

## CANADA'S MICHAEL SMITH GENORE SCENENCES CENTERE

# **ESSENTIALS OF DATA VISUALIZATION** THINKING ABOUT DRAWING DATA + COMMUNICATING SCIENCE







## ENCODING

choosing a data representation

When you think of data visualization, the first ideas that come to mind are a scatter plot, or a bar chart, a box plot or a network diagram. These are all data encodings—methods that relate data values to the positions, sizes and shapes of the lines or symbols that appear on the screen or in a figure.

There are many data encodings—which do you choose? First, it must help answer relevant questions about the data that are difficult, or impossible, to answer by staring at the raw data itself. As such, the encoding may be of the data or some transformation of the data that addresses your questions. Just because you have a network doesn't mean you should automatically draw a force-directed hairball. Second, it should accommodate the uncertainty of the data. How do you incorporate uncertainty in a scatter plot? Easy—error bars. How about a pie chart. Ah, well, umm. We'll come back to this later. Third, it should be flexible enough to address questions that you haven't thought of yet. This sounds vague, I know. What I mean is that if the encoding warps the data or doesn't to at least try to limit occlusion (the phenomena where points overlap and hide behind each other), it's likely to be less useful.

### DEPARTURE DEPART

|                          | au départ                |                  | Departures | trains | Hinweise<br>Remarks                         |          | der zi |
|--------------------------|--------------------------|------------------|------------|--------|---|----------|--------|
| Heure 11 <sup>h</sup> 58 | Destination<br>DAMMARTIN | CREPY SOISSONS A | NIZY-PINON |        | Particularités                              | Train nº | 849917 |
| 12 01                    | PERSAN CHA               | MBLY MERU ST-SUL | PICE BEAUV | AIS    |   | ter-     | 847419 |
| 12 01                    | BRUXELLES -              | MIDI LIEGE AACHE | N KOLN ESS | SEN    | THALYS<br>CONFORT 1 ET CONFORT 2            | 2        | 9437   |
| 12 07                    | CREIL COMP               | IEGNE CHAUNY TER | GNIER ST-C | UENTIN | 14re ET 2eme CLASSE                         |          | 12309  |
| 12 10                    | ORRY-LA-VI               | LLE CHANTILLY-GO | UVIEUX CRE |        |   | ter-     | 847609 |
| 12 * 25                  | BRUXELLES                | ROTTERDAM SCHIPH | OL AMSTERC | AM     | THALYS<br>CONFORT 1 ET CONFORT 2            | 2        | 9339   |
| 12 * 28                  | CREIL LONG               | UEAU AMIENS      |            | BEBEBB | 1 <sup>4re</sup> ET 2 <sup>4me</sup> CLASSE |          | 12011  |
| 12 43                    | LONDON ST                | PANCRAS INT      |            |        | HALL LONDRES                                |          | 9029   |
| 12 46                    | LILLE FLAN               | DRES             |            |        | 1ere ET 2ere CL avec RESERVATION            |          | 7043   |
| 12 49                    | CREIL RIEU               | X PONT LONGUEIL  | COMPIEGNE  |        |   | ter-     | 847807 |
| 12 52                    | ARRAS LENS               | BETHUNE HAZEBRO  | UCK DUNKER | QUE    | 1** ET 2*** CL avec RESERVATION             | Tor      | 7321   |
| 12 * 52                  | ARRAS DOUA               | I VALENCIENNES   |            |        | 1** ET 2*** CL. avec RESERVATION            | Tor      | 7121   |
| 12 55                    | BRUXELLES -              | MIDI             |            | EBBBB  | THALYS<br>CONFORT 1 ET CONFORT 2            |          | 9341   |
| 13 01                    | PERSAN CHA               | MBLY MERU ST-SUL | PICE BEAUV | AIS    |   | term     | 847421 |
| 13 13                    | EBBSFLEET                | LONDON ST PANCRA | SINT       |        | HALL LONDRES                                |          | 9031   |
| 13 16                    | LILLE FLAN               | DRES             |            | EBBBBE | 1** ET 2*** CL. avec RESERVATION            | 75~      | 7045   |
| 13 * 26                  | CREPY VILL               | ERS SOISSONS ANI | ZY-PINON L | AON    |   | ter-     | 849921 |
| 13 59                    | LONGUEAU A               | MIENS ABBEVILLE  |            | BBBBB  | 1ere ET 2eme CLASSE                         |          | 2013   |
| 14 " 01                  | PERSAN CHA               | MBLY MERU ST-SUL | PICE BEAUV | AIS    |   | ter-     | 847423 |
| 14 07                    | CREIL PONT               | COMPIEGNE TERGN  | IER ST-QUE | NTIN   |   | ter-     | 847905 |
|                          |                          |                  |            |        |   |          |        |
|                          |                          |                  |            |        |   | 11.03    |        |
|                          |                          |                  |            |        |   | - aa     |        |
| ectori 🗊 cevral          |                          |                  |            |        |   |          |        |

