

Principles of Visualization

Sourav SEN GUPTA
Lecturer, SCSE, NTU





Data Science Data Visualization

Information Presentation

Is there a “story” hidden in your data?
How to use visuals as information?
How to tell the “story” effectively?

**How to present Data in
the most engaging way?**

There are two goals
when presenting data:
convey your story and
establish credibility.



Edward R. Tufte

<https://www.edwardtufte.com/tufte/>

convey your story

Effectiveness

A visualization is more *effective* than another [vis] if the information conveyed by one visualization is more **readily perceived** than the information in the other visualization.

Jock D. Mackinlay

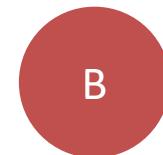
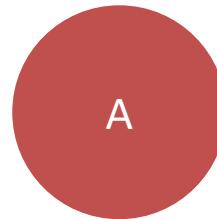
APT : A Presentation Tool (1986)

4



convey your story

readily perceived



Jock D. Mackinlay

APT : A Presentation Tool (1986)

5



establish credibility

Expressiveness

A set of facts is *expressible* in a visual language if the sentences (i.e. the visualizations) in the language express **all the facts** in the set of data, and **only the facts** in the data.

Jock D. Mackinlay

APT : A Presentation Tool (1986)

6

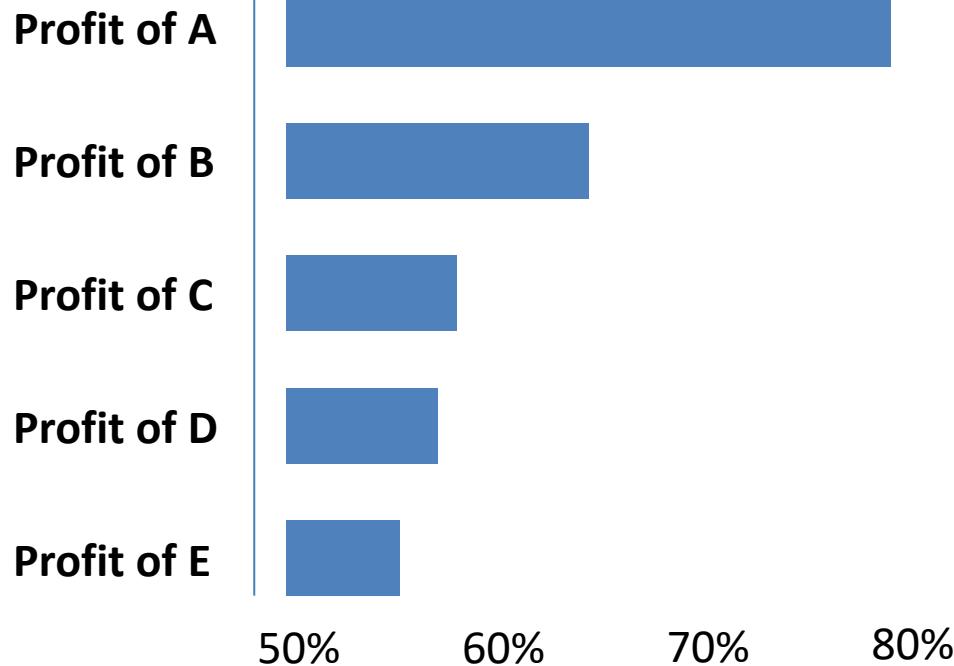


establish credibility

all the facts

Jock D. Mackinlay

APT : A Presentation Tool (1986)



