# Introduction to Data Analysis

# DATA VISUALIZATION BASICS

Patrick Boily

Data Action Lab | uOttawa | Idlewyld Analytics

[pboily@uottawa.ca](mailto:pboily@uottawa.ca)

"Discovery is no longer limited by the collection and processing of data, but rather management, analysis, and visualization."

@DamianMingle

# INFOGRAPHICS

Created for **story-telling** purposes (**subjective**)

Intended for a **specific** audience

**Self-contained** and discrete

Graphic design aspect is key

Cannot usually be re-used with other data

Can incorporate **unquantifiable** information
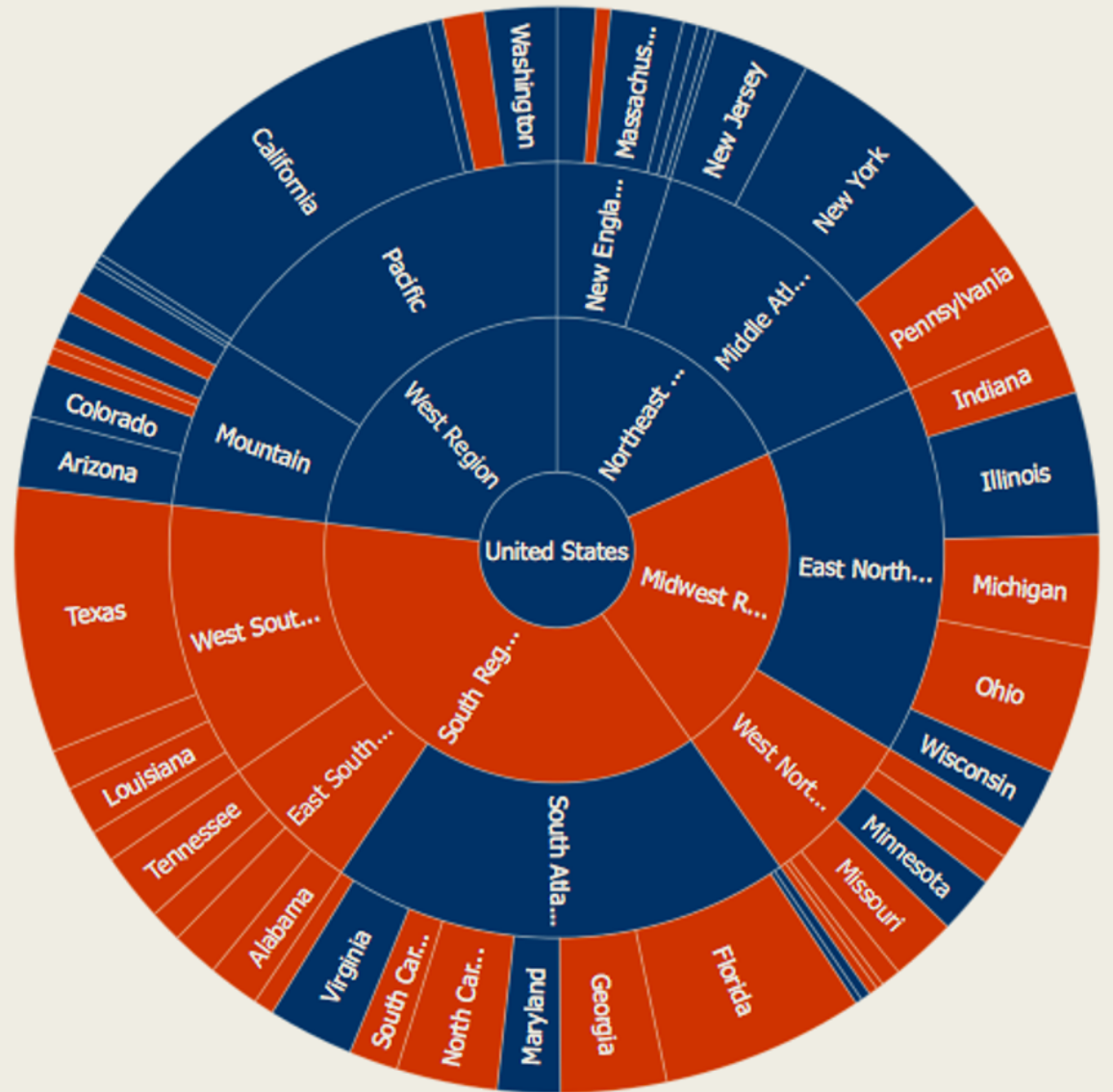
# DATA VISUALIZATION

A **method**, as well as an item (**objective**)

Typically focuses on the **quantifiable**

Used to make sense of the data or to make it **accessible** (datasets can be massive and unwieldy)
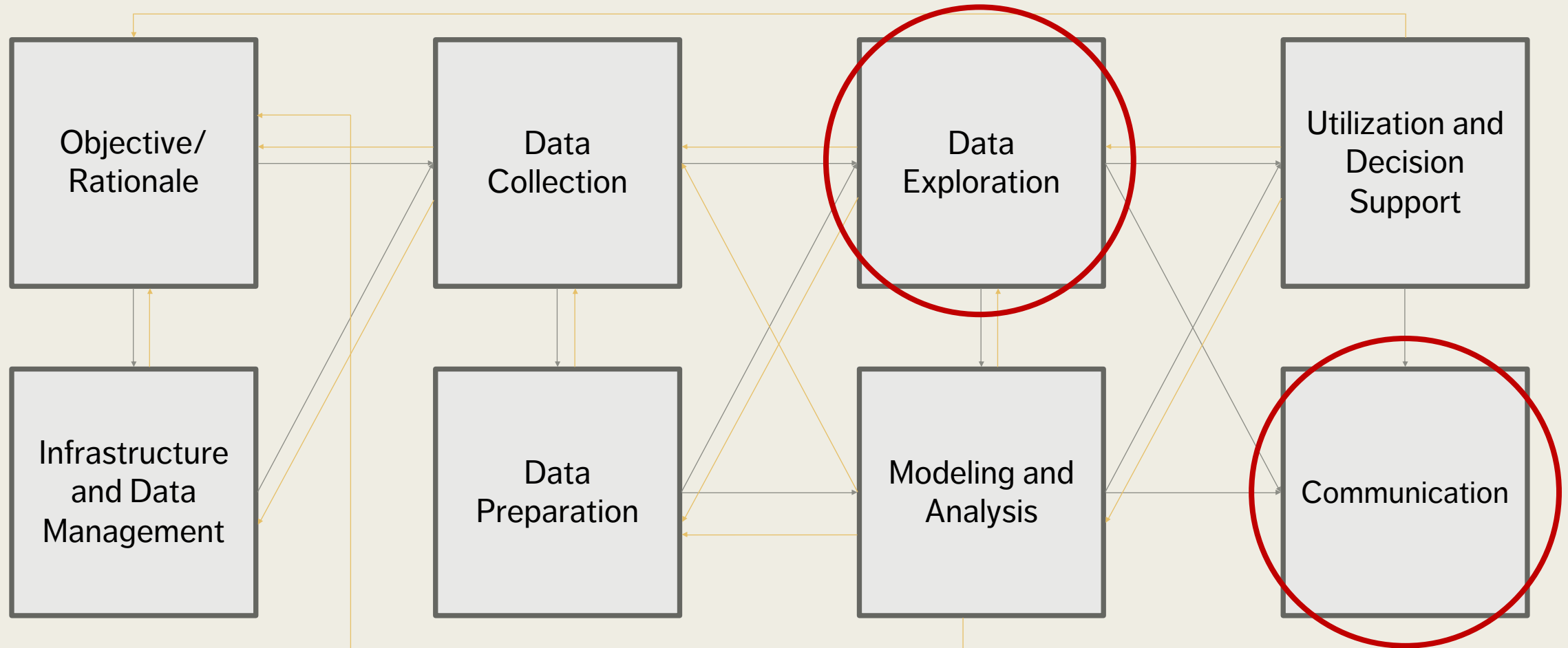
May be generated automatically

The look and feel are less important than the **insights conveyed** by the data
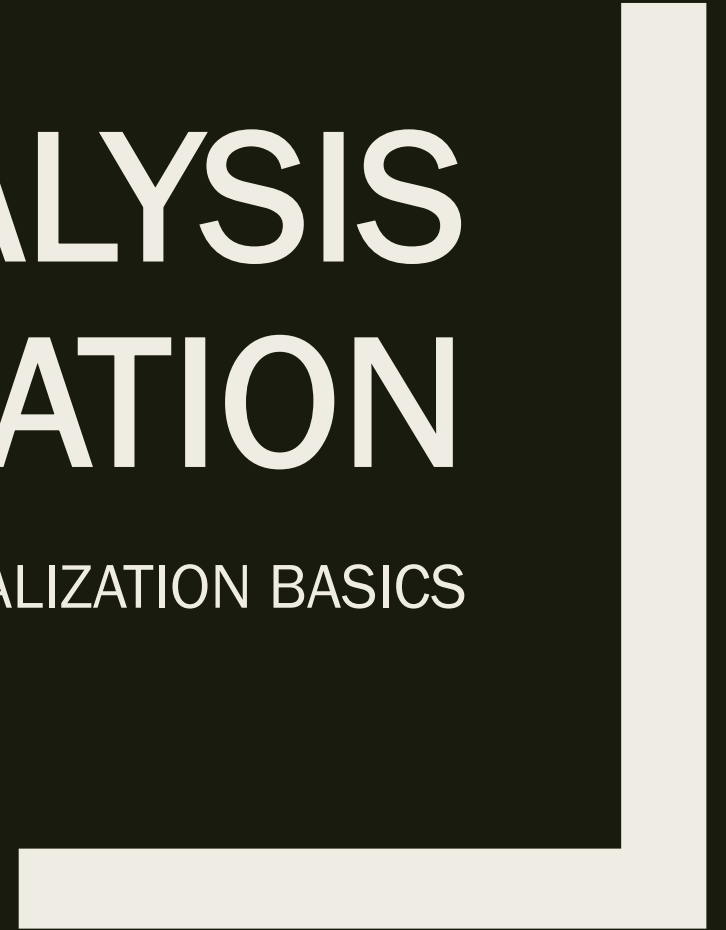
Size Population    Color Median Household Income    ■ Low Income ■ High Income

# THE (MESSY) ANALYSIS PROCESS

# PRE-ANALYSIS DATA VISUALIZATION

## DATA VISUALIZATION BASICS

# SOME BASIC QUESTIONS

What system does your data represent – objects, attributes, relationships?

**How** does it represent this system – i.e. the data model?

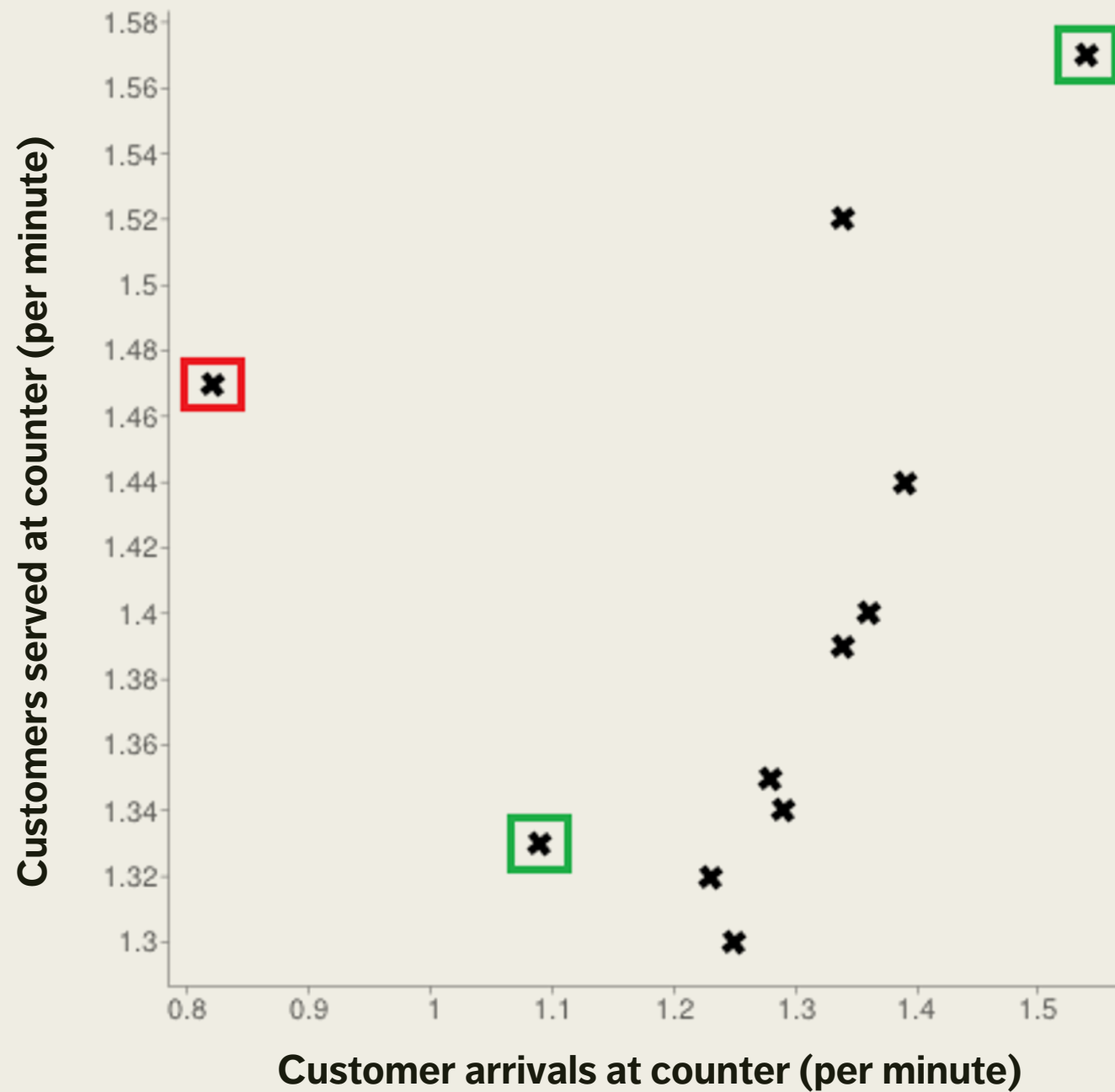Who made this dataset? When? For what purpose?

