

1 Lecture 5: Advanced Pandas

Data Visualization · 1-DAV-105

Lecture by Broňa Brejová

As usual, we start by importing libraries. We also import the country data set from World Bank <https://databank.worldbank.org/home> under CC BY 4.0 license (see Lecture 03b).

```
[1]: import numpy as np
import pandas as pd
from IPython.display import Markdown
import matplotlib.pyplot as plt
import seaborn as sns
pd.options.display.float_format = '{:,.2f}'.format
```

```
[2]: url = 'https://bbrejova.github.io/viz/data/World_bank.csv'
countries = pd.read_csv(url).set_index('Country')
```

1.1 Hierarchical index (MultiIndex)

1.1.1 A small example table

To illustrate a hierarchical index, we first create a very small table consisting of two countries and their population in two years, and convert this table from wide to long format.

```
[3]: example_countries = countries.loc[["Slovak Republic", "Austria"],
                                      ["Population2010", "Population2018"]]
display(Markdown("**A small subset of countries table:**"), example_countries)
# change to long format
example_long = (example_countries.reset_index()
               .melt(id_vars=['Country'],
                    var_name='Year',
                    value_name='Population'))
# change year from string such as "Population2010" to int 2010
example_long.Year = example_long.Year.apply(lambda x : int(x[-4:]))
display(Markdown("**Changed to long format:**"), example_long)
```

A small subset of countries table:

	Population2010	Population2018
Country		
Slovak Republic	5,391,428.00	5,446,771.00
Austria	8,363,404.00	8,840,521.00

Changed to long format:

	Country	Year	Population
0	Slovak Republic	2010	5,391,428.00
1	Austria	2010	8,363,404.00

```
2 Slovak Republic 2018 5,446,771.00
3      Austria    2018 8,840,521.00
```

1.1.2 An index with duplicate labels

The original wide table had country as index, but in the long table, each country can have multiple rows. Pandas still allows us to use country as index with [duplicate values](#). Selecting the name of the country then gives us multiple rows.

```
[4]: # set country name as index in a copy of the table
example_long_indexed = example_long.set_index('Country')
# display table with index
display(Markdown("**Table with country as index:**"), example_long_indexed)
# select Slovakia from this table
display(Markdown("**Selecting multiple rows using `example_long_indexed.`
↳loc['Slovak Republic']`:**"))
display(example_long_indexed.loc['Slovak Republic'])
```

Table with country as index:

	Year	Population
Country		
Slovak Republic	2010	5,391,428.00
Austria	2010	8,363,404.00
Slovak Republic	2018	5,446,771.00
Austria	2018	8,840,521.00

Selecting multiple rows using `example_long_indexed.loc['Slovak Republic']`:

	Year	Population
Country		
Slovak Republic	2010	5,391,428.00
Slovak Republic	2018	5,446,771.00

1.1.3 Finally the hierarchical index

Our table can be more naturally indexed by a pair (country, year), which uniquely specifies a row. An index consisting of two or more levels is called [hierarchical](#) or [multi-level](#).

- MultiIndex can be created by `set_index` with a list of columns to use as index.
- For faster operations, it is a good idea to sort the table by the index using [sort_index](#).
- In `loc` use a tuple with one value per level, or only several initial levels.
- To specify other levels, use [xs](#).

```
[5]: # create MultiIndex by choosing a list of columns
example_multiindexed = example_long.set_index(['Country', 'Year']).sort_index()
display(Markdown("**Table with a multiindex:**"), example_multiindexed)
```

Table with a multiindex:

		Population
Country	Year	
Austria	2010	8,363,404.00
	2018	8,840,521.00
Slovak Republic	2010	5,391,428.00
	2018	5,446,771.00

```
[6]: display(Markdown("**Selecting a row by using a tuple in `loc`:**"))
      display(example_multiindexed.loc[('Slovak Republic', 2010)])
```

Selecting a row by using a tuple in loc:

Population
5,391,428.00

Name: (Slovak Republic, 2010), dtype: float64

```
[7]: display(Markdown("**Selecting all rows for a country using a shorter tuple in_
      ↪ `loc`:**"))
      display(example_multiindexed.loc[('Slovak Republic',)])
```

Selecting all rows for a country using a shorter tuple in loc:

	Population
Year	
2010	5,391,428.00
2018	5,446,771.00

```
[8]: display(Markdown("**Selecting all rows for a year using `xs`:**"))
      display(example_multiindexed.xs(2010, level='Year'))
```

Selecting all rows for a year using xs:

	Population
Country	
Austria	8,363,404.00
Slovak Republic	5,391,428.00

```
[9]: display(Markdown("**Names of index levels can be used in `query`:**"))
      display(example_multiindexed.query('Year > 2015'))
```

Names of index levels can be used in query:

		Population
Country	Year	
Austria	2018	8,840,521.00
Slovak Republic	2018	5,446,771.00

1.2 Combining tables

1.2.1 Concatenating tables using concat

- Function `concat` can be used to concatenate several tables.

- At the default settings, it combines along axis 0, meaning that the rows of second table are added after the rows of the first table.
- We will also use it for `axis=1`, in which case it finds rows with the same index in both tables and combines their columns.
- By default, the result has union of rows of the two tables, but intersection can be obtained by `join='inner'`.

Example Create a second small table of countries and display both tables. Then illustrate various concatenation modes using these tables.

```
[10]: example_countries2 = countries.loc[["Slovak Republic", "Austria", "Hungary"],
                                         ["Area", "Region"]]
display(Markdown("**The first small table:**"), example_countries)
display(Markdown("**The second small table:**"), example_countries2)
```

The first small table:

	Population2010	Population2018
Country		
Slovak Republic	5,391,428.00	5,446,771.00
Austria	8,363,404.00	8,840,521.00

The second small table:

	Area	Region
Country		
Slovak Republic	49,030.00	Europe & Central Asia
Austria	83,879.00	Europe & Central Asia
Hungary	93,030.00	Europe & Central Asia

```
[11]: display(Markdown("**Tables concatenated along axis 0:**"))
display(pd.concat([example_countries, example_countries2]))
```

Tables concatenated along axis 0:

	Population2010	Population2018	Area \
Country			
Slovak Republic	5,391,428.00	5,446,771.00	NaN
Austria	8,363,404.00	8,840,521.00	NaN
Slovak Republic	NaN	NaN	49,030.00
Austria	NaN	NaN	83,879.00
Hungary	NaN	NaN	93,030.00

	Region
Country	
Slovak Republic	NaN
Austria	NaN
Slovak Republic	Europe & Central Asia
Austria	Europe & Central Asia
Hungary	Europe & Central Asia

```
[12]: display(Markdown("**Tables concatenated along axis 1:**"))
display(pd.concat([example_countries, example_countries2], axis=1))
```

Tables concatenated along axis 1:

	Population2010	Population2018	Area \
Country			
Slovak Republic	5,391,428.00	5,446,771.00	49,030.00
Austria	8,363,404.00	8,840,521.00	83,879.00
Hungary	NaN	NaN	93,030.00

	Region
Country	
Slovak Republic	Europe & Central Asia
Austria	Europe & Central Asia
Hungary	Europe & Central Asia

```
[13]: display(Markdown("**Tables concatenated along axis 1 with inner join:**"))
display(pd.concat([example_countries, example_countries2], axis=1,
↪join='inner'))
```

Tables concatenated along axis 1 with inner join:

	Population2010	Population2018	Area \
Country			
Slovak Republic	5,391,428.00	5,446,771.00	49,030.00
Austria	8,363,404.00	8,840,521.00	83,879.00

	Region
Country	
Slovak Republic	Europe & Central Asia
Austria	Europe & Central Asia

1.2.2 Merging tables with merge

- Function `merge` works similarly as `concat` with `axis=1`, but it will match lines of two tables using any specified columns, not necessarily index.
- If values in these columns repeat, it combines all matching pairs of rows.
- Setting `how` in `merge` allows us to include rows that do not have a matching row in the other table.