A



B



C



D



E



0%

10%

20%

30%

40%

50%

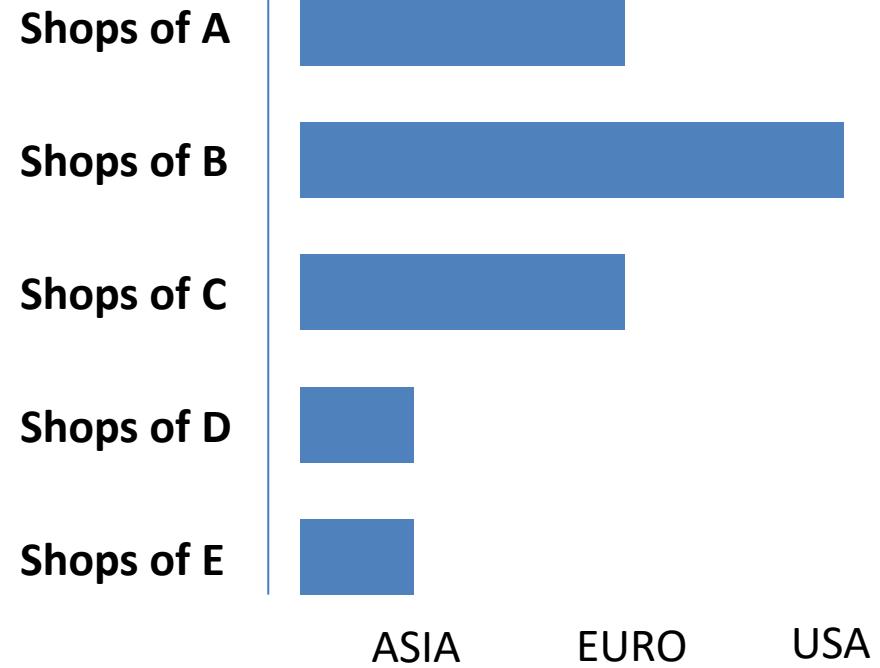
60%

70%

80%

establish credibility

only the facts



Jock D. Mackinlay

APT : A Presentation Tool (1986)

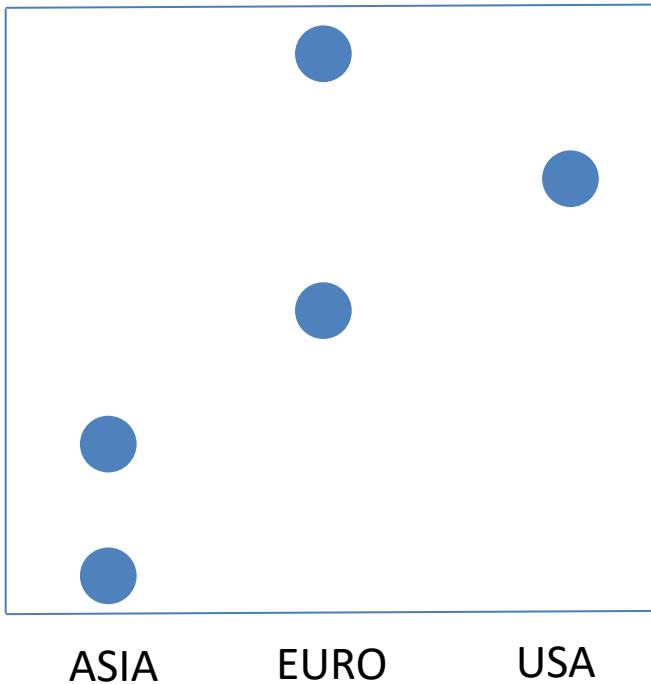
Shops of A

Shops of B

Shops of C

Shops of D

Shops of E



convey your story

Effectiveness

Use encodings that
people decode better.

Better means more accurate and faster.

establish credibility

Expressiveness

Tell the truth and
nothing but the truth.