Assuming a flat file – what do the rows represent? What do the columns represent?
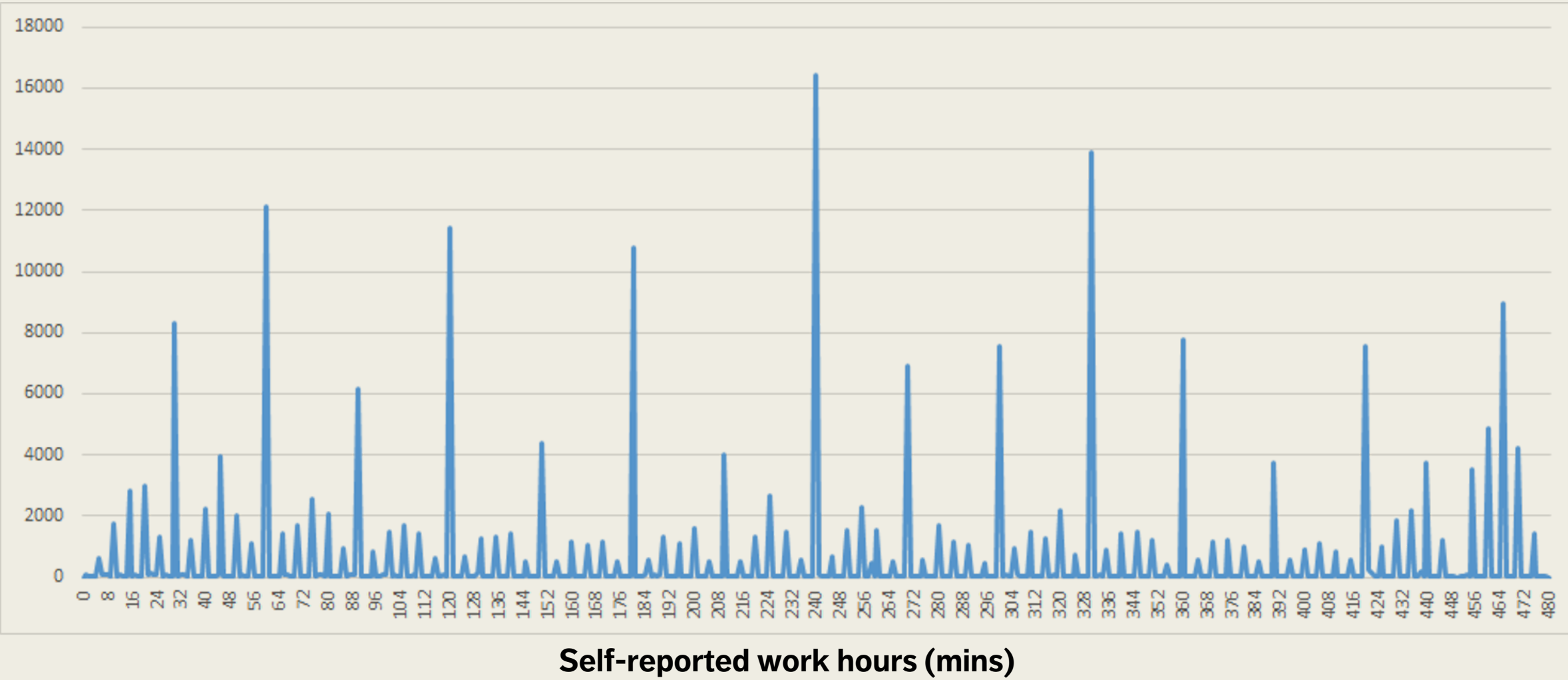
Do you even have enough information (e.g. **metadata**) to answer these questions? Where can you find more information?

# PRE-ANALYSIS USE

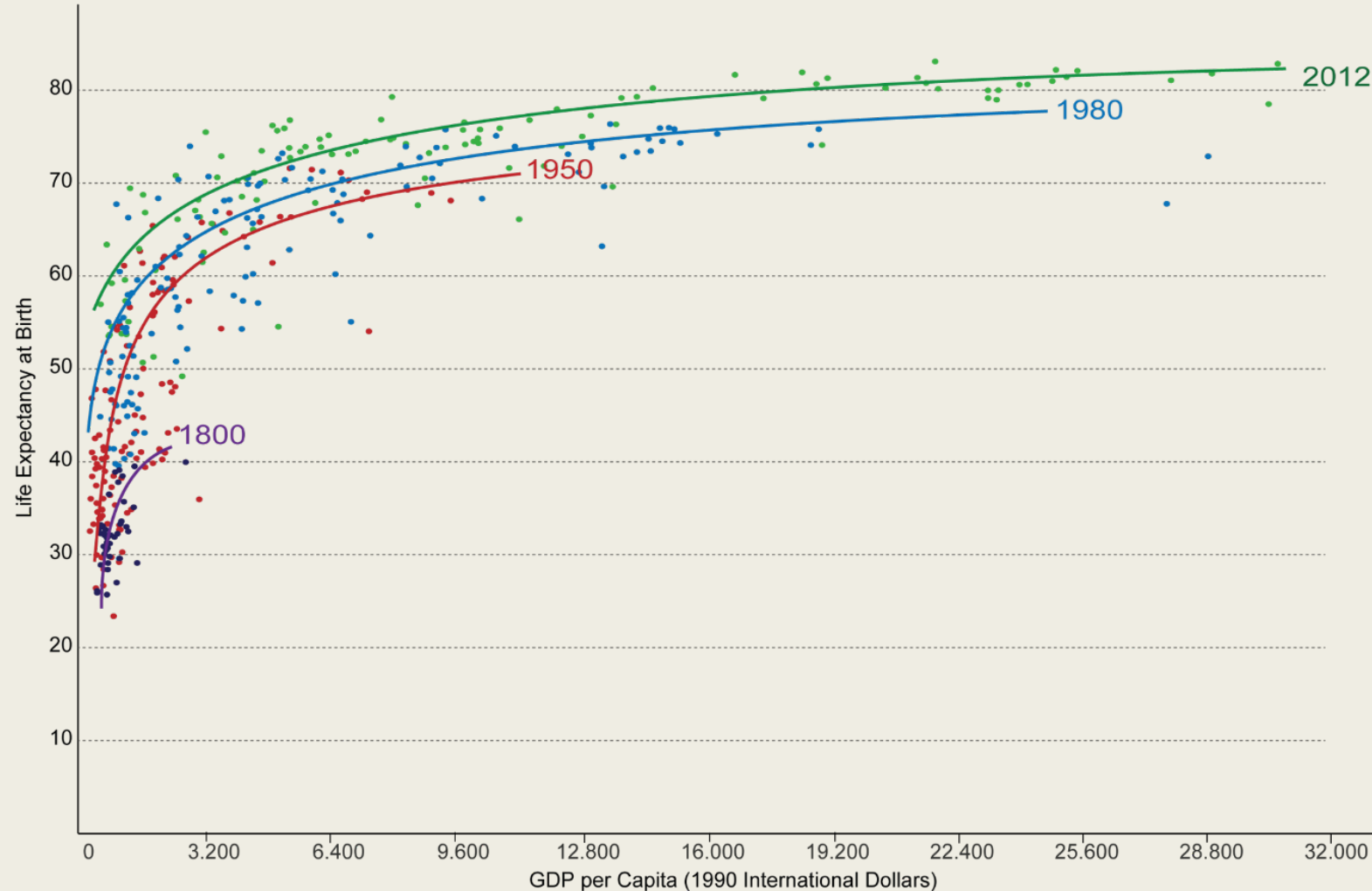Data visualization can be used to set the stage for analysis:

- **detecting anomalous entries**
  *invalid entries, missing values, outliers*

- **shaping the data transformations**
  *binning, standardization, Box-Cox transformations, PCA-like transformations*

- **getting a sense for the data**
  *data analysis as an art form, exploratory analysis*

- **identifying hidden data structure**
  *clustering, associations, patterns informing the next stage of analysis*

**Self-reported work hours (mins)**

**Our World in Data**

## Life Expectancy vs. GDP per Capita from 1800 to 2012 – by Max Roser

GDP per capita is measured in International Dollars. This is a currency that would buy a comparable amount of goods and services a U.S. dollar would buy in the United States in 1990. Therefore incomes are comparable across countries and across time.

*Life Expectancy at Birth (y-axis) vs. GDP per Capita (1990 International Dollars) (x-axis), with trend lines for 1800, 1950, 1980, and 2012.*

Data sources: Data on life expectancy are from Gapminder.org; data on GDP per capita are from the 'New Maddison Project Database'.
The interactive data visualisation is available at OurWorldinData.org. There you find the raw data and more visualisations on this topic.
Licensed under CC-BY-SA by the author Max Roser.

This graph displays the **correlation** between life expectancy and GDP per capita.

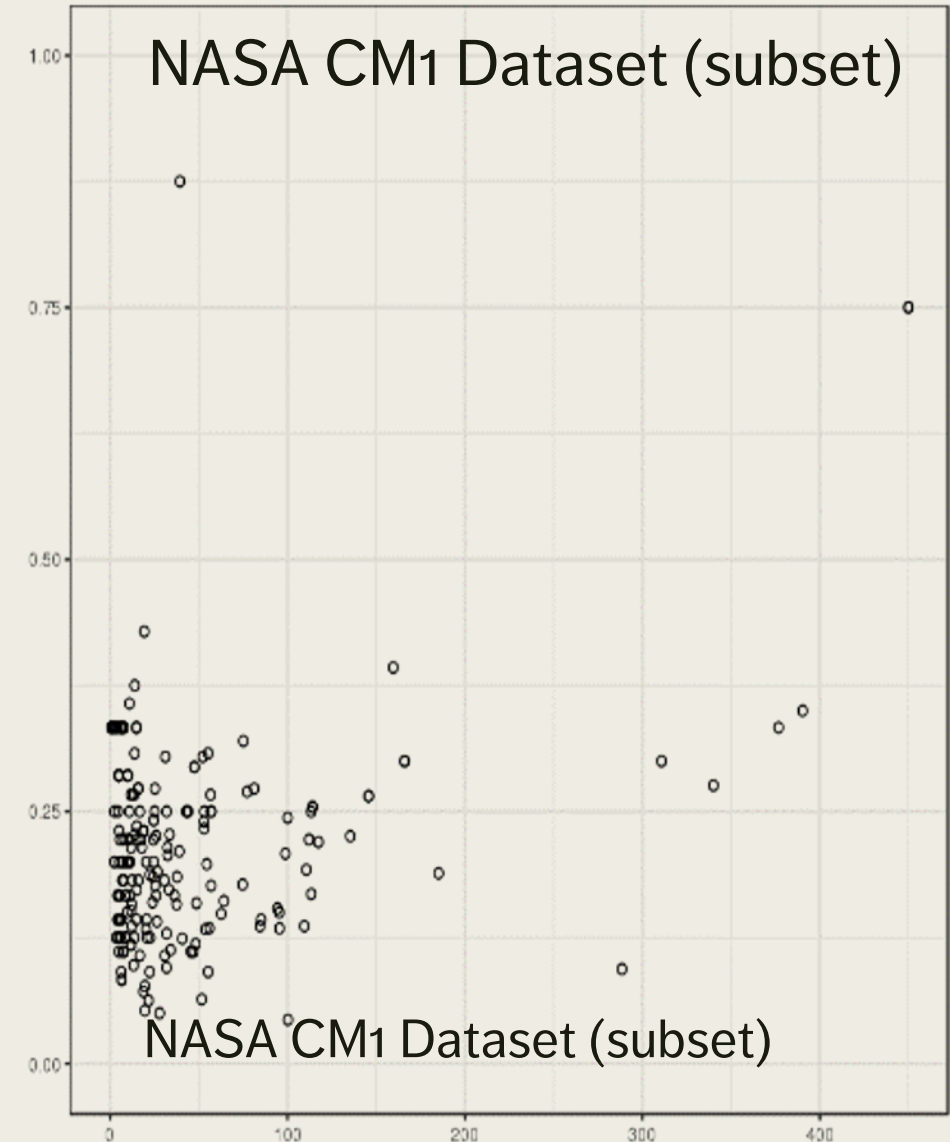Countries with higher GDP have, **in general**, a higher life expectancy.

The relationship seems to follow a **logarithmic trend**: the unit increase in life expectancy per unit increase in GDP decreases as GDP per capita increases.
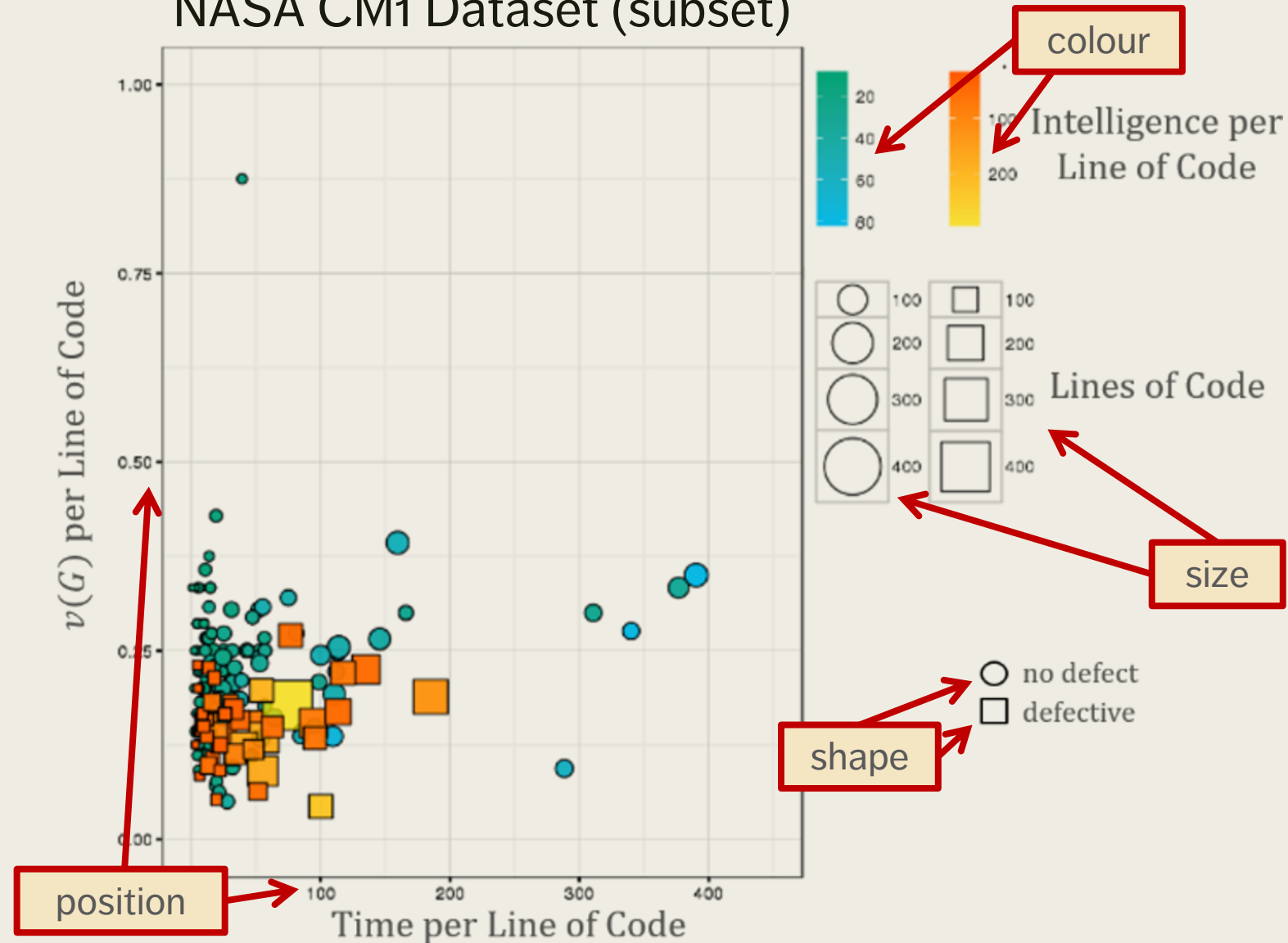
# REPRESENTING OBSERVATIONS

2 variables can be represented by position in the plane.

