

Python

Data & Automation

Data

- We encounter data in many forms, some **well-formed and prepared** and some still **very raw or even faulty**.
- Now that we know how to obtain data from websites, for example, we want to read in data and play it back **fully automatically**.

Scenario: In search of data

The image displays three separate web interfaces side-by-side:

- kaggle**: A sidebar menu on the left with options like Create, Home, Competitions, Datasets, and Models. The main content area shows a "Datasets" section with a "Search" bar and a "+ New Dataset" button.
- Stadt Zürich Open Data**: A header with the city logo and "Stadt Zürich Open Data". Below it is a navigation bar with links to Startseite, Datensätze, Kategorien, and Showcases. The main content area features a large image of a mountain peak.
- opendata.swiss**: A teal-colored landing page with the text "Finden Sie Schweizer Open Government Data" and a button "Erfahren Sie mehr über opendata.swiss".

Scenario: In search of data

Standorte ZüriVelo

Offizielle Standorte der Veloverleihstationen von Züri Velo (PubliBike) Zweck: Dieser Datensatz dient zur Übersicht der Standorte. Die IDs können zur Verknüpfung von...

[csv](#) [dxf](#) [gpkg](#) [json](#) [shp](#) [wfs](#) [wms](#) [wmts](#)

- CSV:
Comma-Separated Values.
- DXF: Drawing Interchange File Format.
- GPKG: GeoPackage.
- JSON: GeoJSON / JSON.
- SHP: Shapefile.
- WFS: Web Feature Service.
- WMS: Web Map Service.
- WMTS: Web Map Tile Service.

Scenario: In search of data

- CSV: Comma-Separated Values.
- DXF: Drawing Interchange File Format.
- GPKG: GeoPackage.
- JSON: GeoJSON / JSON.
- SHP: Shapefile.
- WFS: Web Feature Service.
- WMS: Web Map Service.
- WMTS: Web Map Tile Service.

Standorte ZüriVelo

Offizielle Standorte der Veloverleihstationen von Züri Velo (PubliBike) Zweck: Dieser Datensatz dient zur Übersicht der Standorte. Die IDs können zur Verknüpfung von...

[csv](#) [dxf](#) [gpkg](#) [json](#) [shp](#) [wfs](#) [wms](#) [wmts](#)

Unfortunately, time is too short for GEO data

Scenario: In search of data

We receive the following link from GeoJSON:

[https://www.ogd.stadt-zuerich.ch/wfs/geoportal/Standorte_ZueriVelo?
service=WFS&version=1.1.0&request=GetFeature&outputFormat=GeoJSON&typ
ename=view_zuerivelocommon_publibike](https://www.ogd.stadt-zuerich.ch/wfs/geoportal/Standorte_ZueriVelo?service=WFS&version=1.1.0&request=GetFeature&outputFormat=GeoJSON&typename=view_zuerivelocommon_publibike)

```
// 20240117162054
// https://www.ogd.stadt-zuerich.ch/wfs/geoportal/Standorte_ZueriVelo?
// service=WFS&version=1.1.0&request=GetFeature&outputFormat=GeoJSON&typename=view_zuerivelocommon_publibike

{
  "type": "FeatureCollection",
  "bbox": [
    8.46781,
    47.325159,
    8.650765,
    47.452387
  ],
  "features": [
    {
      "geometry": {
        "coordinates": [
          8.547903,
          47.384902
        ],
        "type": "Point"
      },
      "id": "view_zuerivelocommon_publibike.1",
      "properties": {
        "adresse": "Culmannstrasse 100",
        "datum": "2024-01-17T00:00:03.000",
        "name": "Publibike"
      }
    }
  ]
}
```

Scenario: In search of data

Pretty complicated... Let's look for metadata:



<https://www.geocat.ch/geonetwork/srv/ger/catalog.search#/metadata/a7e477f1-75ce-46d2-a32a-c547e37ccfaf/formatters/xsl-view?root=div&view=advanced>

Scenario: In search of data



```
1      "id": "view_zuerivelopublibike.9",
2      "properties": {
3          "adresse": "Bahnhof Stettbach",
4          "datum": "2024-01-17T00:00:03.000",
5          "id1": 353643.0,
6          "id_publibike": 514.0,
7          "lat": 47.397885,
8          "lon": 8.596202,
9          "name": "Bahnhof Stettbach",
10         "objectid": 9.0,
11         "plz": 8051,
12         "stadt": "Zürich",
13         "status": "Active"
14     },
15     "type": "Feature"
16 },
17 {
18     "geometry": {
19         "coordinates": [
20             8.549831,
21             47.421914
22         ],
23         "type": "Point"
24     },
```

Questions:

- Where is the bike located?
- When is the data from?
- What ID does the object have?

Scenario: In search of data

The bottom line is that finding data is not particularly difficult.

Now it's time to **import and edit**.

Data formats

We have already learned how to query web resources.
In fact, we have used Soup & request to request XML files (HTML).

But now we want to deal with csv and json.

CSV

CSV stands for comma-separated-values and is a file format. The extension of the file name is .csv.

Roughly speaking, each "data record" in the table corresponds to a text line in the text file.

Each field is separated from the next by a comma.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

The code contains a lot to discuss, please take a look at it for 1-2 min.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {" ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.'
```

Question 0:

Why do we not need to install CSV?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {" ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 0:

Why do we not need to install CSV?

Answer:

Because it is part of the Python standard library.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 1:

Why is "open, edit and close" not used here instead of "with"?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 1:

Why is "open, edit and close" not used here instead of "with"?

Response:

The spelling is simpler and it is safely closed.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {" ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.'
```

Question 2:

What does the delimiter mean?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.'
```

Question 2:

What does the delimiter mean?

Answer:

CSV files do not necessarily have to be separated with , - other separators are also possible.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {" ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 3:

What does f' and \t mean?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 3:

What does f' and \t mean?

Answer:

f' is the formatting function of Python \t is a tab or indent.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {" ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 4:

How exactly is the file processed?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 4:

How exactly is the file processed?

Response:

The file is processed line by line, with special treatment for the first line.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {" ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.')
```

Question 5:

Aren't `row[0]` and `row[2]` a bit error-prone?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 00_employee_birthday.py
2
3 import csv
4
5 with open('00_employee_birthday.csv') as csv_file:
6     csv_reader = csv.reader(csv_file, delimiter=',')
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12         else:
13             print(f'\t{row[0]} works in the {row[1]} department, '
14                   f'and was born in {row[2]}.')
15             line_count += 1
16     print(f'Processed {line_count} lines.'
```

Question 5:

Aren't `row[0]` and `row[2]` a bit error-prone?

Answer:

Yes, for this we need to know **exactly** where our information is located.

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 01_employee_birthday_dict.py
2
3 import csv
4
5 with open('00_employee_birthday.csv', mode='r') as csv_file:
6     csv_reader = csv.DictReader(csv_file)
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12             print(f'\t{row["name"]} works in the {row["department"]} '
13                   f'department, and was born in {row["birthday month"]}.')
14             line_count += 1
15             print(f'Processed {line_count} lines.'
```

In response to the question from earlier, we now want to improve a few things, what has been adjusted?

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 01_employee_birthday_dict.py
2
3 import csv
4
5 with open('00_employee_birthday.csv', mode='r') as csv_file:
6     csv_reader = csv.DictReader(csv_file)
7     line_count = 0
8     for row in csv_reader:
9         if line_count == 0:
10             print(f'Column names are {", ".join(row)}')
11             line_count += 1
12             print(f'\t{row["name"]} works in the {row["department"]} '
13                   f'department, and was born in {row["birthday month"]}.')
14             line_count += 1
15             print(f'Processed {line_count} lines.'
```

In response to the question from earlier, we now want to improve a few things, what has been adjusted?

DictReader was used, which treats the csv file like a dictionary.

Accordingly, row[0] also became the speaking row[name].

CSV

Additional information:

- **delimiter** specifies the character used to separate the individual fields.
 - The default is the comma (',').
- **quotechar** specifies the character used to enclose fields containing the delimiter.
 - The default is a double quotation mark ("").
- **escapechar** specifies the character used to enclose the delimiter if no quotation marks are used.
 - The default is not an escape character.

CSV

Mode:

Specify File Mode

Here are five different modes you can use to open the file:

Character	Mode	Description
'r'	Read (default)	Open a file for read only
'w'	Write	Open a file for write only (overwrite)
'a'	Append	Open a file for write only (append)
'r+'	Read+Write	open a file for both reading and writing
'x'	Create	Create a new file

Python File Modes

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 02_employee_birthday_writing.py
2
3 import csv
4
5 with open('02_employee_birthday.csv', mode='w') as employee_file:
6     employee_writer = csv.writer(employee_file, delimiter=',', quotechar='"', quoting=csv.QUOTE_MINIMAL)
7
8     employee_writer.writerow(['John Smith', 'Accounting', 'November'])
9     employee_writer.writerow(['Erica Meyers', 'IT', 'March'])
```

Now we also want to write to the file. Here it is important to use a **write mode** - see last slide.

CSV - Exercise



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```

We will now deal with 03-exercises/00_own_birthday_writing.

Please copy the .csv and .py into your userfolder.

Questions:

- What happens when the file is executed (with the data)?
- Which **mode** could have been used to prevent this? (Simply copy the file again)
- Please add 3 more people

CSV



```
1 name,department,birthday month
2 John Smith,Accounting,November
3 Erica Meyers,IT,March
```



```
1 # 03_employee_birthday_writing_dict.py
2
3 import csv
4
5 with open('employee_file2.csv', mode='w') as csv_file:
6     fieldnames = ['emp_name', 'dept', 'birth_month']
7     writer = csv.DictWriter(csv_file, fieldnames=fieldnames)
8
9     writer.writeheader()
10    writer.writerow({'emp_name': 'John Smith', 'dept': 'Accounting', 'birth_month': 'November'})
11    writer.writerow({'emp_name': 'Erica Meyers', 'dept': 'IT', 'birth_month': 'March'})
```

Of course, this also works with DictWriter (the counterpart to DictReader) and the nice thing is that dictionaries can now be transferred.

CSV



```
1 with open('test.csv', 'w') as csvfile:
2     fieldnames = test_dict.keys()
3     fieldvalues = zip(*test_dict.values())
4
5     writer = csv.writer(csvfile)
6     writer.writerow(fieldnames)
7     writer.writerows(fieldvalues)
```

Multiple lines are only possible here with detours; we will use pandas for this later.

JSON

Although JSON is called **JavaScript Object Notation**, it is a completely independent data format that is not tied to the JavaScript language.

At the same time, with its very simple structure and encoding in the Unicode character set, JSON is the ideal format for exchanging data between systems.

JSON

```
1. {
2.     "customers": [
3.         {
4.             "firstName": "Max",
5.             "middleInitial": "M",
6.             "lastName": "Mustermann",
7.             "gender": "M",
8.             "age": 24,
9.             "address": {
10.                 "streetAddress": "Nordring 98",
11.                 "city": "Nürnberg",
12.                 "state": "Bayern",
13.                 "postalCode": 90409
14.             },
15.             "phoneNumbers": [
16.                 { "type": "home", "number": "0911/161803399" }
17.             ]
18.         },
19.         ,
20.         {
21.             "firstName": "Erika",
22.             "lastName": "Gabler",
23.             "gender": "F",
24.             "age": 57,
25.             "address": {
26.                 "streetAddress": "Neuer Wall 72",
27.                 "city": "Hamburg",
28.                 "state": "Hamburg",
29.                 "postalCode": 20354
30.             },
31.             "phoneNumbers": [
32.                 { "type": "home", "number": "040/27182818" }
33.                 , { "type": "mobil", "number": "0176/31415926" }
34.             ]
35.         }
36.     ]
37. }
```

JSON



```
1  {
2      "firstName": "Jane",
3      "lastName": "Doe",
4      "hobbies": ["running", "sky diving", "singing"],
5      "age": 35,
6      "children": [
7          {
8              "firstName": "Alice",
9              "age": 6
10         },
11         {
12             "firstName": "Bob",
13             "age": 8
14         }
15     ]
16 }
```

As we can see, JSON supports primitive types such as strings and numbers as well as nested lists and objects.

JSON

```
1  {
2      "firstName": "Jane",
3      "lastName": "Doe",
4      "hobbies": [ "running", "sky diving", "singing"],
5      "age": 35,
6      "children": [
7          {
8              "firstName": "Alice",
9              "age": 6
10         },
11         {
12             "firstName": "Bob",
13             "age": 8
14         }
15     ]
16 }
```

Wait, that looks like a Python dictionary! I know, right?

That's because it's a form of UON (Universal Object Notation), so they are very similar to each other.

