat", data=df)
8 plt.subplot(323)
9 plt.plot("cat", "value", data=df)
10 plt.subplot(324)
11 plt.plot("value", "cat", data=df)
12 plt.subplot(325)
13 plt.bar("cat", "value", data=df)
14 plt.subplot(326)
15 plt.bar("value", "cat", data=df)
16 plt.show()
17
```

## `matplotlib`: histograms

```python
df = pd.DataFrame({
  "x1": np.random.normal(size=100),
  "x2": np.random.normal(1,2, size=100)
})
plt.figure(figsize=(4, 6), layout="constrained")
plt.subplot(311)
plt.hist("x1", bins=10, data=df, alpha=0.5)
plt.hist("x2", bins=10, data=df, alpha=0.5)
plt.subplot(312)
plt.hist(df, alpha=0.5)
plt.subplot(313)
plt.hist(df, stacked=True, alpha=0.5)
plt.show()

```
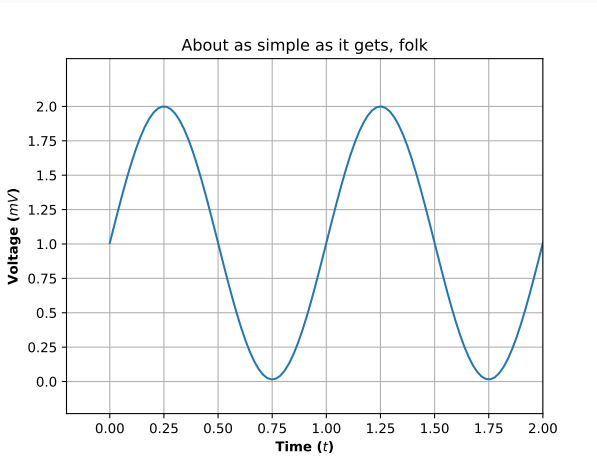
# matplotlib: boxplot

```
1  df = pd.DataFrame({
2    "x1": np.random.normal(size=100),
3    "x2": np.random.normal(1,2, size=100),
4    "x3": np.random.normal(-1,3, size=100)
5  }).melt()
6  plt.figure(figsize=(4, 4), layout="constrained")
7  plt.boxplot("value", positions="variable", data=df)
8
9  ##ValueError: List of boxplot statistics and `positions` values
       must have same the length
10
```

```
1  plt.boxplot(df.value, positions=df.variable)
2  ##ValueError: List of boxplot statistics and `positions` values
       must have same the length
3
```

https://matplotlib.org/stable/plot_types/index.html

Using what we just learnt, recreate the following plot,



From matplotlib examples

# Plotting with **pandas**

## pandas: plotting methods

Both Series and DataFrame objects have a plot method which can be used to create visualizations - `dtypes` determine the type of plot produced. Note these are just pyplot plots and can be formatted as such.

```
1 s = pd.Series(np.cumsum( np.random.normal(size=100) ),
2 index = pd.date_range("2022-01-01", periods=100, freq="D"))
3 plt.figure(figsize=(3,3), layout="constrained")
4 s.plot()
5 plt.show()
6
```

DataFrame plotting.

```
1 df = pd.DataFrame( np.cumsum( np.random.normal(size=(100,4)),
      axis=0),
2 index = pd.date_range("2022-01-01", periods=100, freq="D"),
3 columns = list("ABCD"))
4
5 plt.figure(layout="constrained")
6 df.plot(figsize=(5,3))
7 plt.show()
```