UNIVERSITY OF WARSAW

MIM

FACULTY
OF MATHEMATICS, INFORMATICS
AND MECHANICS

UNIWERSYTET WARSZAWSKI

# Data analysis and visualization (DAV)

*Lecture 05*

Łukasz P. Kozłowski

Warsaw, 2025

l.kozlowski@mimuw.edu.pl

# Data analysis and visualization (DAV)

*Lecture 05*
***Plots part 2***

Łukasz P. Kozłowski

Warsaw, 2025

l.kozlowski@mimuw.edu.pl

**Choose The Right Tool for the Job**
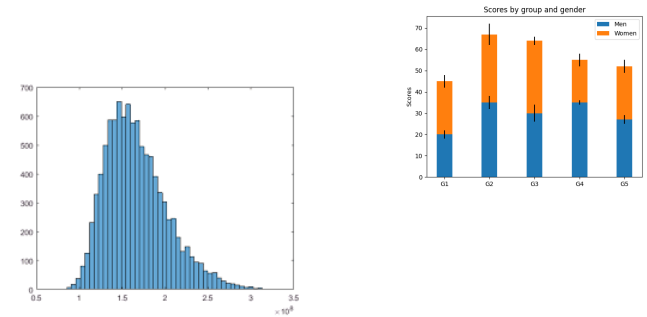


# Welcome in the Visualization ZOO

**Covered so far:**

**- bar charts (normal & stacked)**

**Covered so far:**

**- bar charts (normal & stacked)**

**- histograms**

**Covered so far:**

**- bar charts (normal & stacked)**

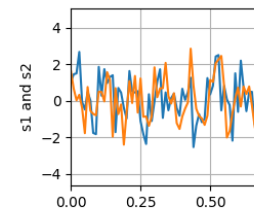**- histograms**

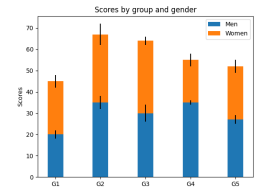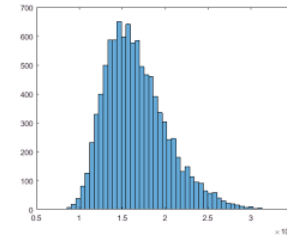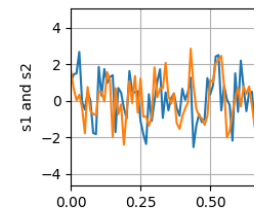**- line charts**

**Covered so far:**

**- bar charts (normal & stacked)**

**- histograms**

**- line charts**

**- scater plots**

Bubble chart show the relationship between **three** variables

With **three** variables you can either create 3d scatter plot or

With **three** variables you can either create 3d scatter plot or

embed 3rd dimension inside the bubble

you can say that the bubble plot is a scatter plot where 3rd dimension is added. The value of an additional variable is represented by the size of the dots (so called bubbles)

Caution: Hans Rosling frequently used 4D data in

(gapfinder) bubble plots

Aantal inwoners
- 30.000
- 30 mln.

A series of bubbles on a map is called a cartogram or sometimes "bubble map"

## Choosing bubble sizes correctly

The human visual system naturally experiences a **disk's size** in terms of its **area**. And the area of a disk — unlike its diameter or circumference — is not **proportional** to its radius, but **to** the **square of the radius**. So if one chooses to scale the disks' radii to the third data values directly, then the apparent size differences among the disks will be non-linear (quadratic) and misleading. To get a properly weighted scale, one must scale each disk's radius to the square root of the corresponding data value v3. This scaling issue can lead to extreme misinterpretations, especially where the range of the data has a large spread. And because many people are unfamiliar with—or do not stop to consider—the issue and its impact on perception, those who are aware of it o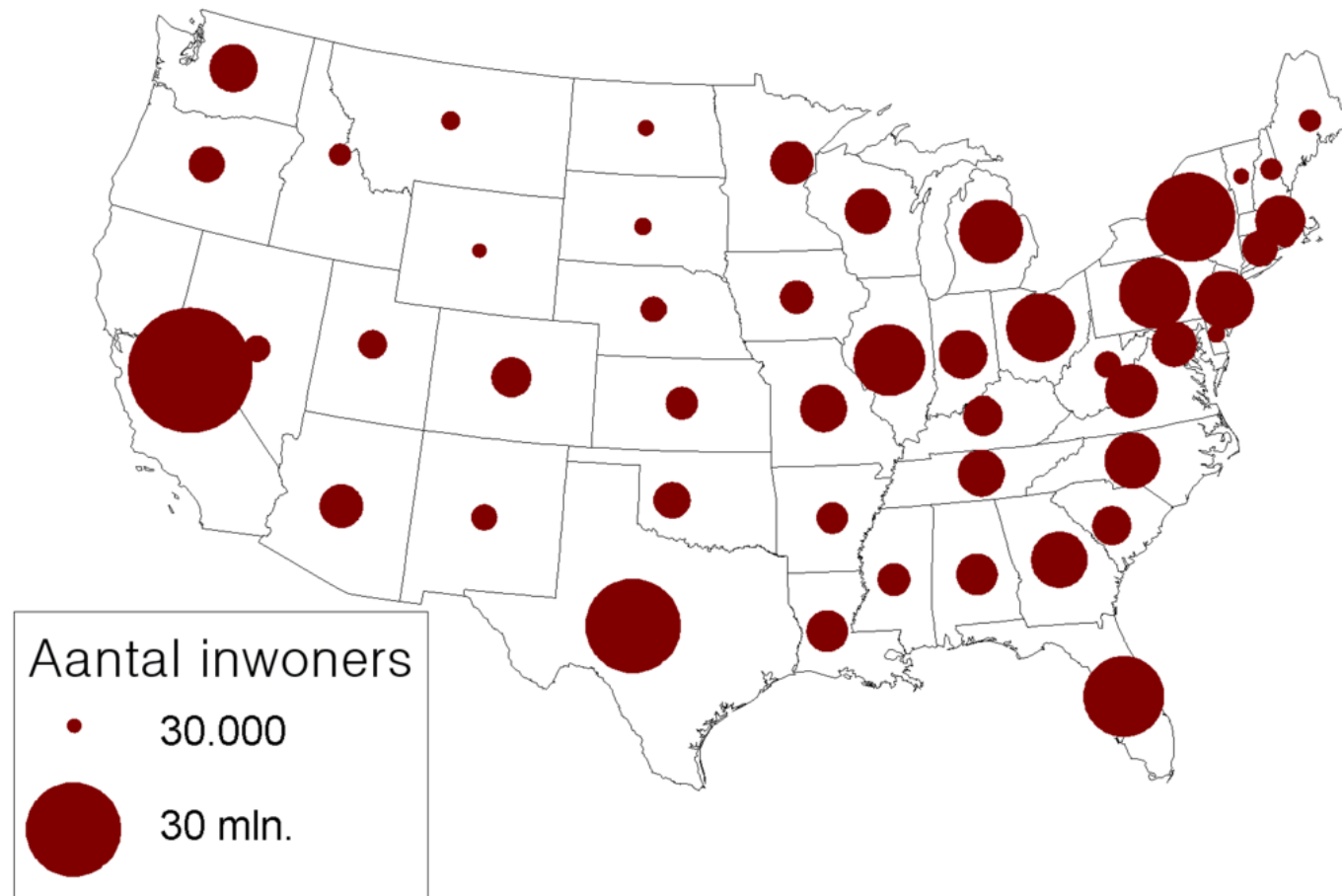ften have to hesitate in interpreting a bubble chart because they cannot assume that **the scaling correction was indeed made**. So it is important that bubble charts not only be scaled in this way, but also be clearly labeled to document that it is area, rather than radius or diameter, that conveys the data.
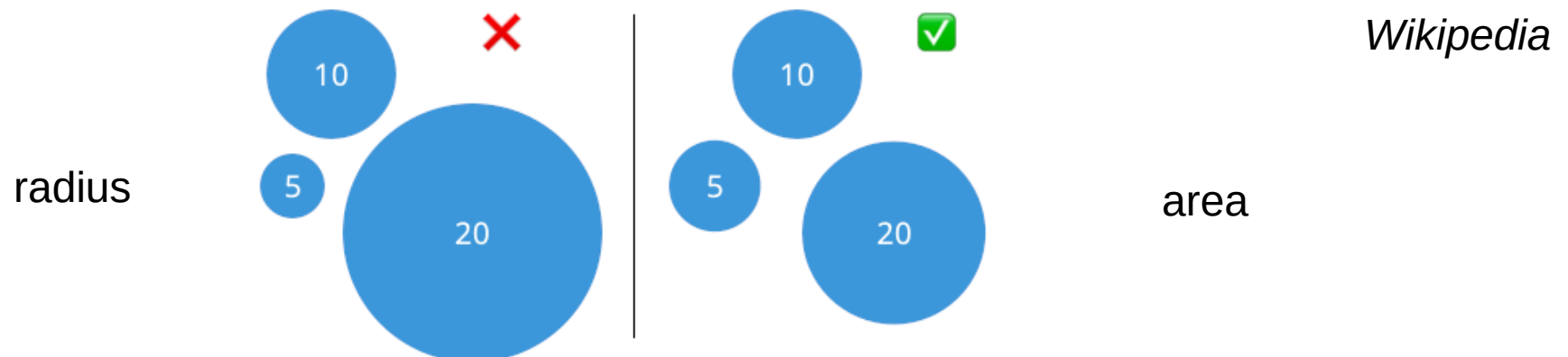
*Wikipedia*