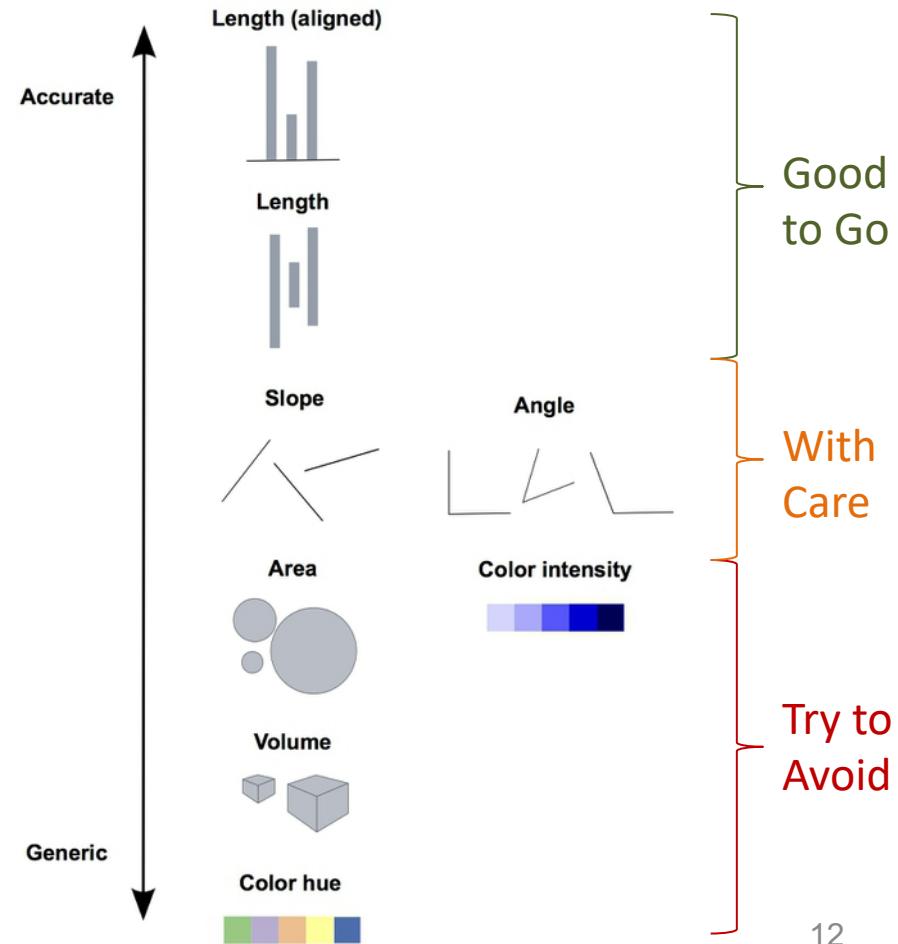
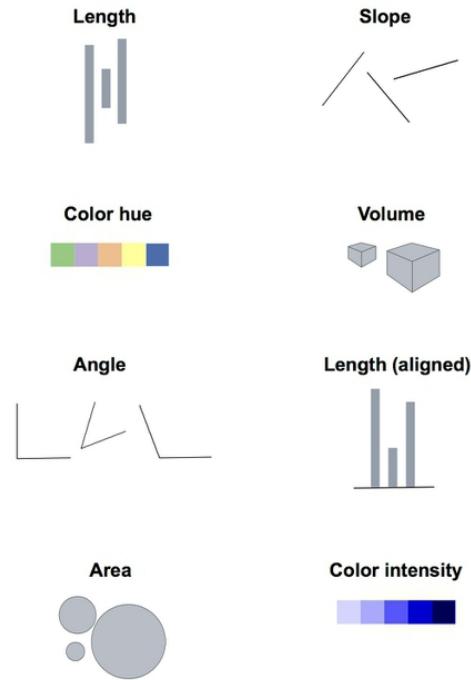
Do not lie, and do not lie by omission.

