# Analysing data

- Currently 33 zettabytes
  - 2025: 173 zettabytes [IDC](https://www.emc.com/collateral/analyst-reports/idc-the-digital-universe-in-2020.pdf (https://www.emc.com/collateral/analyst-reports/idc-the-digital-universe-in-2020.pdf)
  - 173'000'000'000'000'000'000'000 (https://en.wikipedia.org/wiki/Zettabyte)

---

- Your future company will bathe in data about everything
- You need to learn how to use data and extract valuable information from it

---

# Analysing data with models

> All models are wrong, but some are useful

Source: Several (https://en.wikipedia.org/wiki/All_models_are_wrong)

---

# Today: Three relevant questions + geographical show-off

1. Are young people really getting poorer?
2. How green is Denmark really?
3. How many from the Ivory Coast lives in Denmark?

- Plotting with Python

---

# Pandas

Pandas (https://pandas.pydata.org/) is a library for Python that helps you do this.

Install it now by typing:

```
$ pip install pandas
```

---

# Pandas dataframes

Pandas works a bit like the `openpyxl` module: you get a 'sheet' of data and read it in columns.

Let's do that together now.

```python
import pandas as pd
import matplotlib.pyplot as plt
```

But we need a data source!

https://www.dr.dk/nyheder/indland/aeldre-bliver-rigere-unge-fattigere
(https://www.dr.dk/nyheder/indland/aeldre-bliver-rigere-unge-fattigere)

https://www.dst.dk/da/Statistik/emner/arbejde-indkomst-og-formue/indkomster
(https://www.dst.dk/da/Statistik/emner/arbejde-indkomst-og-formue/indkomster)

https://www.dst.dk/da/Statistik/nyt/NytHtml?cid=29483 (https://www.dst.dk/da/Statistik/nyt/NytHtml?
cid=29483)

```python

```

```python
df = pd.read_csv('data/2019721183544253670048AINDK267296307071.csv', header=N
```

```python
print(df)
```

```python
print(df.head())
```

```python
df.head()
```

```python
df.columns = ['','', 'region', 'sex', 'age', '2018', '2017', '2016', '2015',
```

```python
df.head()
```

```python
df['age']
```

```python
both = pd.DataFrame()
```

```
In [ ]:
1  both['age'] = df['age'][1:14]
```

```
In [ ]:
1  both
```

```
In [ ]:
1  # Crop columns to only contain 1 to 13
```

```
In [ ]:
1  for i in range(2012, 2019):
2      name = str(i)
3      both[name] = df[name][1:14]
```

```
In [ ]:
1  both
```

```
In [ ]:
1  both.plot()
```

```
In [ ]:
1  both.T[1:].plot()
```

# Plotting with pandas

Pandas `DataFrame` s have a simple `.plot()` method, which plots **columns as x values** and **rows as y values**.

https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.plot.line.html#pandas.DataFrame.plot.line (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.plot.line.html#pandas.DataFrame.plot.line)

You can also plot barcharts, histograms etc.

```
In [ ]:
1  both_index = pd.DataFrame()
```

```
In [ ]:
1  both_index['age'] = both['age']
```

```
1  for i in range(2012, 2019):
2      name = str(i)
3      both_index[name] = df[name][1:14]
```

```
1  for i in range(2012, 2019):
2      name = str(i)
3      both_index[name] = df[name][1:14] - df['2012'][1:14]
```

```
1  both_index
```

```
1  both_index.T[1:].plot()
```

# Exercise!

Use the dataset available from GitHub: `session-10/data`

- Can you do the same thing for men and women exclusively?
    - Try to turn the problem into a function. The code will almost be the same, but what will change?
- Bonus question: Normalise to percentages
    - The numbers in the graph are pretty big, can you divide by the maximum number in all the columns?
        - This should give you values that are at maximum 1

# How green is Denmark?!

http://databank.worldbank.org/data/reports.aspx?source=&series=EG.ELC.RNEW.ZS
(http://databank.worldbank.org/data/reports.aspx?source=&series=EG.ELC.RNEW.ZS)

https://databank.worldbank.org/source/world-development-indicators
(https://databank.worldbank.org/source/world-development-indicators)

```
1  pd.read_csv('data/abcc94ae-70a2-4ca3-bfc8-5d6ddc6803d8_Data.csv')
```

```
1  wb = pd.read_csv('data/abcc94ae-70a2-4ca3-bfc8-5d6ddc6803d8_Data.csv')
```

wb

## Pandas `.loc()`

Slices the dataset using both row and column indices:

```
df.loc[rows, columns]
```