## Choosing bubble sizes correctly

The human visual system naturally experiences a **disk's size** in terms of its **area**. And the area of a disk — unlike its diameter or circumference — is not **proportional** to its radius, but **to** the **square of the radius**. So if one chooses to scale the disks' radii to the third data values directly, then the apparent size differences among the disks will be non-linear (quadratic) and misleading. To get a properly weighted scale, one must scale each disk's radius to the square root of the corresponding data value v3. This scaling issue can lead to extreme misinterpretations, especially where the range of the data has a large spread. And because many people are unfamiliar with—or do not stop to consider—the issue and its impact on perception, those who are aware of it often have to hesitate in interpreting a bubble chart because they cannot assume that **the scaling correction was indeed made**. So it is important that bubble charts not only be scaled in this way, but also be clearly labeled to document that it is area, rather than radius or diameter, that conveys the data.
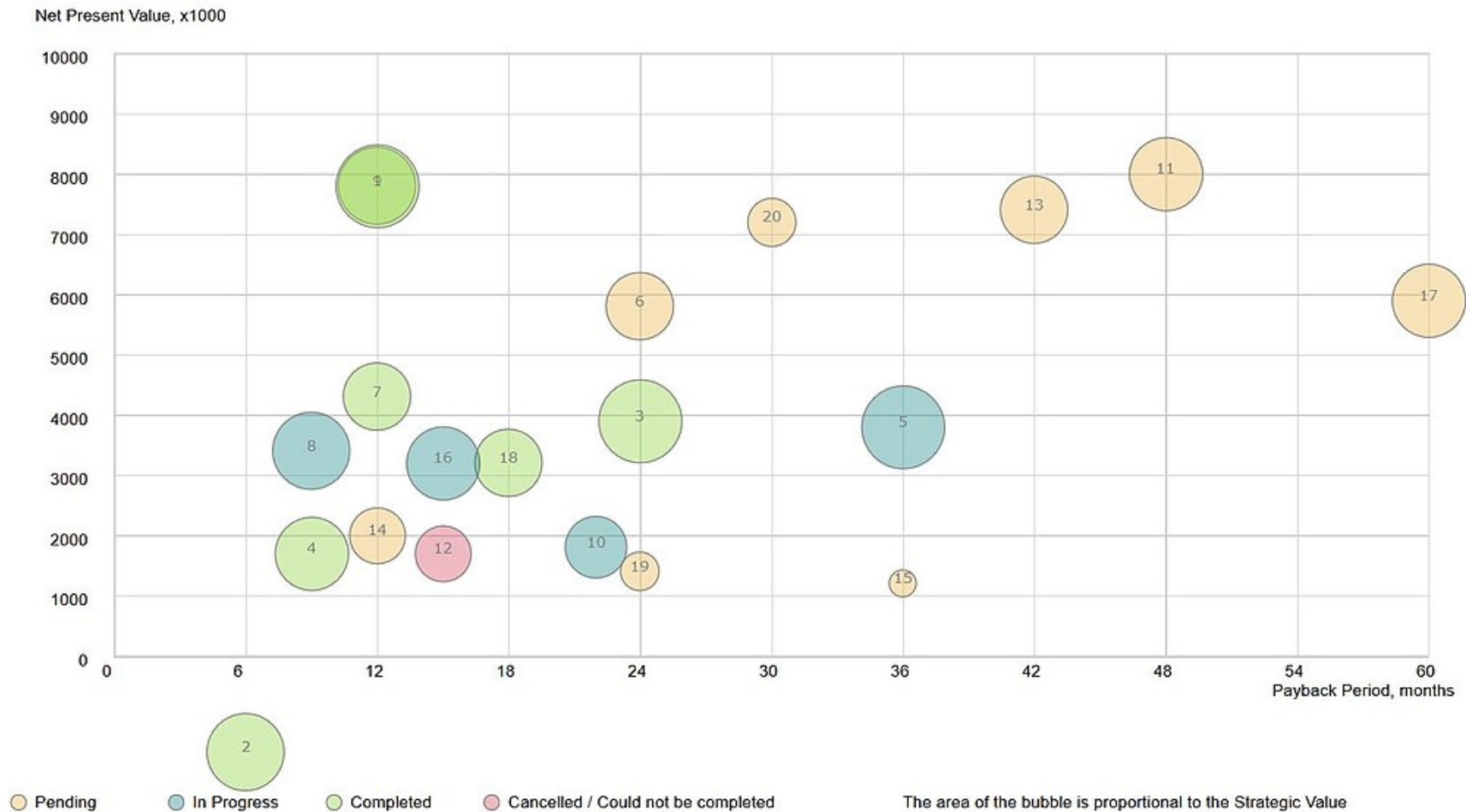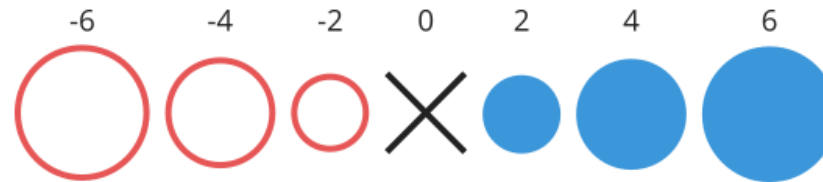


radius ✕    ✅ area    *Wikipedia*

## Displaying zero or negative data values in bubble charts

The metaphoric representation of data values as disk areas cannot be extended for displaying values that are negative or zero. As a fallback, some users of bubble charts resort to graphic symbology to express nonpositive data values. As an example, a negative value v < 0 can be represented by a disk of area v in which is centered some chosen symbol like "×" to indicate that the size of the bubble represents the absolute value of a negative data value. And this approach can be reasonably effective in situations where data values' magnitudes (absolute values) are themselves somewhat important—in other words, where values of v and −v  are similar in some context-specific way—so that their being represented by congruent disks makes sense.