**Additional factors** can be depicted through:

- size
- color
- value
- texture
- line orientation
- shape
- (motion?)

NASA CM1 Dataset (subset)

NASA CM1 Dataset (subset)

NASA CM1 Dataset (subset)

# WORKHORSE VISUALIZATIONS

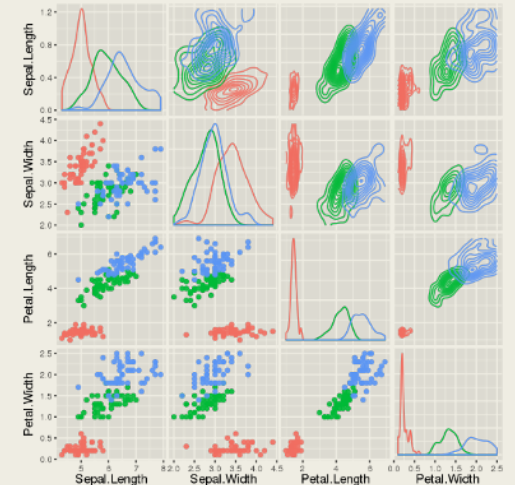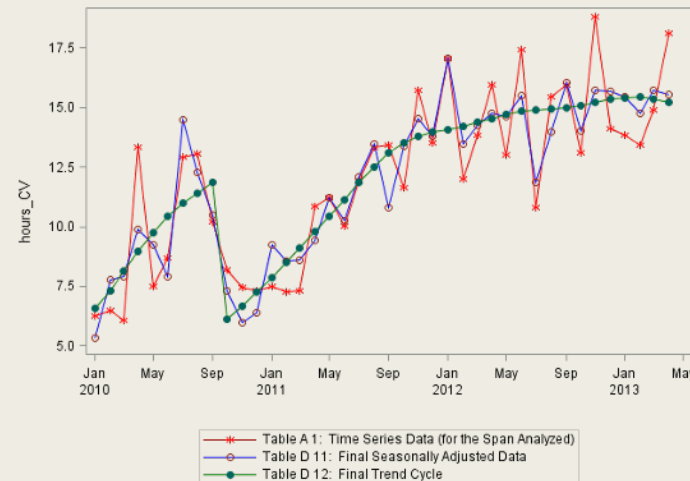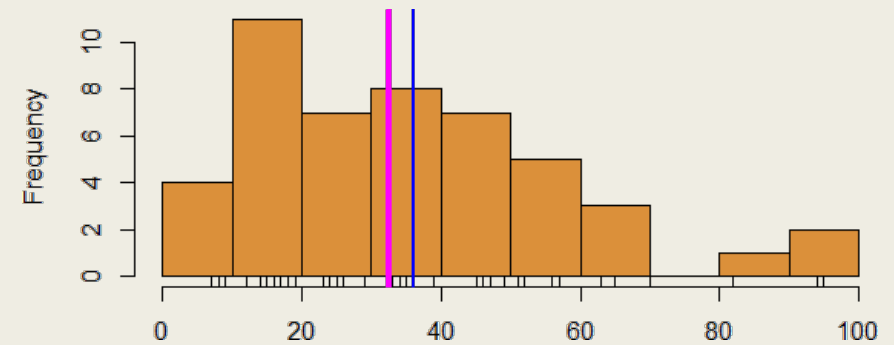Line Chart/Rug Chart/Number Line
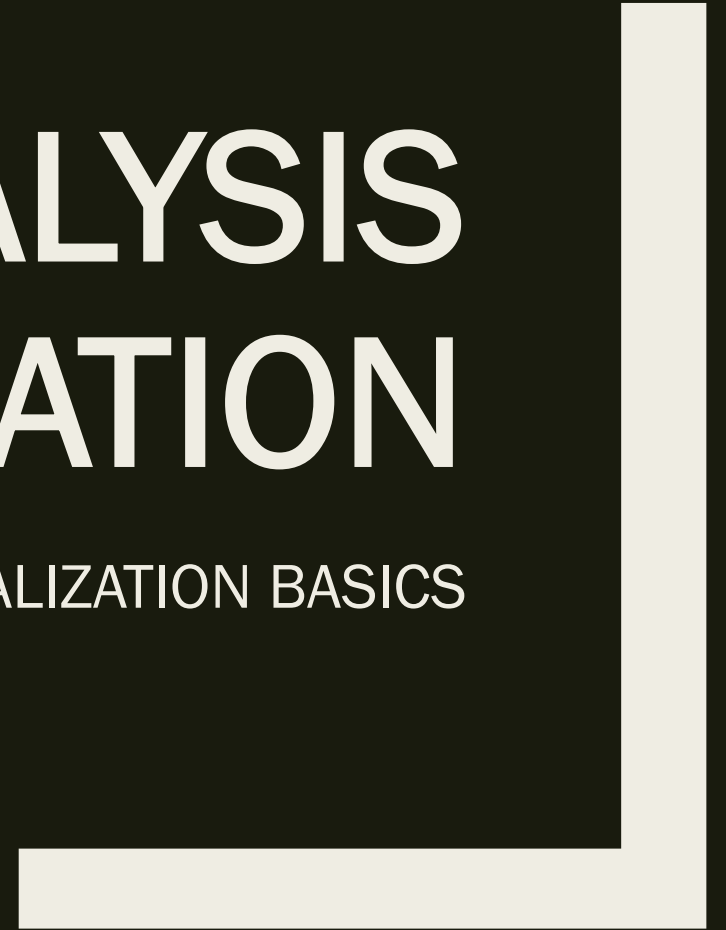
Histogram

Line Graph

Boxplot

Bar Chart

Scatterplot

# POST-ANALYSIS DATA VISUALIZATION
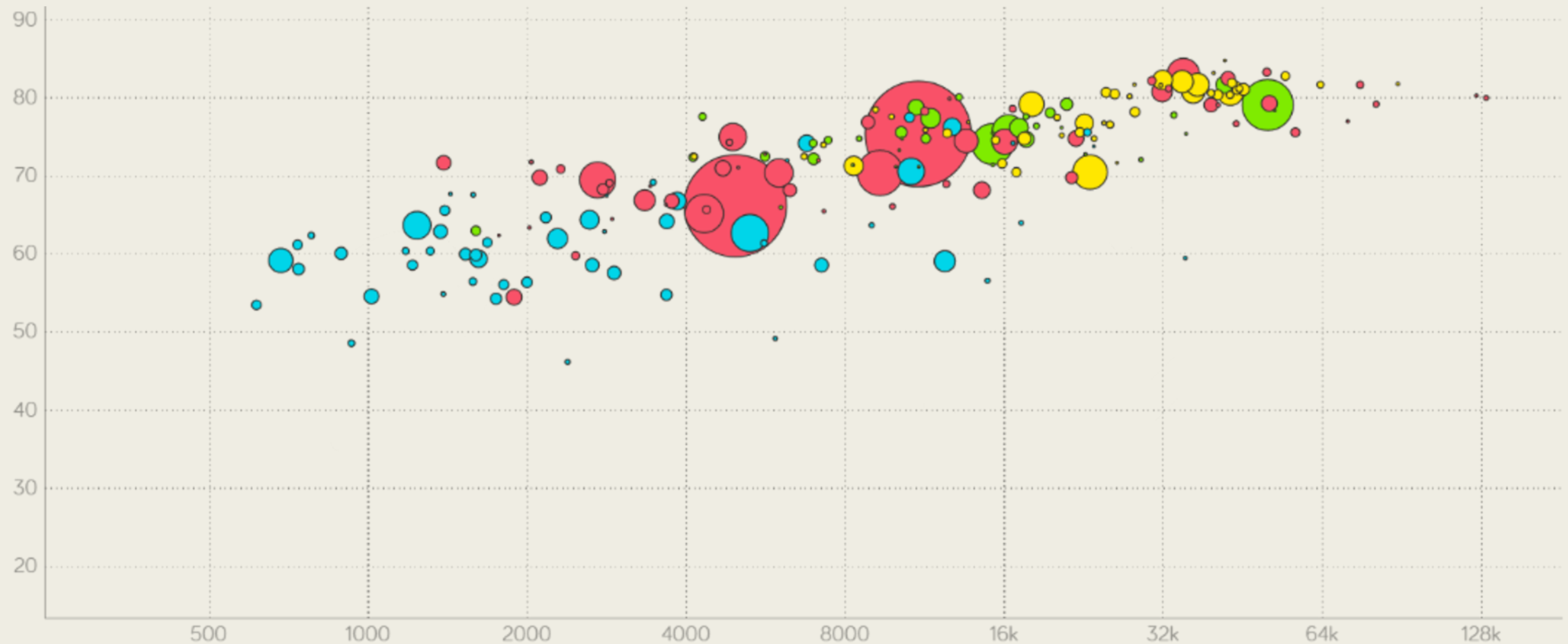
## DATA VISUALIZATION BASICS
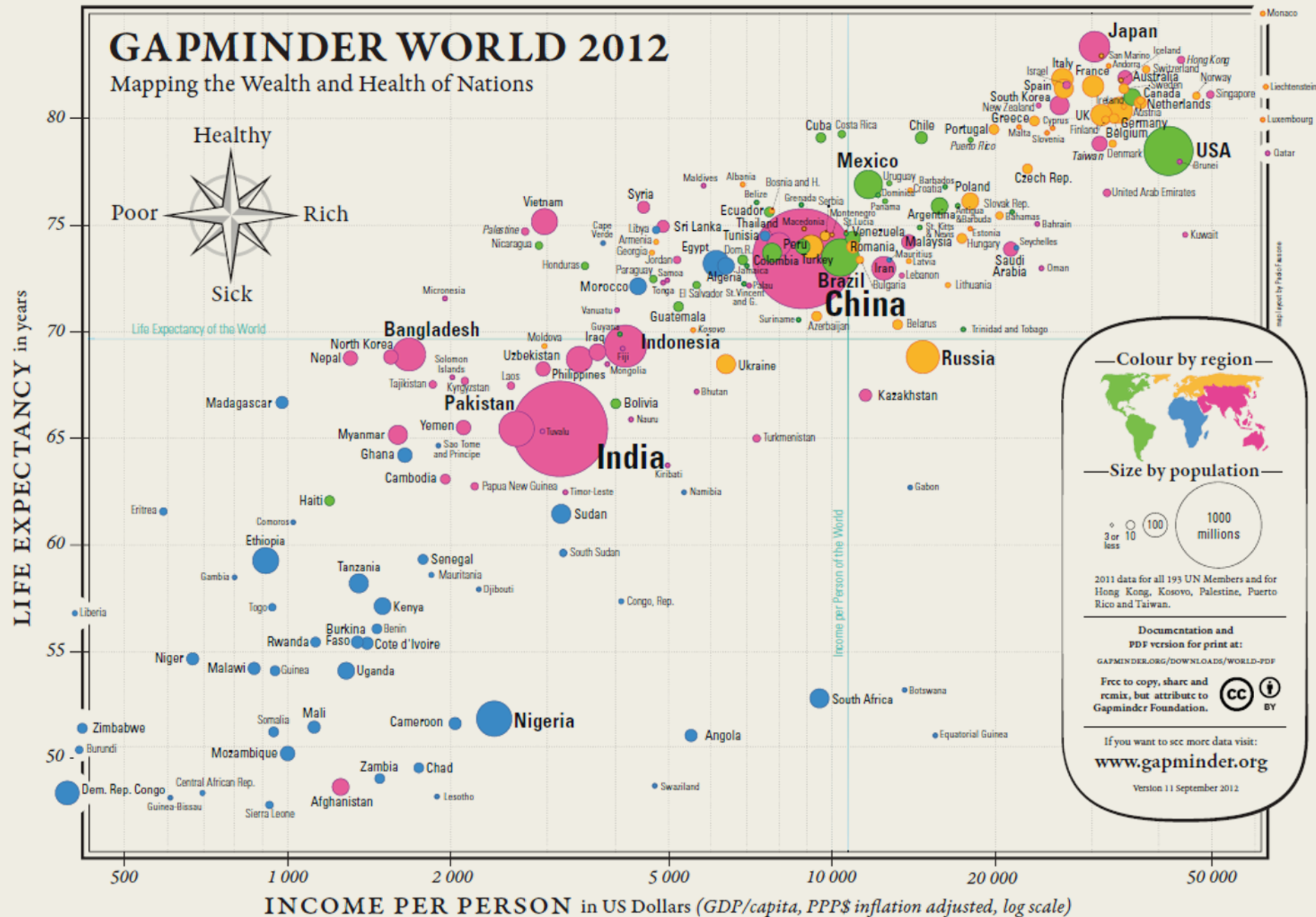
# FUNDAMENTAL PRINCIPLES OF ANALYTICAL DESIGN