JSON

Python

dict

list, tuple

str

int, long, float

True

False

None

JSON

object

array

string

number

true

false

null

JSON

If we now have a Python object:

```
● ● ●  
1 # 04_json_write.py  
2  
3 data = {  
4     "president": {  
5         "name": "Zaphod Beeblebrox",  
6         "species": "Betelgeusian"  
7     }  
8 }
```

We are able to simply serialize it:

```
● ● ●  
1 with open("json_write.json", "w") as write_file:  
2     json.dump(data, write_file)
```

JSON

If we now have a Python object:

```
● ● ●  
1 data = {  
2     "president": {  
3         "name": "Zaphod Beeblebrox",  
4         "species": "Betelgeusian"  
5     }  
6 }
```

and this is to be used directly:

```
● ● ●  
1 json_string = json.dumps(data)
```

JSON

Of course, we can also deserialize the JSON file again:

```
1 # 05_json_read.py
2
3 import json
4
5 with open("04_json_write.json", "r") as read_file:
6     data = json.load(read_file)
7
8 print(data)
9 print(data['president']['species'])
```

JSON - Exercise

Now it's time for a real example, please create a new file in the userfolder (name for example: json_example.py)



```
1 import json
2 import requests
```



```
1 response = requests.get("https://jsonplaceholder.typicode.com/todos")
2 todos = json.loads(response.text)
```



```
1 type(todos)
2 todos[:10]
```

JSON - Exercise Part 2

Now extend json_example.py so that the first 20 todos are written to a file (json_example.json).

Pandas

As soon as we have to deal with larger amounts of data, or have to switch very flexibly from csv to json to excel, the Pandas package comes into play.

Note: Pandas is very generously explained in APPD, I only use it here to read data easier.

Pandas

This time it's about an HR data set:



```
1 Name,HireDate,Salary,SickDaysremaining
2 Graham Chapman,03/15/14,50000.00,10
3 John Cleese,06/01/15,65000.00,8
4 Eric Idle,05/12/14,45000.00,10
5 Terry Jones,11/01/13,70000.00,3
6 Terry Gilliam,08/12/14,48000.00,7
7 Michael Palin,05/23/13,66000.00,8
```



```
1 # 06_read_pandas.py
2
3 import pandas
4
5
6 df = pandas.read_csv('06_hrdata.csv')
7 print(df)
```

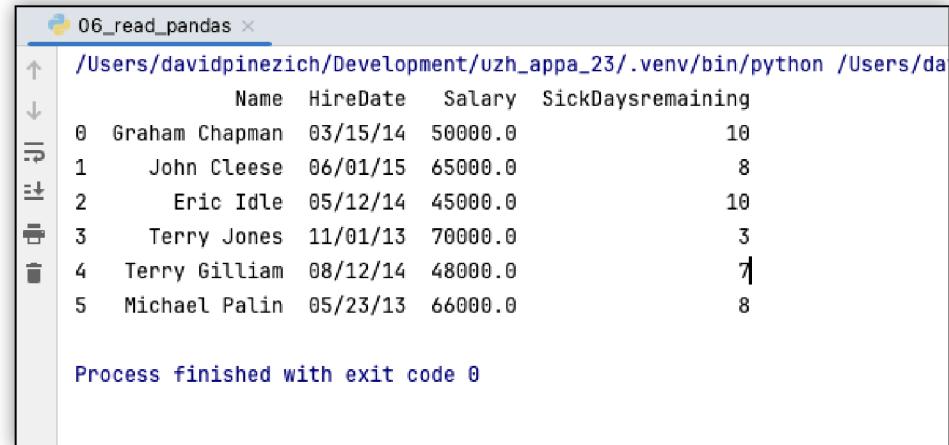
Pandas

Now we get a so-called DataFrame

```
06_read_pandas ×  
/Users/davidpinezich/Development/uZH_APPA_23/.venv/bin/python /Users/davidpinezich/Development/uZH_APPA_23/06_read_pandas.py  
      Name    HireDate   Salary  SickDaysremaining  
0  Graham Chapman  03/15/14  50000.0                  10  
1  John Cleese    06/01/15  65000.0                  8  
2  Eric Idle     05/12/14  45000.0                 10  
3  Terry Jones   11/01/13  70000.0                  3  
4  Terry Gilliam  08/12/14  48000.0                  1  
5  Michael Palin  05/23/13  66000.0                  8  
  
Process finished with exit code 0
```

This is both **highly efficient** and **easy to operate**.

Pandas



The screenshot shows a Jupyter Notebook cell with the title "06_read_pandas". The code executed is "/Users/davidpinezich/Development/uzh_appa_23/.venv/bin/python /Users/da...". The output displays a Pandas DataFrame with 6 rows and 5 columns: Name, HireDate, Salary, SickDays, and remaining. The data is as follows:

	Name	HireDate	Salary	SickDays	remaining
0	Graham Chapman	03/15/14	50000.0	10	
1	John Cleese	06/01/15	65000.0	8	
2	Eric Idle	05/12/14	45000.0	10	
3	Terry Jones	11/01/13	70000.0	3	
4	Terry Gilliam	08/12/14	48000.0	7	
5	Michael Palin	05/23/13	66000.0	8	

At the bottom of the output, it says "Process finished with exit code 0".

- Initially, Pandas recognized that the first line of the CSV file contains column names and used them automatically.
- However, Pandas also uses zero-based integer indices in the DataFrame. This is because we did not tell it **what our index should be called**.
- If you look at the data types of our columns, you will notice that Pandas has properly converted the Salary and Remaining Sick Days columns to numbers, but the Hiring Date column is still a string.

Pandas

Now we want to query something:



```
1 print(type(df['Hire Date'][0]))
```

Or structure it with a different index:



```
1 # 07_hrdata_name.csv
2
3 import pandas
4
5
6 df = pandas.read_csv('06_hrdata.csv', index_col='Name')
7 print(df)
```

Pandas

And also parse data:

```
1 # 08_hrdata_parsed.py
2
3 import pandas
4
5
6 df = pandas.read_csv('06_hrdata.csv',
7                     index_col='Employee',
8                     parse_dates=['Hired'],
9                     header=0,
10                    names=['Employee', 'Hired', 'Salary', 'Sick Days'])
11 print(df)
```

Pandas

And finally, we simply create a JSON from it:

```
● ● ●  
1 # 09_hrdata_to_json.py  
2  
3 import pandas  
4  
5 df = pandas.read_csv('06_hrdata.csv',  
6                      index_col='Employee',  
7                      parse_dates=['Hired'],  
8                      header=0,  
9                      names=['Employee', 'Hired', 'Salary', 'Sick Days'])  
10  
11 df.to_json('09_hrdata_modified.json')
```

Pandas

But that's not all, please copy the file 09_hrdata_to_json into your userfolder.

Try to create an Excel file and an HTML file file.

For Excel you need the package openpyxl (pip install openpyxl).



```
df.to_('09_hrdata_modified.json') You, A minute ago • Uncommitted changes
m to_csv(self, path_or_buf, sep, na_rep, float_form... NDFrame
m to_gbq(self, destination_table, project_id, chu... DataFrame
m to_hdf(self, path_or_buf, key, mode, complevel, c... NDFrame
m to_orc(self, path, engine, index, engine_kwargs) DataFrame
m to_sql
m to_clipboard(self, excel, sep, kwargs)
m to_dict(self, orient, into)
m to_excel
m to_feather(self, path, kwargs)
m to_html(self, buf, columns, col_space, header, ...
m to_json
m to_latex(self, buf, columns, header, index, no_no NDFrame
Press ⇧ to insert, ⇩ to replace Next Tip
```

Flask

We received a wonderful HTML from the last task.

```
<table border="1" class="dataframe">
  <thead>
    <tr style="...">
      <th></th>
      <th>Hired</th>
      <th>Salary</th>
      <th>Sick Days</th>
    </tr>
    <tr>
      <th>Employee</th>
      <th></th>
      <th></th>
      <th></th>
    </tr>
  </thead>
  <tbody>
    <tr>
      <th>Graham Chapman</th>
      <td>2014-03-15</td>
      <td>50000.0</td>
      <td>10</td>
    </tr>
    <tr> You, 3 minutes ago • Uncommitted changes
  </tbody>
</table>
```

Flask

We received a wonderful HTML from the last task.

We can also simply view this in the browser:



Employee	Hired	Salary	Sick Days
Graham Chapman	2014-03-15	50000.0	10
John Cleese	2015-06-01	65000.0	8
Eric Idle	2014-05-12	45000.0	10
Terry Jones	2013-11-01	70000.0	3
Terry Gilliam	2014-08-12	48000.0	7
Michael Palin	2013-05-23	66000.0	8

```
<table border="1" class="dataframe">
  <thead>
    <tr style="...">
      <th></th>
      <th>Hired</th>
      <th>Salary</th>
      <th>Sick Days</th>
    </tr>
    <tr>
      <th>Employee</th>
      <th></th>
      <th></th>
      <th></th>
    </tr>
  </thead>
  <tbody>
    <tr>
      <td>Graham Chapman</td>
      <td>2014-03-15</td>
      <td>50000.0</td>
      <td>10</td>
    </tr>
    <tr> You, 3 minutes ago • Uncommitted changes
```

Flask

We will use this example to familiarize ourselves with Flask.

Flask is a web framework. This means that Flask provides us with tools, libraries, and technologies that can be used to create a web application.

This web application can consist of a few web pages, a blog or a wiki, or even be as large as a web-based calendar application or a commercial website.

Flask



```
1 # 10_flask_starter.py
2
3 import flask
4
5 # Create the application.
6 APP = flask.Flask(__name__)
7
8 @APP.route('/')
9 def index():
10     """ Displays the index page accessible at '/' found in templates
11     """
12     return flask.render_template('10_hrdata_modified.html')
13
14
15 if __name__ == '__main__':
16     APP.debug=True
17     APP.run()
```

Flask

```
P [WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.  
* Running on http://127.0.0.1:5000  
Press CTRL+C to quit  
* Restarting with stat  
* Debugger is active!  
* Debugger PIN: 342-442-388  
127.0.0.1 - - [18/Jan/2024 00:36:27] "GET / HTTP/1.1" 200 -  
127.0.0.1 - - [18/Jan/2024 00:36:27] "GET /favicon.ico HTTP/1.1" 404 -  
* Detected change in '/Users/davidpinezich/Development/uzh_appa_23/data-automation/02-examples/10_flask_starter.py', reloading  
* Restarting with stat  
* Debugger is active!  
* Debugger PIN: 342-442-388  
* Detected change in '/Users/davidpinezich/Development/uzh_appa_23/data-automation/02-examples/10_flask_starter.py', reloading  
* Restarting with stat  
* Debugger is active!  
* Debugger PIN: 342-442-388
```

Trial and error is the order of the day on our local server.

Flask

This website could now easily be placed on the Internet using DigitalOcean or PythonAnywhere for around USD 5 per month.

<https://www.pythonanywhere.com/>

<https://www.digitalocean.com/>

Flask

This way we can easily post the results of our automation on the Internet if desired.

Excel

Next, we also want to look at Excel, because a lot of data is (and will continue to be) in Excel.

Excel - Exercise

With the openpyxl package, working with Excel is surprisingly easy.



```
1 # 11_data_from_excel.py
2
3 import openpyxl
4
5 wb = openpyxl.load_workbook('ch-bevoelkerung.xlsx')
6 wb.sheetnames # The workbook's sheets' names.
7 ['VZ RFP 1850', 'VZ RFP 1860', 'VZ RFP 1870']
8
9 sheet = wb['VZ RFP 1850'] # Get a sheet from the workbook.
10
11 type(sheet)
12 # <class 'openpyxl.worksheet.worksheet.Worksheet'>
13 print(sheet.title) # Get the sheet's title as a string.
14
15 print(sheet['C15'].value)
16 print(sheet['D15'].value)
```

Excel - Exercise

Try displaying all the municipality names:

Excel - Exercise

Try displaying all the municipality names:

```
● ● ●  
1 # 12_data_from_excel_iterate.py  
2  
3 import openpyxl  
4  
5 wb = openpyxl.load_workbook('ch-bevoelkerung.xlsx')  
6 wb.sheetnames # The workbook's sheets' names.  
7 ['VZ RFP 1850', 'VZ RFP 1860', 'VZ RFP 1870']  
8  
9 sheet = wb['VZ RFP 1850'] # Get a sheet from the workbook.  
10  
11 type(sheet)  
12 # <class 'openpyxl.worksheet.worksheet.Worksheet'>  
13 print(sheet.title) # Get the sheet's title as a string.  
14  
15 print(sheet['C15'].value)  
16 print(sheet['D15'].value)  
17  
18 for cell in sheet['C']:  
19     print(cell.value)
```

Excel - Exercise

Now try the same thing but with pandas.



```
1 import pandas  
2  
3 df = pandas.read_excel(open('ch-bevoelkerung.xlsx', 'rb'),  
4                         sheet_name='VZ RFP 1850', index_col=2, header=2)
```

Visualize data

Visualizing data is often something that follows from our previous work.