To represent zero-valued data, some users dispense with disks altogether, using, say, a square centered at the appropriate location. Others use full circles for positive, and empty circles for negative values.

*Wikipedia*

Displaying zero or negative data values in bubble charts

Net Present Value, x1000

Payback Period, months

○ Pending    ○ In Progress    ○ Completed    ○ Cancelled / Could not be completed    The area of the bubble is proportional to the Strategic Value

## Use transparency

## Use area

## Use color

A series of bubbles on a map is called a cartogram or sometimes "bubble map"

Packed circle charts (aka circular packing, bubble cloud)

**Use Bubble charts to:**

- Present relationships between three numerical variables (fourth can be encoded in color or/and animation)

- Plot three sets of variables on one x-y coordinate plane
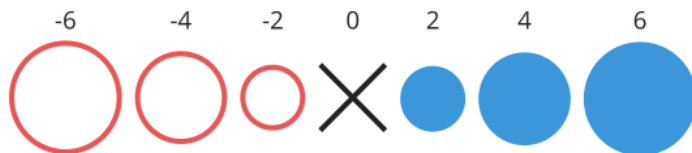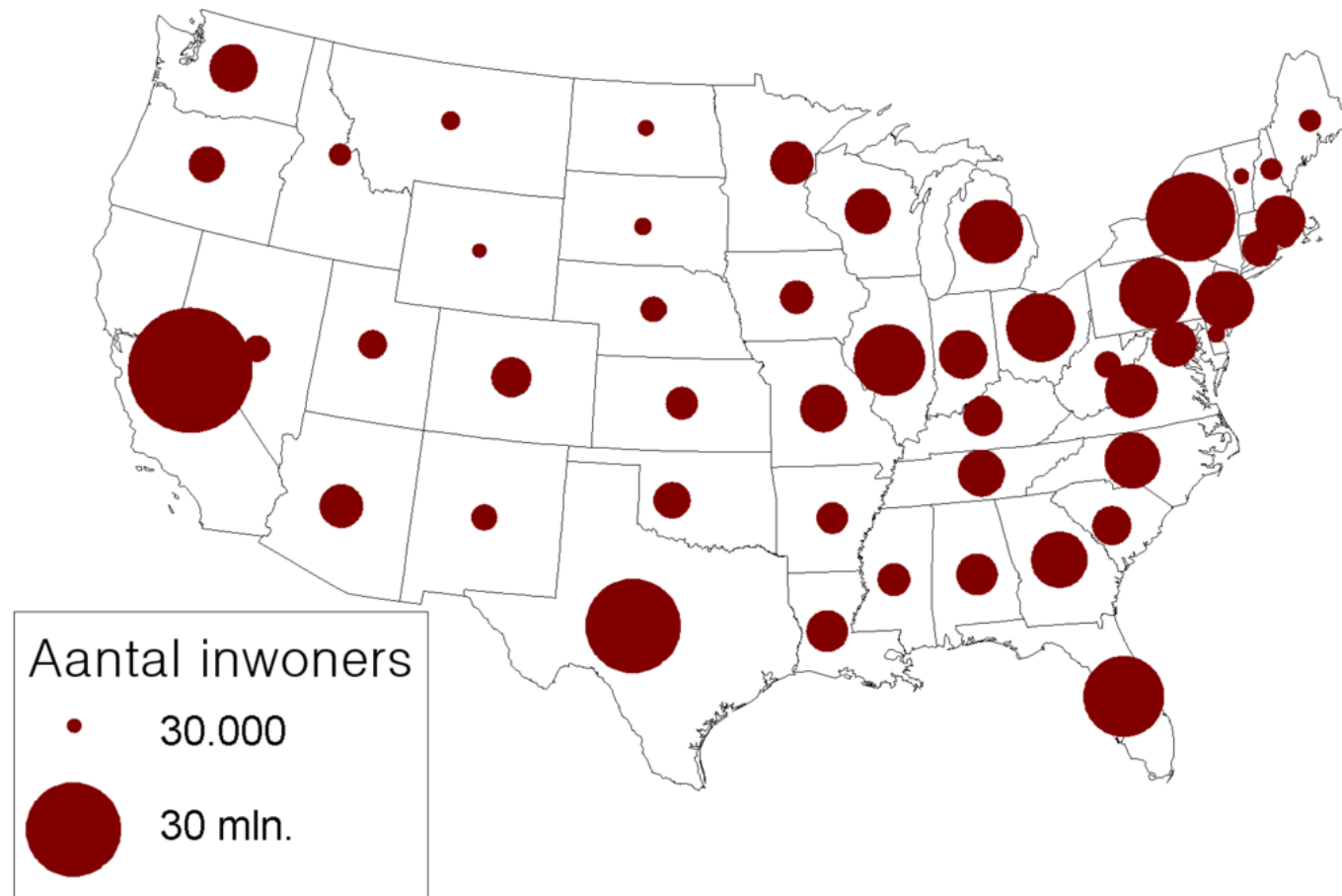
- Turn the horizontal axis into a logarithmic scale when showing the relationships between more widely distributed elements

- Present patterns in large sets of data, linear or non-linear trends, correlations, clusters, or outliers

- Present relationships, but not exact values for comparisons (this may be good or bad)

**More information:**

https://en.wikipedia.org/wiki/Bubble_chart

## More information:

https://en.wikipedia.org/wiki/Bubble_chart

**When queering for new information always start from Wikipedia**

## More information:

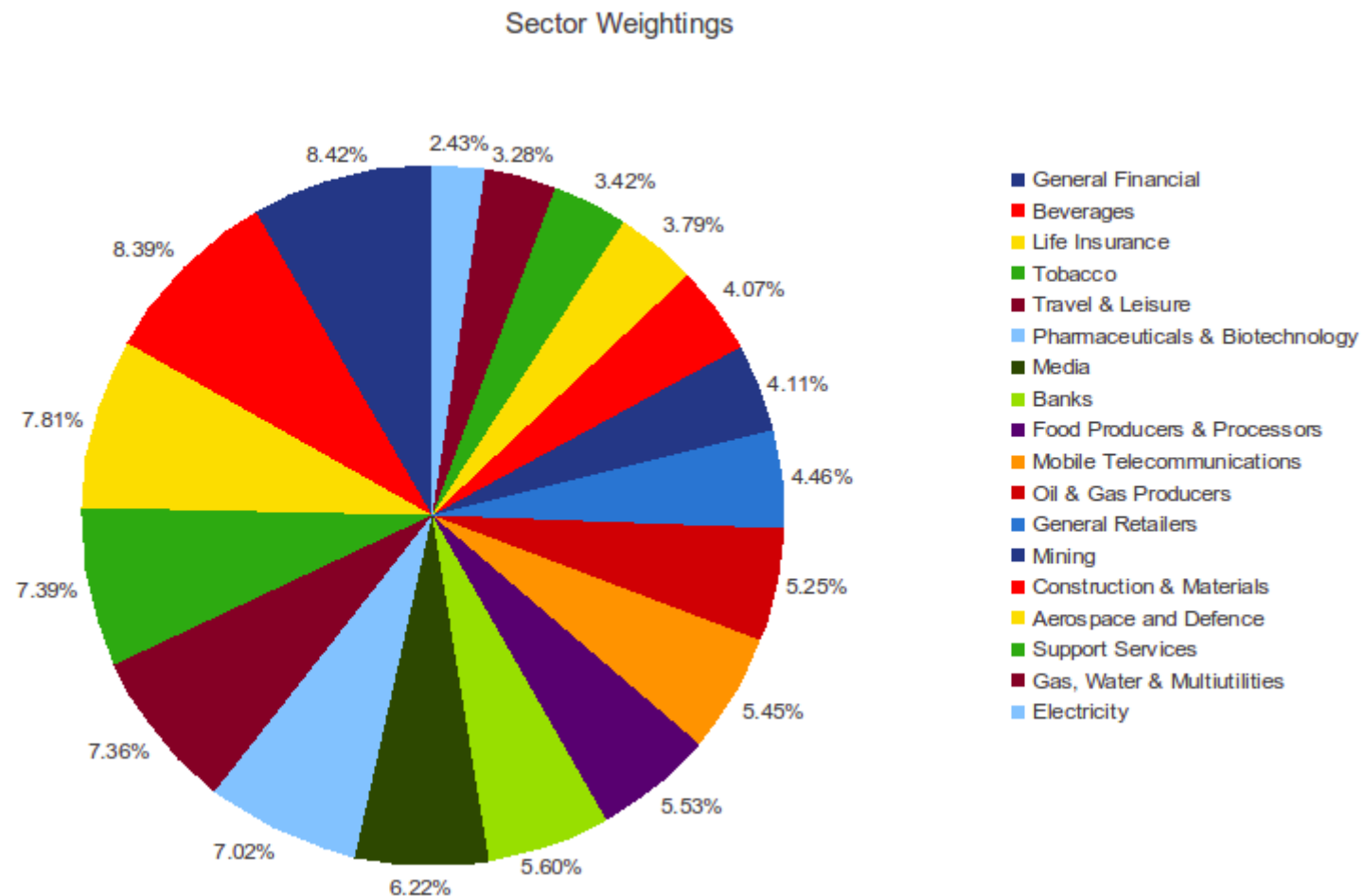https://en.wikipedia.org/wiki/Bubble_chart