Jeffrey Heer

Principles of Data Visualization

11





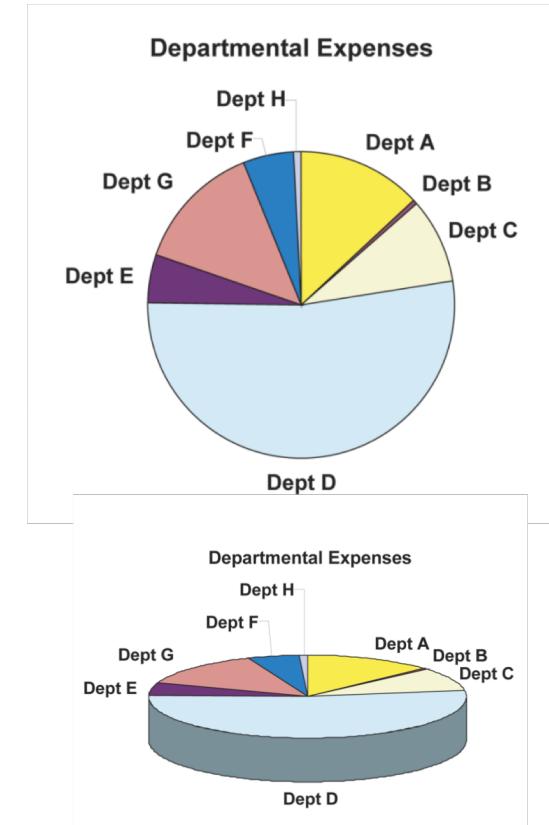
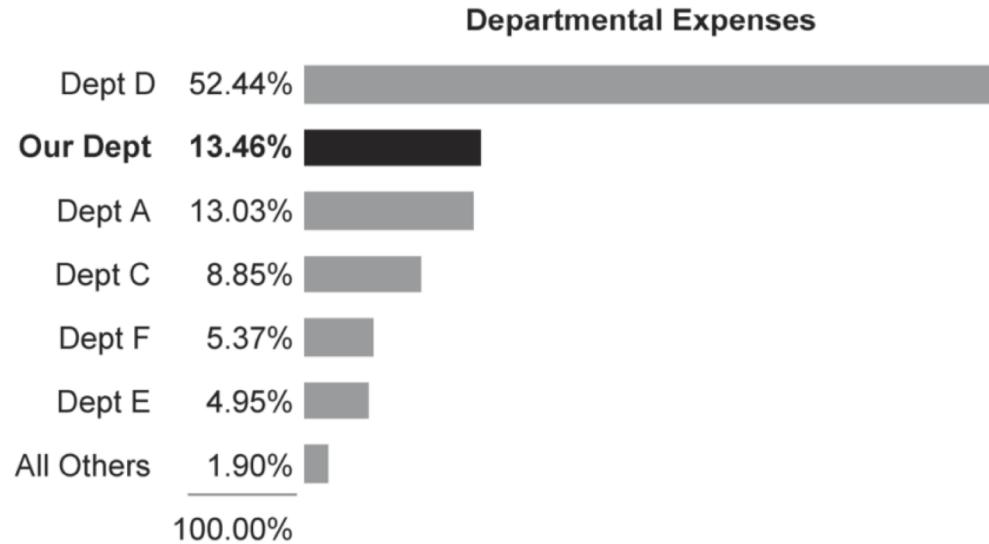
Peter Aldhous

<http://paldhous.github.io/ucb/2018/dataviz/index.html>



Stephen Few

Show me the Numbers



2005 Sales Revenue (USD)

Sales Channel	Q1	Q2	Q3	Q4
Direct	383,383	403,939	437,373	538,583
Indirect	283,733	283,833	257,474	258,474
Total	667,116	687,772	694,847	797,057

2005 Sales Revenue (USD)

Sales Channel	Q1	Q2	Q3	Q4
Direct	383,383	403,939	437,373	538,583
Indirect	283,733	283,833	257,474	258,474
Total	667,116	687,772	694,847	797,057



Data Ink vs. Non-Data Ink

The data-ink ratio must be as high as possible.

