# Data Visualization for Data Science Practitioners

**Olaf Menzer**

Senior Data Scientist

[1] *Pacific Life, Newport Beach*
[2] *University of California, Santa Barbara*

*Advanced Analytics Meetup, Irvine, April 17th 2019*

# Bio



2006: Research Assistant at Max Planck Institute in Jena, Germany… started working with climate data and maps

2011: M.Sc. in Bioinformatics in Jena, Germany

2011: Moved to California, UC Santa Barbara; Research in Ecosystem Science, Statistics

2015: Ph.D. in Geographic Information Science, UCSB

2015-2018: Predictive Analyst, Ingram Micro, Irvine

Since 2018: Data Scientist, Pacific Life, Newport Beach



*Jena, Germany*



*UC Santa Barbara*



*Pacific Life HQ*

2

# Quick Commercial

PyData Socal Meetup in Venice, tomorrow Thursday, 6.30pm

https://www.meetup.com/PyData-SoCal/events/259853926/

# Data Visualization Credits

*Tamara Munzner,*
*Visualization Analysis & Design, 2014*

DataVis @CalTech

VisLab @UCSB

1) Data Visualization – Why do we care?

2) Data Visualization Elements & Tools

3) Advanced Graph Examples (with R code)

4) Emerging Data Visualizations

5) R Exercises

1) **Data Visualization – Why do we care?**

2) Data Visualization Elements & Tools

3) Advanced Graph Examples (with R code)

4) Emerging Data Visualizations

5) R Exercises

# The Data Science "Skillgram"

# Data Visualization in Data Science

# Data Visualization – as a means to...

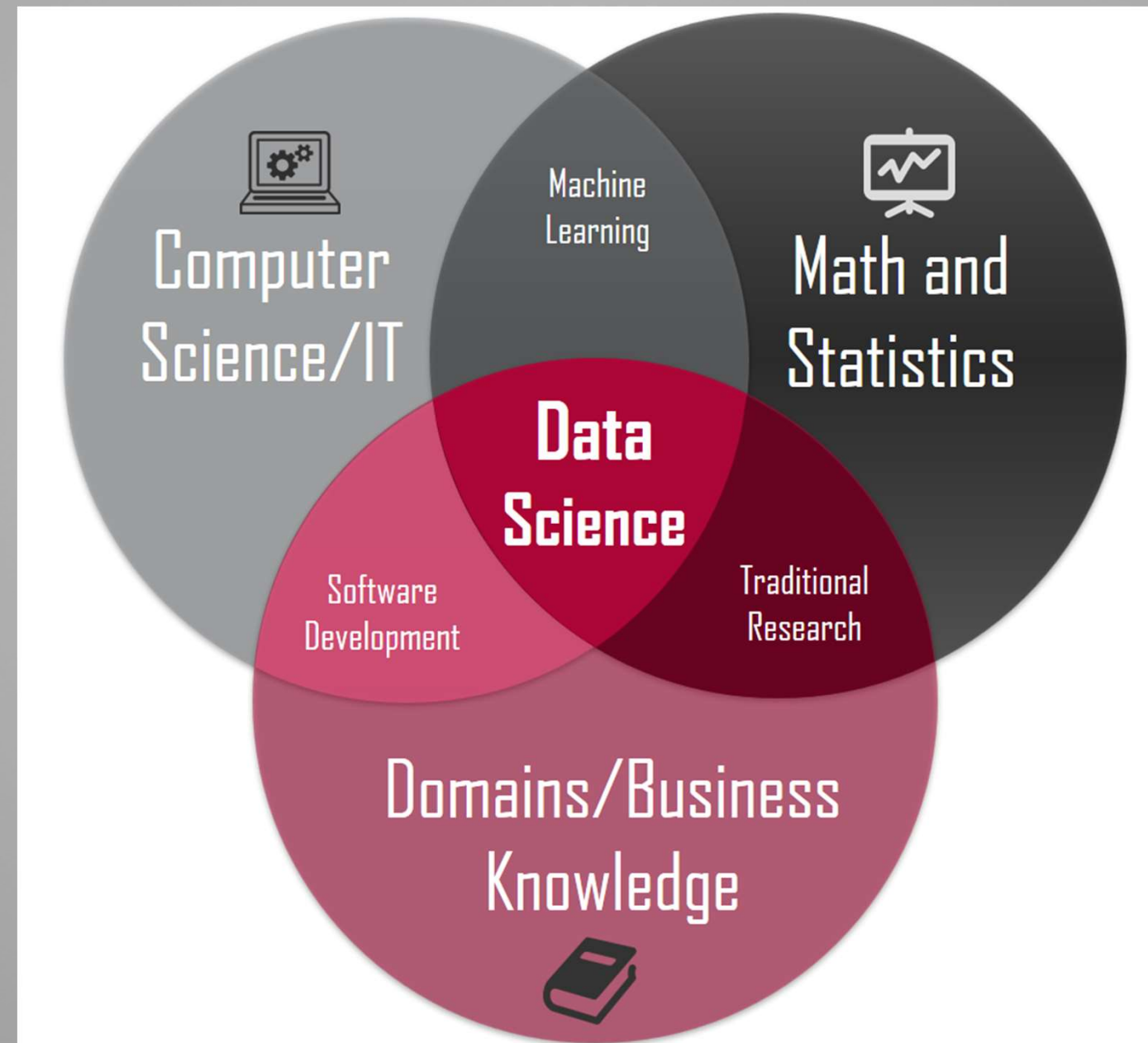| | |
|---|---|
| (1) The need to augment human capabilities to view & interpret data | |
| (2) Understanding relationships between factors | |
| (3) Recognizing patterns | |
| (4) Breaking down complexity / dimensionality reduction | |
| (5) Exploratory Analysis | |
| (6) Communicating results to collaborators, clients, friends, ... | |