# ABFAHRT





The visual display of quantitative information. R. Tufte (2001) 2nd ed. Text in original modified.



The visual display of quantitative information. R. Tufte (2001) 2nd ed. Text in original modified.



The visual display of quantitative information. R. Tufte (2001) 2nd ed. Text in original modified.



The visual display of quantitative information. R. Tufte (2001) 2nd ed. Text in original modified.



The visual display of quantitative information. R. Tufte (2001) 2nd ed. Text in original modified.



J Neurosci (2015) 35:10899-10910.



J Neurosci (2015) 35:10899-10910. Redesign using UpSet encoding (IEEE Trans Vis Comput Graph (2014) 20:1983–1992.)



J Neurosci (2015) 35:10899-10910. Redesign using UpSet encoding (IEEE Trans Vis Comput Graph (2014) 20:1983–1992.)



Submitted to Skills Beyond Science, 2015 Bactory Summer School, Helsingor, Denmark. Redesign published in Environmental Microbiology, dx.doi.org/10.1111/1462-2920.13326



- deletion 105
- duplication 76
  - **SNP** 75
  - indel 46
  - insertion 31



One concept to always keep in mind is the so-called data-to-ink ratio.

Ask yourself: what ink on the page is directly related to data values and what ink is used for labels, grids, navigational components and design elements.

Then consider how to maximize ink used for data and minimize ink for everything else.

At the same time, always keep in mind ways to use less ink to tell the same story and show the same relationships without loss of accuracy and precision—and parsability!

Sometimes you need to use ink to clarify or avoid confusion. That's fine do this. The speed at which your readers understand your message and the depth of this understanding is part of the model.

So, think "data-and-its-understanding-to-ink" ratio.

But not all the data is relevant. If you can figure out which is—that's the holy grail. In fact, if you knew this you might just compute on the data and bypass the visualization. But, I encourage you to think about "actionabledata-to-ink" ratio, too.



- deletion duplication SNP indel
  - insertion

## FREQUENCY OF VARIATION BY TYPE







- deletion duplication SNP indel
  - insertion

## FREQUENCY OF VARIATION BY TYPE







- deletion 105
- duplication 76
  - SNP 75
  - indel 46
  - insertion 31

## FREQUENCY OF VARIATION BY TYPE







- deletion 105
- duplication 76
  - SNP 75
  - indel 46
  - insertion 31

## FREQUENCY OF VARIATION BY TYPE







- deletion 105
- duplication 76
  - SNP 75
  - indel 46
  - insertion 31

## FREQUENCY OF VARIATION BY TYPE



## FREQUENCY OF

# deletion 1 duplication SNP indel insertion

| - VAR | IATION | ΒY | ΤΥΡΕ |
|-------|--------|----|------|
| sam   | ple    |    |      |
| A     | B      |    |      |
| 05    | 51     |    |      |
| 76    | 38     |    |      |
| 75    | 98     |    |      |
| 46    | 5      |    |      |
| 31    | 32     |    |      |

## FREQUENCY OF VARIATION BY TYPE

- deletion
- duplication
  - SNP
  - indel
  - insertion
- - deletion 105
- duplication 76
  - SNP 75
  - indel 46
  - insertion 31

## sample

| Α   | В  |  |
|-----|----|--|
| 105 | 51 |  |
| 76  | 38 |  |
| 75  | 98 |  |
| 46  | 5  |  |
| 31  | 32 |  |

## FREQUENCY OF VARIATION BY TYPE

## sample



## FREQUENCY OF VARIATION BY TYPE

- deletion 1
- duplication
  - SNP
  - indel
  - insertion
- - deletion 105 •
  - duplication 76 SNP 75
    - indel 46
    - insertion 31 •