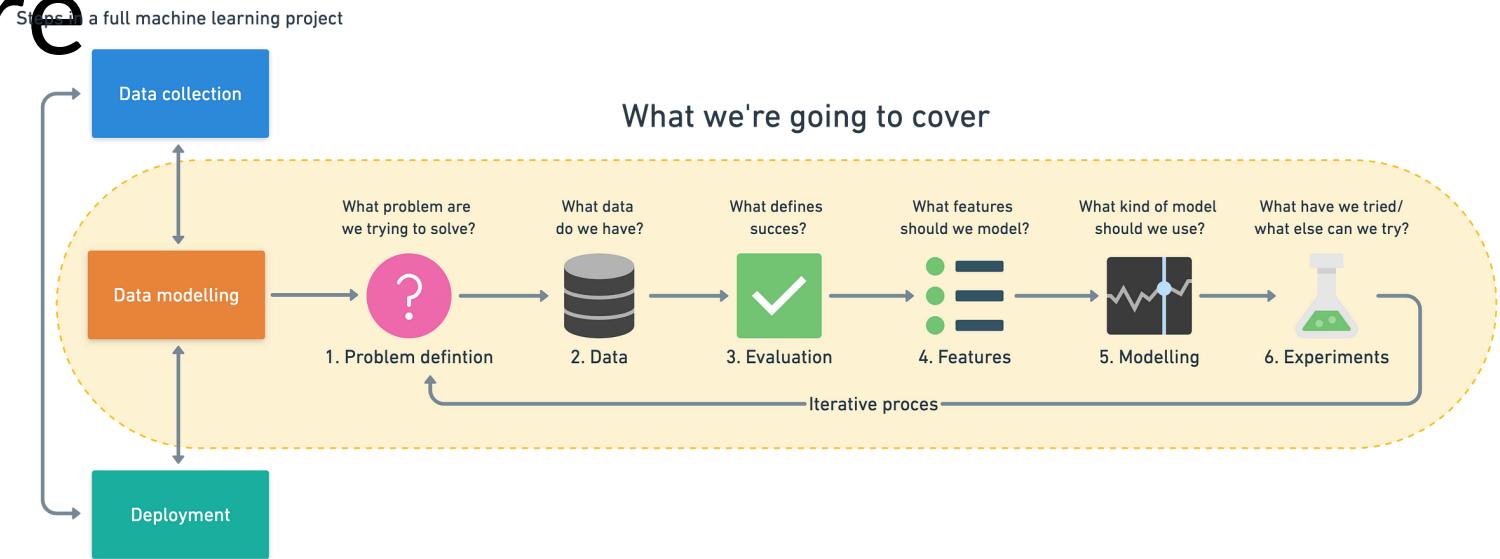
The topic is very broad, and the following is just an overview.

Visualization should be an iterative process

1. problem definition
2. data
3. evaluation (how do I define success)
4. features (what exactly do I need)
5. visualization
6. testing (experimenting)

A step back is possible at any time due to new findings (and not necessarily a bad thing)!

(Rough) procedure



This can also be used for a visualization:

1. problem definition
2. data
3. evaluation (how do I define success)
4. features (what exactly do I need)
5. **visualization**
6. testing (experimenting)

Dog-Science

The trendy brand for your dog



Problem definition (problem definition)

- Dog-Science is a trendy brand from Asia and would like to expand to Zurich.
- In Dog-Science's home country, it is well known where most dog owners live and which products are preferred. But Switzerland, and Zurich in particular, is unknown.
- Dog-Science's budget is limited, but the city of Zurich has (notionally) promised data and support for a pop-up store
- Can we improve Dog-Science's go-to-market position?

A first "look" at the data

The screenshot shows the homepage of the Stadt Zürich Open Data portal. The header features the logo 'Stadt Zürich Open Data' and navigation links for 'Startseite', 'Datensätze', and 'Kategorien'. A search bar is on the right. The main content area displays the dataset 'Hundebestand der Stadt Zürich'. It includes a 'Lizenz' section with 'Creative Commons CCZero' and an 'OPEN DATA' button. Below this, there are tabs for 'Datensatz', 'Kategorien', and 'Showcases'. A detailed description of the dataset is provided, stating it contains information about dog owners and their dogs from the city's dog register. The 'Daten und Ressourcen' section lists a CSV file named '20200306_hundehalter.csv' and a 'Entdecke' button.

Hundebestand der Stadt Zürich

In diesem Datensatz finden Sie Angaben zu Hunden und deren Besitzerinnen und Besitzern aus dem aktuellen Bestand des städtischen Hunderegisters. Bei den hundehaltenden Personen sind Informationen zur Altersgruppe, dem Geschlecht und dem statistischen Quartier des Wohnorts angegeben. Zu jedem Hund ist die Rasse, der Rassetyp, das Geschlecht, das Geburtsjahr und die Farbe erfasst. Das Hunderegister wird von der Abteilung Hundekontrolle der Stadtpolizei Zürich geführt.

Daten und Ressourcen

20200306_hundehalter.csv
Comma-Separated Values. Weitere Informationen zu CSV finden Sie in unserer...

Entdecke

A first "look" at the data

HALTER_ID	ALTER	GESCHLECHT	STADTKREIS	STADTQUARTIER	RASSE1	GEBURTSJAHR_HUND	GESCHLECHT_HUND	HUNDEFARBE
574	61-70	w		2	23 Mischling gross	2013	w	schwarz
695	41-50	m		6	63 Labrador Retriever	2012	w	braun
893	71-80	w		7	71 Mittelschnauzer	2010	w	schwarz
916	41-50	m		3	34 Mischling klein	2015	w	hellbraun
1177	51-60	m		10	102 Shih Tzu	2011	m	schwarz/weiss
4054	51-60	w		11	111 Lagotto Romagnolo	2016	w	weiss/beige
4135	41-50	w		9	91 Mischling klein	2016	w	schwarz
4206	71-80	w		8	82 Havaneser	2016	w	hellbraun/weiss
4281	61-70	w		9	91 Chihuahua	2011	w	hellbraun
4388	61-70	w		11	115 Mops	2006	m	beige
4726	51-60	m		5	52 Mischling gross	2007	m	schwarz
4726	51-60	m		5	52 Golden Retriever	2013	w	creme
4747	61-70	m		2	24 Chihuahua	2013	m	weiss/braun
4850	51-60	m		4	42 Chihuahua	2013	w	beige
4862	51-60	m		4	42 Mops	2006	m	braun
5040	61-70	m		10	102 Labrador Retriever	2016	w	gelb
5088	61-70	m		7	72 Labrador Retriever	2014	w	schwarz
5113	61-70	m		11	119 Beagle	2010	w	tricolor
5113	61-70	m		11	119 Chihuahua	2017	w	beige
5225	71-80	m		3	34 Lagotto Romagnolo	2007	m	braun
5227	71-80	m		10	101 Border Terrier	2011	w	tricolor

A first "look" at the data

- The data is of good, but not perfect quality
 - The breeds were unfortunately defined very imprecisely
- Often defined as "mixed breed large" or "mixed breed small", but not further defined
- Often represented
 - Chihuahua 573
 - Labrador Retriever 426
- Rarely represented
 - Austrian Pinscher 1
 - Daisy-Dog 1

A first "look" at the data

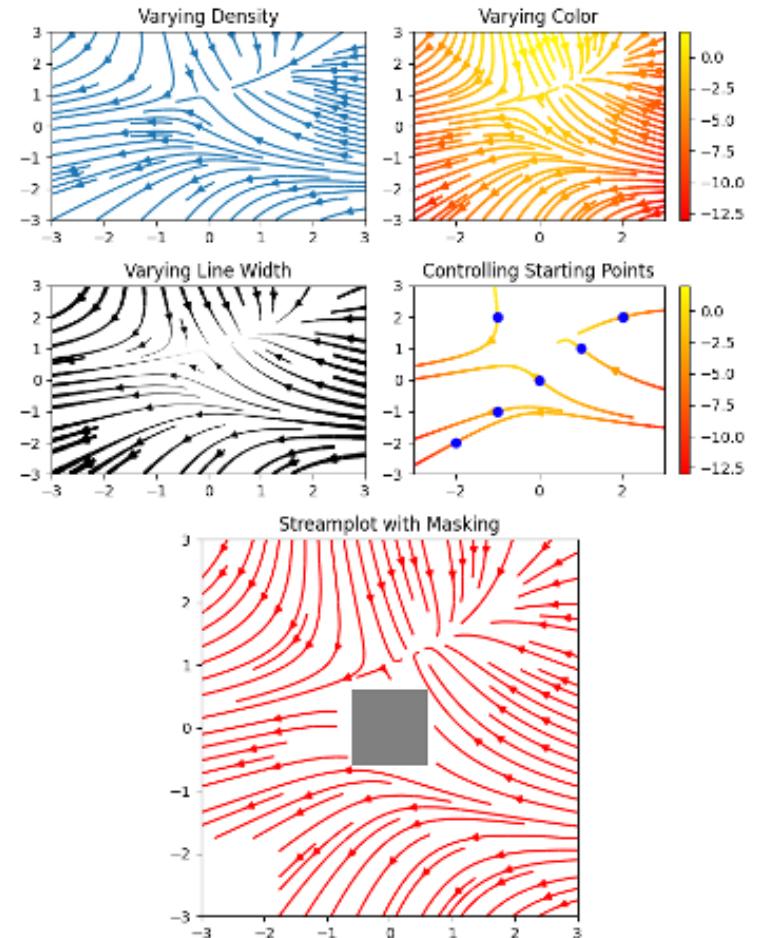
- Significantly more female than male dogs (5402 to 2439)
Black is the most common coat color (800), followed by tricolor (725) and white (634), less often "black mottled (melliert)" or "yellow / black" once each.
- We take a look at the age of the dogs in the visualizations (please note the outliers)
- What is the dog owner / dog ratio?
 - 6562 dog owners have one registered dog
 - 581 dog owners have 2 or more registered dogs
 - The top dog owner has 14 registered dogs

Visualization

- Many libraries with different specialties
 - Matplotlib: For static, animated and interactive visualizations (one of the oldest libraries)
 - Pygal: Dynamic SVG charting library
 - Seaborn: Based on Matplotlib and offers a high-level interface for statistical graphics
 - Altair: Based on Vega/Vega Lite and is a declarative statistical visualization library
 - Ggplot2: System for the creation of declarative graphics
 - Plotly: Interactive and analyzable by the user
 - Bokeh: Library for interactive visualizations
 - Geoplotlib: Mainly for maps

Matplotlib

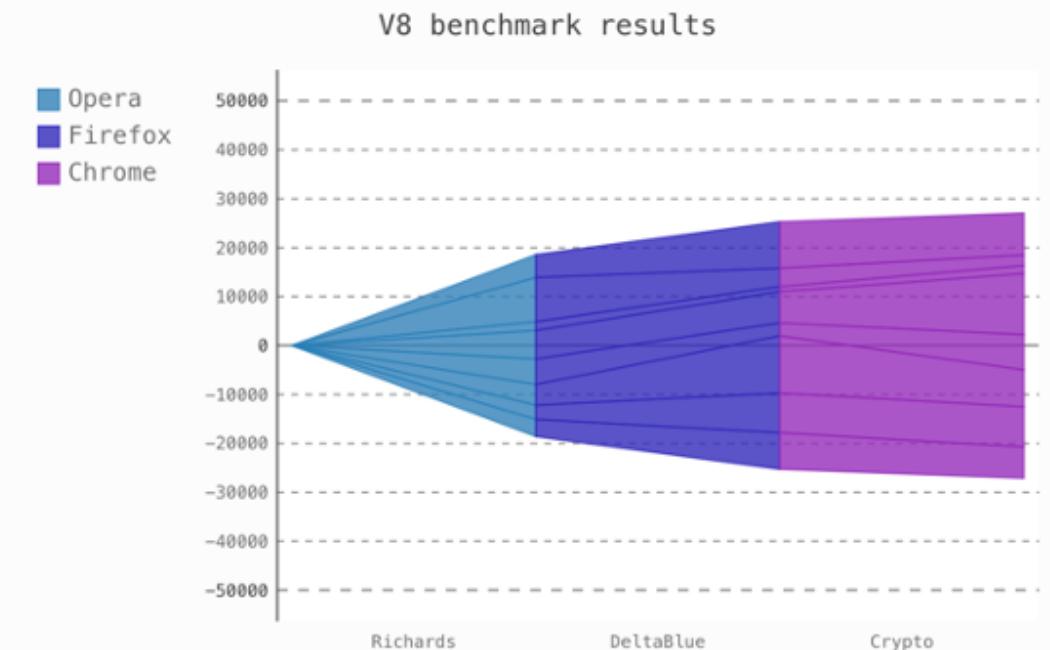
- Suitable for simple and complex representations
- Multiple representations possible via subplots
- Offers the largest and most general spectrum



Streamplot with various plotting options.