**Reasoning and communicating** our thoughts are intertwined with our lives in a causal and dynamic multivariate Universe.
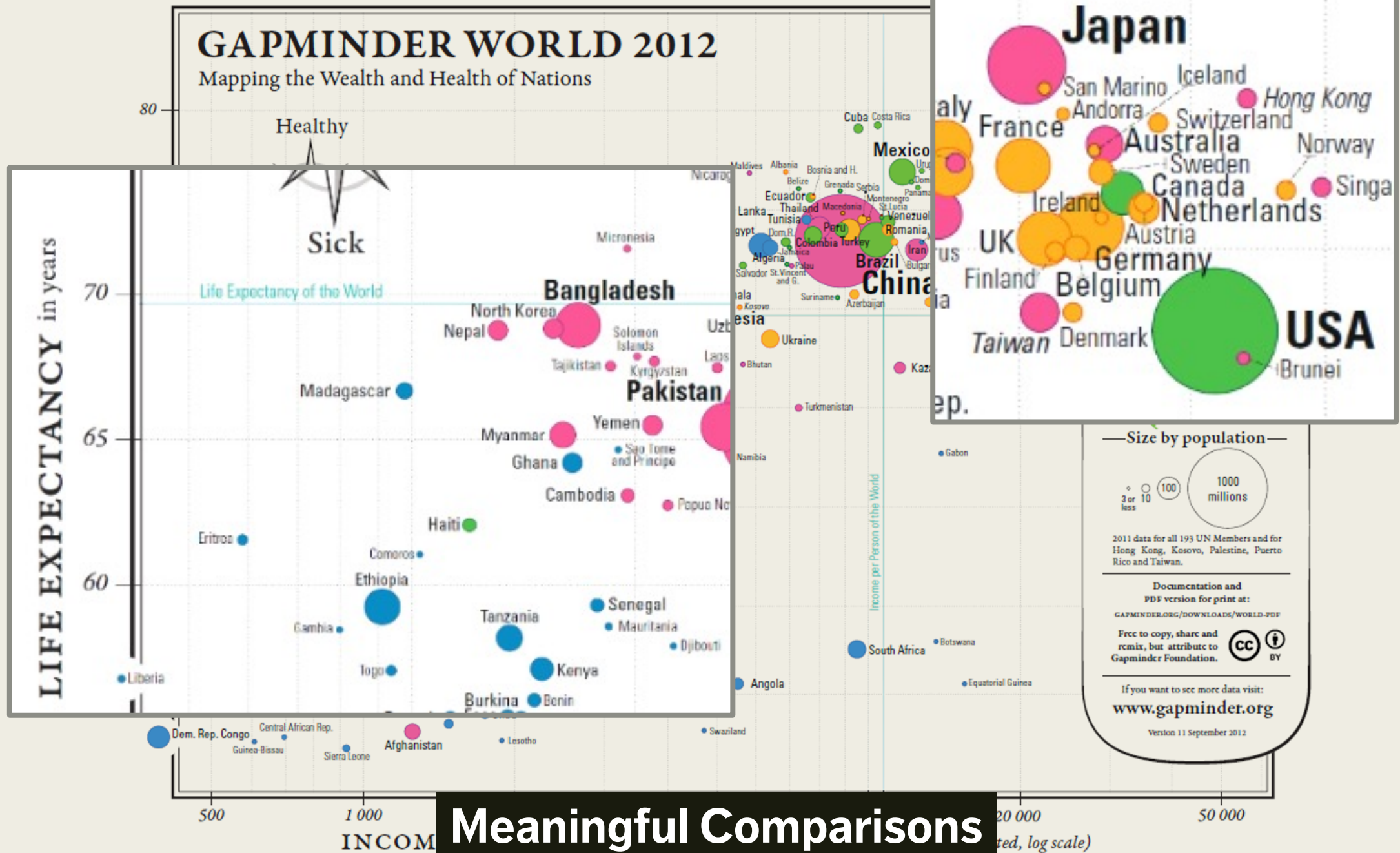
**Symmetry** to visual displays of evidence: consumers should be seeking exactly what producers should be providing, namely
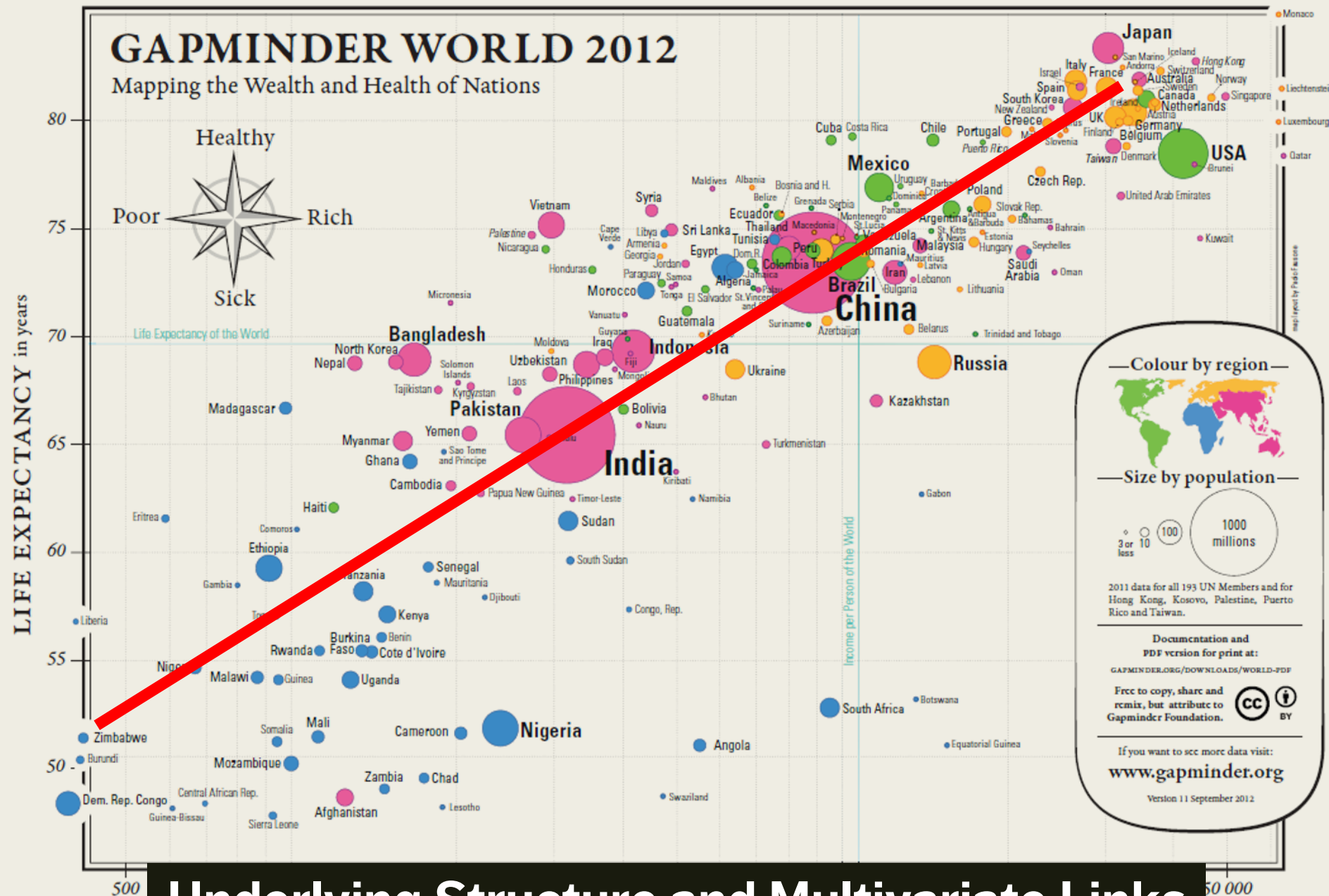
- meaningful comparisons
- causal networks and underlying structure
- multivariate links
- integrated and relevant data
- honest documentation
- primary focus on content