# Anscombe's quartet



| Property | Value |
|---|---|
| Mean of x | 9 |
| Variance of x | 11 |
| Mean of y | 7.5 |
| Variance of y | 4.125 |
| Correlation | 0.816 |
| Regression Line | $y = 3 + 0.5x$ |
| Coefficient of determination | 0.67 |

# Why can Data Viz be difficult?



A Nested Model for Visualization Design and Validation. Tamara Munzner. IEEE TVCG (Proc. InfoVis 2009), 15(6):921-928, 2009

# What makes a good graph successful?



A Nested Model for Visualization Design and Validation. Tamara Munzner. IEEE TVCG (Proc. InfoVis 2009), 15(6):921-928, 2009

# What makes a good graph successful?



*Good in theory,*

*But... what next?*

A Nested Model for Visualization Design and Validation. Tamara Munzner. IEEE TVCG (Proc. InfoVis 2009), 15(6):921-928, 2009

# Practical considerations



To make good graphs, practitioners are faced with:

1) Tasks that require skills in Computer Science, Design, Psychology, etc.

2) **Too many choices** in graph types and platforms

3) Searching for the **template that is most suitable** for the data

4) Having to test if the graph represents information correctly

5) Finally, **developing a narrative** / story (at the least, in form of a headline)

1) Data Visualization – Why do we care?

2) **Data Visualization Elements & Tools**

3) Advanced Graph Examples (with R code)

4) Emerging Data Visualizations

5) R Exercises

# I. Data Representation and Abstraction ("*What?*")



Data and Dataset Types

| Tables | Networks & Trees | Fields | Geometry | Clusters, Sets, Lists |
|--------|------------------|--------|----------|----------------------|
| Items | Items (nodes) | Grids | Items | Items |
| Attributes | Links | Positions | Positions | |
| | Attributes | Attributes | | |

(Munzner (2018) Visualization Analysis & Design)

16

# II. Visualization Actions ("*Why?*")

# III. Most Widely Used Visualization Tools ("*How?*")



(visual.ly; online-behavior.com)

# Time Series example #1: 250 years of Sunspots data



(Wilkinson, 2005,
"The Grammar of Graphics)



"Stretch transformation" shows periodicity of 11 years much more clearly

# III. How to improve your Graph?



Isaac Reyes, DataSeer, as seen at ODSC West in San Francisco 2018

# III. How to improve your Graph?



Airbnb Bookings now outnumber Hotel Bookings 3 to 1

Isaac Reyes, DataSeer, as seen at ODSC West in San Francisco 2018

# Visual Design Aspects (~Tufte)

Graphical displays must:

1) Show the data

2) Present many numbers in a small space (cf. data – ink ratio)

3) Make large data sets coherent

4) Encourage the eye to compare different pieces of data

5) Reveal the data at several levels of detail

6) Serve a clear purpose: description, exploration, tabulation or decoration

7) Be closely integrated with statistical and verbal descriptions of a data set.

8) ...

   (*Tufte E. - The Visual Display of Quantitative Information*)

# Other considerations

1) Who is the audience of your graph?
   → Always know your audience.

2) What is the target medium? A laptop, a large projector, a magazine, a print poster or a handout?
   → Consider resolution, color depth, format, etc.

3) Gestalt principles from Psychology Research of human perception

1) Data Visualization – Why do we care?