Example

- Countries belong to various international organizations and a single country can belong to many. We will represent this as a table having one row for each pair of country and an organization it belongs to.
- To combine this with other country data, we apply `merge` to get a table in which each country is copied for each organization it is in.
- Then we can for example aggregate value for individual organizations.

```
[14]: # we create a small membership table by parsing a CSV-format string
import io
membership_str = io.StringIO("""Country,Member
Slovak Republic,NATO
Slovak Republic,EU
Slovak Republic,UN
Austria,UN
Austria,EU
""")
membership = pd.read_csv(membership_str)
display(Markdown("**A small country membership table:**"), membership)
```

A small country membership table:

	Country	Member
0	Slovak Republic	NATO
1	Slovak Republic	EU
2	Slovak Republic	UN
3	Austria	UN
4	Austria	EU

```
[15]: # merging tables using column Country in both
example_membership = pd.merge(example_countries, membership, on='Country')
display(Markdown("**Merged table:**"), example_membership)
```

Merged table:

	Country	Population2010	Population2018	Member
0	Slovak Republic	5,391,428.00	5,446,771.00	NATO
1	Slovak Republic	5,391,428.00	5,446,771.00	EU
2	Slovak Republic	5,391,428.00	5,446,771.00	UN
3	Austria	8,363,404.00	8,840,521.00	UN
4	Austria	8,363,404.00	8,840,521.00	EU

```
[16]: #aggregating organisations
display(Markdown("**The sum of country populations for each organization**",
↳(only for our two countries)))
display(example_membership.groupby('Member')['Population2018'].sum())
```

The sum of country populations for each organization (only for our two countries)

```
Member
EU      14,287,292.00
NATO     5,446,771.00
UN       14,287,292.00
Name: Population2018, dtype: float64
```

Similar operations are often done in relational databases, where `merge` is called `join`. Aggregation is also frequently used. More in a specialized database course in the third year.

1.3 Aggregation, split-apply-combine (groupby)

We have already seen simple examples of aggregation by `groupby` in Lecture 04. Here we discuss it in more detail.

Pandas follow the [split-apply-combine strategy](#) introduced in R by [Hadley Wickham](#).

Split: split data into groups, often by values in some column, e.g. `Region` in the `countries` table

Apply: apply some computation on each group, obtaining some result (single value, Series, DataFrame)

Combine: concatenate results for all groups together to a new table

Typical operations in the apply step:

- **aggregation:** e.g. compute group size, mean, median etc.
- **transformation:** e.g. compute percentage or rank of each item within group
- **filtering:** e.g. include only groups that are large enough

In Pandas, this is done by combination of `groupby` for the split step and additional functions for the apply step. The combine step is done implicitly. Pandas library provides many options, we will cover only basics.

1.3.1 Simple aggregation in the apply step

Apply functions such as `sum`, `mean`, `median`, `min`, `max`, `size`, `count`, `describe` after `groupby`.

- `size` gives the number of rows in the group.
- `count` gives the number of non-missing values in each column.

```
[42]: display(Markdown("**The number of countries in each region:**"))
display(countries.groupby('Region').size())
```

The number of countries in each region:

```
Region
East Asia & Pacific      37
Europe & Central Asia    58
Latin America & Caribbean 42
Middle East & North Africa 21
North America            3
South Asia               8
Sub-Saharan Africa       48
dtype: int64
```

```
[43]: display(Markdown("**Sums of country indicators in each region**"))
display(Markdown(" (including nonsense sums such as life expectation or GDP per_
↳capita)"))
display(countries.groupby('Region').sum(numeric_only=True))
```

Sums of country indicators in each region

(including nonsense sums such as life expectation or GDP per capita)

	Population2000	Population2010	Population2018	\
Region				
East Asia & Pacific	2,025,467,590.00	2,183,746,243.00	2,304,563,792.00	
Europe & Central Asia	861,278,548.00	887,926,820.00	917,922,618.50	
Latin America & Caribbean	520,903,450.00	589,932,554.00	640,467,174.00	
Middle East & North Africa	315,326,801.00	385,917,886.00	448,912,962.00	
North America	312,909,974.00	343,391,679.00	363,809,186.00	
South Asia	1,390,946,064.00	1,638,792,934.00	1,814,388,744.00	
Sub-Saharan Africa	665,327,581.00	869,025,106.00	1,074,853,734.00	