**Non-Integrated Data**

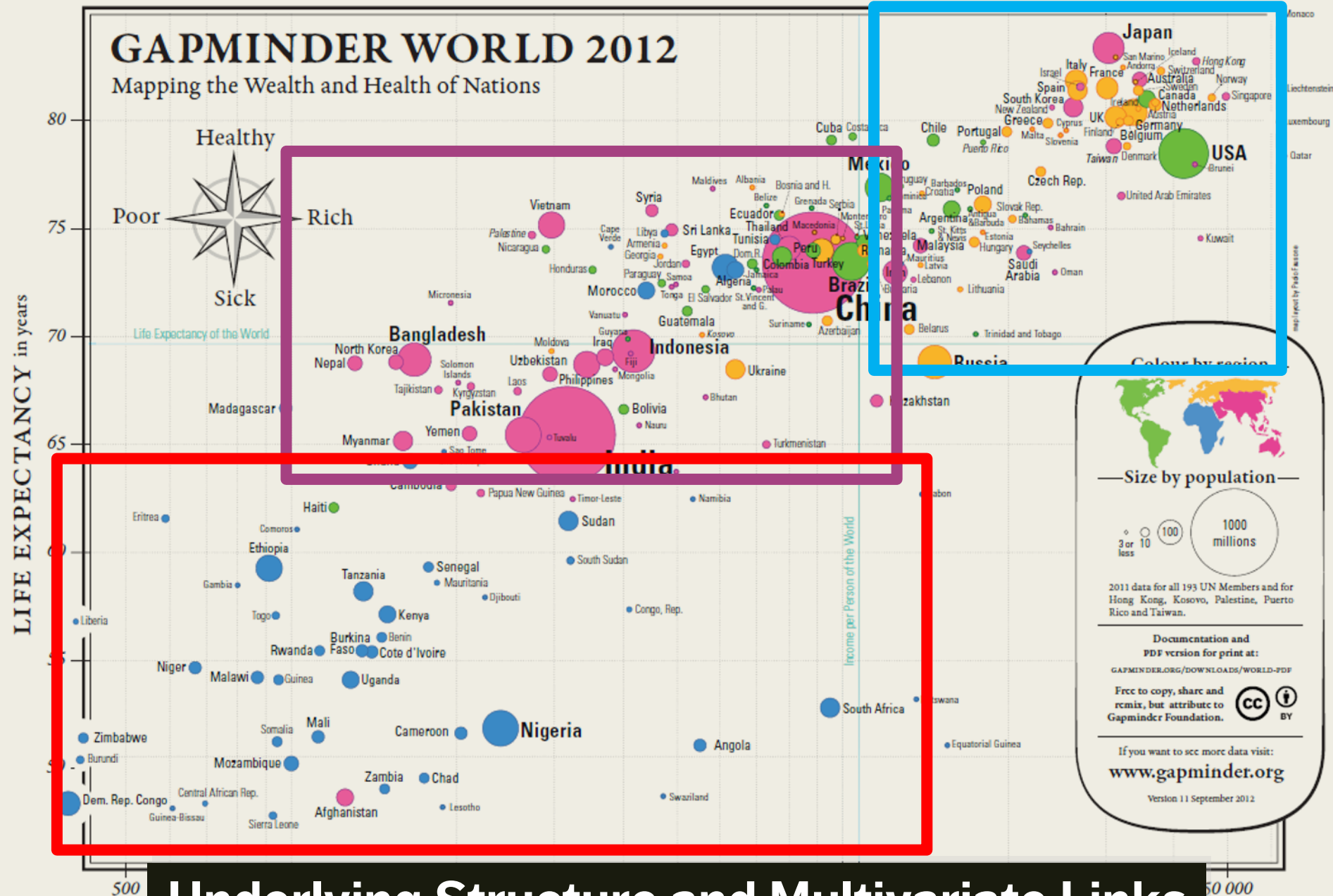**Meaningful Comparisons**

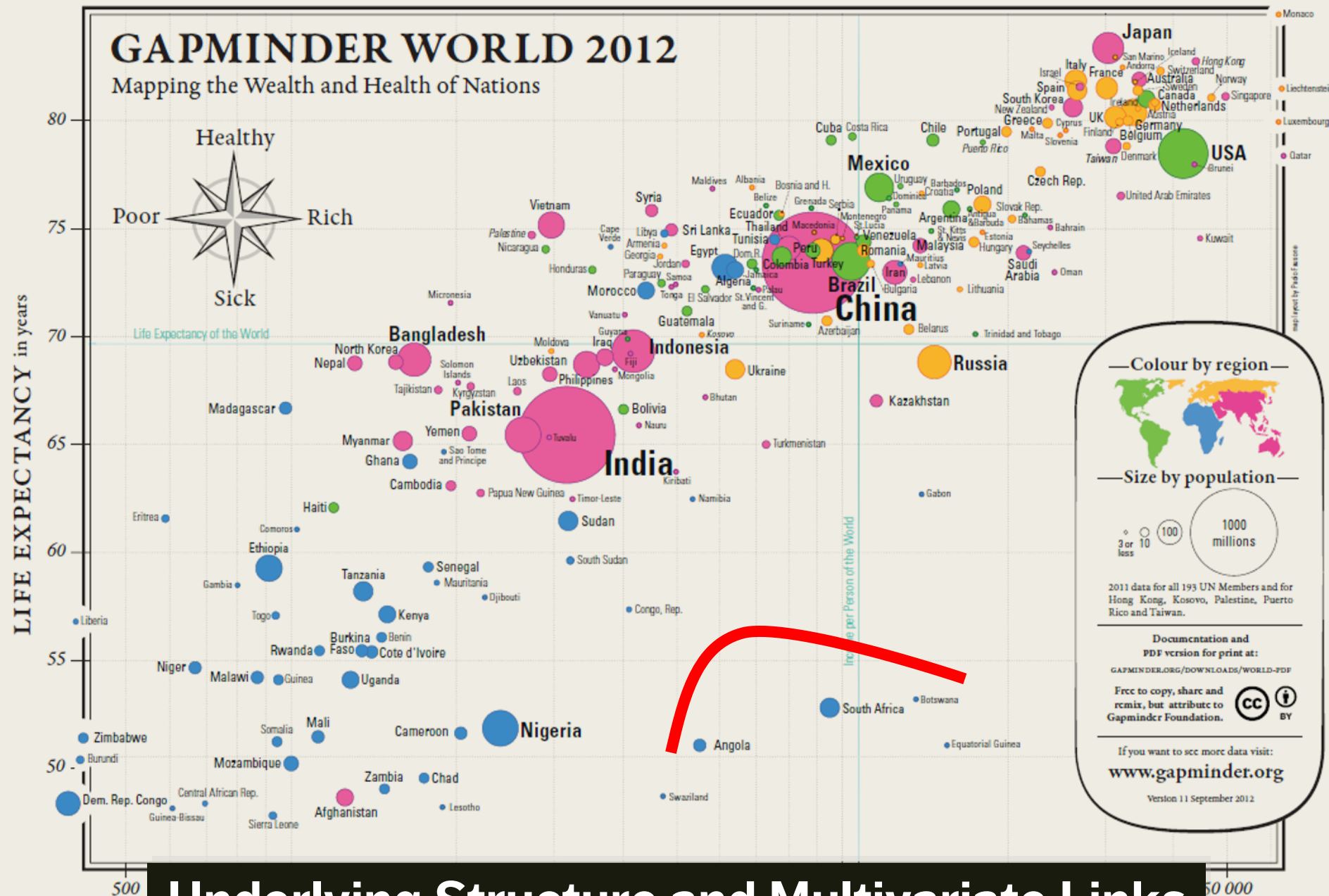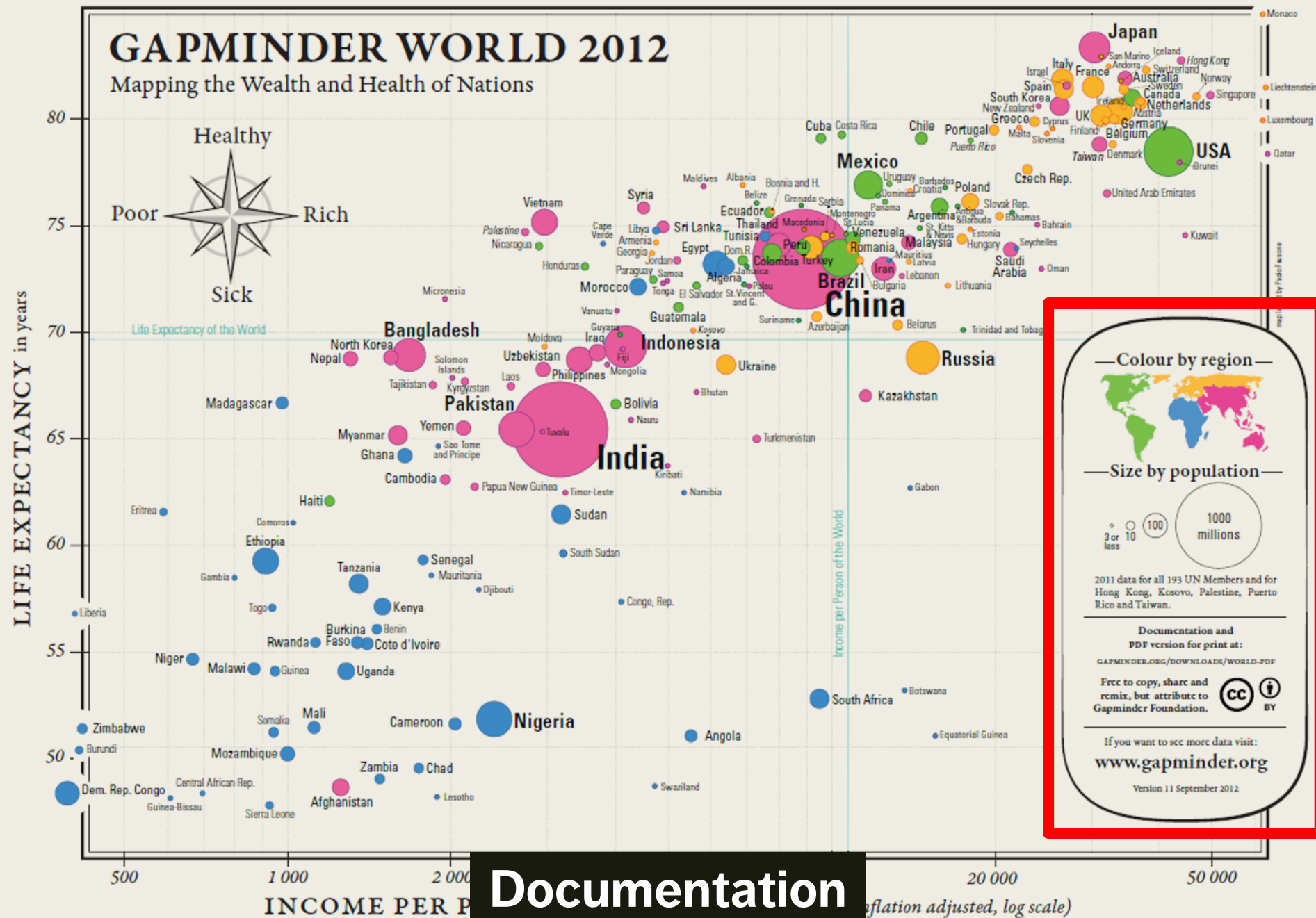**Underlying Structure and Multivariate Links**

**Underlying Structure and Multivariate Links**

**Underlying Structure and Multivariate Links**

# BASIC RULES

**1. Check the data**
outliers, spikes, anomalies

**2. Explain encoding**
don't assume the reader knows what everything means



**3. Label axes**
knowing the scale is important

# BASIC RULES

**4. Include units**
eliminate the need for guesswork

**5. Keep your geometry in check**
circles and 2D shape are sized by area, bar charts by length

**6. Include your sources**
protect yourself, and let those who want to dig deeper do so

**7. Consider your audience**
a poster can be wordy, a presentation should be minimalist

3.5 radius

3.5 area

Cancer cases

Incidence
Mortality

Eggplant consumption, per week

# PRESENTING ANALYSIS RESULTS

Graphics should be clear and engaging.

Not every pretty picture tells a story, but if a story can't be told with pretty pictures, perhaps it's time to re-think the story…

Graphical representation techniques appear regularly – it's too early to tell which ones will stand the test of time.

Don't be afraid to try something new if it helps **convey the message**.

# A WORD ABOUT ACCESSIBILITY

Charts cannot usually be translated to Braille. Describing the features and emerging structures in a visualization is a possible solution... **if they can be spotted.**

Analysts must produce clear and meaningful visualizations, but they must also describe them and their features in a fashion that allows all to "see" the insights. This requires analysts to have "seen" all the insights, which is not always possible.

**Conditions:** colourblindness, low vision, motor impairment, cognitive disability, ADHD, etc.

**Best Practices:** high contrast text/elements, zoom/magnifications, keyboard navigation, assistive design, short summaries, undo/redo functionality, etc. [F. Elavsky]

# A WORD ABOUT ACCESSIBILITY

**Data Perception:**

- texture-based representations
- text-to-speech
- sound/music
- odor-based or taste-based representations (?!?)

**Sonifications:**

- [TRAPPIST Sounds : TRAPPIST-1 Planetary System Translated Directly Into Music](#)
- [Listening to data from the Large Hadron Collider, L. Asquith](#)

# TAKE-AWAYS

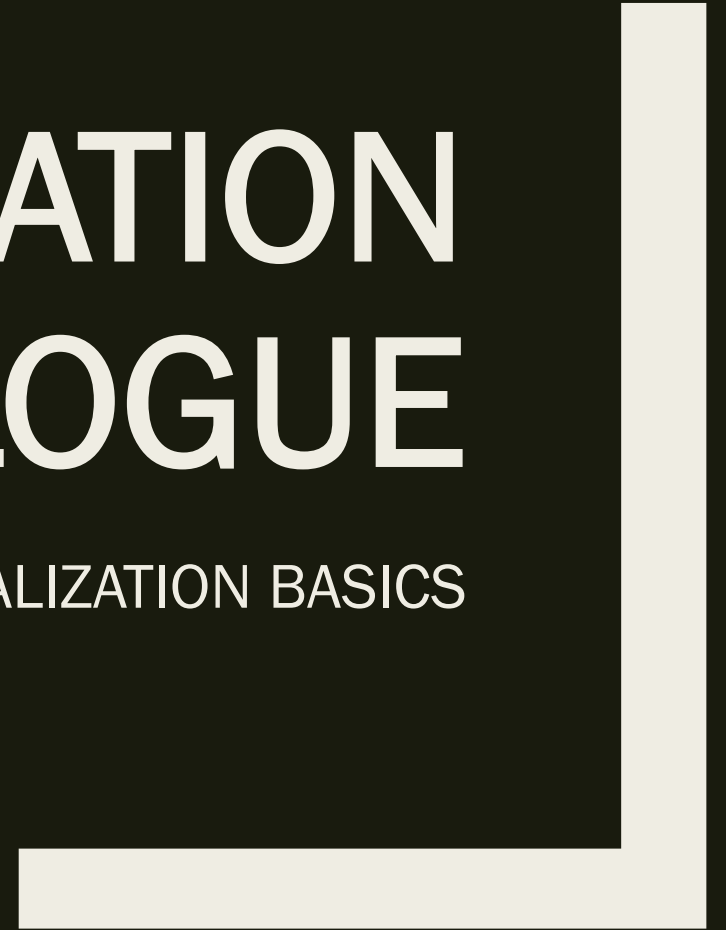**Is the point getting across?** Integrated data helps convey the message.

Adding design elements can enhance our understanding of the data.

How we spot patterns affect what we get out of data presentations.

Data displays are not just about picking a random visualization method. The result varies depending on the structure of the data and the (combinations of) questions.

# VISUALIZATION CATALOGUE

DATA VISUALIZATION BASICS

# DATA DISPLAYS

With data displays, we try to highlight:

1. a **relationship** (show a connection or correlation between two or more variables);

2. a **comparison** (set some variables apart from others, and display how those two variables interact);

3. a **composition** (collect different types of information that make up a whole and display them together), and

4. a **distribution** (lay out a collection of related or unrelated information to see how it correlates, if at all, and to understand if there's any interaction between the variables).

# A CLASSIFICATION OF CHART TYPES

## Data comparison charts

### Comparison

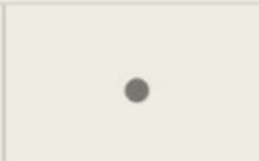**Bars**



**Dot plot**



**Bullet**



**ID Scatterplot**



**Heat map**



**Slope**



**Alert**



### Composition

**Pie**



**Pareto**



**Multidimensional Pie**



### Distribution

**Histogram**



**ID Scatterplot**



**Boxplot**



## Data reduction charts

### Evolution

**Line**



**Horizon**



**Step**



**Connected Scatterplot**



### Relationship

**Scatterplot**



**Connected Scatterplot**



**Bubble**



### Profiling

**Grouped bars**



**Cycle plot**



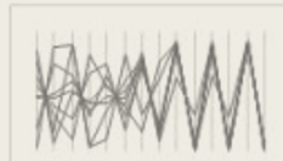**Scatterplot matrix**



**Reorderable matrix**
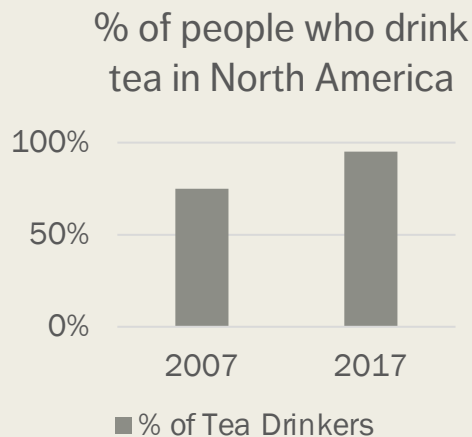


**Horizon**



**Parallel Plot**



**Trellis**



v 0.9

# SIMPLE TEXT AND TABLES

One or two numbers to focus on may help "set the scene" and **draw focus** to an area of the report.

% of people who drink tea in North America

95% of the population drinks tea today compared to 75% in 2007

Tables interact with our **verbal** system (we **read** them):
- used to **compare** values
- audiences will look for **their** rows

Table design needs to blend into background:
- the data should stand out, not the borders
- dense table: use **alternating** row colour

Leverage colour to convey magnitude:
- use **single colour saturation**
- use a legend to remove values

100%

50%

0%

2007  2017

■ % of Tea Drinkers

# TABLES AND TABLE HEATMAPS

| Name | Last Year | This Year |
|------|-----------|-----------|
| Ron | 20 | 30 |
| Fred | 30 | 40 |
| George | 10 | 15 |

| Name | Last Year | This Year |
|------|-----------|-----------|
| Ron | 20 | 30 |
| Fred | 30 | 40 |
| George | 10 | 15 |

| | Last Year | This Year | Next Year | Optimum |
|--------|-----------|-----------|-----------|---------|
| George | 20 | 20 | 20 | 20 |
| Peter | 40 | 35 | 30 | 25 |
| John | 10 | 10 | 5 | 5 |
| Sandra | 25 | 30 | 35 | 40 |

| | Last Year | This Year | Next Year | Optimum |
|--------|-----------|-----------|-----------|---------|
| George | 20 | 20 | 20 | 20 |
| Peter | 40 | 35 | 30 | 25 |
| John | 10 | 10 | 5 | 5 |
| Sandra | 25 | 30 | 35 | 40 |

| | Last Year | This Year | Next Year | Optimum |
|--------|-----------|-----------|-----------|---------|
| George | | | | |
| Peter | | | | |
| John | | | | |
| Sandra | | | | |

# SCATTERPLOTS

Show relationship between 2 variables (**scatterplot**) or 3 variables (**bubble plot**):

- use average lines (dotted lines) to provide context

- far fewer options in Power BI than in R or Excel

- consider using groupings to add clarity (e.g. **colour gradients**)

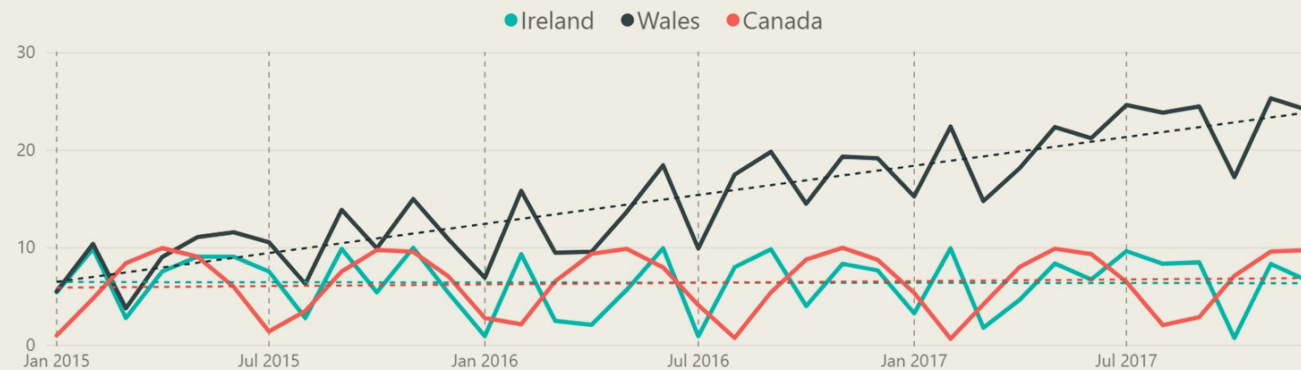How long should the perfect cup of tea be steeped?