2) Data Visualization Elements & Tools

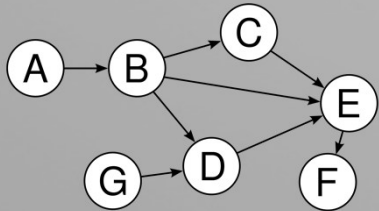3) **Advanced Graph Examples (with R code)**

4) Emerging Data Visualizations

5) R Exercises

# Beyond line charts, scatterplots and histograms …

References

https://www.r-graph-gallery.com/

https://flowingdata.com/

https://plot.ly/r/shiny-gallery/

# [1] Streamgraph

# [2] Sunburst

*... or how to tell the difference between 66 sorts of cheese.*

Source

# [3] Sankey Chart



Source

# [4] Circular Plot / Chord Diagram



- Ability to show large amount of detail
- Heavily used in genomics research
- Interactivity can be key

- Source: circos.ca

# [4] Circular Plot / Chord Diagram [R]



- shows and differentiates inbound & outbound migration patterns per country,
- magnitudes are easier to compare than on a map (due to fixed reference)

| country | to_Aus | to_Chin | to_Ind | to_Jap | to_Thai | to_Mal |
|---------|--------|---------|--------|--------|---------|--------|
| Australia | 1 | 35000 | 10000 | 7000 | 70000 | 60000 |
| India | 150000 | 1 | 7000 | 8000 | 30000 | 90000 |
| China | 90000 | 10000 | 1 | 175000 | 22000 | 110000 |
| Japan | 180000 | 12000 | 40000 | 1 | 120000 | 14000 |
| Thailand | 15000 | 25000 | 5000 | 11000 | 1 | 30000 |
| Malaysia | 10000 | 8000 | 4000 | 18000 | 40000 | 1 |

Source

# Advanced Graphs in R

In summary:

1) **tradeoff** between showing enough detail and cognitive load

2) Tendency to adding more complexity and dimensions

3) Complex graphs in this section are a good example for visualizing more data rich relationships between variables...

4) ... but need a lot of preparation, crafting, and often a detailed description to convey a story

1) Data Visualization – Why do we care?

2) Data Visualization Elements & Tools

3) Advanced Graphs in R

4) **Emerging Data Visualizations**

5) R Exercises

32

# "Tomorrow's Earth" (*Science* June 28th 2018)

# "Tomorrow's Earth" (*Science* June 28th 2018)



**Tomorrow's Earth will be**

**More quantitative.**

**More digital.**

**Therefore, more visual.**

# >10 Million data points in one large network graph

When you have, say, 10M data pairs or more...

How to test if your graph
represents the data correctly?

# Practical recommendations

To reduce bias in charting, the following steps can help:

1. Plot multiple samples of the same data set.

2. Plot even more samples. Repeat.

3. Plot various aggregation levels (sum, group by, mean vs median).

4. Validate with external data sets (ideally, obtained from independent sources).

5. Use interactivity

6. Test for visual significance, e.g., with error bars —> "graphing uncertainty".

7. Transform axis (log, log-log, sigmoid, etc.), align, stretch/compress axis ranges.

Dynamic Data Visualization has not been solved, but we're getting a lot closer → H2O.ai's "Auto Visualization" method

What will graphs look like in 5-10 years?

# AlexNet



A visualization of … a computational graph of **18.7 million vertices** and **115.8 million edges**
**Using a good amount of abstraction to show different IPUs (Intelligent Processing Units)**

40

# Take away

(1) Data Visualization is a highly interdisciplinary task

(2) Data Visualization needs are ubiquitous along the Data Science Stack

(3) Data Visualization is not fully automated at all

(4) New techniques offer a way to pack more data into a graph, but come at a cost

(heavy manual editing, detailed descriptions, danger of cognitive overload)


Learning how to make better visualizations, and to fit visualizations to given data sets

for a given audience will help us to improve our means of Visual Communication.

# Further Reading

1) Munzner, T. (2017) Visualization Analysis & Design, *VIS 17 Tutorial*, *Phoenix AZ 10/17* (pdf slides available here http://www.cs.ubc.ca/~tmm/talks.html#vad17fullday

2) Tufte, E. (1990) Envisioning Information. Graphics Press

3) Wilkinson, L. (2005) The Grammar of Graphics. Springer

1) Data Visualization – Why do we care?

2) Data Visualization Elements & Tools

3) Advanced Graphs in R

4) Emerging Data Visualizations

5) **R Exercises**

43

# Exercise 1: Biopharmaceutical App

Github link: https://github.com/olaf-menzer/biopharma-app

Example of an interactive, visual clustering app built in R shiny

It is centered around interactivity in data exploration, through the following:

- Plotting a heatmap
- A simple hierarchical clustering algorithm
- Plotting a phylogenetic tree to visualize the clusters

# Exercise 2: Complex Graphs in R

Github link: https://github.com/olaf-menzer/Rgraphs

We will learn how to plot the charts in the section "Beyond line charts, scatterplots and histograms" in R, some of them with interactivity.

4 example R scripts (mini examples) for plotting:

(1) Streamgraph      (01_Streamgraph.R)
(2) Sunburst Chart      (02_SunburstR.R)
(3) Sankey Chart      (03_Sankey.R)
(4) Circular Chart      (04_Circular_Chord.R)

Example data sets given in file Explore_DataSets.R
(2) – (4) have small, coded mini example data sets built in so we can run them independent of WiFi bandwidth

# *Thank you.*