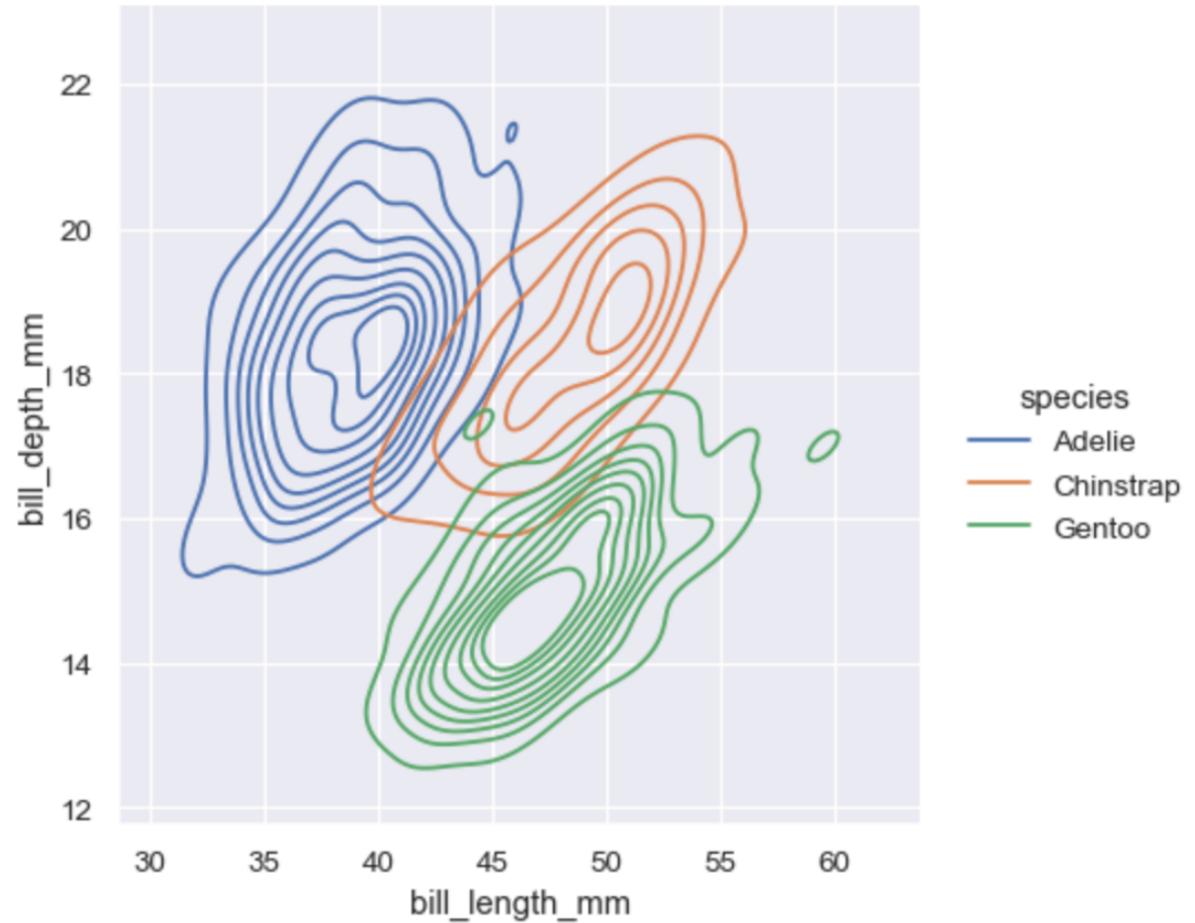
Pygal

- Unfortunately no longer (actively) developed
- Good range of functions and relatively easy to learn
- Easy to output as interactive HTML



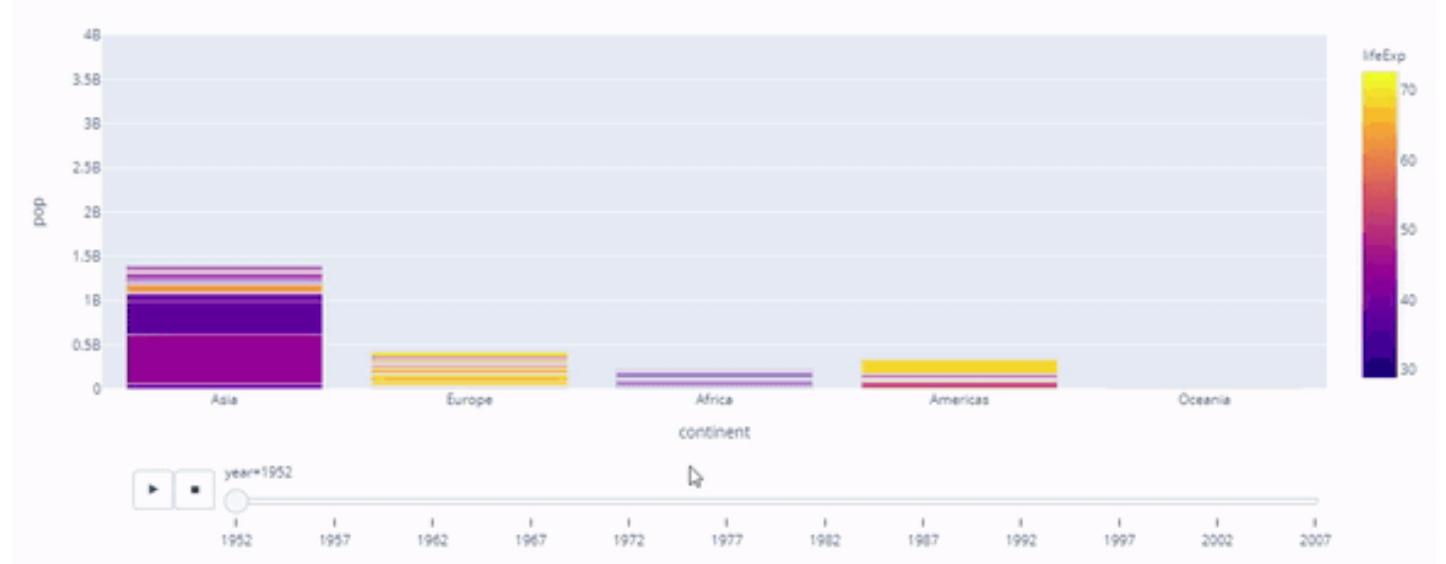
Seaborn

- Based on Matplotlib
- Especially suitable for statistical visualizations
- **Will be our first example**



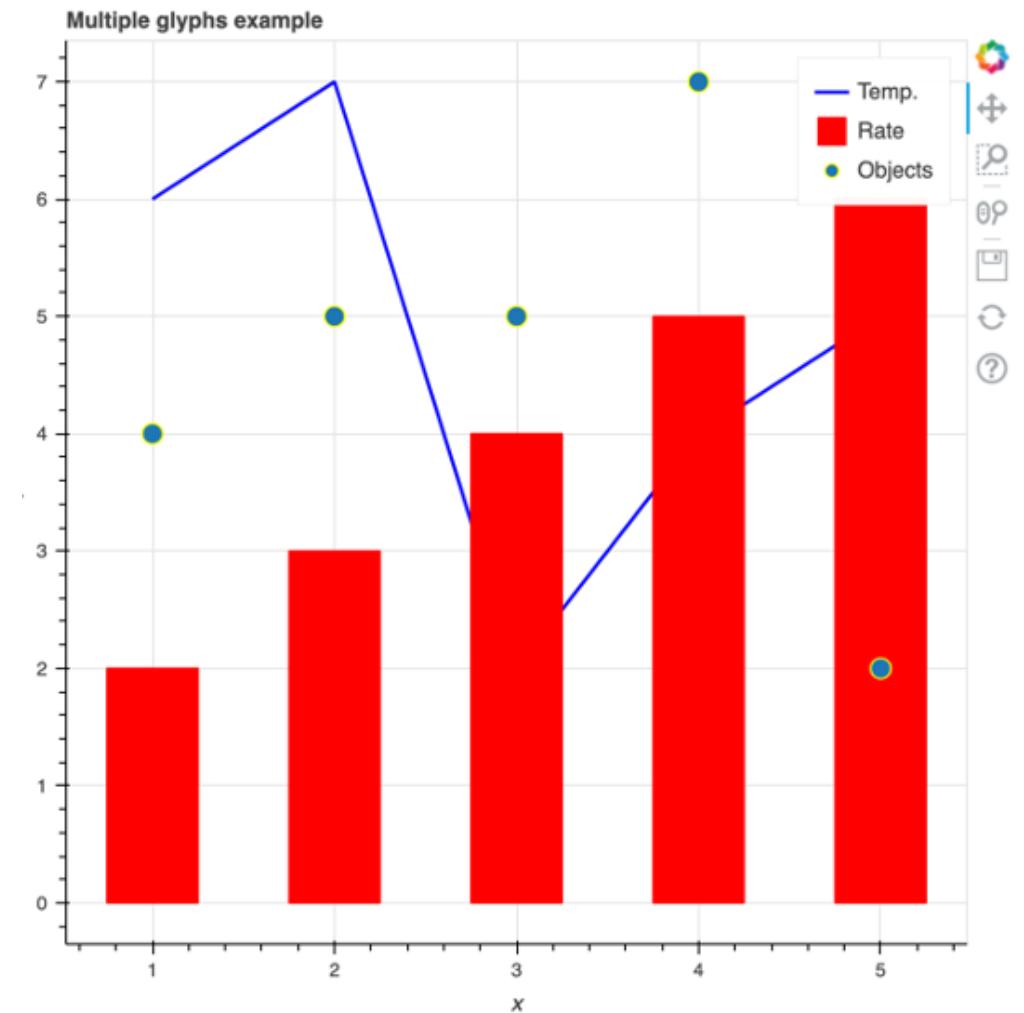
Plotly

- Offers a wide range of charts
- Can be animated
- Has "out of the box" manipulation tools