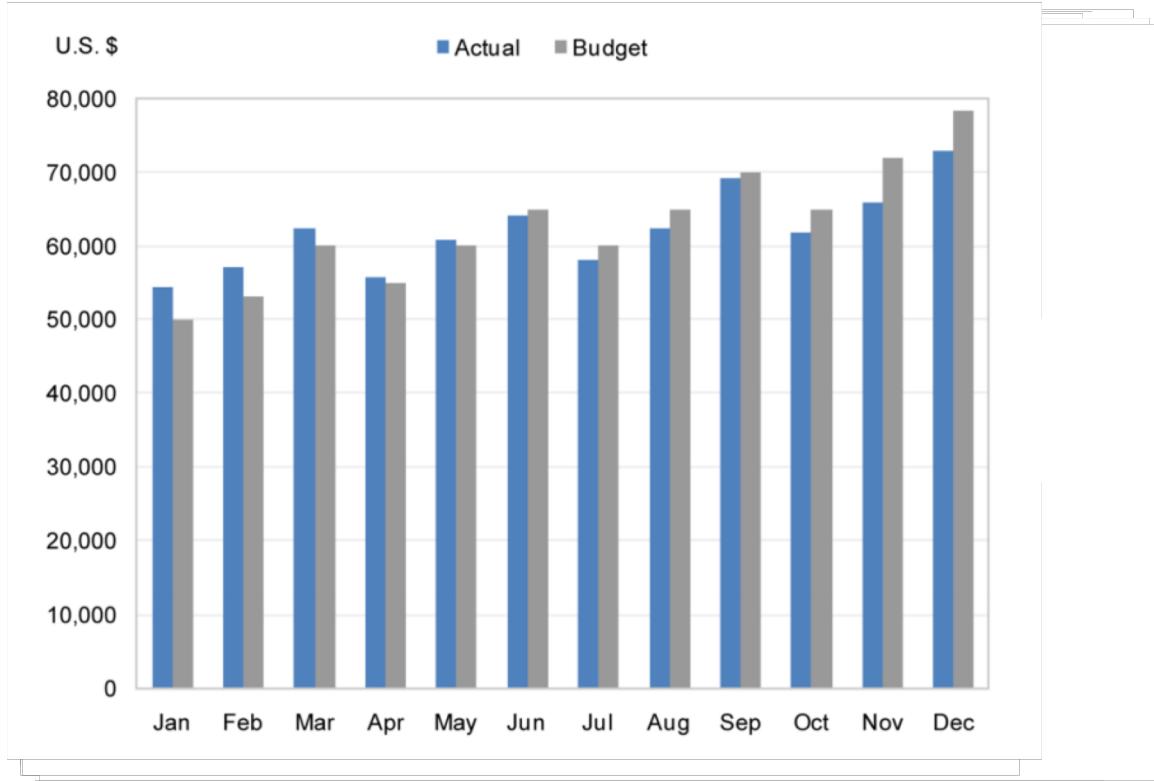
- Tufte

Stephen Few

Show me the Numbers

14





Data Ink vs. Non-Data Ink

The data-ink ratio must be as high as possible.

