# Data Analysis & Visualisation

## Week 2

Raoul Grouls, 26-2-2024

# Recap leerdoelen les 1

- de gestalt principes & five guidelines toepassen op visualisaties

- een virtual environment activeren met pdm

- nieuwe features extraheren met behulp van regular expressions

- Een script vanaf de terminal opstarten

- click gebruiken voor command line arguments bij een script

- begrijpt de opzet van een project (src folder, data/raw en data/processed, pyproject.toml, notebooks) en kan dit zelf opzetten

- kan een eigen git-repo maken

- Regular expressions toepassen:
  - start ^
  - end $
  - or: [Bb]
  - ranges [a-zA-Z]
  - any char .
  - zero or more a*
  - one or more a+
  - not in range [^a-z]
  - shortcuts (\w, \s, \d)
  - lookbehind (?<=...)
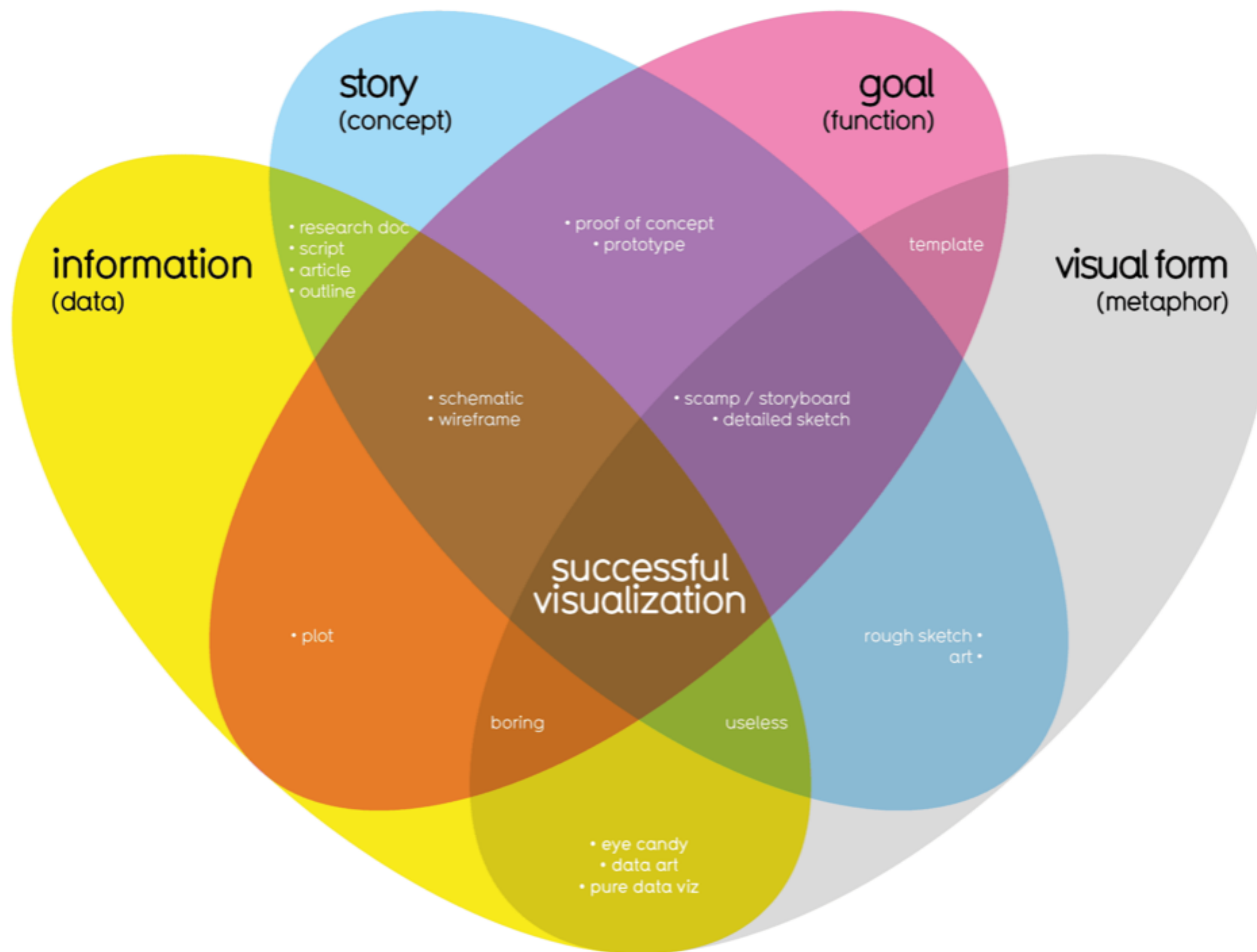  - lookahead (?=...)