https://www.youtube.com/watch?v=Fygi3mQY0Ag

https://www.youtube.com/watch?v=Uu2mfzWajQY

https://chartio.com/learn/charts/bubble-chart-complete-guide/

**When queering for new information always start from Wikipedia**

**One of the most frequently used and … also misused charts**

**General rule:**
**if you can, do not use it**
**(if you think that you need pie chart, think once more)**

Sector Weightings



Legend:
- General Financial
- Beverages
- Life Insurance
- Tobacco
- Travel & Leisure
- Pharmaceuticals & Biotechnology
- Media
- Banks
- Food Producers & Processors
- Mobile Telecommunications
- Oil & Gas Producers
- General Retailers
- Mining
- Construction & Materials
- Aerospace and Defence
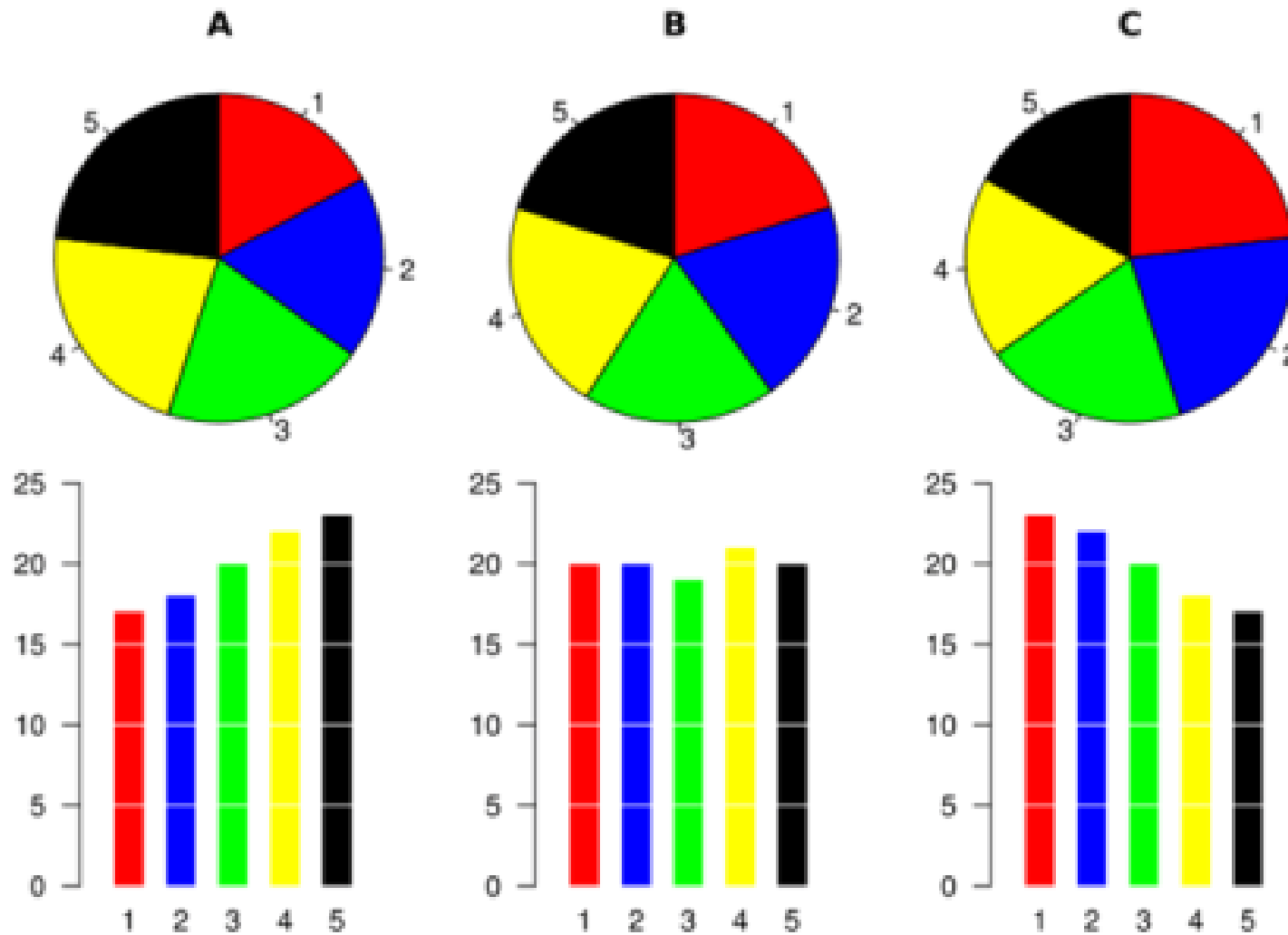- Support Services
- Gas, Water & Multiutilities
- Electricity

Represents numbers in percentages, used to visualize a part to whole relationship or a composition

Not meant to compare individual sections to each other or to represent exact values (use bar plot instead)

When possible, avoid pie charts and donuts. The human mind thinks linearly but, when it comes to angles and areas, most of us can't judge them well.