# LINE CHARTS

Line chart can show a single series or multiple series of data *(particularly useful for **time series**)*.

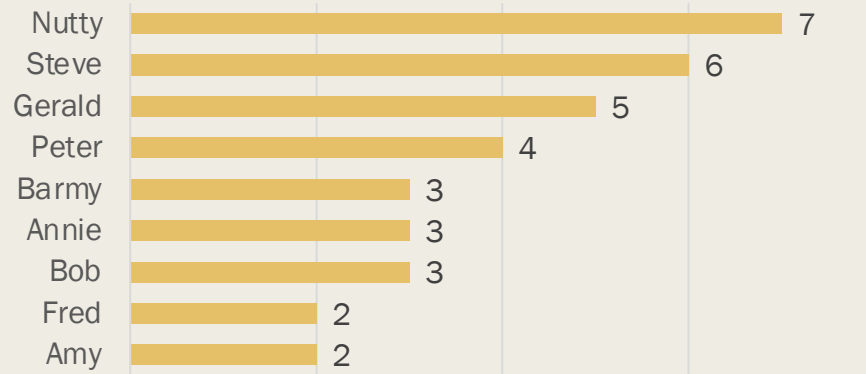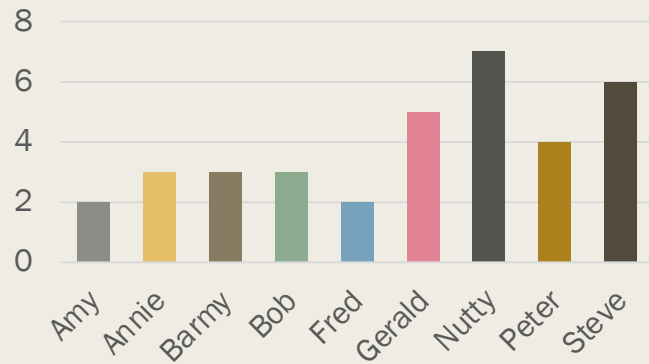Axis scale should be **clear** and **relevant**.

May wish to "**anchor**" $y-$axis if using dynamic filters
- otherwise the graph can jump around as people interact with it



Comparison of Countries – cups of tea drunk per week per person

# BAR CHARTS


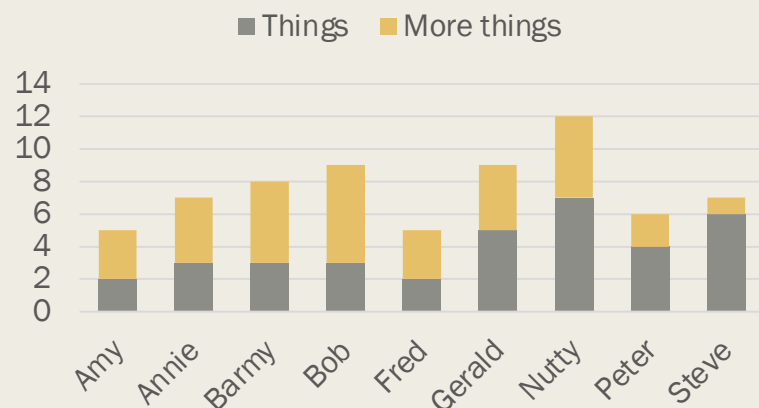
**Versatile** and useful.

ALWAYS (?) use a zero baseline.

Use graph axis OR data labels: axis for broad statements, data labels for details.

Horizontal charts are apparently **easier to read** (according to many studies).
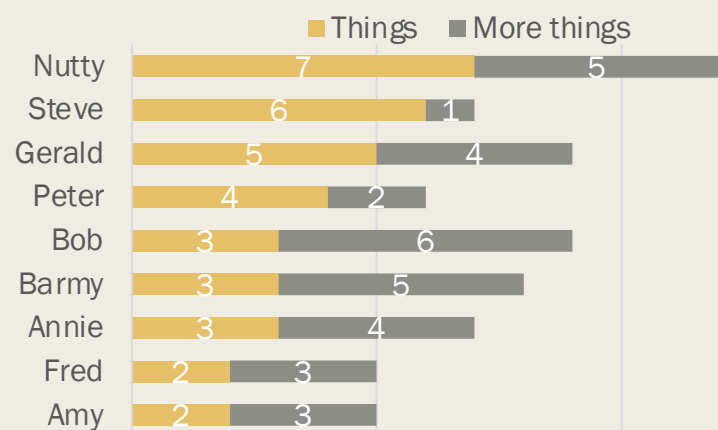
Think about the ordering of categories.
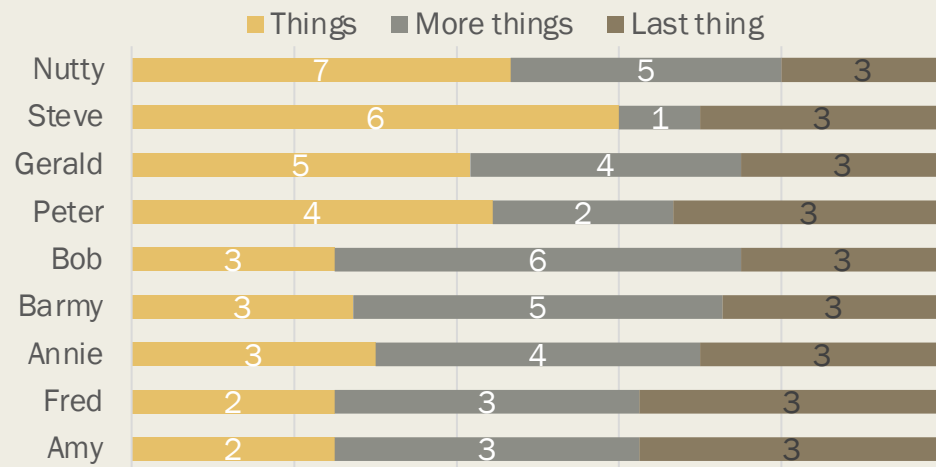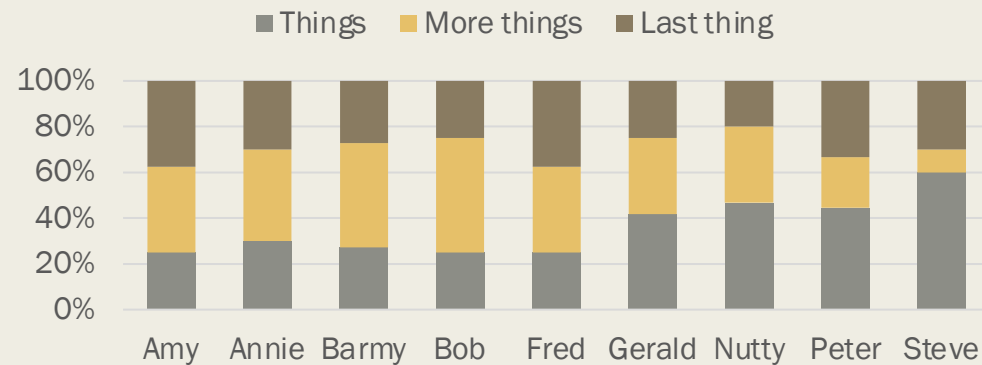
# STACKED BAR CHARTS



Designed for **comparing totals**, but can quickly become **overwhelming**.

Hard to sort / order.

Filtering is complicated in Power BI (what do you click on & how does the chart responds when the filter is clicked on?).

# 100% BAR CHARTS



Work well for visualizing **portions of a whole** on scale from negative to positive.

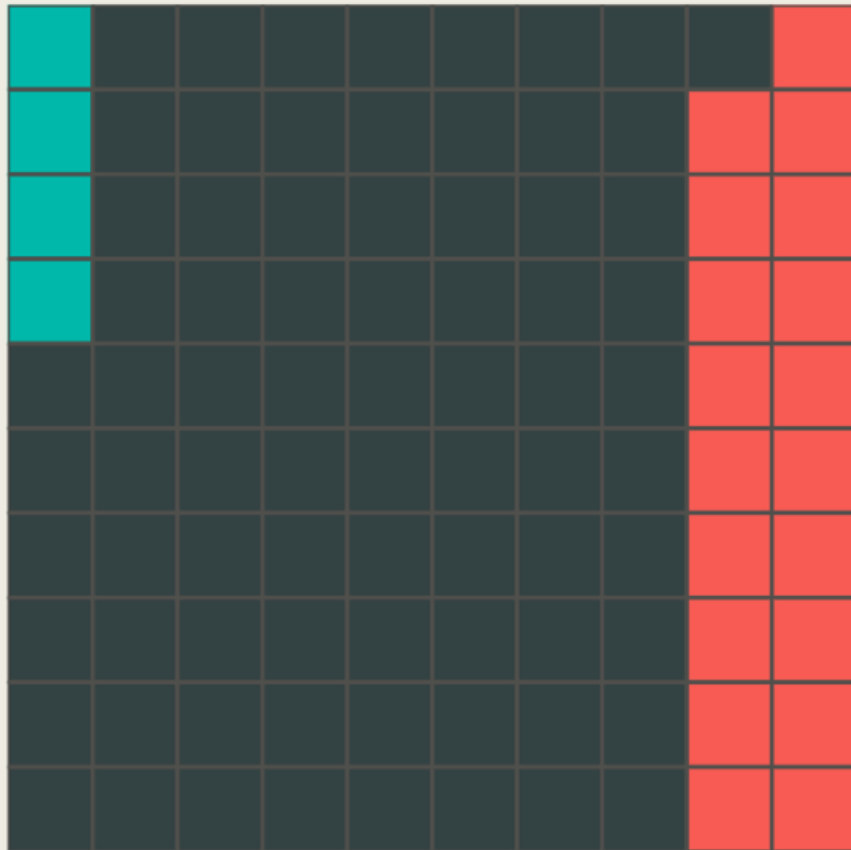Consistent baseline on far left and right.

Easy to compare.

**Problem:** there is no relative measure to magnitude of data (unless labeled).

Again, research shows that horizonal charts are **easier to process**.

# AREA CHARTS



Category ●Customers ●Leads ●Prospects

**Try to avoid:** human brains have a hard time attributing a value to a 2D area…

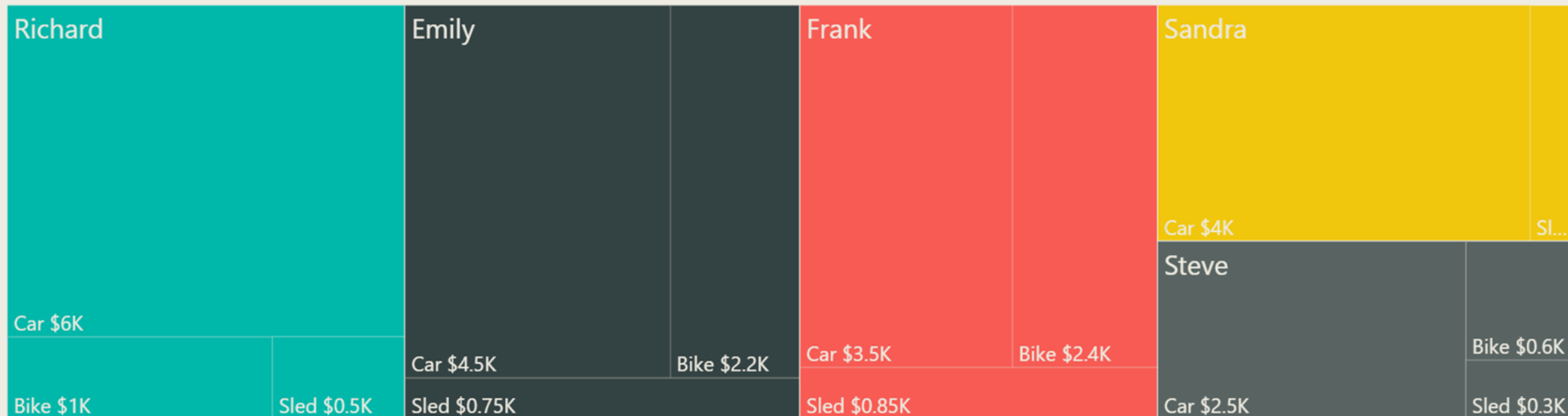… except for numbers with **vastly different** magnitudes.

# TREEMAPS

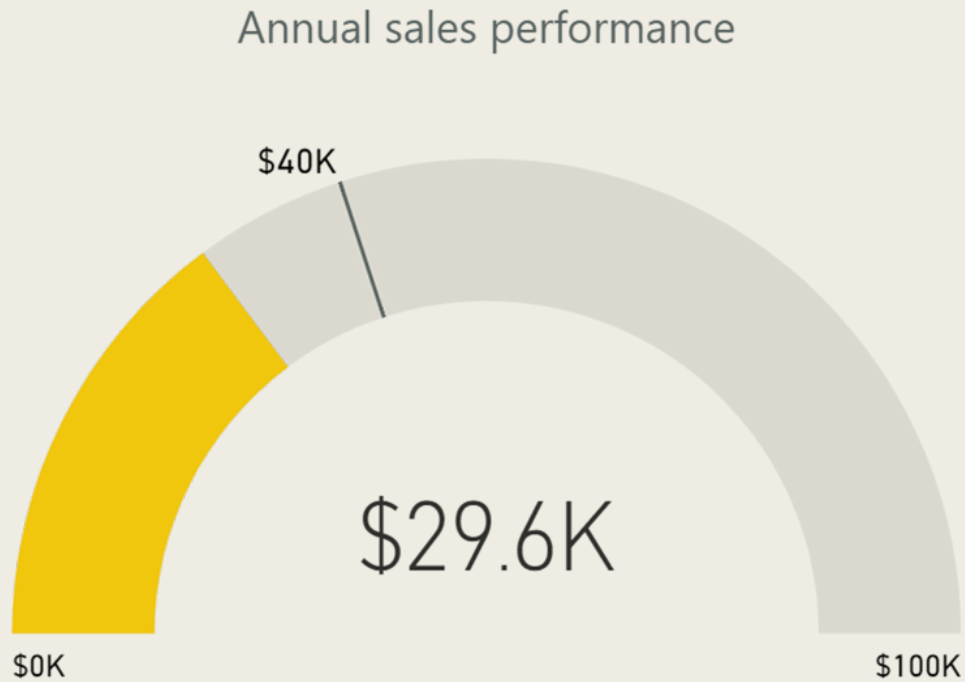Simultaneously show **big picture** and **comparisons** easily.