# leerdoelen les 2

- Leren toepassen van visualisatie principes

- Omgaan met venv, pdm, path, scripts

- Oefenen met nieuwe features extraheren met behulp van regular expressions

- Vergelijken van categorieen met behulp van data visualisaties:
  - Barplots
  - Barbell plot
  - heatmaps

- Werken met palettes (en list comprehensions)

- Pandas
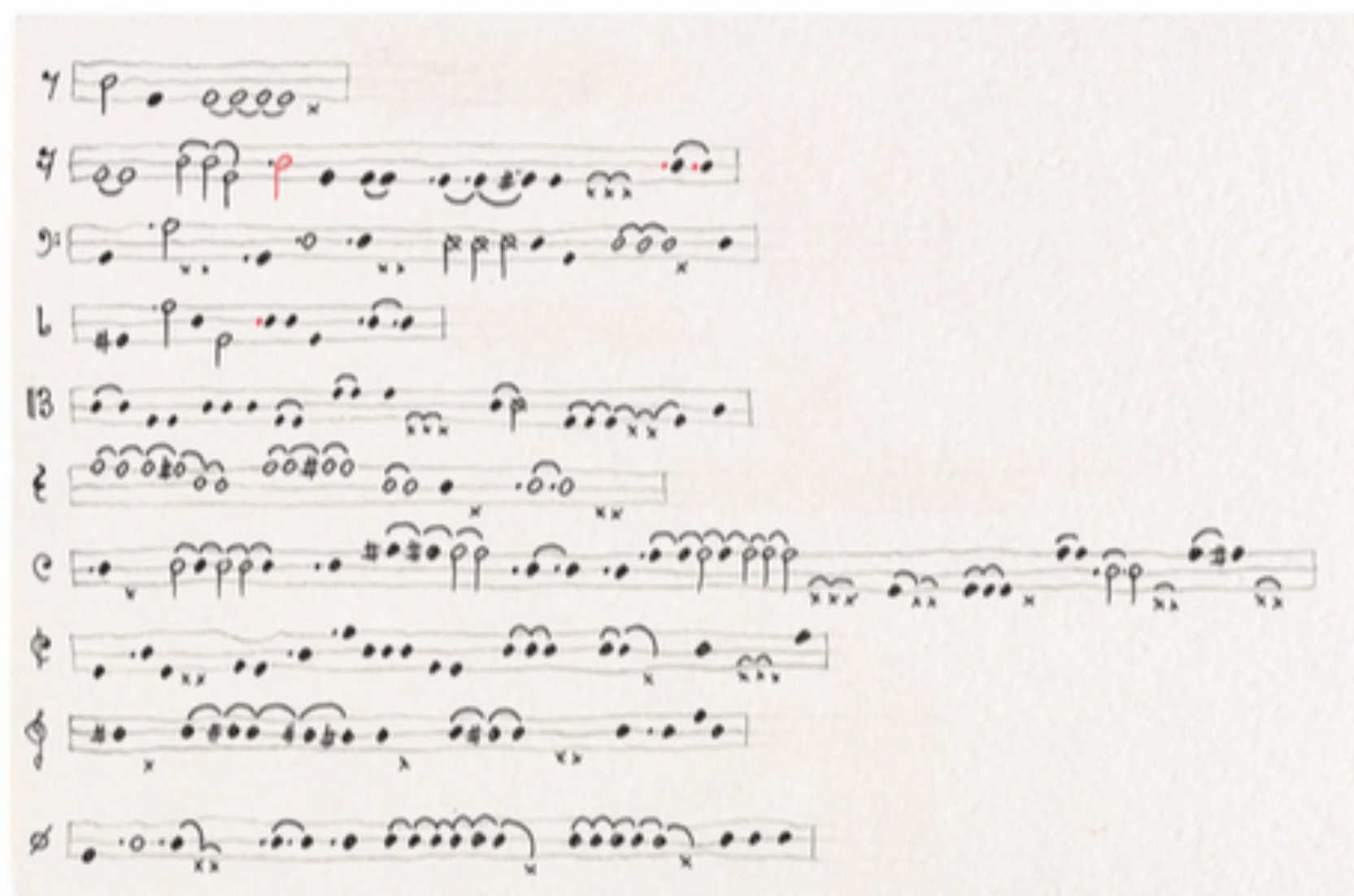  - Pandas groupby & aggregate
  - Pandas cut