*Source: Oracle.com*

Human mind thinks linearly but, not by angles or areas

# Pie Charts and Donut Charts

Donuts charts are like pie charts, but
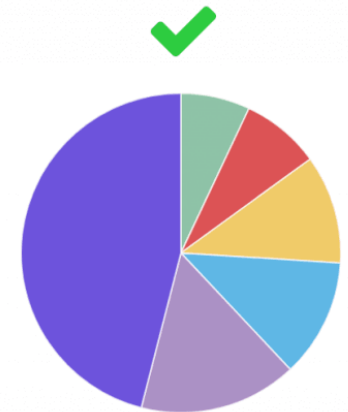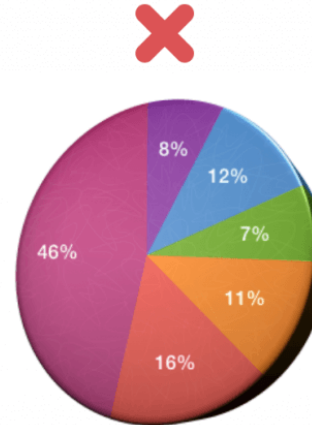


Countries by Area

- Russia — 27.3%
- Canada — 16%
- China — 15.4%
- United States — 15.4%
- Brazil — 13.6%
- Australia — 12.3%

Salary BreakDown

Total Salary: 24900

- Raphaely — 44....
- Khoo — 12.45%
- Baida — 11....
- Tobias — 11....
- Himuro — 10....
- Colmenares — 10.04%

Donuts charts are like pie charts, but worse



**Countries by Area**

| | |
|---|---|
| 27.3% | Russia |
| 16% | Canada |
| 15.4% | China |
| 15.4% | United States |
| 13.6% | Brazil |
| 12.3% | Australia |

**Salary BreakDown**

Total Salary: 24900

44....
12.45%
11....
11....
10....
10.04%

Raphaely  Khoo  Baida  Tobias  Himuro  Colmenares

Donuts charts are like pie charts, but worse



**Countries by Area**

- Russia — 27.3%
- Canada — 16%
- China — 15.4%
- United States — 15.4%
- Brazil — 13.6%
- Australia — 12.3%



**Salary BreakDown**

Total Salary: 24900

- 44....
- 12.45%
- 11....
- 11....
- 10....
- 10.04%

Raphaely ■ Khoo ■ Baida ■ Tobias ■ Himuro ■ Colmenares

**Many variations**:
- 3D pie chart and perspective pie cake
- Doughnut chart
- Exploded pie chart
- Polar area diagram

- Ring chart,
- Sunburst chart
- Multilevel pie chart
- Spie chart
- Square chart / Waffle chart

Standard Pie (year 1801)



Donut



Exploded Pie



Multilevel Pie



Spie Pie

# Pie Charts and Donut Charts

| Atmospheric air content | |
|---|---|
| Gas | percentage |
| Nitrogen | 78.00% |
| Oxygen | 21.00% |
| Others | 1.00% |

Legend:
- ■ Nitrogen
- ■ Oxygen
- ■ Others

UNITS SOLD (2017)

PRODUCT
- ■ A
- ■ B
- ■ C
- ■ D
- ■ E

Exploded and 3d pie in one

Square/Waffle chart

Polar chart by Florence Nightingale, 1858

Standard Pie (year 1801)



Donut



Exploded Pie



Multilevel Pie



Spie Pie

**VISUALIZE NO MORE THAN 5 CATEGORIES PER CHART**

It is difficult to differentiate between small values; depicting too many slices decreases the impact of the visualization. If needed, you can group smaller values into an "other" or "miscellaneous" category, but make sure it does not hide interesting or significant information.

**DON'T USE MULTIPLE PIE CHARTS FOR COMPARISON**

Slice sizes are very difficult to compare side-by-side. Use a stacked bar chart instead.

**MAKE SURE ALL DATA ADDS UP TO 100%**

Verify that values total 100% and that pie slices are sized proportionate to their corresponding value.

# PIE CHART

## DESIGN BEST PRACTICES

**ORDER SLICES CORRECTLY**

There are two ways to order sections, both of which are meant to aid comprehension:

**OPTION 1**

Place the largest section at 12 o'clock, going clockwise. Place the second largest section at 12 o'clock, going counterclockwise. The remaining sections can be placed below, continuing counterclockwise.

**OPTION 2**

Start the largest section at 12 o'clock, going clockwise. Place remaining sections in descending order, going clockwise.

## Dos and Don'ts for Pie chart

- Make sure that the total sum of all segments equals 100 percent

- Use pie charts only if you have **less than six categories**, unless there's a clear winner you want to focus on

- Ideally, there should be **only two categories**, like men and women visiting your website, or only one category, like a market share of your company, compared to the whole market

- **Don't use** a pie chart if the category **values are almost identical or completely different**. You could add labels, but that's a patch, not an improvement

- **Don't use 3D or blow apart effects** — they reduce comprehension and show incorrect proportions.

**More:**

**https://infogram.com/blog/do-this-not-that-pie-charts/**

**https://speakerdeck.com/cherdarchuk/data-looks-better-naked-pie-chart-edition?slide=45**

**Gantt charts** were adapted by Karol Adamiecki in 1896. But the name comes from Henry Gantt who independently adapted this bar chart type much later, in the 1910s



WEEKS: 1 2 3 4 5 6 7 8 9 10 21 22 23

WBS 1 Summary Element 1 — 57% complete

WBS 1.1 Activity A — 75% complete

START-TO-START

WBS 1.2 Activity B — 67% complete

FINISH-TO-START

WBS 1.3 Activity C — 50% complete

FINISH-TO-FINISH

WBS 1.4 Activity D — 0% complete

WBS 2 Summary Element 2 — 0% complete

WBS 2.1 Activity E — 0% complete

WBS 2.2 Activity F — 0% complete

WBS 2.3 Activity G — 0% complete

TODAY

**Gantt charts** were adapted by Karol Adamiecki in 1896. But the name comes from Henry Gantt who independently adapted this bar chart type much later, in the 1910s



**bar chart that illustrates a project schedule**

Also can show the dependency relationships between activities and current schedule status

| ID | Task Name | Predecessors | Duration |
|----|-----------|--------------|----------|
| 1 | Start | | 0 days |
| 2 | a | 1 | 4 days |
| 3 | b | 1 | 5.33 days |
| 4 | c | 2 | 5.17 days |
| 5 | d | 2 | 6.33 days |
| 6 | e | 3,4 | 5.17 days |
| 7 | f | 5 | 4.5 days |
| 8 | g | 6 | 5.17 days |
| 9 | Finish | 7,8 | 0 days |

Can also be used for instance for in rental businesses, displaying a list of items for rent (cars, rooms, apartments) and their rental periods

**Minimum requirement: a start date and an end date of each event**

Frequently, Gantt chart are used in parallel to network diagrams (that add another dimension i.e. the relationship between tasks

- improve the project structure design

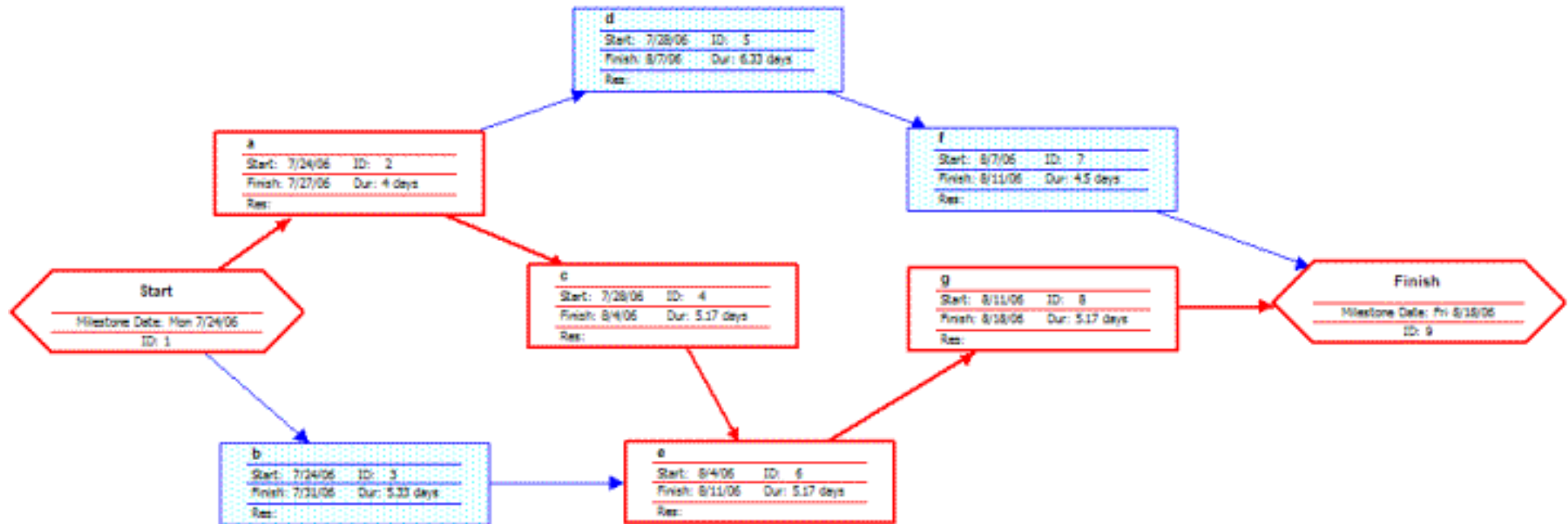- should be updated if needed

- improve the project structure design

- should be updated if needed

Unfortunetelly, many funding bodies / companies use it as definitive plan of the project and do not allow to any extensions or re-scheduling