Easy to process data sub-categories.

Useful to prioritize "**big ticket items**" in dynamic dashboards.

Labeling and colouring are tricky.

# GAUGE CHARTS

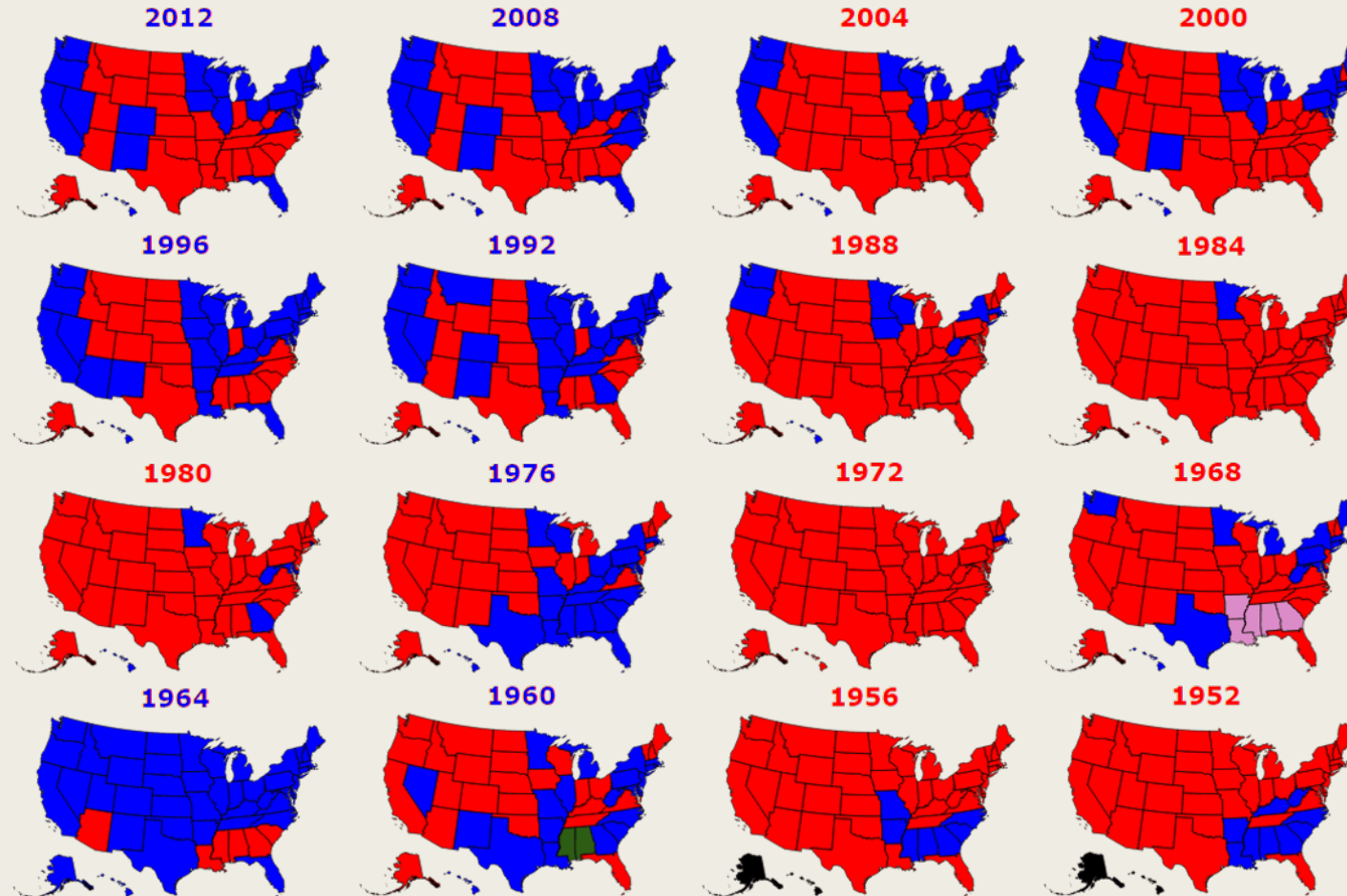Annual sales performance

$40K

$0K

$100K

$29.6K

Often used as a dashboard component (with or without needle).

Display **single value measures** towards goal / KPI.

Great to **show progress** (a bit of a management fad, though…)

Displays information that can be quickly **scanned** and **understood**.

# SMALL MULTIPLES

**U.S. Electoral College Results 1952 – 2012**

# OTHER CHART TYPES

Heatmaps and Choropleth Maps

Geographical Maps

Parallel Coordinates

Chernoff Faces

Word Clouds

Network Diagrams

Dendrograms and Trees

Sparklines

etc.

# INTERACTIVE AND ANIMATED VISUALIZATIONS

Animation **does not always** improve a visualization. What insights can interactivity provide? That depends on the data, and on the visualization method.

**Examples:**

- [The Clubs That Connect the World Cup](#), NY Times, 2014
- [Who Marries Whom](#), Bloomberg, 2016
- [Hipparcos Star Mapper](#), European Space Agency, 2016
- [The Internet of Things – a Primer](#), Information is Beautiful, 2016
- [The Genealogy and History of Popular Music Genres](#), Musicmap, 2016

# INTERACTIVE AND ANIMATED VISUALIZATIONS

**Examples (continued):**

- Sequences Sunburst, Kerry Rodden, 2015

- Health and Wealth of Nations, Gapminder Foundation

- Mobius Transformations Revealed, Arnold D.N, Rogness, J, 2007

- Visualizing the Riemann ζ Function and Analytic Continuation, 3Blue1Brown, 2016

- Small Arms and Ammunition – Imports and Exports, Google, 2012

- The Evolution of the Web, Google, Hyperakt, Vizzuality, 2012

- peoplemovin, Carlo Zapponi, 2012

# DISCUSSION

"There is always a danger that if certain types of visualization techniques take over, the kinds of questions that are particularly well-suited to providing data for these techniques will come to dominate the landscape, which will then affect data collection techniques, data availability, future interest, and so forth." (P. Boily)

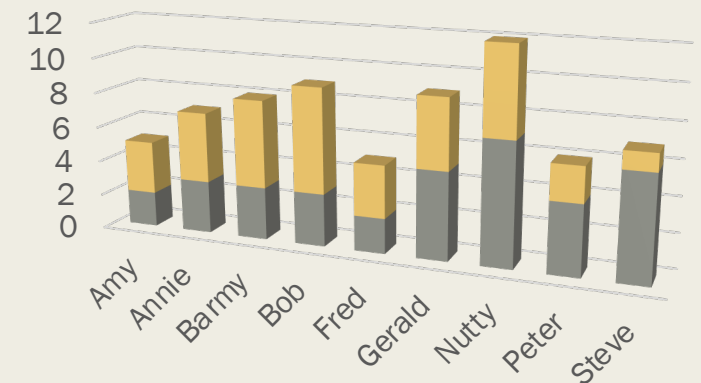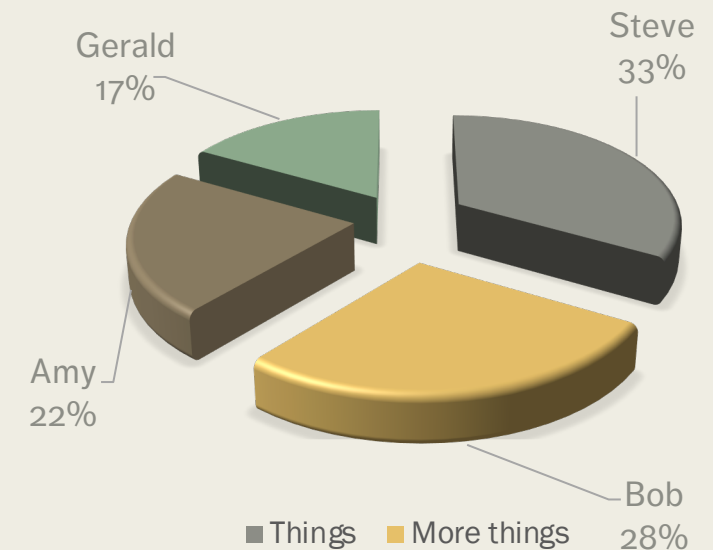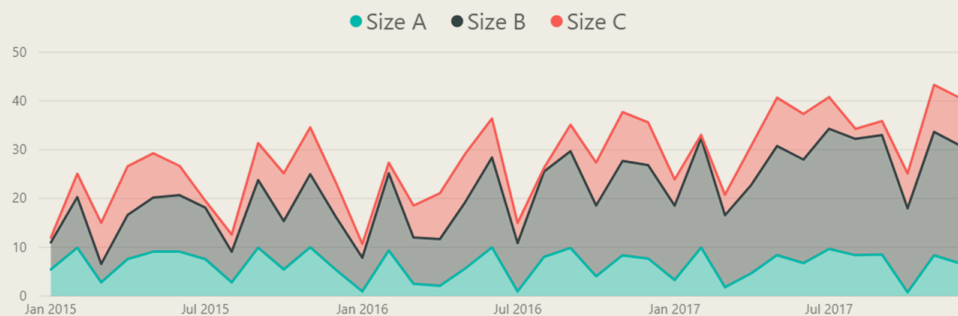Even when done well, 85% of users don't bother with interactive viz (NY Times).

**Take-Away:** explore the data and try different methods
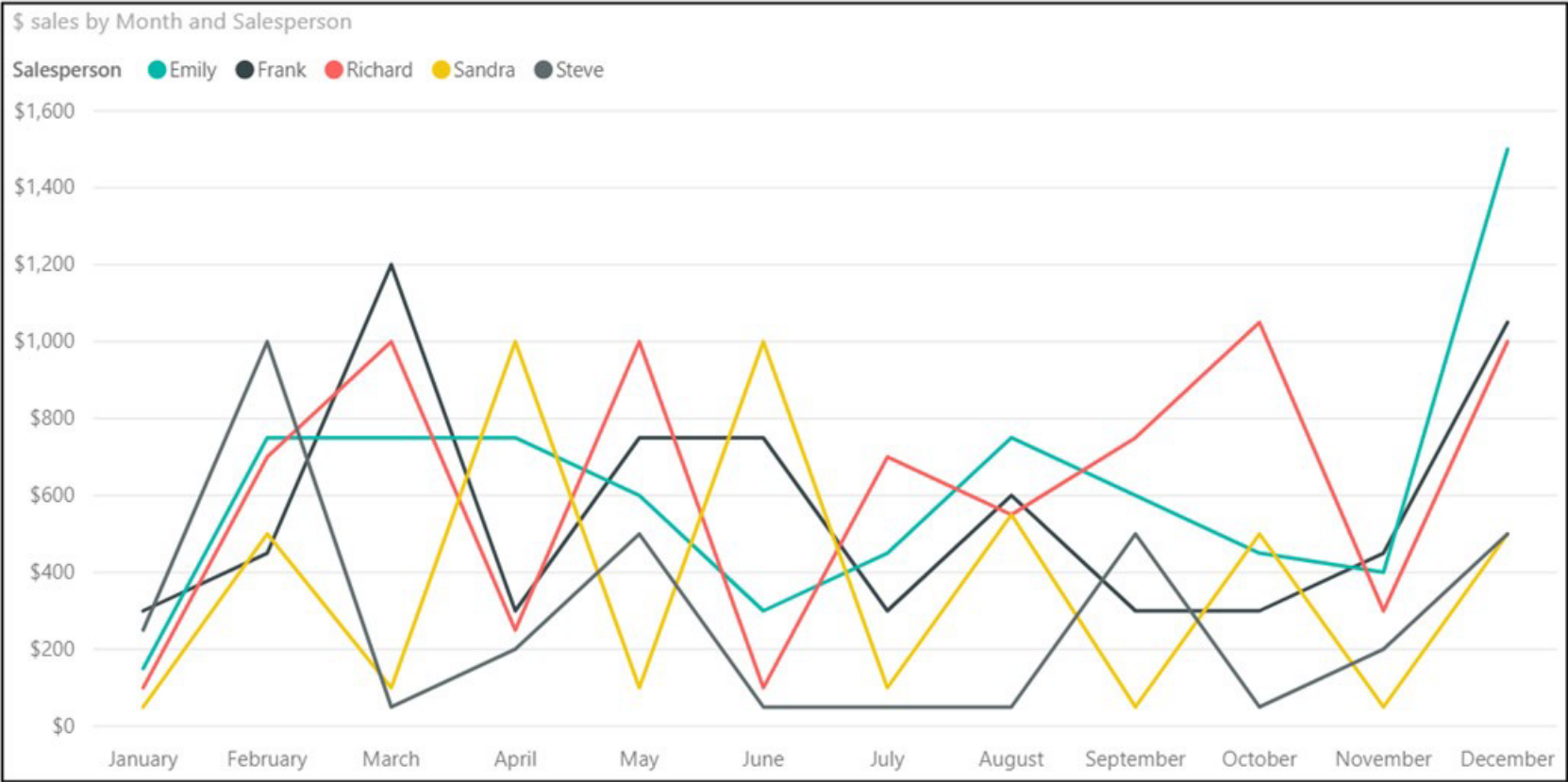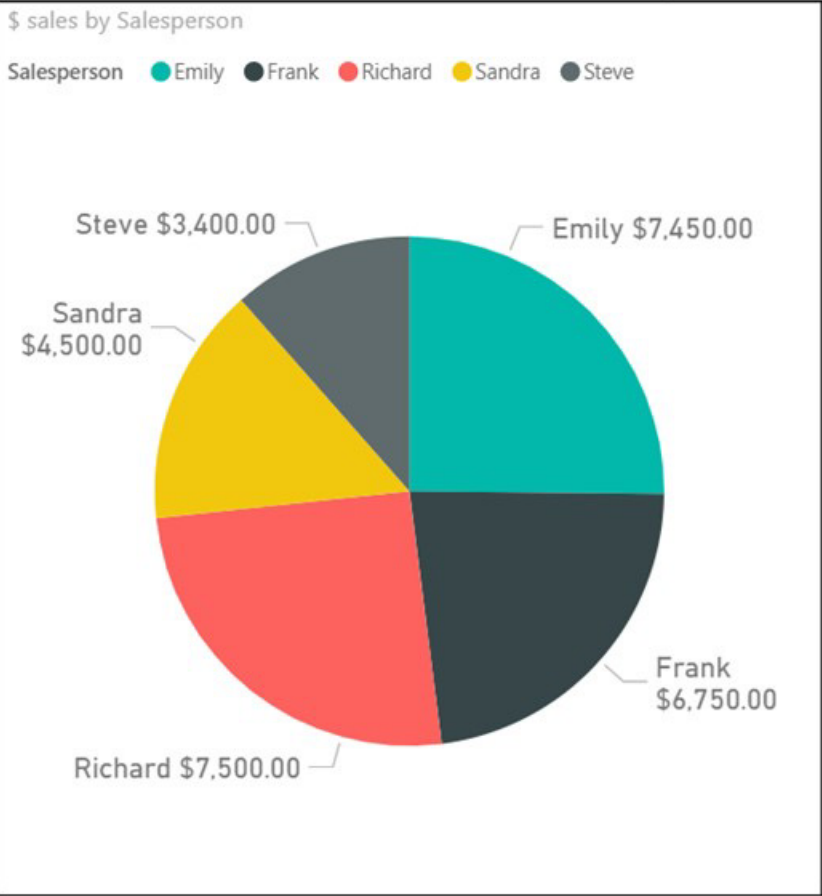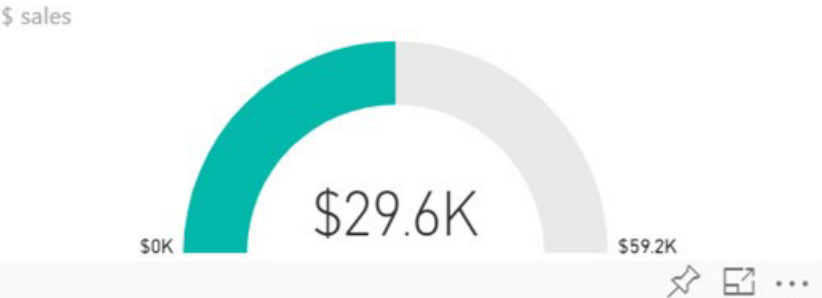
# CHARTS TO AVOID (within reason)