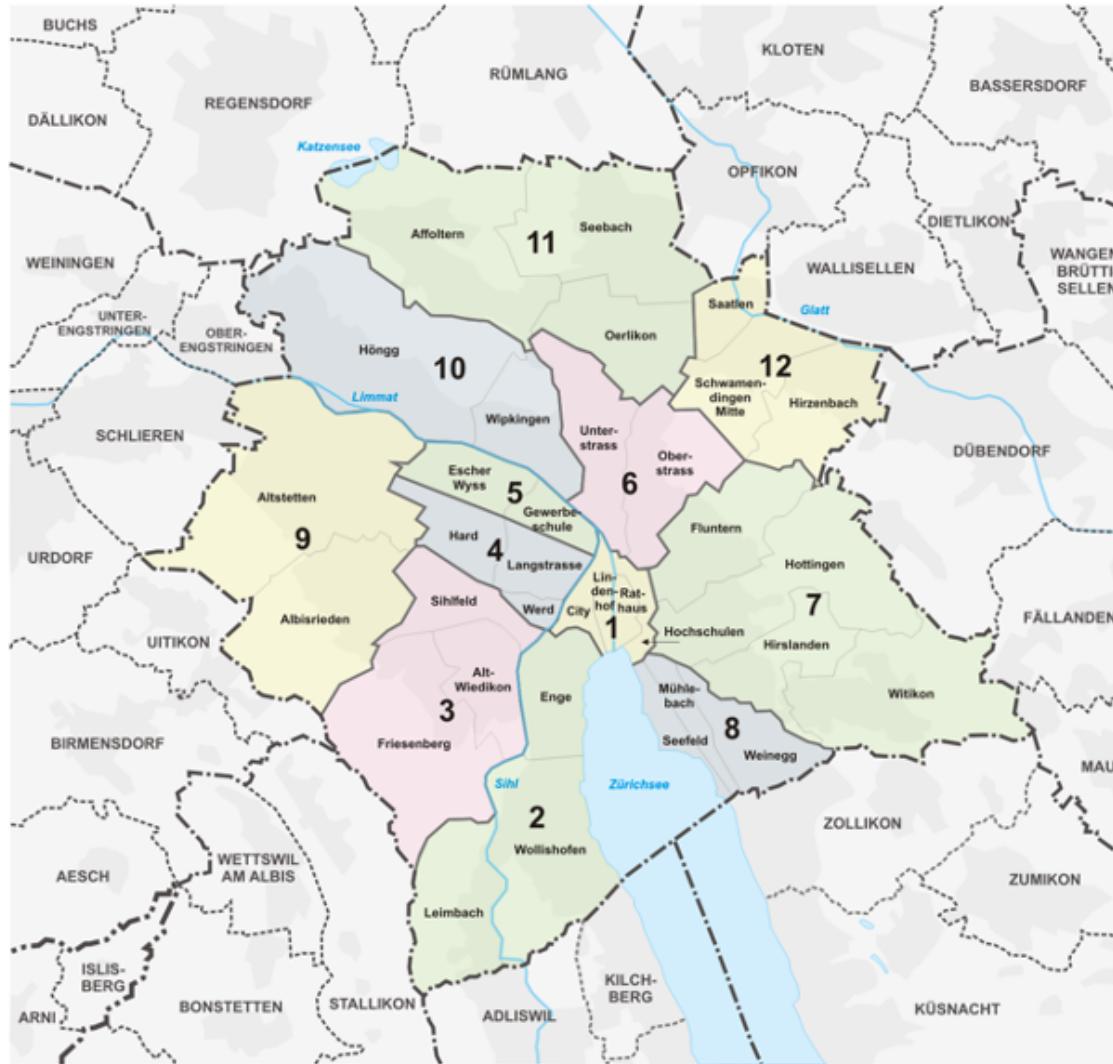


Bokeh

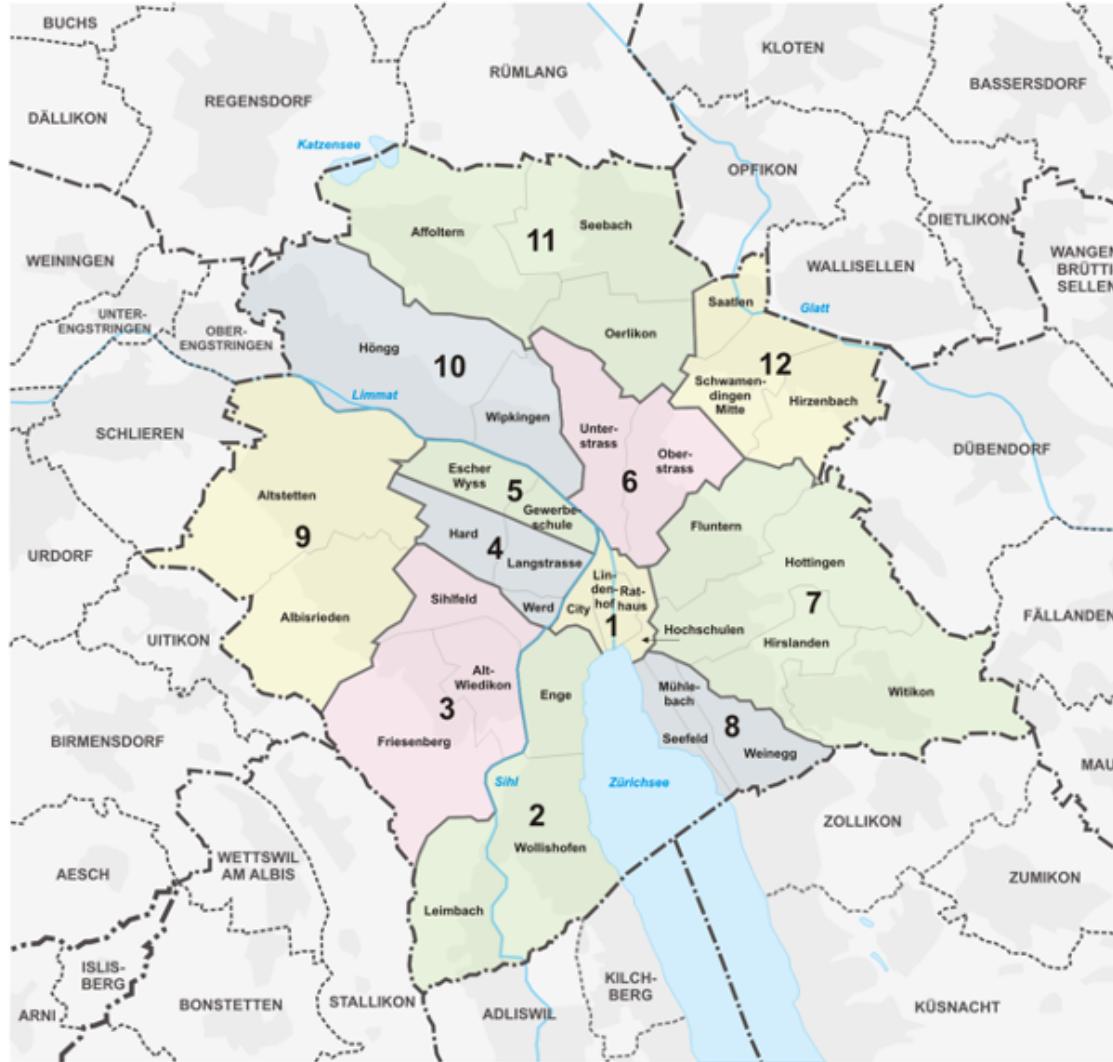
- Offers a wide range of interactive charts
- Has "out of the box" manipulation tools



The place is important...



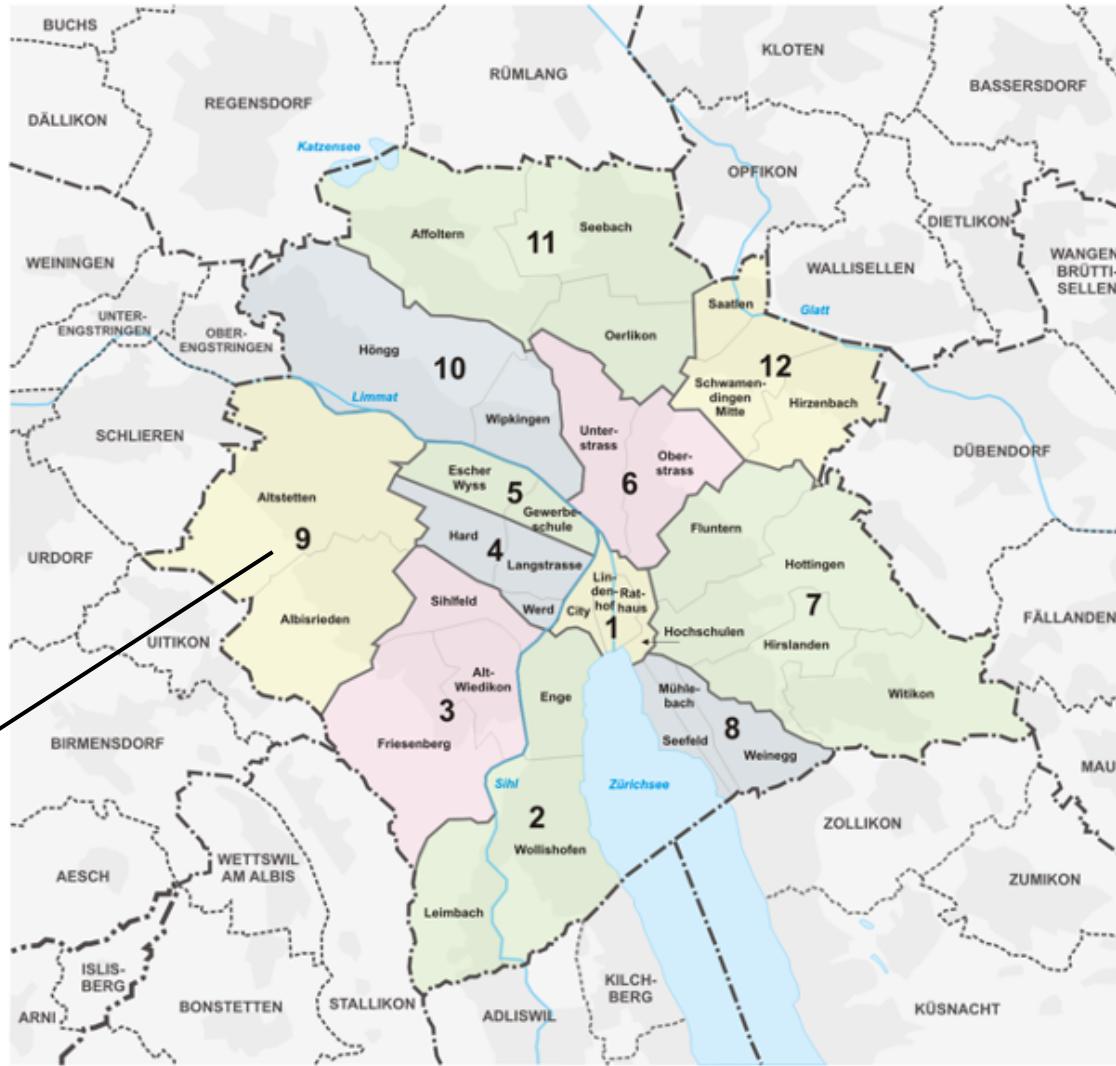
The place is important...



Place 3: 984

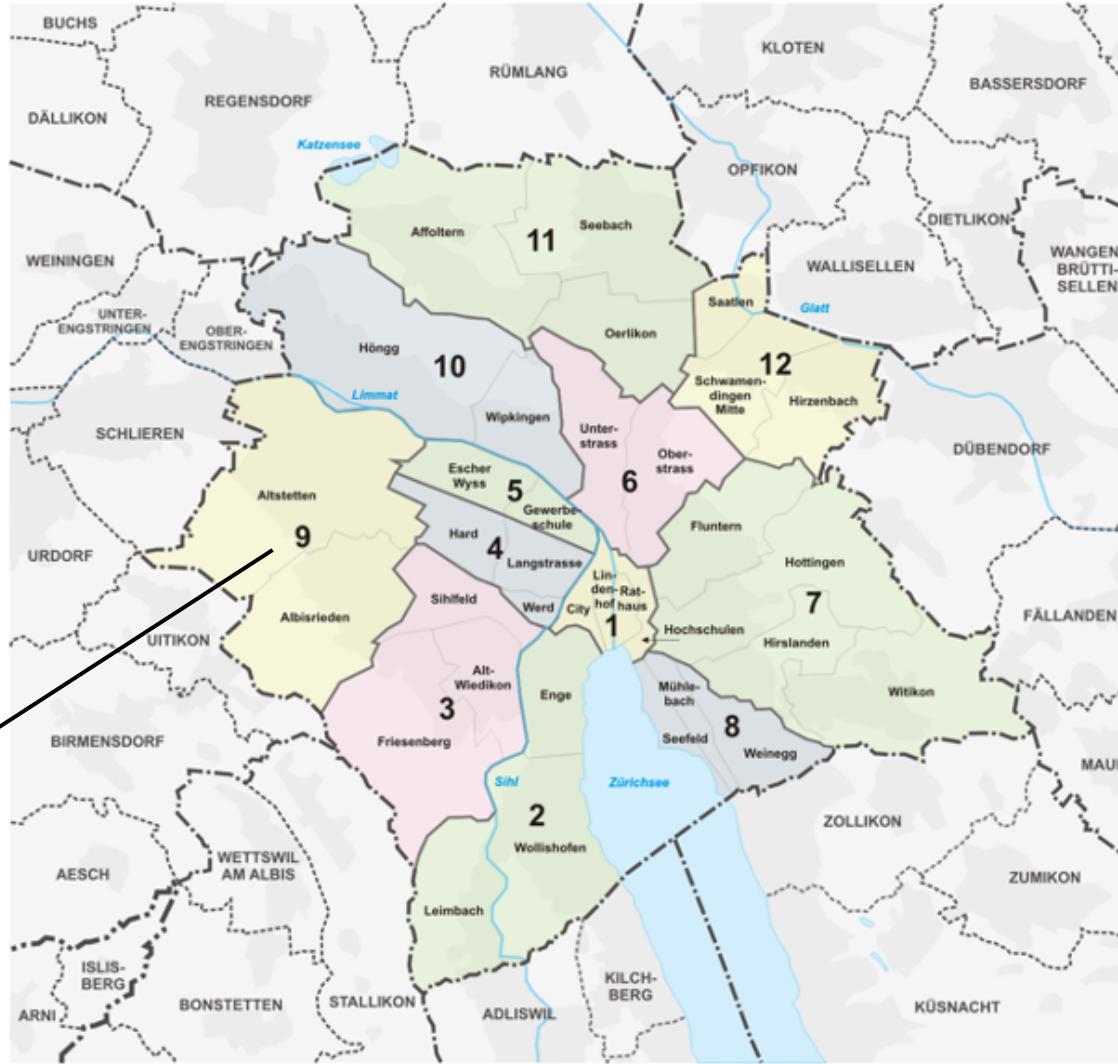
The place is important...