but

if stage „N" relay on stage „N-1" and stage „N-1" is delayed because some unexpected, extra tasks, it is not possible to move to stage „N", no matter how pretty your Gantt plot is*

* this means that you need to think about some extra time between critical milestones of the project (the more risky the project, the more extra task will appear e.g. science projects in novel areas are extreme examples of that)

In python:

matplotlib
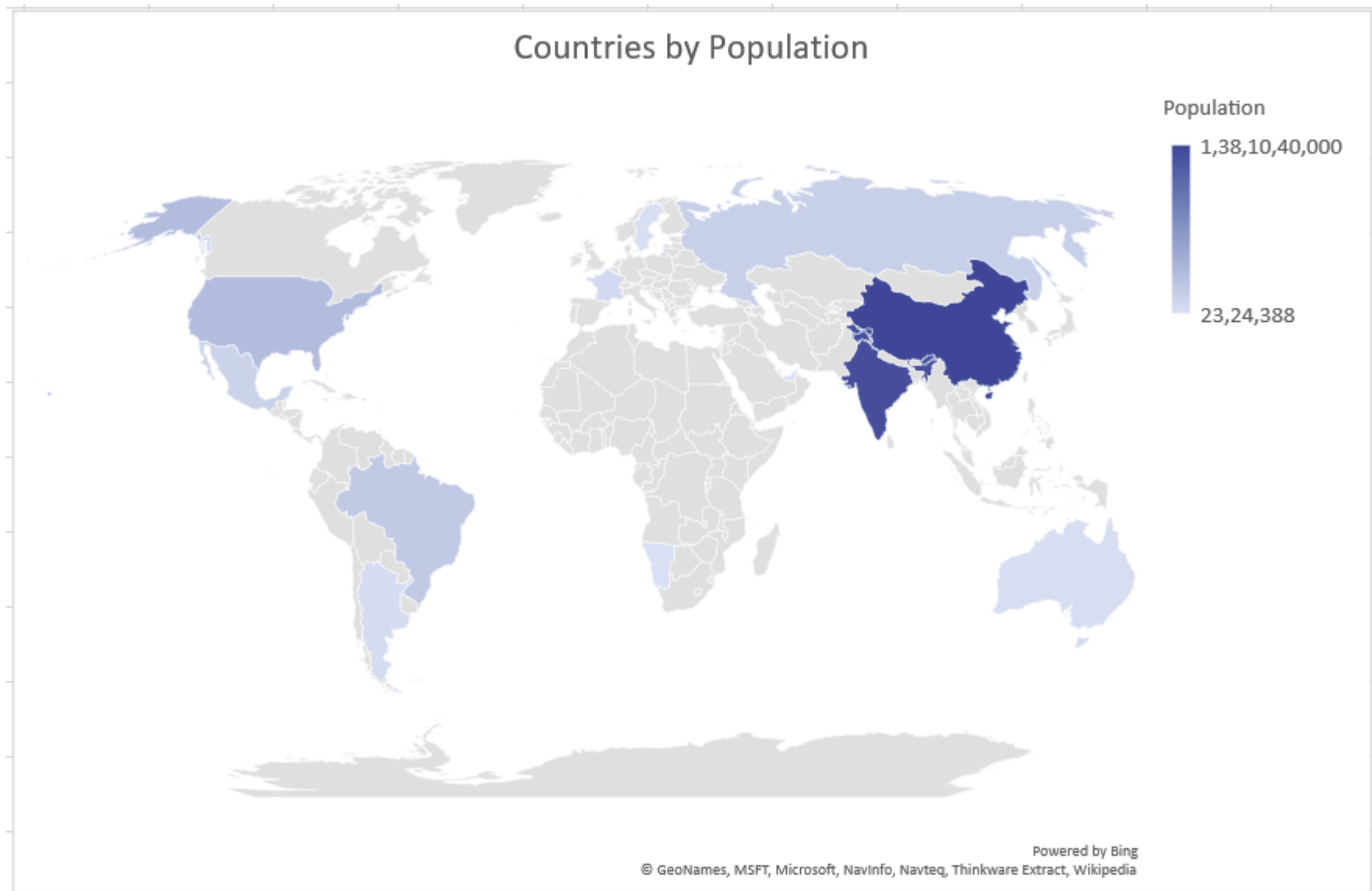https://www.geeksforgeeks.org/python-basic-gantt-chart-using-matplotlib/

Plotly
https://www.geeksforgeeks.org/python-basic-gantt-chart-using-matplotlib/

Additional packages (python-gantt)

http://www.clowersresearch.com/main/gantt-charts-in-matplotlib/

Frequently use color to connect county with the value (e.g. population, product sold, etc.)



Countries by Population

Population

1,38,10,40,000

23,24,388

Powered by Bing
© GeoNames, MSFT, Microsoft, NavInfo, Navteq, Thinkware Extract, Wikipedia

The area can be deformed by the value (area cartograms)