	Area	GDP2000	GDP2010	GDP2018	\
Region					
East Asia & Pacific	24,791,783.40	230,490.22	450,761.32	587,887.86	
Europe & Central Asia	28,691,030.98	880,359.02	1,726,243.88	1,991,497.58	
Latin America & Caribbean	20,425,545.98	195,238.12	454,889.22	511,688.30	
Middle East & North Africa	11,370,619.99	170,287.55	305,224.50	332,228.47	
North America	19,820,470.00	116,809.33	197,790.81	222,331.06	
South Asia	5,135,333.06	5,489.53	16,396.33	24,319.91	
Sub-Saharan Africa	21,754,456.00	43,129.49	109,876.37	115,339.80	

	Expectancy2000	Expectancy2010	Expectancy2018	\
Region				
East Asia & Pacific	2,297.62	2,256.05	2,311.10	
Europe & Central Asia	3,892.03	4,034.44	4,057.23	
Latin America & Caribbean	2,513.09	2,667.81	2,633.67	
Middle East & North Africa	1,500.65	1,549.98	1,581.12	
North America	233.66	239.08	242.14	
South Asia	511.29	550.20	570.71	
Sub-Saharan Africa	2,521.07	2,783.95	3,002.43	

	Fertility2000	Fertility2010	Fertility2018
Region			
East Asia & Pacific	91.76	80.27	74.22
Europe & Central Asia	91.23	95.88	92.16
Latin America & Caribbean	94.52	82.26	74.70
Middle East & North Africa	71.23	59.39	53.70
North America	5.28	5.32	4.83
South Asia	31.48	24.13	20.20
Sub-Saharan Africa	261.73	235.17	209.67

```
[45]: display(Markdown("**Specifically sum only population in 2018 per region:**"))
display(countries.groupby('Region')['Population2018'].sum())
```

Specifically sum only population in 2018 per region:

Region	
East Asia & Pacific	2,304,563,792.00
Europe & Central Asia	917,922,618.50

Latin America & Caribbean	640,467,174.00
Middle East & North Africa	448,912,962.00
North America	363,809,186.00
South Asia	1,814,388,744.00
Sub-Saharan Africa	1,074,853,734.00

Name: Population2018, dtype: float64

1.3.2 Transformation in the apply step

Function `groupby` can be followed by function `apply` which gets a function (e.g. lambda expression) and runs it on each group, producing a transformed version of the group. These are finally combined together.

Here we compute for each country what percentage is its population from the population of the region.

```
[47]: # apply gets a single column with population sizes within a group
# and divides each number by the sum of the group
pop_within_group = (countries.groupby('Region',
    ↪group_keys=False)['Population2018']
    .apply(lambda x : x / x.sum()))
display(Markdown("**For each country, what fraction is its population within_
    ↪region:**"))
display(pop_within_group.head())
```

For each country, what fraction is its population within region:

Country	
Afghanistan	0.02
Albania	0.00
Algeria	0.09
American Samoa	0.00
Andorra	0.00

Name: Population2018, dtype: float64

```
[51]: display(Markdown("**Add back region name using concat:**"))
pop_within_group2 = pd.concat([pop_within_group, countries['Region']], axis=1)
display(pop_within_group2.head())

display(Markdown("**Look up value for Slovakia:**"))
display(pop_within_group2.loc["Slovak Republic"])
```

Add back region name using concat:

	Population2018	Region
Country		
Afghanistan	0.02	South Asia
Albania	0.00	Europe & Central Asia
Algeria	0.09	Middle East & North Africa

American Samoa	0.00	East Asia & Pacific
Andorra	0.00	Europe & Central Asia

Look up value for Slovakia:

Population2018	0.01
Region	Europe & Central Asia

Name: Slovak Republic, dtype: object

```
[52]: display(Markdown("**Check that the sum of each region is 1:**"))
display(pop_within_group2.groupby('Region').sum())
```

Check that the sum of each region is 1:

	Population2018
Region	
East Asia & Pacific	1.00
Europe & Central Asia	1.00
Latin America & Caribbean	1.00
Middle East & North Africa	1.00
North America	1.00
South Asia	1.00
Sub-Saharan Africa	1.00

1.3.3 Filtering in the apply step

Finally, `groupby` can be followed by `filter` to use only some of the groups in the result.

