

Introduction to Data Visualization

Morine Amutorine
Benjamin Akera
Elaine Nsoesie





Instructor introductions ...

Contacts:

Morine - morine.amutorine@one.un.org / Twitter:
@M_moryn

Elaine - onelaine@bu.edu / Twitter: @ensoesie

Ben - akeraben@gmail.com



Resources

Github Repo

https://github.com/ensoesie/DSA_Visualization

Google Trends

<https://trends.google.com>

Twitter

<https://developer.twitter.com>

Why visualize data?

A picture is worth a thousand words


It is easier to remember pictures than text

Useful for understanding data

Can summarize large amounts of complex data

What Makes a Good Visualization?

explicit (implicit)





Visualization in Data Science can be used to:

- Explore data
- Analyze data
- Communicate findings
- Quickly draw attention to key messages



**How to use
visualizations to
communicate
effectively?**

1

Decide on what your visualization should convey

FOCUS ON THE DATA

The style and structure of your visualization will depend on its purpose

Design for a specific audience

Tell a good story with a clear message


2

Use color and size to highlight and suppress information


East Asia and Pacific South Asia Europe and Central Asia Middle East and North Africa Sub-Saharan Africa Latin America and Caribbean North America


WORLD


The average life expectancy in the world in 2009 was 69 years.



Source: The World Bank; Graphic by: Nathan Yau








3

Use length and position to express quantitative information. Use color for categorical information

Scatter plots and **bar charts** allow for **more accurate comparison of information over time** compared to pie charts




4

Think carefully about color selection and usage

Use color to create groupings

Add a **single color** to a black and white image

Use black and white to add contrast to an image with a single color gradient



4

Think carefully about color selection and usage

Some colors have pre-established meanings

Consider those with color blindness

Red

Stop

Dangerous

Hot

Green

Moving

Money

Plants

Blue

Water

Cool

Safe


5

Use all available space and proper scales

Child Mortality in 1980 and 2015

Child mortality is the probability that a newborn will die before reaching the age of 5.

Our World
in Data



Data source: UN Child Mortality Estimates
This data visualization is part of AfricanData.org – an Our World in Data project.

Licensed under CC-BY-SA by the author Max Roser.

Scale does not always have to include zero

Optimize the ratio between plot objects to capture accurate relationships

Transform data to a different scale e.g. use log scale to show percentage change over time

6

Use text and labels to improve interpretation

AVERAGE DAILY CONSUMPTION, PER PERSON

SPEED SCALE BY

YEAR
1985



Use meaningful titles

Label axis, as needed


Add texts directly to the image - do not always rely on legends

Lines should not obstruct points

Use colors (e.g. light grey) and weight that lessen focus on tick marks and grids

7

Balance complexity and clarity



Color World Regions



Select Search...

- Afghanistan
- Albania
- Algeria
- Andorra
- Angola

Size Population

Zoom 100%

DATA DOUBTS






Examples



When to use?

Visualize
correlation/association

Bubbles



Color World Regions ▼




Select Search...

- Afghanistan
- Albania
- Algeria
- Andorra
- Angola


Size Population ▼ ?

Zoom 100%



Scatterplot

- Connected scatter
- Correlogram
- Heatmap






Maps

When to use?

Useful for spatial visualizations


Child Mortality in 1980 and 2015

Child mortality is the probability that a newborn will die before reaching the age of 5.





Data source: UN Child Mortality Estimates
This data visualization is part of AfricanData.org – an Our World in Data project.

Licensed under CC-BY-SA by the author Max Roser.



Urban population
(% of population, 2014)






- Maps with bubbles
- Maps with pins

When to use?

Useful for rankings

Bar plots


Top five themes of hashtags around the world




When to use?

Useful for showing evolution

Area/density plots







- Line plot
- (Stacked) area plot
- Stream chart

WORLD

The average life expectancy in the world in 2009 was 69 years.



When to use?

Useful for information flow

Networks

13 TRAXLER, GAVRIN, and LINDELL

PHYS. REV. PHYS. EDUC. RES. 14, 020107 (2018)




FIG. 4. Forum networks from weeks 7–8 in semester 1 (left) and semester 2 (right). Line opacity is scaled by edge weight, so darker lines indicate more threads in common for a student pair. Nodes are sized by total contributions over the semester and colored by grade (red low, yellow medium, blue high). Nodes without grades (withdrawals and instructor or CN staff accounts) are white, and the instructor's node is labeled "I."


Flows of Global Health Financing

Total for 2016: \$37.5 billion in 2017 US dollars

Source: All

Channel: All

Region: Sub-Saharan Africa



Year: 2016

Target


Region

Health Focus Area

Reset

Global Migration 1960-1965

Estimates of Global Bilateral
Migration Flows by Gender
between 1960 and 2015.
doi.org/10.1111/imre.12327



Code available from:

[https://guyabel.com/post/
animated-directional-
chord-diagrams/](https://guyabel.com/post/animated-directional-chord-diagrams/)

Chord diagram



Bad visualizations




Which of these images has issues?

Have You Ever Liked a Brand on Facebook?




A

Years Experience Required by Employers




C

Which of these images has issues?




B




D

Conflicting polls



What's wrong with these images?





Tools and Resources



Python libraries

- Matplotlib
- ggplot
- Seaborn
- Bokeh
- Pygal
- Plotly
- Geoplotlib
- Gleam
- Missingno
- Leather
- Pydot

Deviation	Correlation	Ranking	Distribution	Change over Time	Competition	Magnitude	Part-to-whole	Spatial	Flow
<p>Dot plot</p> <p>Use when the relationship between two or more things. For example, the number of goals scored by a team can be a target of a player's performance. Can be used to show the number of goals scored by a player in a season.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Scatter plot</p> <p>Use when the relationship between two or more things. For example, the number of goals scored by a team can be a target of a player's performance. Many tables will require the use of a scatter plot to show the relationship between two variables.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Bar chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Area chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Line chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Bar chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Bar chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Stacked bar chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Map</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>	<p>Line chart</p> <p>Use when a single value is being compared across two or more things. For example, the number of goals scored by a team can be a target of a player's performance.</p> <p>Example FT uses</p> <p>Table with 2 columns: player, goals scored.</p>

Visual vocabulary

Designing with data

There are so many ways to visualise data - how do we know which one to pick? Use the categories across the top to decide which data relationship is most important in your story, then look at the different types of chart within the category to form some initial ideas about what might work best. This list is not meant to be exhaustive, nor a wizard, but is a useful starting point for making informative and meaningful data visualisations.



The Chart Doctor





Other tools

- Tableau
- R ggplot2 and others
- D3

Next ... ipython tutorial