For example:

- `df.loc[:, :]`
- `df.loc[0, 1960:2018]`

Great resource: https://medium.com/dunder-data/selecting-subsets-of-data-in-pandas-6fcd0170be9c (https://medium.com/dunder-data/selecting-subsets-of-data-in-pandas-6fcd0170be9c)

In [ ]:
```
1  wb.loc[:, '1960 [YR1960]':'2018 [YR2018]']
```

In [ ]:
```
1  wb_data = wb.loc[:10, '1960 [YR1960]':'2018 [YR2018]'].T
```

In [ ]:
```
1  wb_data
```

In [ ]:
```
1  wb_data = wb.loc[:10, '1990 [YR1990]':'2015 [YR2015]'].T
```

In [ ]:
```
1  wb_data.columns = wb['Country Name'][:11]
```

In [ ]:
```
1  wb_data.plot()
```

In [ ]:
```
1  wb_data.plot(figsize=(20, 20))
```

# Exercise

We only care about our northern bretheren of course. Can you filter our all the other countries? Perhaps using `.loc` ?

Use the dataset available from GitHub: `session-10/data`

```python
import pandas as pd
import matplotlib.pyplot as plt

wb = pd.read_csv('abcc94ae-70a2-4ca3-bfc8-5d6ddc6803d8_Data.csv')
wb_data = wb.loc[:10, '1990 [YR1990]':'2015 [YR2015]'].T
wb_data.columns = wb['Country Name'][:11]
wb_data.plot()
plt.show()
```

# How many from the Ivory Coast lives in Copenhagen?

In [ ]:

```python
1  cph = pd.read_csv('data/befkbhalderstatkode_small.csv')
```

In [ ]:

```python
1   KBH_NEIGHBORHOODS = {
2       '1': 'Indre By',
3       '2': 'Østerbro',
4       '3': 'Nørrebro',
5       '4': 'Vesterbro/Kgs. Enghave',
6       '5': 'Valby',
7       '6': 'Vanløse',
8       '7': 'Brønshøj-Husum',
9       '8': 'Bispebjerg',
10      '9': 'Amager Øst',
11      '10': 'Amager Vest',
12      '99': 'Udenfor inddeling'
13  }
```