## sample

| Α  | В  |  |
|----|----|--|
| 05 | 51 |  |
| 76 | 38 |  |
| 75 | 98 |  |
| 46 | 5  |  |
| 31 | 32 |  |

## FREQUENCY OF VARIATION BY TYPE

## sample





1.5

# log<sub>2</sub>(error) 2





Proceedings of the 28th international conference on Human factors in computing systems (2010) ACM: Atlanta, Georgia, USA. 203–212.



















average line angle





of data visualization. Simple examples embody these.

about the same things:

Am I using ink responsibly?

Am I making comparisons easy to make?

Am I visually emphasizing the things that are relevant?

you might have a matrix of scatter plots or bar plots. Each of those individual small plots must be handled with care.

- We've started with some simple examples. There's a good reason for this.
- Always appreciate and use to your advantage the fundamental principles
- At no point, do you ever eject these principles. I don't care how big your data set is. Sure, they might be more nuanced or interrelated, if you're showing a lot of stuff on the page, but fundamentally you're still thinking

- Am I drawing shapes whose position and size can be accurately judged?
- Once your data size grows, a popular technique in visualization is called small multiples. Here, you break your data down into a large number of smaller sets and represent each set with the same primitive encoding. So,

created by Martin Krzywinski, Kim Bell-Anderson & Philip Poronnik

written and designed by

Martin Krzywinski

production One Ski Digital Media Productions

with financial support by

University of Sydney

University of Sydney, Australia

filmed at

Represent the "frequency of variation by type" two-column table as a bar plot.

Pay careful attention to bar width and the distance within groups (e.g. the deletion bars for A and B) and between groups (the deletion and duplication groups). Start with a ratio of

bar width : within : between = 1 : 0.2 : 1.5

How does this look to you? Incorporate the golden ratio ( $\phi$  = 1.62) and its inverse (1/ $\phi$  =  $\phi$ – 1 = 0.62) and use ratio of

bar width : within : between =  $1 : 1/\phi^2 : \phi$ 

Does this look better?

Draw both vertical and horizontal versions of your bar plot. In which one is the text easier to read? What would make you consider reordering the categories?

|             | sample |    |
|-------------|--------|----|
|             | Α      | В  |
| deletion    | 105    | 51 |
| duplication | 76     | 38 |
| SNP         | 75     | 98 |
| indel       | 46     | 5  |
| insertion   | 31     | 32 |

Differences in data are important—sometimes main important than the data themselves. Think of the differences as the data.

For example, in the "frequency of variation by typ two-column table, the number of deletions drops 54 from 105 in A to 51 in B.

Try to add this value to your bar plot. Is a relative absolute difference meaningful here? Suggest reasons to show one or the other.

| ore    | FREQUENCY   | OF VA | RIATION | BY TYPE |
|--------|-------------|-------|---------|---------|
| )<br>) |             | sar   | nple    |         |
|        |             | Α     | В       |         |
| pe"    | deletion    | 105   | 51      |         |
| s by   | duplication | 76    | 38      |         |
|        | SNP         | 75    | 98      |         |
| e or   | indel       | 46    | 5       |         |
|        | insertion   | 31    | 32      |         |

Draw a scatter plot using the data in the "frequency" of variation by type" two-column table. Use the number of variations in A on the X axis and the relative change on the Y axis.

Does this representation make some questions easier to answer?

Address the challenge of labeling the points—this can be difficult!

What has been gained?

|             | sample |    |
|-------------|--------|----|
|             | Α      | В  |
| deletion    | 105    | 51 |
| duplication | 76     | 38 |
| SNP         | 75     | 98 |
| indel       | 46     | 5  |
| insertion   | 31     | 32 |

You are a terrific visual calculator and comparer.

How much faster can you compare lengths than numbers? Let's check.

Print this slide. Now, with pen and paper, time yourself to see how quickly you can identify the larger of each of the two numbers in each pair. Next, time yourself to see how quickly you can identify the longer of the two lengths in each pair.

Are you surprised at the timing difference? Which task was more tiring?

Notice that as the numbers get larger the longer the comparison takes—your eye has to travel further. Can you think of a way to align the number pairs differently to speed this up?