Here we report all countries in regions that have at least billion inhabitants.

```
[53]: # filter gets a function returning a boolean value for each group
filtered = (countries.groupby("Region")
            .filter(lambda x : x['Population2018'].sum() > 1e9))
display(Markdown("**Filtered data:**"))
display(filtered.head())
display(Markdown("**Check sums in regions for selected countries:**"))
display(filtered.groupby('Region')['Population2018'].sum())
```

Filtered data:

	Region	Income Group	Population2000	\
Country				
Afghanistan	South Asia	Low income	20,779,953.00	
American Samoa	East Asia & Pacific	Upper middle income	57,821.00	
Angola	Sub-Saharan Africa	Lower middle income	16,395,473.00	
Australia	East Asia & Pacific	High income	19,153,000.00	
Bangladesh	South Asia	Lower middle income	127,657,854.00	

	Population2010	Population2018	Area	GDP2000	\
Country					
Afghanistan	29,185,507.00	37,172,386.00	652,860.00	NaN	

American Samoa	56,079.00	55,465.00	200.00	NaN
Angola	23,356,246.00	30,809,762.00	1,246,700.00	556.84
Australia	22,031,750.00	24,982,688.00	7,741,220.00	21,679.25
Bangladesh	147,575,430.00	161,356,039.00	147,630.00	418.07

	GDP2010	GDP2018	Expectancy2000	Expectancy2010	\
Country					
Afghanistan	543.30	493.75	55.84	61.03	
American Samoa	10,271.22	11,466.69	NaN	NaN	
Angola	3,587.88	3,289.65	46.52	55.35	
Australia	52,022.13	57,354.96	79.23	81.70	
Bangladesh	781.15	1,698.35	65.45	69.88	

	Expectancy2018	Fertility2000	Fertility2010	Fertility2018
Country				
Afghanistan	64.49	7.49	5.98	4.47
American Samoa	NaN	NaN	NaN	NaN
Angola	60.78	6.64	6.19	5.52
Australia	82.75	1.76	1.93	1.74
Bangladesh	72.32	3.17	2.32	2.04

Check sums in regions for selected countries:

Region	
East Asia & Pacific	2,304,563,792.00
South Asia	1,814,388,744.00
Sub-Saharan Africa	1,074,853,734.00
Name: Population2018, dtype: float64	

1.3.4 Grouping by multiple values

Function `groupby` can get a single column, but also a list of columns or a Series which will be used as if it was a column of the table.

```
[54]: display(Markdown("**Populations split by both region and income group**"))
display(countries.groupby(['Region', 'Income Group'])['Population2018'].sum())
```

Populations split by both region and income group

Region	Income Group	
East Asia & Pacific	High income	222,924,695.00
	Low income	25,549,819.00
	Lower middle income	293,430,538.00
	Upper middle income	1,762,658,740.00
Europe & Central Asia	High income	520,487,513.00
	Low income	9,100,837.00
	Lower middle income	86,607,467.00
	Upper middle income	301,726,801.50
Latin America & Caribbean	High income	32,341,730.00
	Low income	11,123,176.00

	Lower middle income	33,826,921.00
	Upper middle income	563,175,347.00
Middle East & North Africa	High income	66,016,247.00
	Low income	45,404,970.00
	Lower middle income	193,774,373.00
	Upper middle income	143,717,372.00
North America	High income	363,809,186.00
South Asia	Low income	37,172,386.00
	Lower middle income	1,776,700,662.00
	Upper middle income	515,696.00
Sub-Saharan Africa	High income	1,362,065.00
	Low income	519,523,613.00
	Lower middle income	488,057,804.00
	Upper middle income	65,910,252.00

Name: Population2018, dtype: float64

- Now we create a Series classifying each countries as small, medium and large using cutoff 1 million for small and 100 million for medium.
- We then use this series in `groupby`.
- The classification is created by `pd.cut` function.

```
[60]: bin_ends = [0, 1e6, 1e8, 1e10]
bin_labels = ["small", "medium", "large"]
size_groups = pd.cut(countries['Population2018'],
                      bins=bin_ends, labels=bin_labels).rename("SizeCategory")
display(Markdown("**Country size classification:**"))
display(size_groups.head())
```

Country size classification:

```
Country
Afghanistan      medium
Albania           medium
Algeria           medium
American Samoa   small
Andorra           small
Name: SizeCategory, dtype: category
Categories (3, object): ['small' < 'medium' < 'large']
```

```
[63]: # now use size_groups in groupby
display(Markdown("**The number of countries in each size group:**"))
display(countries.groupby(size_groups).size())
display(Markdown("**The number of countries in each size group and region:**"))
display(countries.groupby(['Region', size_groups]).size())
```