In [ ]:

```python
1    COUNTRY_CODES = {
2        '0': 'Uoplyst (1)',
3        '5001': 'Uoplyst (2)',
4        '5100': 'Danmark',
5        '5101': 'Grønland',
6        '5102': 'Udlandet uoplyst',
7        '5103': 'Statsløs',
8        '5104': 'Finland',
9        '5105': 'Island, ligeret dansk',
10       '5106': 'Island',
11       '5107': 'Liechtenstein',
```

```
12          '5108': 'Luxembourg',
13          '5109': 'Monaco',
14          '5110': 'Norge',
15          '5114': 'Europa uoplyst',
16          '5115': 'Kongelig',
17          '5120': 'Sverige',
18          '5122': 'Albanien',
19          '5124': 'Andorra',
20          '5126': 'Belgien',
21          '5128': 'Bulgarien',
22          '5129': 'Tjekkoslovakiet',
23          '5130': 'Frankrig',
24          '5134': 'Grækenland',
25          '5140': 'Nederlandene',
26          '5142': 'Irland',
27          '5150': 'Italien',
28          '5151': 'Serbien og Montenegro',
29          '5152': 'Jugoslavien',
30          '5153': 'Malta',
31          '5154': 'Polen',
32          '5156': 'Portugal',
33          '5158': 'Rumænien',
34          '5159': 'San Marino',
35          '5160': 'Schweiz',
36          '5162': 'Sovjetunionen',
37          '5164': 'Spanien',
38          '5170': 'Storbritannien',
39          '5172': 'Tyrkiet',
40          '5174': 'Ungarn',
41          '5176': 'Vatikanstaten',
42          '5180': 'Tyskland',
43          '5182': 'Østrig',
44          '5199': 'Europa uoplyst',
45          '5202': 'Algeriet',
46          '5204': 'Angola',
47          '5207': 'Botswana',
48          '5213': 'Burundi',
49          '5214': 'Etiopien',
50          '5215': 'Comorerne',
51          '5216': 'Eritrea',
52          '5222': 'Gambia',
53          '5228': 'Ghana',
54          '5230': 'Ækvatorialguinea',
55          '5231': 'Guinea-Bissau',
56          '5232': 'Guinea',
57          '5233': 'Kap Verde',
58          '5234': 'Kenya',
59          '5235': 'Lesotho',
60          '5236': 'Liberia',
61          '5238': 'Libyen',
62          '5240': 'Mozambique',
63          '5242': 'Madagaskar',
64          '5243': 'Mali',
65          '5244': 'Marokko',
66          '5245': 'Mauritius',
67          '5246': 'Nigeria',
68          '5247': 'Namibia',
```

```
 69          '5248': 'Marshalløerne',
 70          '5255': 'Sierra Leone',
 71          '5258': 'Sudan',
 72          '5259': 'Swaziland',
 73          '5260': 'Sydsudan',
 74          '5262': 'Sydafrika',
 75          '5266': 'Tanzania',
 76          '5268': 'Tunesien',
 77          '5269': 'Uganda',
 78          '5272': 'Egypten',
 79          '5273': 'Tuvalu',
 80          '5274': 'Kiribati',
 81          '5275': 'Vanuatu',
 82          '5276': 'Centralafrikanske Republik',
 83          '5277': 'Cameroun',
 84          '5278': 'Congo, Demokratiske Republik',
 85          '5279': 'Congo, Republikken',
 86          '5281': 'Benin',
 87          '5282': 'Elfenbenskysten',
 88          '5283': 'Gabon',
 89          '5284': 'Mauretanien',
 90          '5285': 'Niger',
 91          '5287': 'Rwanda',
 92          '5288': 'Senegal',
 93          '5289': 'Somalia',
 94          '5292': 'Tchad',
 95          '5293': 'Togo',
 96          '5294': 'Burkina Faso',
 97          '5295': 'Zimbabwe',
 98          '5296': 'Zambia',
 99          '5297': 'Malawi',
100          '5298': 'Seychellerne',
101          '5299': 'Afrika uoplyst',
102          '5302': 'Argentina',
103          '5303': 'Bahamas',
104          '5304': 'Bolivia',
105          '5305': 'Barbados',
106          '5306': 'Brasilien',
107          '5308': 'Guyana',
108          '5309': 'Antigua og Barbuda',
109          '5310': 'Nauru',
110          '5311': 'Skt. Vincent og Grenadinerne',
111          '5314': 'Canada',
112          '5316': 'Chile',
113          '5318': 'Colombia',
114          '5319': 'Syd- og Mellemamerika uoplyst',
115          '5322': 'Costa Rica',
116          '5324': 'Cuba',
117          '5326': 'Dominikanske Republik',
118          '5328': 'Ecuador',
119          '5338': 'Guatemala',
120          '5339': 'Grenada',
121          '5342': 'Haiti',
122          '5344': 'Surinam',
123          '5345': 'Dominica',
124          '5347': 'Skt. Lucia',
125          '5348': 'Honduras',
126          '5352': 'Jamaica',
```

```
127    '5354': 'Mexico',
128    '5356': 'Nicaragua',
129    '5358': 'Panama',
130    '5364': 'Paraguay',
131    '5366': 'Peru',
132    '5372': 'El Salvador',
133    '5374': 'Trinidad og Tobago',
134    '5376': 'Uruguay',
135    '5390': 'USA',
136    '5392': 'Venezuela',
137    '5395': 'Vestindiske Øer',
138    '5397': 'Nordamerika uoplyst',
139    '5398': 'Syd- og Mellemamerika uoplyst',
140    '5402': 'Yemen',
141    '5403': 'Forenede Arabiske Emirater',
142    '5404': 'Afghanistan',
143    '5406': 'Bahrain',
144    '5408': 'Bhutan',
145    '5410': 'Bangladesh',
146    '5412': 'Brunei',
147    '5414': 'Myanmar',
148    '5416': 'Cambodja',
149    '5418': 'Sri Lanka',
150    '5422': 'Cypern',
151    '5424': 'Taiwan',
152    '5432': 'Indien',
153    '5434': 'Indonesien',
154    '5435': 'Østtimor',
155    '5436': 'Irak',
156    '5438': 'Iran',
157    '5442': 'Israel',
158    '5444': 'Japan',
159    '5446': 'Jordan',
160    '5448': 'Kina',
161    '5452': 'Kuwait',
162    '5454': 'Laos',
163    '5456': 'Libanon',
164    '5457': 'Maldiverne',
165    '5458': 'Malaysia',
166    '5459': 'Mongoliet',
167    '5462': 'Oman',
168    '5464': 'Nepal',
169    '5466': 'Nordkorea',
170    '5468': 'Vietnam (1)',
171    '5471': 'Asien uoplyst',
172    '5472': 'Pakistan',
173    '5474': 'Filippinerne',
174    '5478': 'Saudi-Arabien',
175    '5482': 'Singapore',
176    '5484': 'Sydkorea',
177    '5486': 'Syrien',
178    '5487': 'Mellemøsten uoplyst',
179    '5488': 'Vietnam (2)',
180    '5492': 'Thailand',
181    '5496': 'Qatar',
182    '5499': 'Asien uoplyst',
183    '5502': 'Australien',
```

```
184        '5505': 'Tonga',
185        '5508': 'Fiji',
186        '5514': 'New Zealand',
187        '5522': 'Samoa',
188        '5525': 'Djibouti',
189        '5526': 'Belize',
190        '5534': 'Papua Ny Guinea',
191        '5599': 'Øer i Stillehavet',
192        '5607': 'Estland',
193        '5609': 'Letland',
194        '5611': 'Litauen',
195        '5621': 'Sao Tome og Principe',
196        '5623': 'Salomonøerne',
197        '5625': 'Skt. Kitts og Nevis',
198        '5700': 'Rusland',
199        '5704': 'Ukraine',
200        '5706': 'Hviderusland',
201        '5708': 'Armenien',
202        '5710': 'Aserbajdsjan',
203        '5712': 'Moldova',
204        '5714': 'Usbekistan',
205        '5716': 'Kasakhstan',
206        '5718': 'Turkmenistan',
207        '5720': 'Kirgisistan',
208        '5722': 'Tadsjikistan',
209        '5724': 'Georgien',
210        '5750': 'Kroatien',
211        '5752': 'Slovenien',
212        '5754': 'Bosnien-Hercegovina',
213        '5756': 'Makedonien',
214        '5757': 'Serbien',
215        '5758': 'Jugoslavien, Forbundsrepublikken',
216        '5759': 'Montenegro',
217        '5761': 'Kosovo',
218        '5776': 'Tjekkiet',
219        '5778': 'Slovakiet',
220        '5779': 'Cookøerne',
221        '5800': 'Land ukendt (2)',
222        '5901': 'Færøerne uoplyst',
223        '5902': 'Færøerne',
224        '5999': 'Land ukendt (1)'
225  }
```

