Announcements

| |

1

HW1

• Due Tue 2/8

Plan

Wrap up Relational Data

- SQLite examples
- DB joins

Visualization and Data Exploration

15-388/688 - Practical Data Science: Visualization and Data Exploration

Pat Virtue Carnegie Mellon University Spring 2022

Slide credits: CMU AI, Zico Kolter

Outline

Basics of visualization

Data types and visualization types

Outline

Basics of visualization

Data types and visualization types

Software plotting libraries

Two types of visualization

Data exploration visualization: figuring out what is true

Data presentation visualization: convincing other people it is true

This lecture will mostly be focused on the first, some later lectures will touch on the second

"Data exploration" is much broader than just visualization (most of the analysis techniques we will cover fit into it)

Importance of visualization

Before you run any analysis, build any machine learning system, etc, always visualize your data

If you can't identify a trend or make a prediction for your dataset, it's unlikely that an automated algorithm will

This is especially important to keep in mind as you hear stories of "superhuman" performance of AI methods (it is possible, but takes a long time, and is not the norm)

Visualization vs. statistics

Visualization almost always presents a more informative (though less quantitative) view of your data than statistics (the noun, not the field)



[Source: https://twitter.com/JustinMatejka/status/770682771656368128 Credit: @JustinMatejka, @albertocairo]

This is a mathematical property: n data points and m equations to satisfy, with n > m

Outline

Basics of visualization

Data types and visualization types

Software plotting libraries

Data types

Nominal: categorical data, no ordering Example – Pet: {dog, cat, rabbit, …} Operations: =, ≠

Ordinal: categorical data, with ordering Example – Rating: {1,2,3,4,5}
Operations: =, ≠, ≥, ≤, >, <</p>

Interval: numerical data, zero doesn't mean zero "quantity" Example – Temperature Fahrenheit Operations: $=, \neq, \geq, \leq, >, <, +, -$

Ratio: numerical data, zero has meaning related to zero "quantity" Example – Temperature Kelvin Operations: $=, \neq, \geq, \leq, >, <, +, -, \div$

Poll 1

What data type is temperature data in Celsius

- A. Nominal
- B. Ordinal
- C. Interval

D. Ratio

Poll 2

What data type is data for number of children?



✓B. Ordinal

 χ C. Interval



Visualization Types

Most discussion of visualization types emphasizes what elements the chart is trying to convey

Instead, we are going to focus on the type and dimensionality of the underlying data

Visualization types (not an exhaustive list):

- 1D: bar chart, pie chart, histogram
- 2D: scatter plot, line plot, box and whisker plot, heatmap
- 3D+: scatter matrix, bubble chart

1D DATA

Bar chart



Bar chart (bad)

Don't use lines within a bar chart for categorial or ordinal features!



Pie chart



Histogram H 1400 -1200 · Data 1000 Nominal X 800 Ordinal X Interval \checkmark 600 Ratio \checkmark 400 200 -0 0 10 20 30 40 0 Data \rightarrow Count per bin (range of values) \rightarrow Plot with bar chart

Histogram

OK to use lines within a histogram (but not very informative)

	Dim 1	Dim 2
Nominal	X	X
Ordinal	X	X
Interval	\checkmark	~
Ratio	\checkmark	\checkmark

Why not ordinal data in first dimension?

Scatter plot

	Dim 1	Dim 2
Nominal	X	X
Ordinal	X	X
Interval	\checkmark	\checkmark
Ratio	\checkmark	\checkmark

Why not ordinal data in first dimension?

Heatmap (density, or 2D histogram)

Several factors to consider. For example:

- Accessibility
- Printing in grayscale
- Unintentional boundaries
- Intentional boundaries
- Color semantics

Several factors to consider. For example:

- Accessibility
- Printing in grayscale
- Unintentional boundaries
- Intentional boundaries
- Color semantics

Image: https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0199239

Several factors to consider. For example:

- Accessibility
- Printing in grayscale
- Unintentional boundaries
- Intentional boundaries
- Color semantics

Image: https://jakevdp.github.io/blog/2014/10/16/how-bad-is-your-colormap/

Several factors to consider. For example:

- Accessibility
- Printing in grayscale
- Unintentional boundaries
- Intentional boundaries
- Color semantics

Image: https://eagereyes.org/basics/rainbow-color-map

Several factors to consider. For example:

- Accessibility
- Printing in grayscale
- Unintentional boundaries
- Intentional boundaries
- Color semantics

Image: https://weather.com/maps/currentusweather

Several factors to consider. For example:

- Accessibility
- Printing in grayscale
- Unintentional boundaries
- Intentional boundaries
- Color semantics

Community Transmission in US by County

Image: https://covid.cdc.gov/covid-data-tracker/#county-view

Scatter plot (bad)

Box and whiskers

Violin plot

Line plot

Why not ordinal data in first dimension?

Heatmap (matrix)

	Dim 1	Dim 2	ho
Nominal	\checkmark	\checkmark	
Ordinal	\checkmark	\checkmark	
Interval	X	X	
Ratio	X	X	

Bubble plot

3D+ DATA

3D scatter plot

	Dim 1	Dim 2	Dim 3
Nominal	X	X	X
Ordinal	X	X	X
Interval	X	X	X
Ratio	X	X	X

Scatter plot matrix

Nominal

Ordinal

Interval

Ratio

Bubble plot

Color scatter plot

t-SNE (be careful!)

t-SNE (T-distributed Stochastic Neighbor Embedding) works to embed high dimensional data in smaller (often 2 or 3) dimensions. Van der Maaten, L., & Hinton, G. (2008).

Easy to misread : <u>https://distill.pub/2016/misread-tsne/</u>

Outline

Basics of visualization

Data types and visualization types

Software plotting libraries

Matplotlib

Matplotlib is the standard for plotting in Python / Jupyter Notebook

Matplotlib used to generate fairly ugly plots by default, but in recent versions this is no longer the case, so minimal need for additional libraries

It is aimed at generating static plots, not very good for interacting with data (with a few exceptions)

A number of additional libraries provide some level of interactive plot (and static plots), but matplotlib is enough of a standard that we'll use it here