- Tufte

Stephen Few

Show me the Numbers

15



Expenses Percentage Variance from Budget

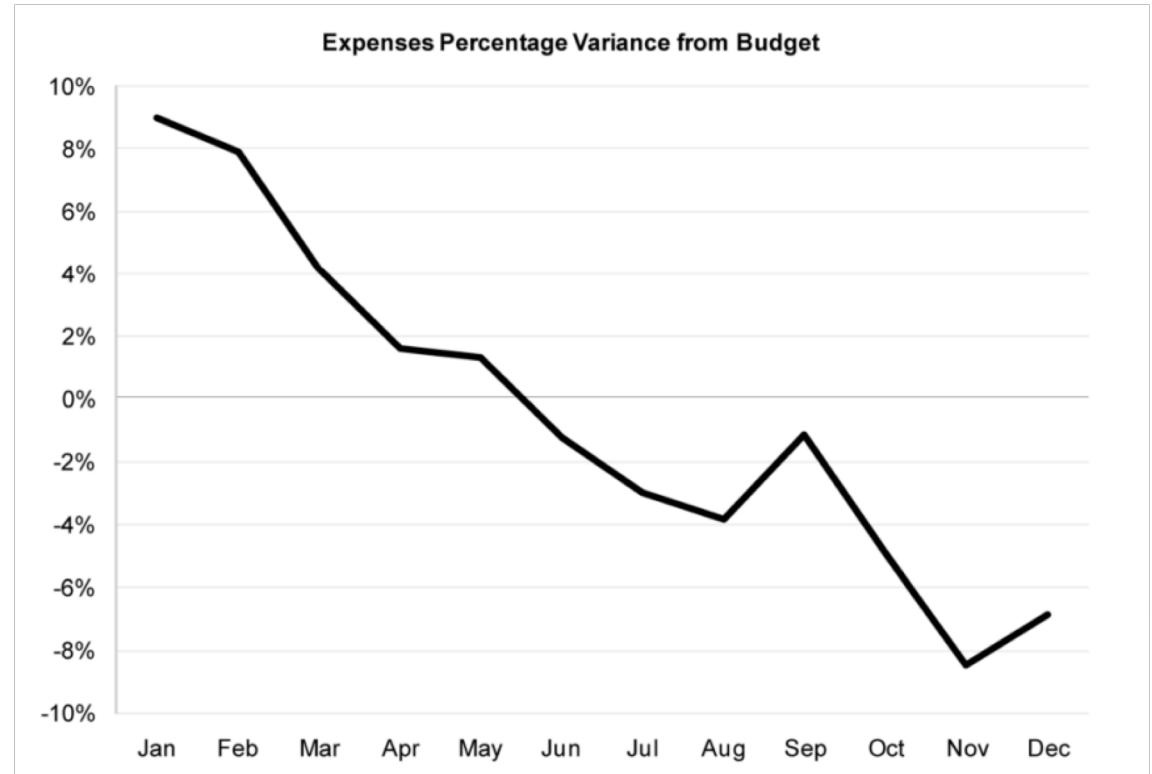
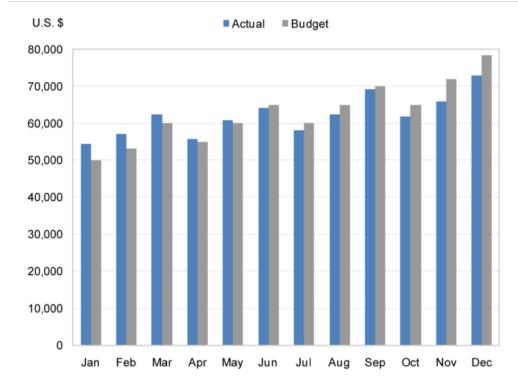
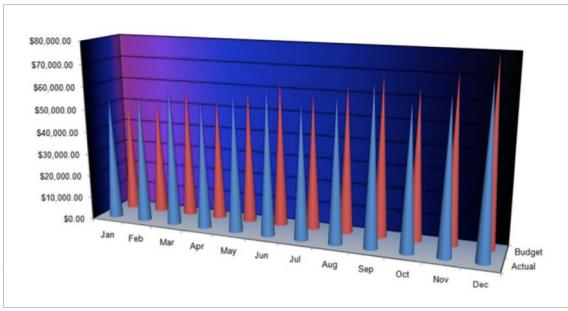


Stephen Few

Show me the Numbers

Interpretation

Realizing the context and the most effective way to present your data.



Above all else, show the Data

Edward R. Tufte

17



Data Type

Numerical
Categorical
Mixed Type

Map Data
Network
Time Series

Distribution

Relationship

Comparison

What do you want to show?

Connection

Composition
(parts of the whole)

Location

Peter Aldhous

<http://paldhous.github.io/ucb/2018/dataviz/index.html>

18



Deviation

Emphasise variations (+/-) from a fixed reference point. Typically the reference point is zero but it can also be a target or a known average. Can also be used to show if something is skewed (positive/heuristic/negative).

Example FT uses
Trade surplus/deficit, climate change

Diverging bar

A simple standard bar chart that can handle negative and positive magnitude values.

Diverging stacked bar

Perfect for presenting survey results which have two opposing sides (disagree/neural/agree).

Spine

Spins a single value into two contrasting components (eg male/female).

Surplus/deficit line

The shaded area of these charts allows a baseline to be shown – either against a baseline or between two series.

XY heatmap

A good way of showing relationships between 2 categories of data (less effective at showing fine differences in amounts).

Bubble

Like a scatterplot but adds additional detail by using bubbles or circles according to a third variable.

Lollipop

Lollipop draws more attention to the data value than standard bubble charts can also show rank and value effectively.

Bump

Effective for showing changing proportions over time or dates. For large datasets, consider grouping lines using colour.

Cumulative curve

A good way of showing frequency distributions, i.e. a curve is always cumulative (frequency x axis leaves a measure).

Frequency polygon

For displaying multiple distributions of data. Like a regular line chart but limited to a maximum of 3 or 4 datasets.

Beeswarm

Use to emphasise individual points in a dataset. Points can be scaled on an additional variable. Best suited to medium-sized datasets.