In [ ]:

```
1  cph.groupby('STATKODE').sum().loc[5100]
```

# Exercise

How many from the Ivory Coast lives in Copenhagen?

1. Open the `data/befkbhalderstatkode_small.csv` dataset using `pandas`
2. Group by the `'STATKODE'` column
3. Sum all the values
4. Find the row corresponding to the Ivory Coast (`5282`)

# Plotting maps with folium (showcase)

https://python-visualization.github.io/folium/ (https://python-visualization.github.io/folium/)

In [ ]:

```
1  cph_lat, cph_lon = 55.6867243, 12.5700724
```

In [ ]:

```
1  import requests
2
3  # Copenhagen map data: http://wfs-kbhkort.kk.dk/web/
4  url = 'http://wfs-kbhkort.kk.dk/k101/ows?service=WFS&version=1.0.0&request=Ge
5  geo_json = requests.get(url).json()
```

In [ ]:

```
1  import folium
2
3  # Create a map on a specific location (tuple)
4  map_osm = folium.Map(location=(cph_lat, cph_lon), zoom_start=10)
5  # Using geospatial data formatted in JSON, add the points from the dataset to
6  folium.GeoJson(geo_json, name='geojson').add_to(map_osm)
7  # Show the map
8  map_osm
```

# Data sources

All ripe for the taking!

- World Bank (https://www.worldbank.org/)
- WTO (https://data.wto.org/)
- WHO (https://www.who.int/hiv/data/en/)
- Twitter (http://www.tweepy.org/)
- Kaggle (https://www.kaggle.com/datasets)
- Københavns datasæt (https://data.kk.dk/dataset)

- Your future company will bathe in data about everything
- You need to learn how to use data and extract valuable information from it

# Analysing data with models

> All models are wrong, but some are useful

Source: [Several (https://en.wikipedia.org/wiki/All_models_are_wrong)](https://en.wikipedia.org/wiki/All_models_are_wrong)

# Today: Three relevant questions + geographical show-off

1. Are young people really getting poorer?
2. How green is Denmark really?
3. How many from the Ivory Coast lives in Denmark?