The number of countries in each size group:

```
SizeCategory
small      58
```

```
medium    145
large      13
dtype: int64
```

The number of countries in each size group and region:

Region	SizeCategory	
East Asia & Pacific	small	18
	medium	15
	large	4
Europe & Central Asia	small	12
	medium	45
	large	1
Latin America & Caribbean	small	19
	medium	21
	large	2
Middle East & North Africa	small	2
	medium	19
	large	0
North America	small	1
	medium	1
	large	1
South Asia	small	2
	medium	3
	large	3
Sub-Saharan Africa	small	4
	medium	41
	large	2

```
dtype: int64
```

1.4 Categorical variables

Categorical variables have values from a small set, such as region and income group in the table of countries. So far we have represented them only as strings, but we can explicitly convert them to a [categorical data type](#) in Pandas.

This has several advantages: * Strings are internally replaced by numerical IDs within the table, potentially saving memory. * Categories can be ordered and then sorting, minimum, maximum etc works as desired, not alphabetically. * Pandas is aware of the full set of possible values. For example categories without members can appear in the `groupby` results.

Example Income groups in our table are strings, we will convert them to an ordered categorical variable.

```
[68]: # creating a categorical type
cat_type = pd.api.types.CategoricalDtype(categories=["Low income",
                                                    "Lower middle income",
                                                    "Upper middle income",
                                                    "High income"],
```

```

                                ordered=True)
# converting Income Group column to cat_type in a new DataFrame
countries_cat = countries.astype({'Income Group': cat_type})

display(Markdown("**Income Group column in the old table:**"),
        countries['Income Group'].head(3))
display(Markdown("**Income Group column in the new table:**"),
        countries_cat['Income Group'].head(3))

```

Income Group column in the old table:

```

Country
Afghanistan          Low income
Albania              Upper middle income
Algeria              Lower middle income
Name: Income Group, dtype: object

```

Income Group column in the new table:

```

Country
Afghanistan          Low income
Albania              Upper middle income
Algeria              Lower middle income
Name: Income Group, dtype: category
Categories (4, object): ['Low income' < 'Lower middle income' < 'Upper middle_
income' < 'High income']

```

```

[73]: display(Markdown("**Minimum and maximum income group in the table with_
      ↪categorical values:**"
                    " (manually fixed order):"))
display(countries_cat['Income Group'].min())
display(countries_cat['Income Group'].max())

display(Markdown("**Minimum and maximum income group in the table with_
      ↪strings**"
                    " (alphabetical order):"))
display(countries['Income Group'].min())
display(countries['Income Group'].max())

```

Minimum and maximum income group in the table with categorical values: (manually fixed order):

```

'Low income'
'High income'

```

Minimum and maximum income group in the table with strings (alphabetical order):

```

'High income'
'Upper middle income'

```

- Note that if categories do not need a fixed order, they can be created directly in the `astype` function as in the code below.
- Notice that `groupby` creates even empty groups which would not happen with strings.

```
[74]: # convert region to an unordered category
countries_cat2 = countries_cat.astype({'Region': 'category'})
# count the number of countries for each combination of income group and region
countries_cat2.groupby(['Income Group', 'Region']).size()
```

```
[74]: Income Group      Region
Low income            East Asia & Pacific      1
                   Europe & Central Asia      1
                   Latin America & Caribbean  1
                   Middle East & North Africa  2
                   North America              0
                   South Asia                 1
                   Sub-Saharan Africa        23
Lower middle income  East Asia & Pacific      12
                   Europe & Central Asia      4
                   Latin America & Caribbean  4
                   Middle East & North Africa  6
                   North America              0
                   South Asia                 6
                   Sub-Saharan Africa        18
Upper middle income  East Asia & Pacific      10
                   Europe & Central Asia      15
                   Latin America & Caribbean  20
                   Middle East & North Africa  5
                   North America              0
                   South Asia                 1
                   Sub-Saharan Africa         5
High income          East Asia & Pacific      14
                   Europe & Central Asia      38
                   Latin America & Caribbean  17
                   Middle East & North Africa  8
                   North America              3
                   South Asia                 0
                   Sub-Saharan Africa         2

dtype: int64
```

1.5 Dates and times

An important type of data sets are time series, where some variables are measured repeatedly over time. Pandas has an extensive support [for work with times and dates](#). Here we show only a small example.