Correlation

Show the relationship between two or more variables. Its strength may tell you if they tell them otherwise, many readers will assume the relationships there must be causal (i.e. one causes the other).

Example FT uses
Inflation and unemployment, income and life expectancy

Scatterplot

The standard way to show the relationship between two continuous variables, each of which has its own axis.

Column + line timeline

A good way of showing the relationship between an amount (columns) and a rate (line).

Connected scatterplot

Usually used to show how the relationship between 2 variables changes over time.

Dot strip plot

Dots placed in order on a strip are a quick method of laying out ranks across multiple categories.

Barcode plot

Like a dot strip plot, good for displaying all the data at once, they work best when highlighting individual values.

Slope

Perfect for showing how ranks have changed over time or vary between categories.

Boxplot

Summarise multiple dimensions of the data, showing the median (Centre) and range of the data.

Population pyramid

A standard way for showing the age and sex distribution of a population.

For large datasets, consider grouping lines using colour.

Candlestick

Usually focused on day-to-day activity, these charts show the price movement and high/low points of each day.

Facet chart (projection)

Used to show the uncertainty in future projections. If daily, this gives the further forward to projection.

Ranking

Use where an item's position in an ordered list is more important than its absolute or relative value. Don't be afraid to highlight the points of interest.

Example FT uses
Wealth, deprivation, legislature, constituency election results

Ordered bar

Standard bar charts display the ranks of variables more easily when sorted into order.

Histogram

The standard way to show a statistical distribution between the gaps between bins small to highlight the shape of the data.

Line

The standard way to show a changing time series. Consider markers to represent data points.

Dot plot

A simple way of showing the count or range (minimum or maximum) of data across multiple categories.

Dot strip plot

Good for showing individual values in a distribution, can be a problem if there are too many dots have the same value.

Column + line timeline

A good way of showing the relationship between an amount (columns) and a rate (line).

Barcode plot

Like a dot strip plot, good for displaying all the data at once, they work best when highlighting individual values.

Slope

Good for showing as long as the data can be simplified into 2 or 3 points without missing a key part of story.

Area chart

Summarise multiple dimensions of the data, showing the median (Centre) and range of the data.

Violin plot

Similar to a box plot but more effective with complex distributions (and that cannot be summarised with simple averages).

Lollipop

Lollipop draws more attention to the data value than standard bubble charts can also show rank and value effectively.

Population pyramid

A standard way for showing the age and sex distribution of a population.

Distribution

Show values in a dataset or frequency of occurrence. Its mode that tells you if the data is more important than its absolute or relative value. Don't be afraid to highlight the points of interest.

Example FT uses
Income distribution, population, constituency election results

Example FT uses

Income distribution, population, constituency election results

Change over Time

Give emphasis to changing trends. These can be short (intra-day), movements or extended series (e.g. annual). Choosing the correct time period is important to provide suitable context for the reader.

Example FT uses
Corporate acquisition, market capitalisation, volumes in general

Example FT uses

Corporate acquisition, market capitalisation, volumes in general

Magnitude

Show up comparisons. These can be relative (first having to see larger/bigger) or absolute to see the raw numbers. Usually these are number of examples, bars, dollars or people rather than a calculated rate or per cent.

Example FT uses
Fiscal budgets, company structures, national election results

Example FT uses

Fiscal budgets, company structures, national election results

Part-to-whole

Show how one entity can be broken down into its component elements. If the reader's interest is solely in the size of the components, consider a magnitude-type chart instead.

Example FT uses
Population density, natural resource locations, natural disaster risk/impact, settlement areas, variation in election results

Example FT uses

Population density, natural resource locations, natural disaster risk/impact, settlement areas, variation in election results

Spatial

Above from location maps only used when precise locations or geographical patterns in the data are more important to the reader than anything else.

Example FT uses
Movement of funds, trade, migrations, lawsuits, information, relationship graphs.

Example FT uses

Movement of funds, trade, migrations, lawsuits, information, relationship graphs.

Sankey

Shows changes in flows from one condition to another, allowing for tracing the eventual outcome of a complex process.

Example FT uses
Waterfall

Example FT uses

Waterfall

Chord

Designed to show the flow of data through a process, typically budget. Chord (red winnner) in a matrix.