# What Makes a Good Visualization?

rollover for more detail

explicit (implicit)

**story** (concept)

**goal** (function)

**information** (data)

**visual form** (metaphor)

• research doc
• script
• article
• outline

• proof of concept
• prototype

template

• schematic
• wireframe

• scamp / storyboard
• detailed sketch

**successful visualization**

• plot

rough sketch •
art •

boring

useless

• eye candy
• data art
• pure data viz

David McCandless
InformationisBeautiful.net

taken from new book
**Knowledge is Beautiful**

find out more
bit.ly/KIB_Books

What better visual reference than a musical score to show the repetitiveness of Giorgia's protests and the "level" of complaint: whether they are justified or totally out of place.



Note the hand-drawn stamps: these postcards were delivered in person in New York!

# What is Consciousness?

Make up your own mind



A field that exists in its own parallel "realm" of existence outside reality so can't be seen.
(Substance Dualism)

A sensation that "grows" inevitably out of complicated brain states.
(Emergent Dualism)

A physical property of all matter, like electromagnetism, just not one the scientists know about.
(Property Dualism)

All matter has a psychic part. Consciousness is just the psychic part of our brain.
(Pan Psychism)

Simply, mental states are physical events that we can see in brain scans.
(Identity Theory)

Consciousness and its states (belief, desire, pain) are simply functions the brain performs.
(Functionalism)

Literally just behaviour. When we behave in a certain way, we appear conscious.
(Behaviourism)

An accidental side-effect of complex physical processes in the brain.
(Epiphenomenalism)

# THE COUNTRIES THAT USE THE MOST TOILET PAPER

If you laid out all of the rolls of toilet paper used in each country in one year, what would that look like? For China it would be an incredible 4 billion miles long, which is further than the distance from Earth to Neptune. For the USA it's 2.65 billion miles, while both Japan and Germany use enough toilet paper to stretch out beyond Jupiter.

► TOILET PAPER USAGE PER YEAR PER COUNTRY

| Country | KM | Miles |
|---|---|---|
| CHINA | 4,287,630,419 KM | 4,007,208,240 MILES |
| UNITED STATES | 1,052,415,691 KM | 653,940,793 MILES / 2,651,782,602 MILES |
| JAPAN | 1,026,601,290 KM | 637,900,468 MILES |
| GERMANY | 788,352,099 KM | 489,859,284 MILES |
| UNITED KINGDOM | 738,584,457 KM | 458,935,105 MILES |
| BRAZIL | 423,771,299 KM | 263,319,277 MILES |
| FRANCE | 387,004,058 KM | 240,473,172 MILES |
| ITALY | 346,295,809 KM | 215,178,240 MILES |
| SPAIN | 328,163,799 KM | 203,911,531 MILES |
| SOUTH KOREA | 286,444,833 KM | 177,988,567 MILES |
| CANADA | 205,190,427 KM | 127,499,420 MILES |
| AUSTRALIA | 200,423,386 KM | 124,537,316 MILES |
| MEXICO | 152,446,413 KM | 94,725,810 MILES |
| TAIWAN | 137,809,685 KM | 85,630,958 MILES |
| VENEZUELA | 127,737,029 KM | 79,372,110 MILES |
| PORTUGAL | 127,170,615 KM | 79,020,157 MILES |
| BELGIUM | 115,715,644 KM | 71,902,368 MILES |
| ARGENTINA | 106,753,978 KM | 66,333,846 MILES |
| RUSSIA | 104,879,125 KM | 65,168,867 MILES |
| CHILE | 100,043,997 KM | 62,164,458 MILES |
| INDONESIA | 97,922,922 KM | 60,846,483 MILES |
| CZECH REPUBLIC | 93,054,631 KM | 57,821,467 MILES |
| COLOMBIA | 82,194,172 KM | 57,286,809 MILES |
| SOUTH AFRICA | 81,285,958 KM | 50,496,325 MILES |
| SWEDEN | | |