Place 3: 984



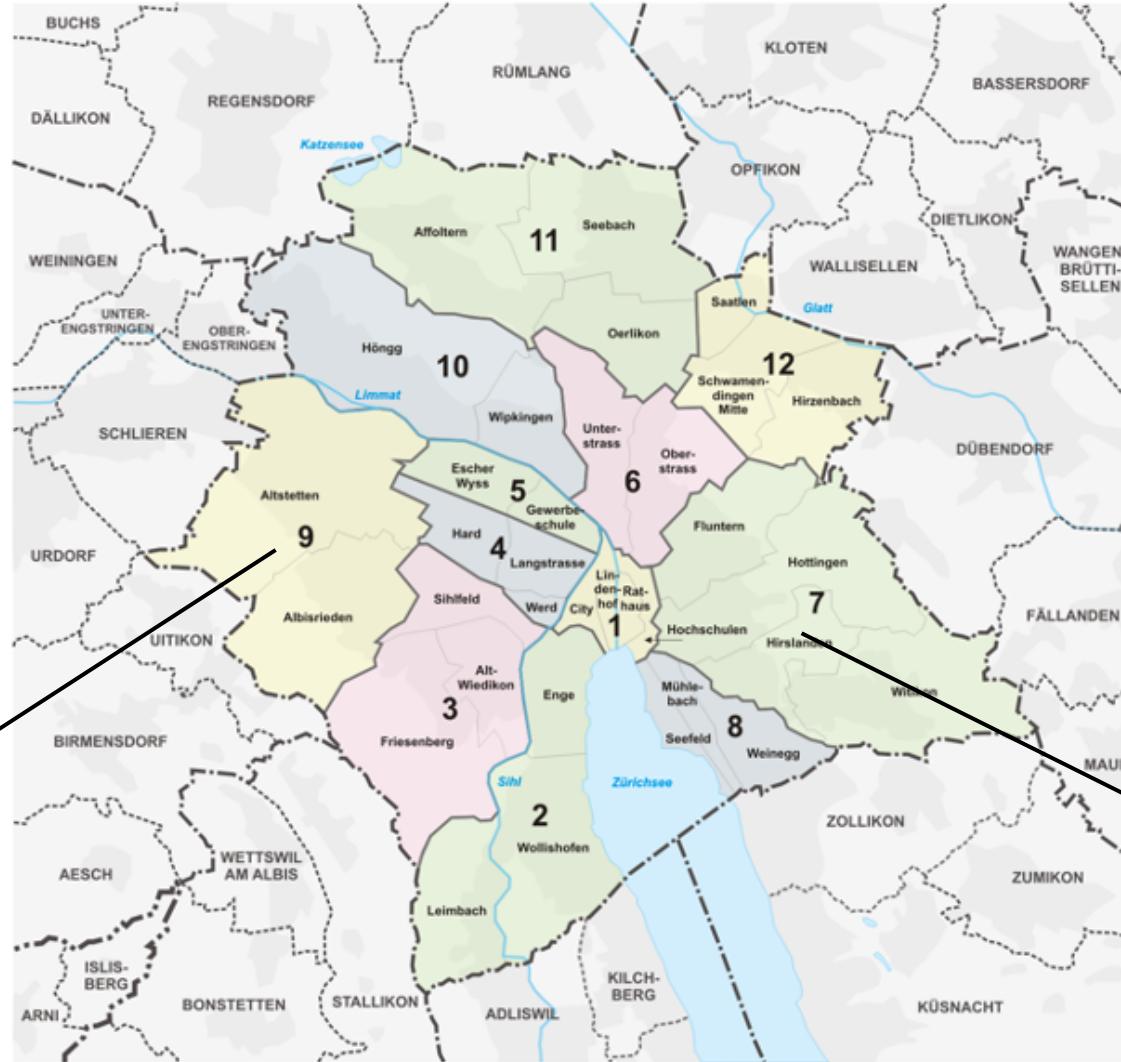
The place is important...

Place 3: 984



Place 2: 1087

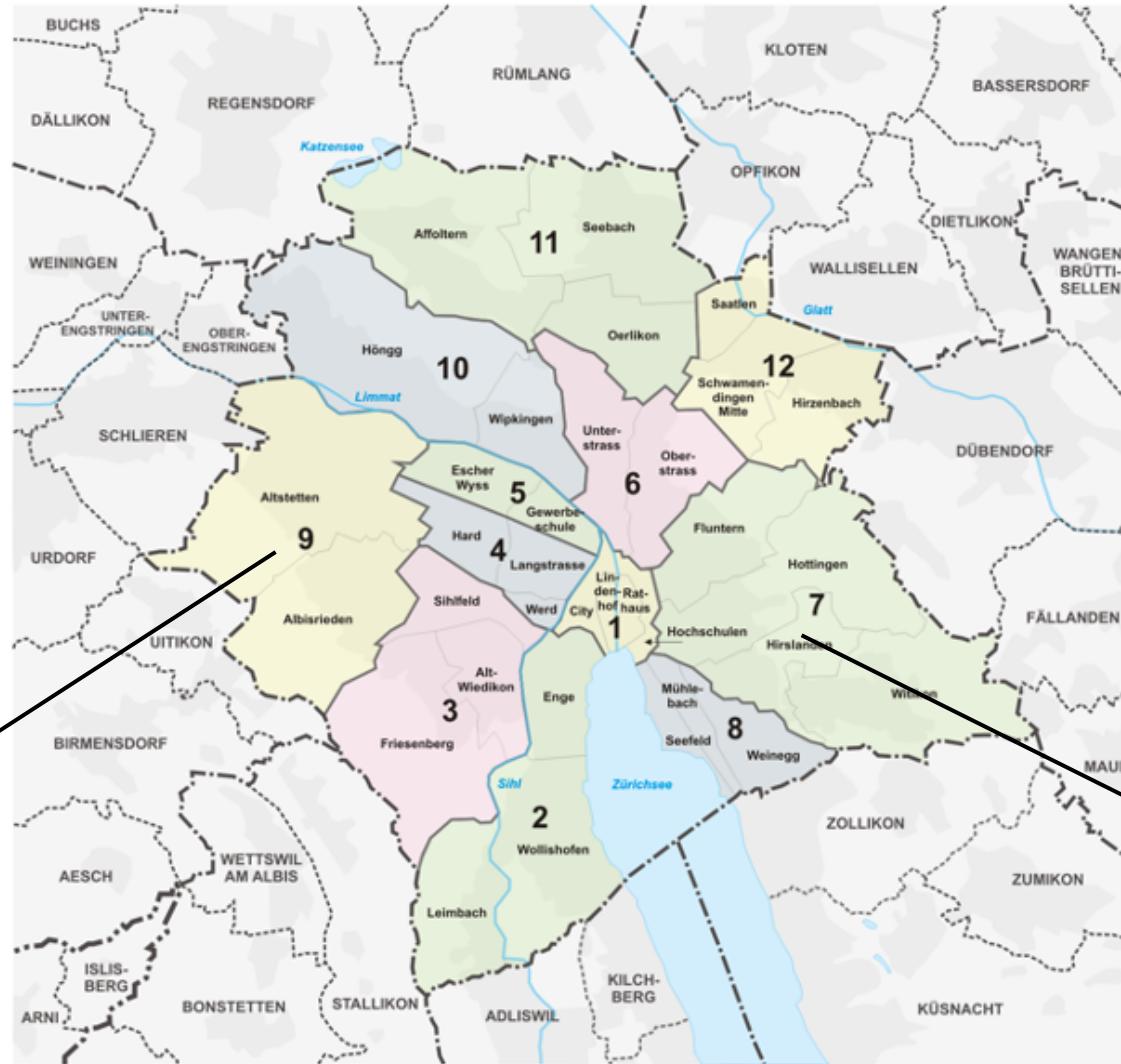
The place is important...



Place 3: 984

Place 2: 1087

The place is important...



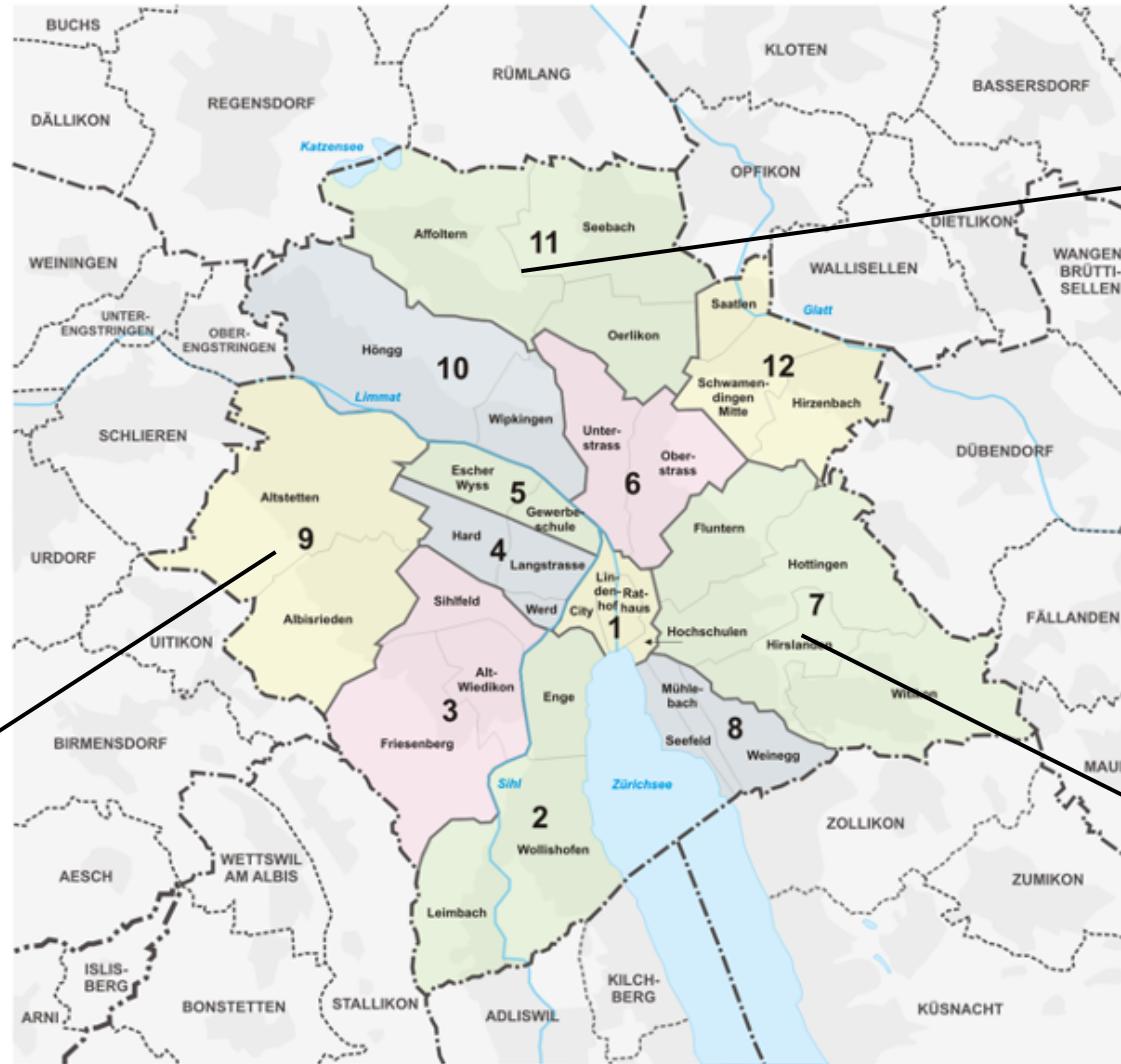
Place 3: 984

Place 1: 1352

Place 2: 1087

The place is important...

Place 3: 984



Place 1: 1352

Place 2: 1087

The place is important...

 Stadt Zürich
Open Data

Startseite Datensätze Kategorien 🔍

🏠 / Datensätze / Hundebestand der Stadt Zürich

🔒 Lizenz Creative Commons CCZero OPEN DATA

CSV **Datensatz** Kategorien Showcases

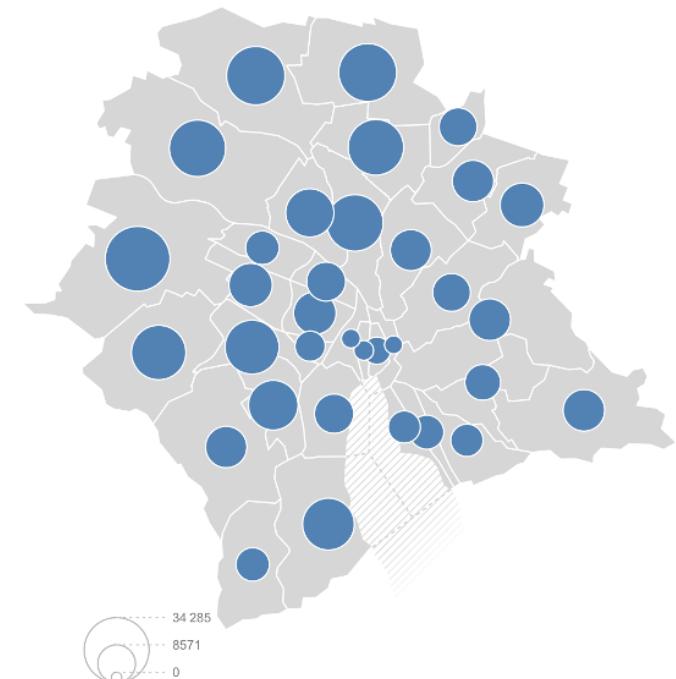
Hundebestand der Stadt Zürich

In diesem Datensatz finden Sie Angaben zu Hunden und deren Besitzerinnen und Besitzern aus dem aktuellen Bestand des städtischen Hunderegisters. Bei den hundehaltenden Personen sind Informationen zur Altersgruppe, dem Geschlecht und dem statistischen Quartier des Wohnorts angegeben. Zu jedem Hund ist die Rasse, der Rassetyp, das Geschlecht, das Geburtsjahr und die Farbe erfasst. Das Hunderegister wird von der Abteilung Hundekontrolle der Stadtpolizei Zürich geführt.