- Plotting with Python

# Pandas

[Pandas (https://pandas.pydata.org/)](https://pandas.pydata.org/) is a library for Python that helps you do this.

Install it now by typing:

```
$ pip install pandas
```

# Pandas dataframes

Pandas works a bit like the `openpyxl` module: you get a 'sheet' of data and read it in columns.

Let's do that together now.

```
In [ ]:
```

But we need a data source!

[https://www.dr.dk/nyheder/indland/aeldre-bliver-rigere-unge-fattigere (https://www.dr.dk/nyheder/indland/aeldre-bliver-rigere-unge-fattigere)](https://www.dr.dk/nyheder/indland/aeldre-bliver-rigere-unge-fattigere)

https://www.dst.dk/da/Statistik/emner/arbejde-indkomst-og-formue/indkomster (https://www.dst.dk/da/Statistik/emner/arbejde-indkomst-og-formue/indkomster)

https://www.dst.dk/da/Statistik/nyt/NytHtml?cid=29483 (https://www.dst.dk/da/Statistik/nyt/NytHtml?cid=29483)

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

# Plotting with pandas

Pandas `DataFrame`s have a simple `.plot()` method, which plots **columns as x values** and **rows as y values**.

https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.plot.line.html#pandas.DataFrame.plot.line (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.plot.line.html#pandas.DataFrame.plot.line)

You can also plot barcharts, histograms etc.

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

# Exercise!

Use the dataset available from GitHub: `session-10/data`

- Can you do the same thing for men and women exclusively?
  - Try to turn the problem into a function. The code will almost be the same, but what will change?
- Bonus question: Normalise to percentages
  - The numbers in the graph are pretty big, can you divide by the maximum number in all the columns?
    - This should give you values that are at maximum 1

# How green is Denmark?!

http://databank.worldbank.org/data/reports.aspx?source=&series=EG.ELC.RNEW.ZS (http://databank.worldbank.org/data/reports.aspx?source=&series=EG.ELC.RNEW.ZS)

https://databank.worldbank.org/source/world-development-indicators (https://databank.worldbank.org/source/world-development-indicators)

In [ ]:

In [ ]:

wb

## Pandas `.loc()`

Slices the dataset using both row and column indices:

 `df.loc[rows, columns]`

For example:

- `df.loc[:, :]`
- `df.loc[0, 1960:2018]`

Great resource: https://medium.com/dunder-data/selecting-subsets-of-data-in-pandas-6fcd0170be9c (https://medium.com/dunder-data/selecting-subsets-of-data-in-pandas-6fcd0170be9c)

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

```
In [ ]:
```

## Exercise

We only care about our northern bretheren of course. Can you filter our all the other countries? Perhaps using `.loc` ?

Use the dataset available from GitHub: `session-10/data`

```python
import pandas as pd
import matplotlib.pyplot as plt

wb = pd.read_csv('abcc94ae-70a2-4ca3-bfc8-5d6ddc6803d8_Data.csv')
wb_data = wb.loc[:10, '1990 [YR1990]':'2015 [YR2015]'].T
wb_data.columns = wb['Country Name'][:11]
wb_data.plot()
plt.show()
```

## How many from the Ivory Coast lives in Copenhagen?

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

## Exercise

How many from the Ivory Coast lives in Copenhagen?

1. Open the `data/befkbhalderstatkode_small.csv` dataset using `pandas`
2. Group by the `'STATKODE'` column
3. Sum all the values
4. Find the row corresponding to the Ivory Coast (`5282`)

## Plotting maps with folium (showcase)

https://python-visualization.github.io/folium/ (https://python-visualization.github.io/folium/)

In [ ]:

In [ ]:

In [ ]:

# Data sources

All ripe for the taking!

- World Bank (https://www.worldbank.org/)
- WTO (https://data.wto.org/)
- WHO (https://www.who.int/hiv/data/en/)
- Twitter (http://www.tweepy.org/)
- Kaggle (https://www.kaggle.com/datasets)
- Københavns datasæt (https://data.kk.dk/dataset)