Example FT uses
Chord

Example FT uses

Chord

Network

A complex but powerful diagram which can illustrate connections (red winnner) in a matrix.

Example FT uses
Network

Example FT uses

Network

Equalized cartogram

Converting each unit on a map into an equal area and equal-sized shapes, good for representing voting regions with varying populations.

Example FT uses
Equalized cartogram

Equalized cartogram

Equalized cartogram

Scaled cartogram

Scaling and stretching a map so that each area is sized according to a particular value.

Example FT uses
Scaled cartogram

Scaled cartogram

Scaled cartogram

Dot density

Used to show the location of individual events. Make sure to annotate any patterns the reader should see.

Example FT uses
Dot density

Dot density

Dot density

Heat map

Grid-based data values mapped with an intensity scale. As choropleth map – but not shaped to an administrative unit.

Example FT uses
Heat map

Heat map

Heat map

Venn

Generally only used for schematic representation.

Example FT uses
Venn

Venn

Venn

Parallel coordinates

An alternative to radar charts – again, the variables whose values are important are highlighted.

Example FT uses
Parallel coordinates

Parallel coordinates

Parallel coordinates

Waterfall

Can be useful for showing part-to-whole relationships where some of the components are negative.

Example FT uses
Waterfall

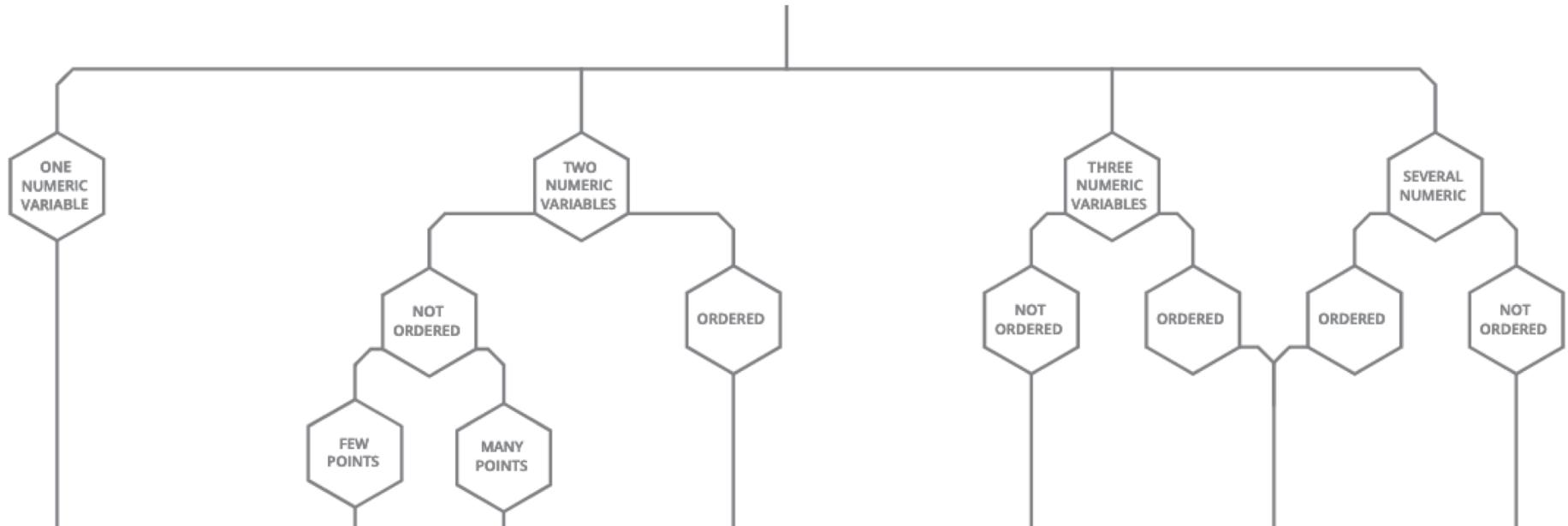
Waterfall

Waterfall

Page 19

<http://ft.com/vocabulary>

Visual vocabulary



from Data to Viz

<https://www.data-to-viz.com/>

20