Daten und Ressourcen

 **20200306_hundehalter.csv**
Comma-Separated Values. Weitere Informationen zu CSV finden Sie in unserer... Entdecke

Bevölkerung nach Stadtquartier

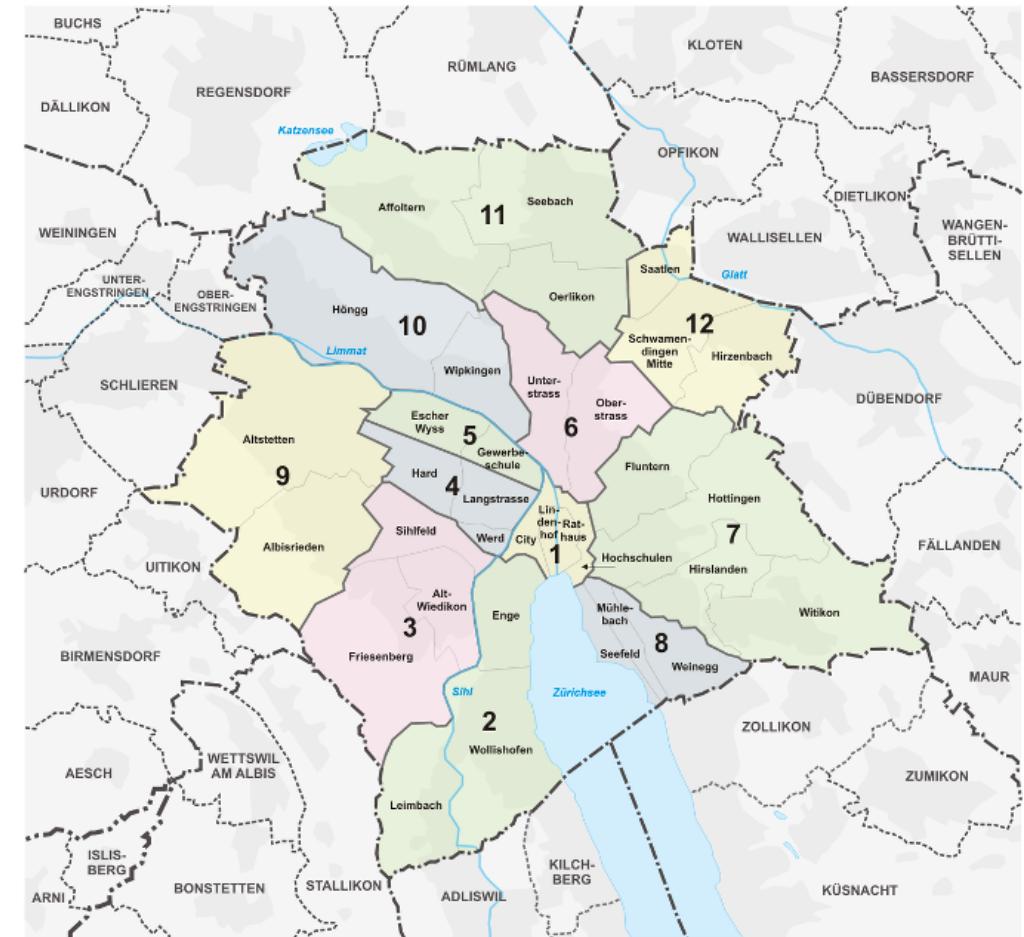
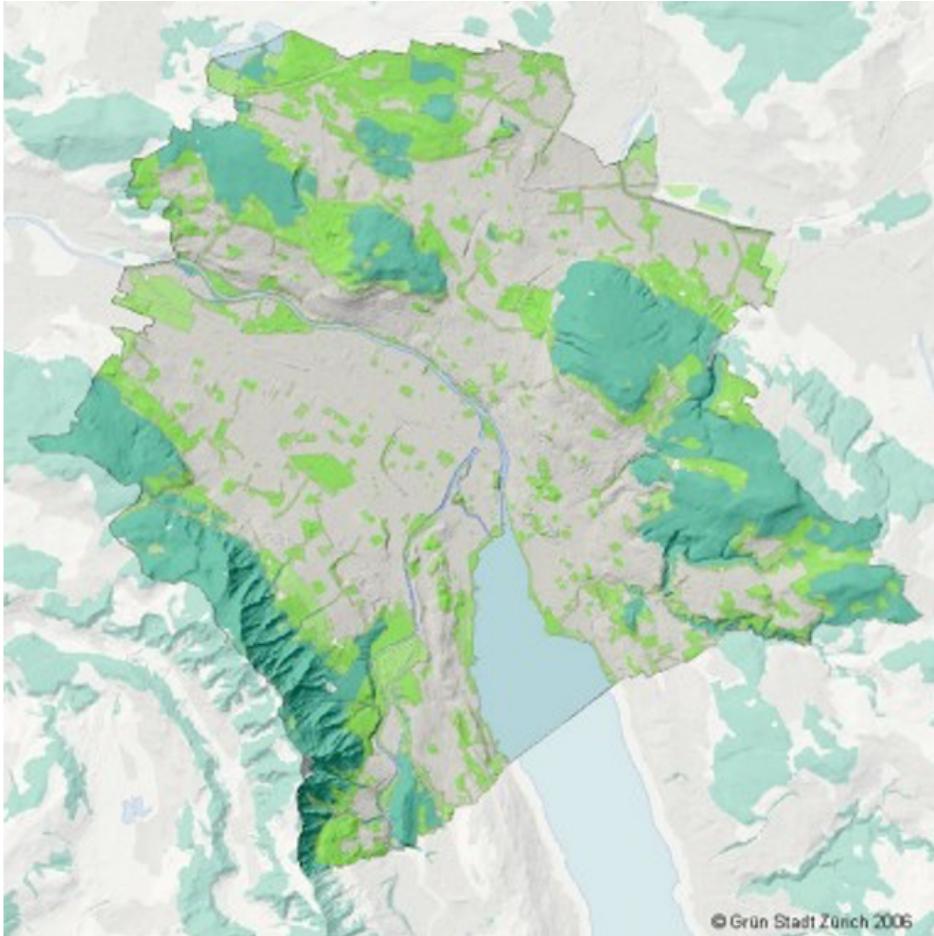


Pay particular attention to temporal differences:
Dog data: March 2020 / Population data: 1993-2020

Is the number of people important?

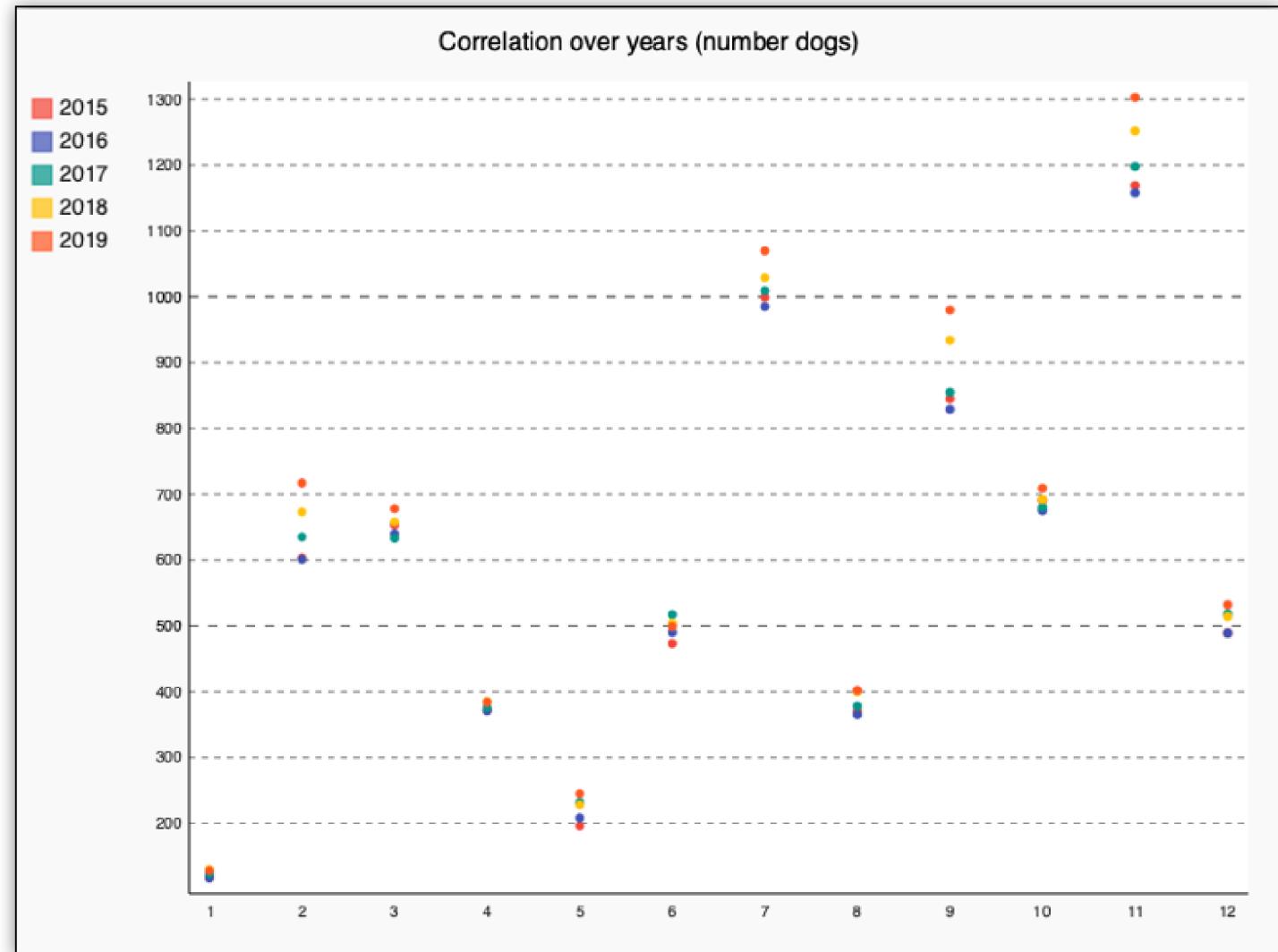
- Kreis 1; 5 831
- Kreis 2; 35 552
- Kreis 3; 50 756 => Many people, fewer dogs (697)
- Kreis 4; 29 034
- Kreis 5; 15 622
- Kreis 6; 35 317
- Kreis 7; 38 629
- Kreis 8; 17 456
- Kreis 9; 56 462
- Kreis 10; 41 044
- Kreis 11; 76 188
- Kreis 12; 32 845

Are Green Spaces important?