100%
Democrat

50/50

100%
Republican

## Linear cartograms



Time to Travel from High Barnet

0  5  10
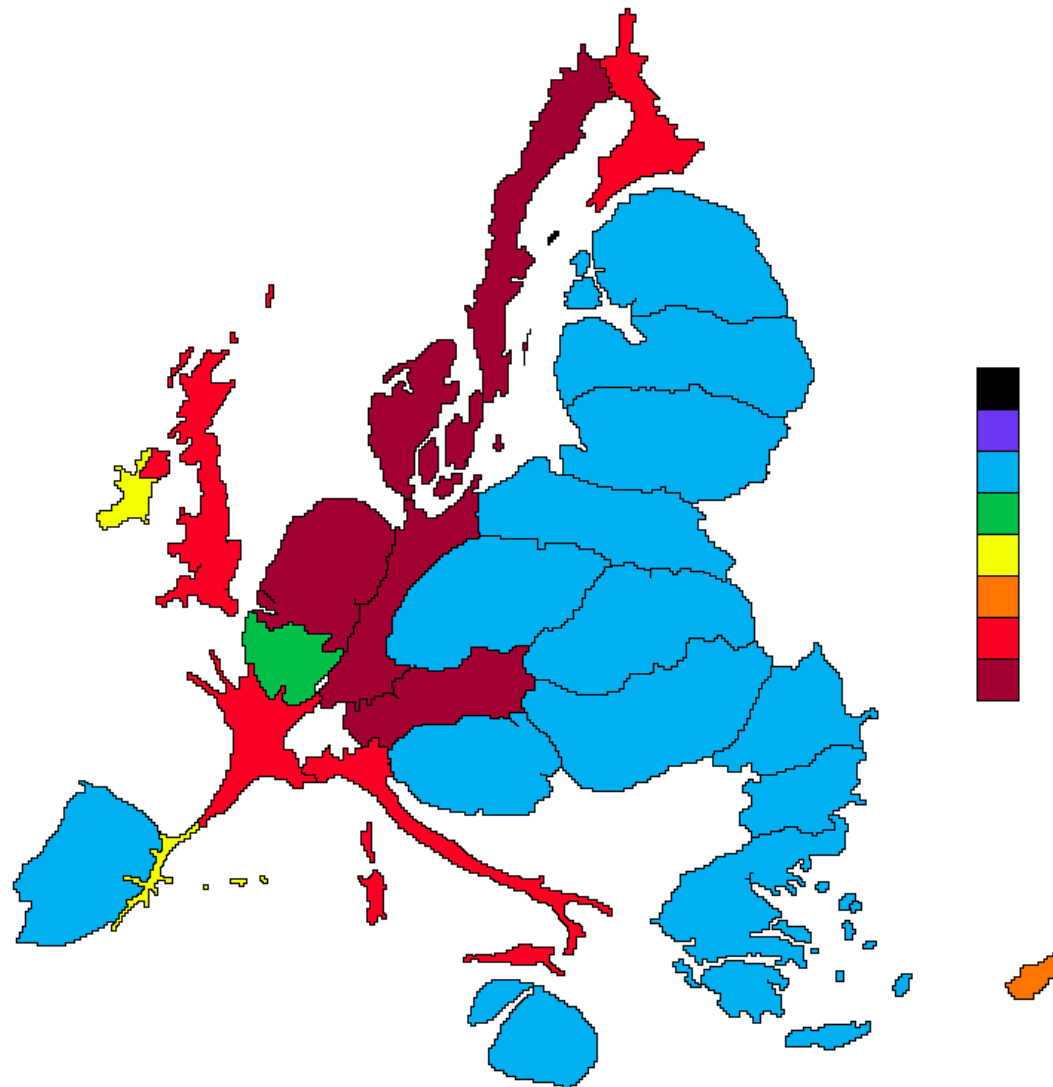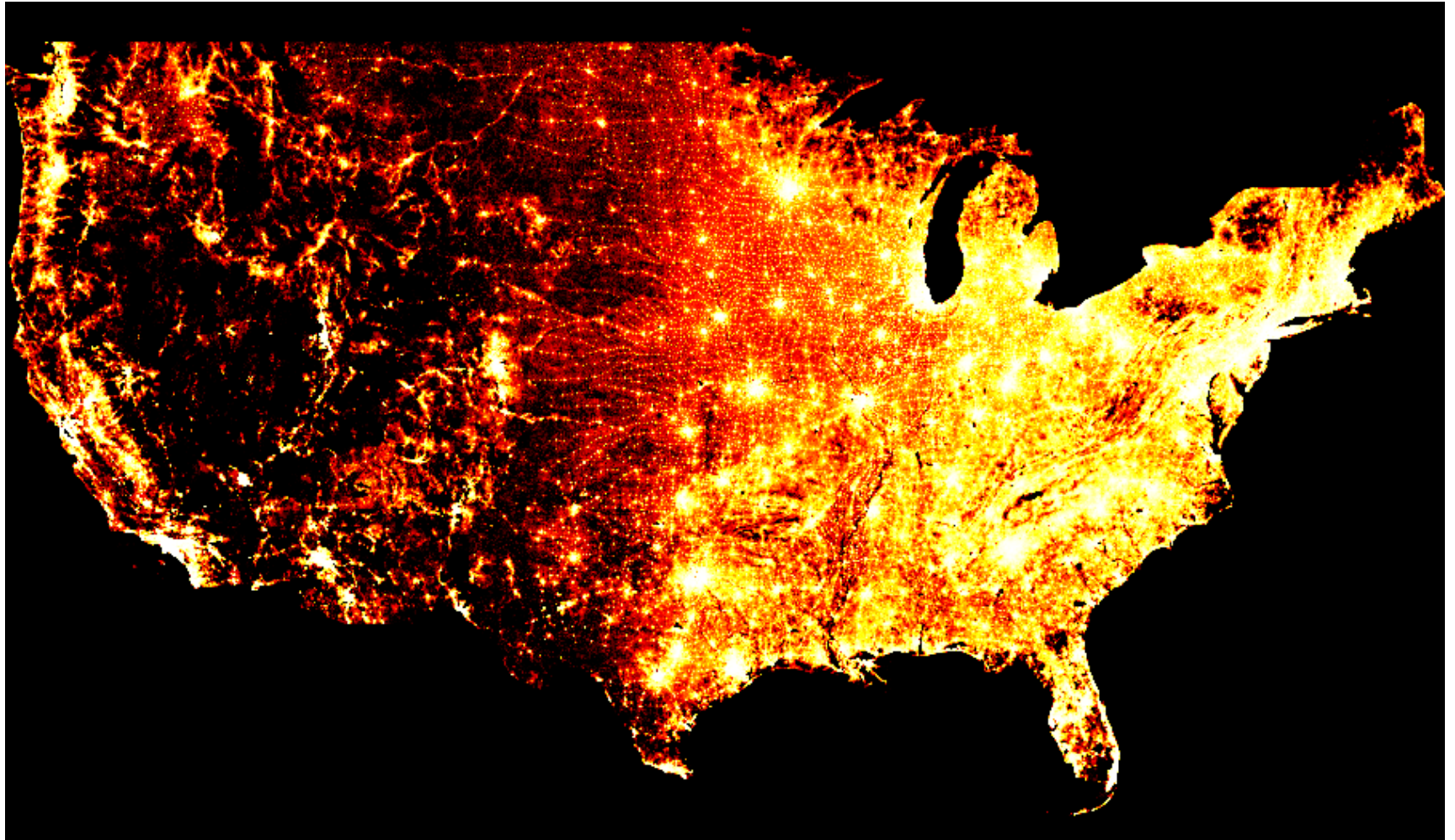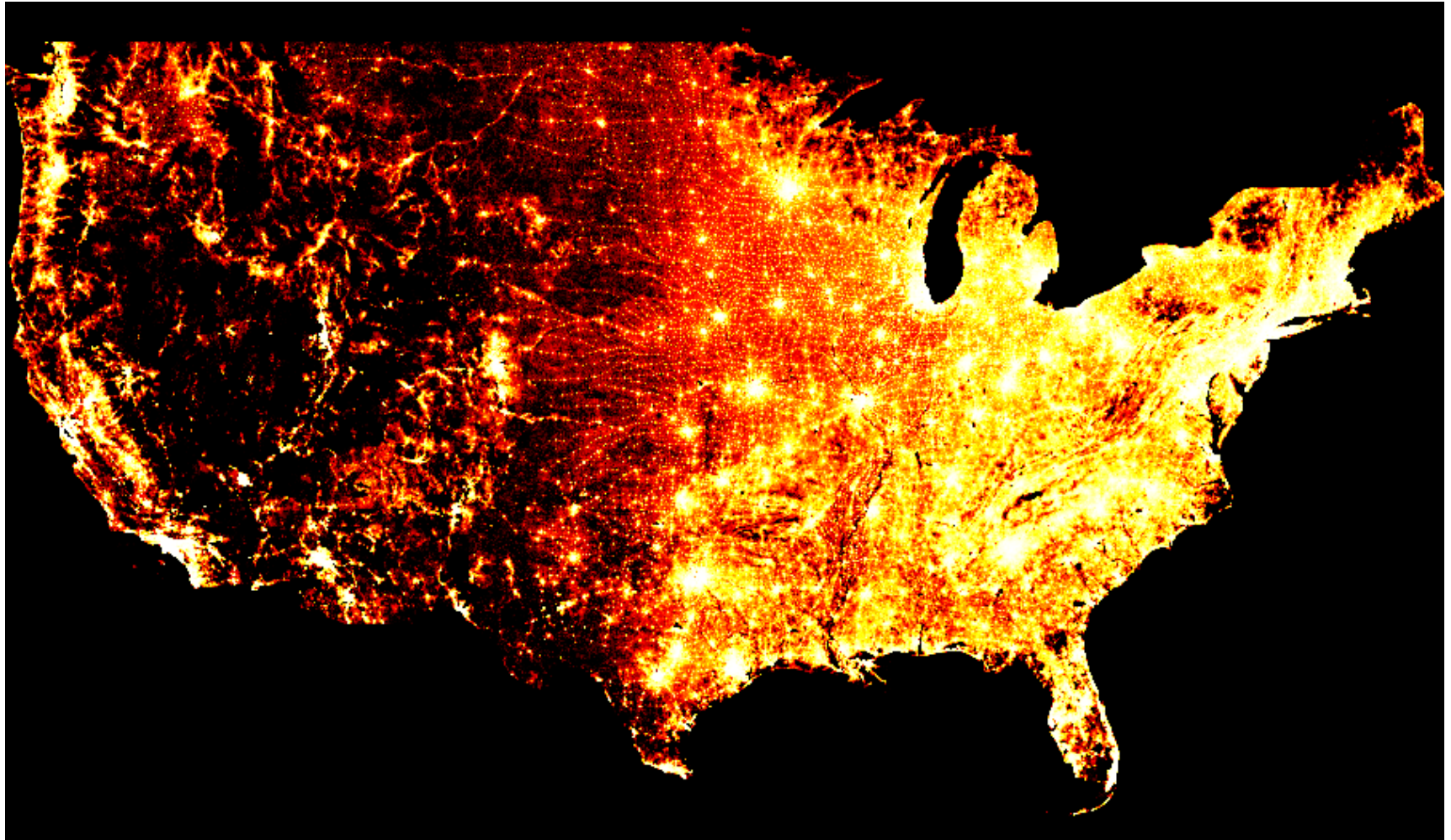minutes