**Distance to the Sun**
93,000,000 MILES

**Furthest distance to Mars**
250,000,000 MILES

**Furthest distance to Jupiter**
601,000,000 MILES

**Furthest distance to Neptune**
2,700,000,000 MILES

# DEATH TOLL
[HIGHEST TO LOWEST]

the disease is new to medicine, and data is still coming in.

**200M**
**Black Death (Bubonic Plague)**
1347-1351

**56M**
**Small Pox**
1520

**40-50M**
**Spanish Flu**
1918-1919

**30-50M**
**Plague of Justinian**
541-542

The plague originated in rats and spread to humans via infected fleas.

The outbreak wiped out 30-50% of Europe's population. It took more than 200 years for the continent's population to recover.

**Smallpox** killed an estimated 90% of Native Americans. In Europe during the 1800s, an estimated 400,000 people were being killed by smallpox annually. The first ever vaccine was created to ward off smallpox.

The death toll of this plague is still under debate as new evidence is uncovered, but many think it may have helped hasten the fall of the Roman Empire.

A series of **Cholera** outbreaks spread around the world in the 1800s killing millions of people. There is no solid consensus on death tolls.

**25-35M**
**HIV/AIDS**
1981-PRESENT

**12M**
**The Third Plague**
1855

**5M**
**Antonine Plague**
165-180

**3M**
**17th Century Great Plagues**
1665

**1.1M**
**Asian Flu**
1957-1958

**1M**
**Russian Flu**
1889-1890

**1M**
**Hong Kong Flu**
1968-1970

**1M**
**Cholera 6 outbreak**
1817-1923

**1M**
**Japanese Smallpox Epidemic**
735-737

**600K**
**18th Century Great Plagues**
1817-1923

**200K**
**Swine Flu**
2009-2010

**100-150K**
**Yellow Fever**
LATE 1800s

**11.3K**
**Ebola**
2014-2016

**850**
**MERS**
2015-PRESENT

**770**
**SARS**
2002-2003

**4.7K***
**Novel Coronavirus (COVID-19)**
2019-PRESENT

# Causes of death in children under 5

(2013)

14% die from non-communicable diseases.

6% die from injuries.

80% die from preventable causes.

Source: IHME

**gates**notes

# Which Animals Could You Beat in a Fight?

"Which of the following animals, if any, do you think you could beat in a fight if you were unarmed?"

● Brits ● Americans

| Animal | Brits | Americans |
|---|---|---|
| Rat | ~68% | ~72% |
| House cat | ~68% | ~70% |
| Goose | ~45% | ~60% |
| Medium sized dog | ~40% | ~49% |
| Eagle | ~23% | ~30% |
| Large dog | ~15% | ~23% |
| Chimpanzee | ~12% | ~17% |
| King cobra | ~9% | ~13% |
| Kangaroo | ~9% | ~14% |
| Wolf | ~8% | ~12% |
| Crocodile | ~6% | ~9% |
| Gorilla | ~5% | ~8% |
| Lion | ~5% | ~8% |
| Elephant | ~5% | ~9% |
| Grizzly bear | ~5% | ~8% |

0%  20%  40%  60%  80%

# Trend – publicatie jaar vs Date Read



Jaar van uitgave versus leesdatum

# Color palettes

- **Qualitative**: representing categorical data

- **Sequential**: perceptually uniform

- **Diverging**: both large low and high values are interesting and span a midpoint value

# Introduction regular expressions

- Regular expressions (regex) are a powerful tool for matching patterns in text.