- We illustrate this on the movie dataset from [Kaggle](#) (see lecture 04).
- The column labeled `release_date` is recognized as date by passing `parse_dates` parameter

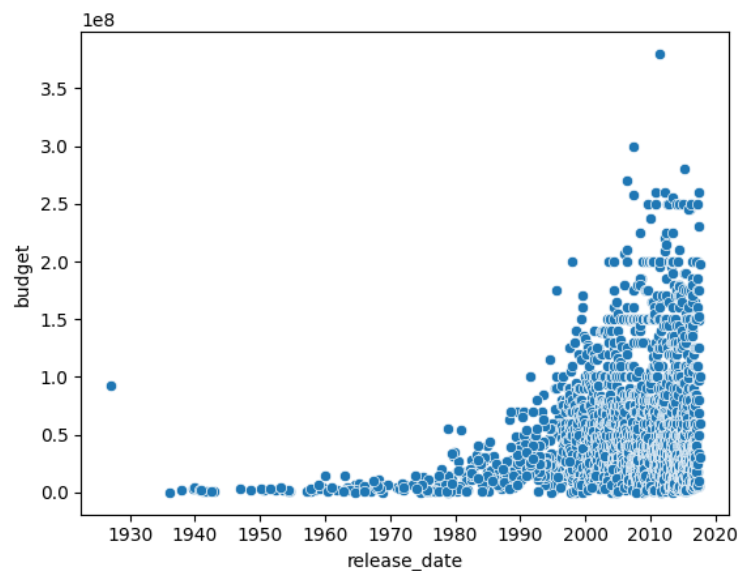
to read_csv.

- Then we call function `day_name()` to get the day of week for each release day and use `value_counts` to see which days are most frequent as movie release dates.
- We also use the release date as the x-coordinate in a scatterplot.

```
[77]: # import data, including parsing of dates
url = 'https://bbrejova.github.io/viz/data/Movies_small.csv'
movies = pd.read_csv(url, parse_dates=['release_date'])
# get days of week for realse dates
days = movies['release_date'].apply(lambda x : x.day_name())
days.value_counts()
```

```
[77]: Friday          639
Thursday         515
Wednesday        474
Tuesday          175
Saturday          94
Monday           87
Sunday           65
Name: release_date, dtype: int64
```

```
[78]: # use release date is x-coordinate
sns.scatterplot(data=movies, x='release_date', y='budget')
pass
```



1.6 Missing values

Data sets are often incomplete, and Pandas provides techniques for [working with missing data](#).

- Missing data are typically imported as `np.nan` (not-a-number).
- These cannot occur in int-type columns, so ints are converted to floats, but can be handled in a [special way](#).

Bellow we show a small example what happens when working with missing data, including functions `isna`, `dropna`, `fillna`.

```
[86]: # create a small series with one missing value
a = pd.Series([1, 2, np.nan, 3])
display(Markdown("**`a.sum()` skips missing values:**"),
        a.sum())
display(Markdown("**`a.count()` counts non-missing values:**"),
        a.count())
display(Markdown("**`a.mean()` also considers only non-missing:**"),
        a.mean())
display(Markdown("**`a > 2` evaluates missing values as `False`, similarly `<`,  
↪ `==`:**"),
        a > 2)
display(Markdown("**`a == np.nan` also evaluates as `False`:**"),
        a == np.nan)
display(Markdown("**`a.isna()` can be used to detect missing values:**"),
        a.isna())
display(Markdown("**`a.dropna()` omits missing values:**"),
        a.dropna())
display(Markdown("**`a.fillna(-1)` replaces them with a specified value:**"),
        a.fillna(-1))
```

`a.sum()` skips missing values:

6.0

`a.count()` counts non-missing values:

3

`a.mean()` also considers only non-missing:

2.0

`a > 2` evaluates missing values as `False`, similarly `<`, `==`:

0 False

1 False

2 False

3 True

dtype: bool

`a == np.nan` also evaluates as `False`:

0 False

1 False

2 False

```
3    False
dtype: bool
```

a.isna() can be used to detect missing values:

```
0    False
1    False
2     True
3    False
dtype: bool
```

a.dropna() omits missing values:

```
0    1.00
1    2.00
3    3.00
dtype: float64
```

a.fillna(-1) replaces them with a specified value:

```
0    1.00
1    2.00
2   -1.00
3    3.00
dtype: float64
```

1.7 Pandas efficiency

Below we show several examples how different ways of implementing the same operation can have very different running time on large data. Pandas functions are usually much faster than manual iteration. However, if you do not work on huge data sets, the difference is not so important.

To measure time, we use a special Jupyter command `%timeit`. * It runs the code several times to estimate the time per one repeat.

```
[87]: # generate a Series of million random numbers and also convert it to Python list
length = int(1e6)
xs = pd.Series(np.random.uniform(0,100, length))
xl = list(xs)
```

Below we see that method `sum()` on Series is faster than built-in `sum` on a Python list, but built-in `sum` on Series is much slower, because it iterates over elements of Series.

```
[90]: display(Markdown("**Method `sum` on `Series` `xs.sum()`:**"))
%timeit result = xs.sum()
display(Markdown("**Python `sum` on Python list `sum(xl)`:**"))
%timeit result = sum(xl)
display(Markdown("**Python `sum` on Series `sum(xs)`:**"))
%timeit result = sum(xs)
```

Method `sum` on Series `xs.sum()`:

473 μ s \pm 26.3 μ s per loop (mean \pm std. dev. of 7 runs, 1,000 loops each)

Python sum on Python list sum(x1):

4.33 ms \pm 68.4 μ s per loop (mean \pm std. dev. of 7 runs, 100 loops each)

Python sum on Series sum(xs):

36.5 ms \pm 1.07 ms per loop (mean \pm std. dev. of 7 runs, 10 loops each)

Below we compare three ways of generating a sequence of squared values. Multiplying Series with * is the fastest, Python list comprehension is much slower and apply function from Pandas is even slower.

```
[92]: display(Markdown("**Pandas `Series` multiplication `x2s = xs * xs`**"))
      %timeit x2s = xs * xs
      display(Markdown("**Python list comprehension on a list `x2l = [x * x for x in_
      ↪x1]``**"))
      %timeit x2l = [x * x for x in x1]
      display(Markdown("**Pandas `apply` function `x2s = xs.apply(lambda x : x *_
      ↪x)`**"))
      %timeit x2s = xs.apply(lambda x : x * x)
```

Pandas Series multiplication x2s = xs * xs:

621 μ s \pm 27.9 μ s per loop (mean \pm std. dev. of 7 runs, 1,000 loops each)

Python list comprehension on a list x2l = [x * x for x in x1]:

34.5 ms \pm 200 μ s per loop (mean \pm std. dev. of 7 runs, 10 loops each)

Pandas apply function x2s = xs.apply(lambda x : x * x)

123 ms \pm 1.29 ms per loop (mean \pm std. dev. of 7 runs, 10 loops each)

The code below creates the Series of squares by creating a Series filled with zeroes and then assigning individual values using for-loop. This is again much slower than all methods above, so to make the code reasonably fast, we run it on data which is 100 times smaller than above.

```
[36]: length2 = 10000
      xs_small = xs.iloc[0:length2]
      def assignments(len, x):
          x2 = pd.Series([0] * len)
          for i in range(len):
              x2[i] = x[i] * x[i]
          return x2
      %timeit x2s_small = assignments(length2, xs_small)
```

79.8 ms \pm 418 μ s per loop (mean \pm std. dev. of 7 runs, 10 loops each)

Finally the code below is even worse. It appends individual squares to a Series which starts with size 1. We run it on even smaller list of size 1000.

```
[40]: length3 = 1000
      xs_tiny = xs.iloc[0:length3]
      def assignments(len, x):
```

```
x2 = pd.Series([0])
for i in range(len):
    x2[i] = x[i] * x[i]
return x2
%timeit x2s_tiny = assignments(length3, xs_tiny)
```

203 ms \pm 8.1 ms per loop (mean \pm std. dev. of 7 runs, 1 loop each)