A linear cartogram of the London Underground, with distance distorted to represent travel time from High Barnet station

Total European Union net budget expenditure in euros for the whole period 2007–2013, per capita

300M points
1 point = 1 person

https://datashader.org/
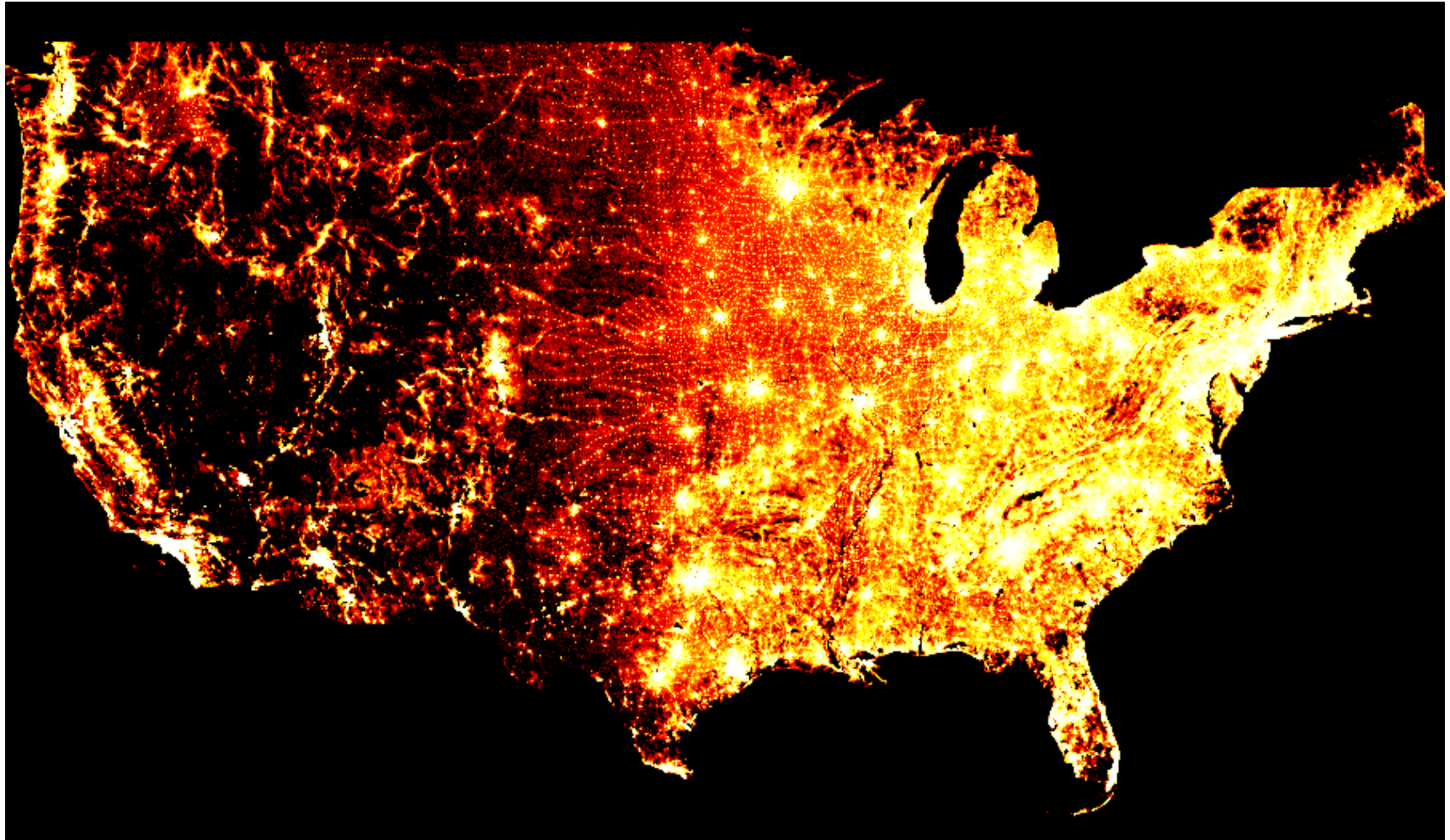
# The biggest problems of cartogram maps:

- the color is the worst to differentiate the values

## The biggest problems of cartogram maps:

- the color is the worst to differentiate the values

- frequently all you get is population density map (most of the features will be in one or another way corellated with the population density and wealth)

**When to use map charts?**

- If you want to display quantitative information on a map

- To present spatial relationships and patterns

- When a regional context for your data is important

- To get an overview of the distribution across geographic locations

- Only if your data is standardized (that is, it has the same data format and scale for the whole set)

**When to use map charts?**

- If you want to display quantitative information on a map

- To present spatial relationships and patterns

- When a regional context for your data is important

- To get an overview of the distribution across geographic locations

- Only if your data is standardized (that is, it has the same data format and scale for the whole set)

**More: https://en.wikipedia.org/wiki/Cartogram**

# Thank you for your time
# and
# See you at the next lecture

Any other
questions & comments

**l.kozlowski@mimuw.edu.pl**

# Credits:

**Infographics: The Power of Visual Storytelling by Ross Crooks, Jason Lankow and Josh Ritchie (Wiley 2012);**

**The Wall Street Journal Guide to Information Graphics by Dona Wong (Dow Jones & Company 2010)**

**Visualize This by Nathan Yau (Wiley 2011)**

**Michael Grossberg (2015) Data Visualization Basics Tools, Principles and Pitfalls**

**https://www.kaggle.com/timoboz/data-science-cheat-sheets**

*visage*

HubSpot