- Used for searching, editing, or manipulating text and data.

- Use https://regex101.com/ to develop and test your regexes

- chatGPT is pretty good at creating and explaining regexes.

# Basic symbols

- **^ (Start):** Matches the start of a line.

- **$ (End):** Matches the end of a line.

- **. (Any Char):** Matches any character except a newline.


- Example: To match any line that starts with "A", we use `^A`

# The OR operator and Ranges

- **[Bb] (Or):** Matches either "B" or "b".

- **[a-zA-Z] (Ranges):** Matches any letter, regardless of case.

- **[0-9]:** matches any number from 0 to 9

Example:

- To find any line that starts with a lowercase letter or number, use `^[a-z0-9]`

# Quantifiers

- **a\* (Zero or More):** Matches zero or more occurrences of "a".

- **a+ (One or More):** Matches one or more occurrences of "a".

- **a{3} (Exactly Three):** Matches exactly three occurrences of "a".

- **a{2,5} (Two to Five):** Matches between two and five occurrences of "a".


- **Example**: Combined with ranges: `[a-z]+`

- To match a string that contains four to six a's in a row, use `a{4,6}`

# Negation and Shortcuts

- **[^a-z] (Not in Range):** Matches any character not in the range "a" to "z".

- Shortcuts:

  - **\w (Word Char):** Matches any word character (letter, number, underscore).

  - **\s (Whitespace):** Matches any whitespace character (space, tab, newline).

  - **\d (Digit):** Matches any digit.

- **Question:** How to find lines <u>not</u> starting with any lowercase letter?

# Lookaround

- **Lookahead (?=...):** Matches a group after the main expression without including it in the result.

- **Lookbehind (?<=...):** Matches a group before the main expression without including it in the result.

- **Example:** To find words directly after a colon without including the colon, use `(?<=:)\s*\w+`

- Question: what does the \s* do?

# Combining everything

Message format:

- [11-01-2022 08:05:25] ~Raoul: Ah

- Can you explain what the other formats should be, as defined by the regexes?

```python
iosRegexes = BaseRegexes(
    timestamp=r"(?<=\[)\d{2}-\d{2}-\d{4} \d{2}:\d{2}:\d{2}(?=\])",
    author=r"(?<=\]\s)(.*?)(?=:)",
    clear=r"\[\d{2}-\d{2}-\d{4} \d{2}:\d{2}:\d{2}\]\s[~a-zA-Z\s]+:",
    fmt="%d-%m-%Y %H:%M:%S",
)


androidRegexes = BaseRegexes(
    timestamp=r"^\d{2}-\d{2}-\d{4} \d{2}:\d{2}",
    author=r"(?<=\s-\s)(.*?)(?=:)",
    clear=r"^\d{2}-\d{2}-\d{4} \d{2}:\d{2}[-~a-zA-Z\s]+:",
    fmt="%d-%m-%Y %H:%M",
)


oldRegexes = BaseRegexes(
    timestamp=r"^\d{1,2}/\d{1,2}/\d{2}, \d{2}:\d{2}",
    author=r"(?<=\s-\s)(.*?)(?=:)",
    clear=r"^\d{1,2}/\d{1,2}/\d{2}, \d{2}:\d{2}[-~a-zA-Z0-9\s]+:",
    fmt="%m/%d/%y, %H:%M",
)
```

# Improving with chatGPT

[…] this is the dataformat, […] this the regex, can you improve it?

`\[([0-2]?\d|3[01])-(0?\d|1[012])-(\d{4}) ([01]?\d|2[0-3]):([0-5]?\d):([0-5]?\d)\]`

This regex breaks down as follows:

- `([0-2]?\d|3[01])`: Matches days from `1-31`, allowing for a leading zero.

- `(0?\d|1[012])`: Matches months from `1-12`, also allowing for a leading zero.

- `(\d{4})`: Matches any four digits for the year.

- `([01]?\d|2[0-3])`: Matches hours from `0-23`, with optional leading zero.

- `([0-5]?\d)`: Matches minutes and seconds from `0-59`, with optional leading zero.