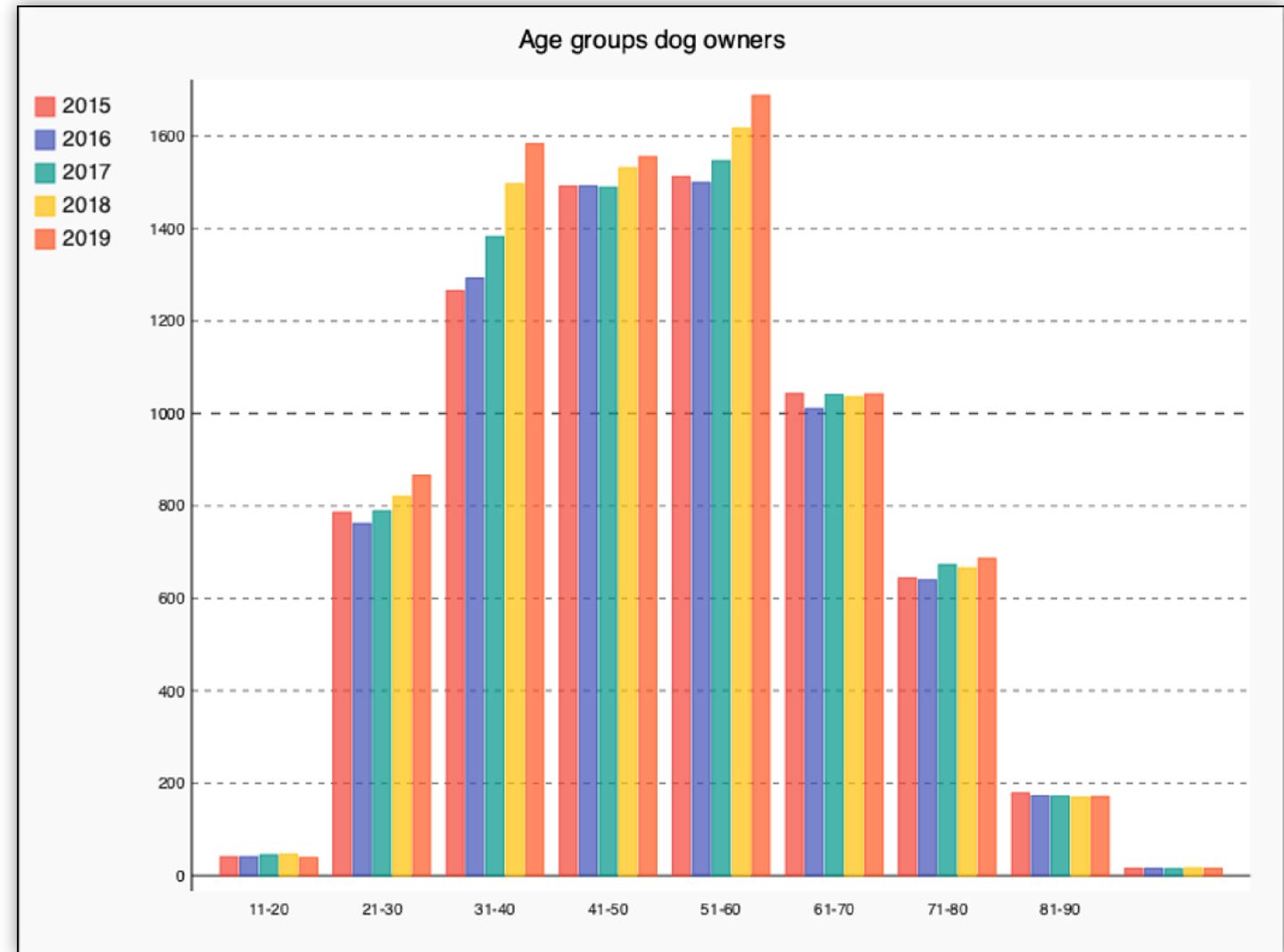
Calculate the correlation

See Step 2 in the documents



Age groups

The most pronounced group is 31-60



What next?

- Of course, many more methods could be used here:
 - K-Nearest Neighbor
 - Naive Bayes
 - Linear Regression
 - Multiple Regression
 - Logistic regression
 - Decision Trees
 - Neural Networks / Deep Learning
 - clustering
 - etc.

Conclusion: What do we recommend Dog- Science?

- We recommend Dog-Science to apply for a place **in Kreis 11 near a green area.**
- We recommend designing the product for 31-61 year olds.
- It may cost something "more" (see average salary in Zurich).
- It should not be too playful and should be easy to understand (interpretation of the age group).

Capstone Project

We build a dashboard together with Dash.

Capstone Project

Take a look at the following first:

<https://tracking-dashboard-app.herokuapp.com/dashboard>

Capstone Project

Our example is simpler, but combines dash (plotly) with our data skills and visualization ideas.

Capstone Project

Dash Hello World

Capstone Project

Dash Hello World



```
1 # 13_hello_world_dash.py
2
3 import openpyxl
4 from dash import Dash, html
5
6 app = Dash(__name__)
7
8 app.layout = html.Div([
9     html.Div(children='Hello World')
10])
11
12 if __name__ == '__main__':
13     app.run(debug=True)
```

Capstone Project

Connect with data

Capstone Project

Connect with data



```
1 # 14_dash_with_data.py
2
3 # Import packages
4 from dash import Dash, html, dash_table
5 import pandas as pd
6
7 # Incorporate data
8 df = pd.read_csv('gapminder2007.csv')
9
10 # Initialize the app
11 app = Dash(__name__)
12
13 # App layout
14 app.layout = html.Div([
15     html.Div(children='My First App with Data'),
16     dash_table.DataTable(data=df.to_dict('records'), page_size=10)
17 ])
18
19 # Run the app
20 if __name__ == '__main__':
21     app.run(debug=True)
```

Capstone Project

And a visualization

Capstone Project

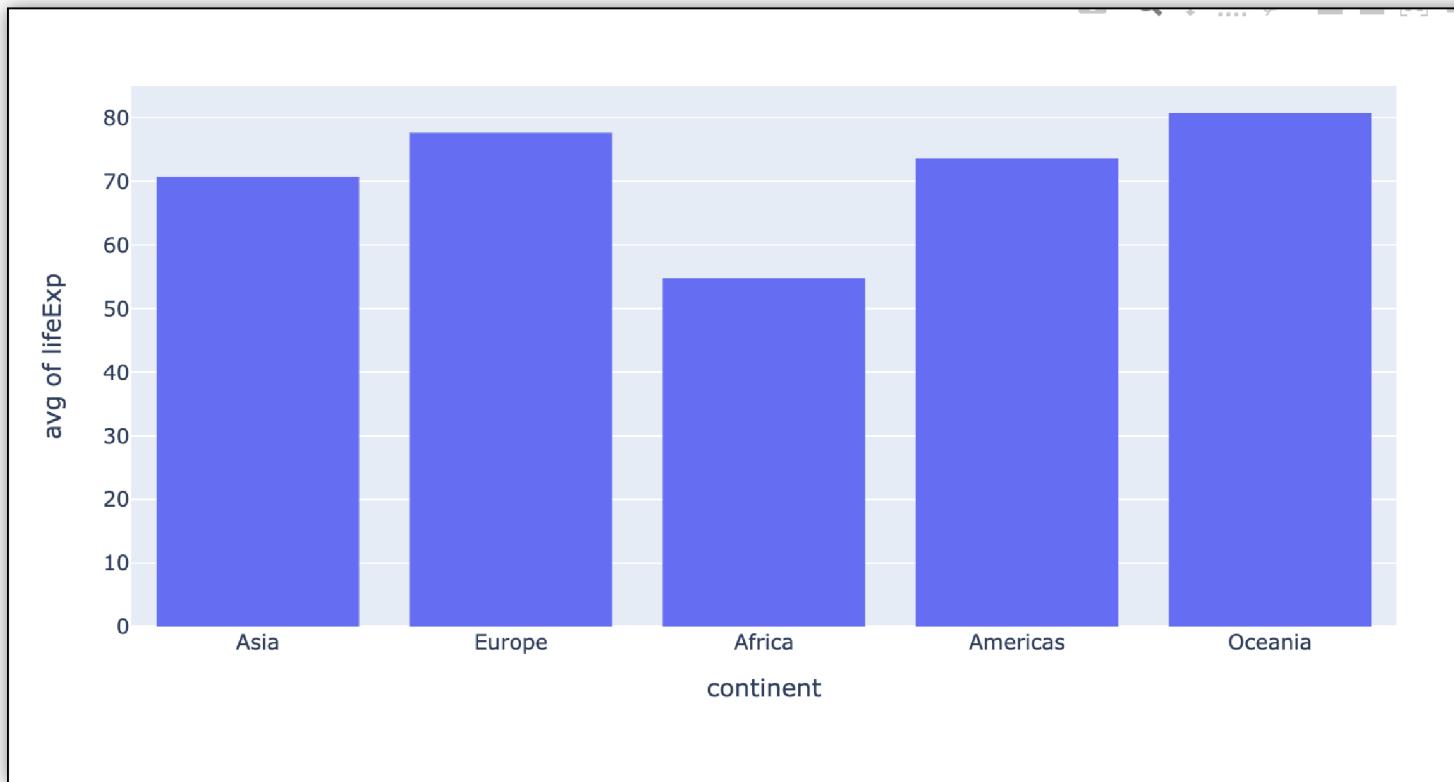
And a visualization



```
1 # 15_dash_with_data_and_vis.py
2
3 # Import packages
4 from dash import Dash, html, dash_table, dcc
5 import pandas as pd
6 import plotly.express as px
7
8 # Incorporate data
9 df = pd.read_csv('gapminder2007.csv')
10
11 # Initialize the app
12 app = Dash(__name__)
13
14 # App layout
15 app.layout = html.Div([
16     html.Div(children='My First App with Data and a Graph'),
17     dash_table.DataTable(data=df.to_dict('records'), page_size=10),
18     dcc.Graph(figure=px.histogram(df, x='continent', y='lifeExp', histfunc='avg'))
19 ])
20
21 # Run the app
22 if __name__ == '__main__':
23     app.run(host='0.0.0.0', port=5005)
```

Capstone Project

And a visualization



Capstone Project

And Controls

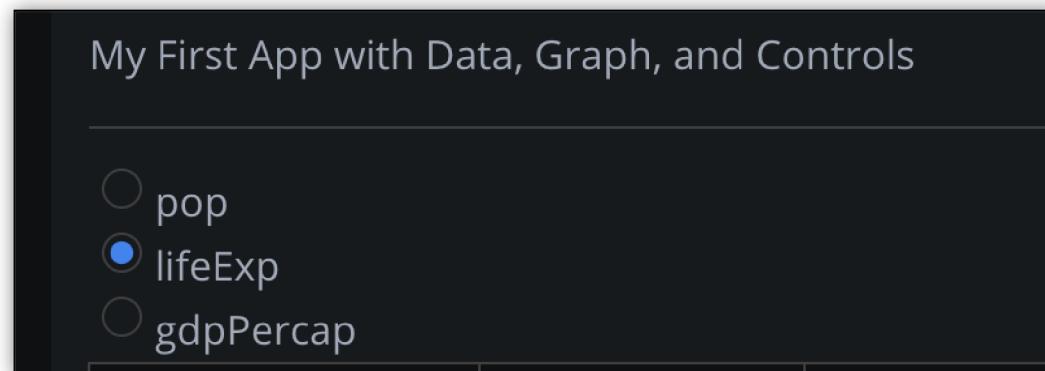
Capstone Project

And Controls

```
● ● ●  
1 # 16_dash_with_data_and_vis_cont.py  
2  
3 # Add controls to build the interaction  
4 @callback(  
5     Output(component_id='controls-and-graph', component_property='figure'),  
6     Input(component_id='controls-and-radio-item', component_property='value')  
7 )  
8 def update_graph(col_chosen):  
9     fig = px.histogram(df, x='continent', y=col_chosen, histfunc='avg')  
10    return fig
```

Capstone Project

As well as the effective radio buttons

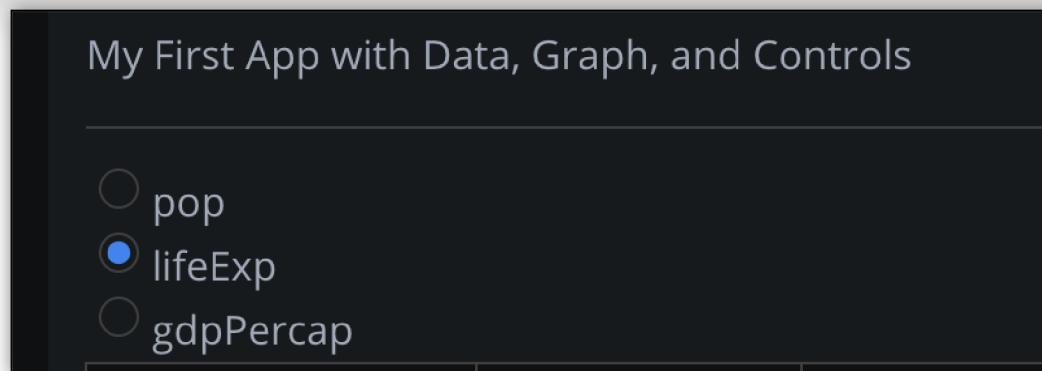


Capstone Project

As well as the effective radio buttons



```
1 # 16_dash_with_data_and_vis_cont.py
2
3 # App layout
4 app.layout = html.Div([
5     html.Div(children='My First App with Data, Graph, and Controls'),
6     html.Hr(),
7     dcc.RadioItems(options=['pop', 'lifeExp', 'gdpPercap'], value='lifeExp', id='controls-and-radio-item'),
8     dash_table.DataTable(data=df.to_dict('records'), page_size=6),
9     dcc.Graph(figure={}, id='controls-and-graph')
10 ])
```



Capstone Project

Let's take a look at the final - externally styled variant in the code.
06-capstone.