**AVOID (?) anything with an arc** (except gauge charts): pie, donut, etc: human brains have a hard time **comparing arcs** -- without labels, how different are Steve & Bob?

**AVOID 3D charts:** it is difficult to compare them visually (and they add **too much** clutter).
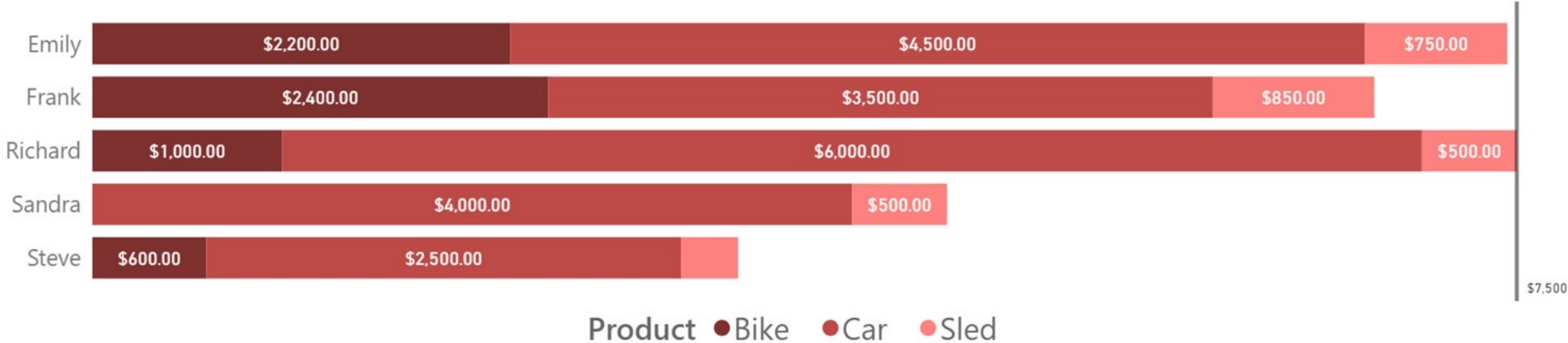
**AVOID stacked area charts:** way too confusing.

# Sales Dashboard

## $ sales



$0K    $29.6K    $59.2K

## $ sales by Salesperson

Salesperson  ● Emily  ● Frank  ● Richard  ● Sandra  ● Steve



Steve $3,400.00
Sandra $4,500.00
Emily $7,450.00
Richard $7,500.00
Frank $6,750.00

## $ sales by Month and Salesperson

Salesperson  ● Emily  ● Frank  ● Richard  ● Sandra  ● Steve



| | January | February | March | April | May | June | July | August | September | October | November | December |

## $ sales by Product and Salesperson

Product  ● Car  ● Bike  ● Sled



| Car | | | | Bike | | Sled | |
|-----|-----|-----|-----|------|------|------|------|
| | Emily $4.5K | | | | Emily $2.2K | Frank... | Emil... |
| | | | | Frank $2.4K | | | |
| | | | | | | Sand... | Rich... |
| Richard $6K | Sandra $4K | Frank $3.5K | Steve $2.5K | Richard $1K | Steve $0.6K | Steve $0.3K | |

# TAKE-AWAYS

Effective data visualizations **provide insights** and **facilitate understanding**.

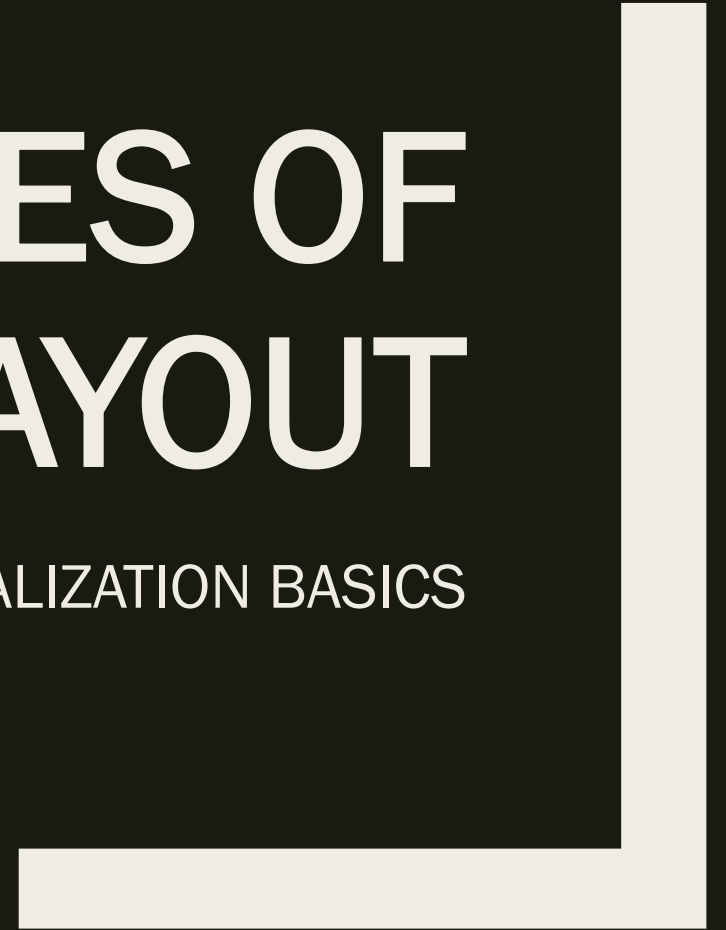The basic principles can guide your visualization design and consumption.

Be **creative** but keep your data and your representations **honest**.

Be mindful of attempts to distort trends and conclusions with flashy visuals.

Data and code should be made available along with the displays.

# BASIC RULES OF DESIGN AND LAYOUT
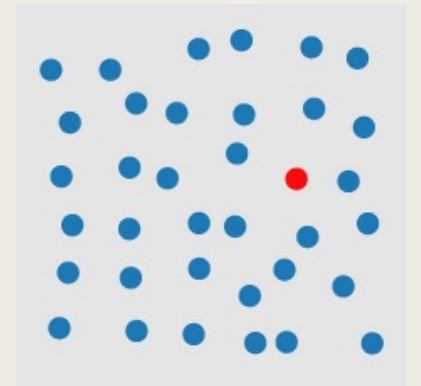
DATA VISUALIZATION BASICS

# VISUAL PROCESSING

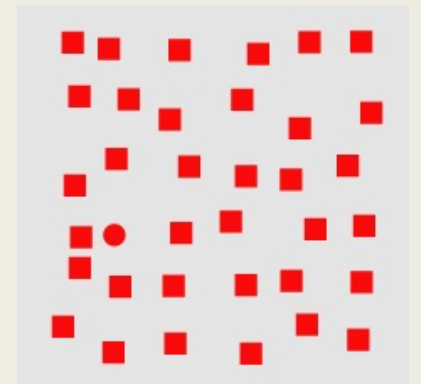Perception is **fragmented** − eyes are continuously scanning.

Visual thinking seeks patterns

- **Pre-attentive processes:** fast, instinctive, efficient, multitasking
  *gather information and build patterns:*

  features → patterns → objects

- **Attentive process:** slow, deliberate, focused
  *discover features in the patterns:*

  objects → patterns → features

**pre-attentive**



**attentive**

# PRE-ATTENTIVE ATTRIBUTES

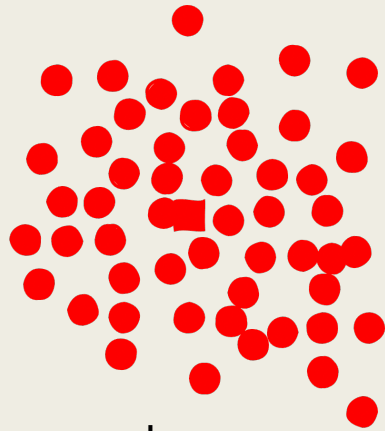Pre-attentive attributes are the domain of iconic memory (brief): they

- help to define a hierarchy of focus
- push non-message impacting components into the background

Use pre-attentive attributes to help **emphasize the story** (but don't overdo them):
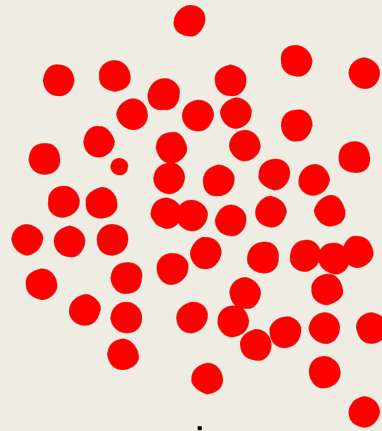
- easier to do in Excel and R, harder in Power BI

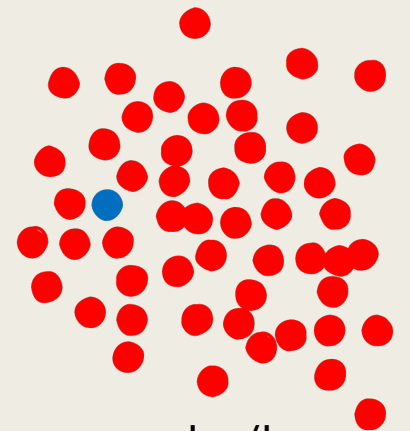**Challenge:** highlighting one aspect of a chart can make other aspects harder to see.
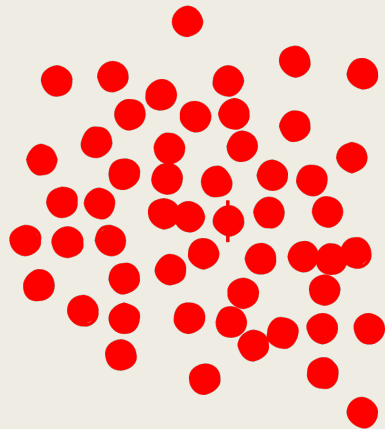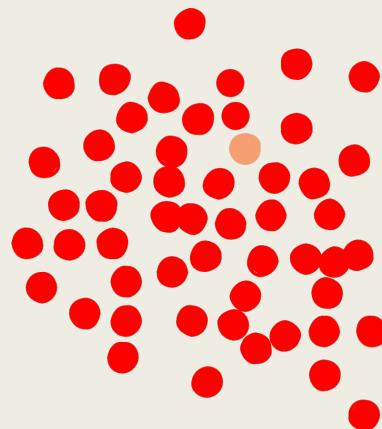
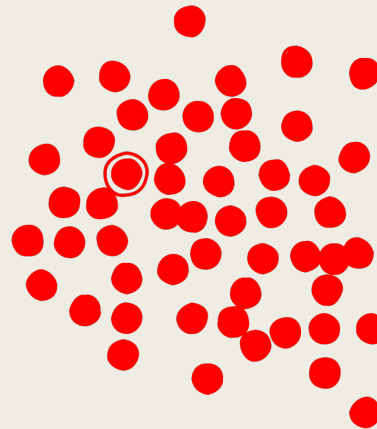# PRE-ATTENTIVE FEATURES

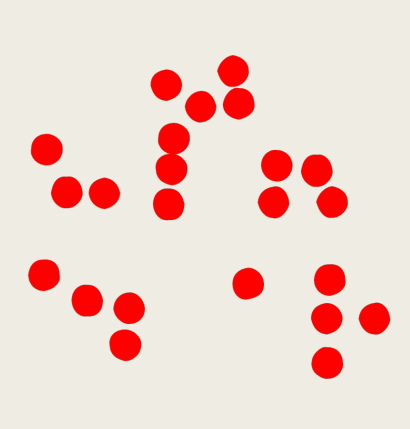

shape

size

sharpness

color/hue

markings

intensity/value

enclosure

numerosity

# PRE-ATTENTIVE ATTRIBUTES

How many 6's are there on the next slide?

28694O8609876

934858674 8676

2967303986739

3967496749674

286940860987 6

93485867486 76

2967303986739

396749674967 4

28694086609876
9348586748676
2967303986739
3967496749674

286940860998766

9348586748676

296730398674939

39674967496674

2869408609876
93485867486**76**
29**67**30398**67**39
39**67**49**67**49**67**4

2869408606098 76
93485867486 76
296730398673 9
396749674967 4

# DECLUTTERING (or, Less is More)

**CLUTTER IS THE ENEMY!**

- every element on a page adds **cognitive load**

- identify anything that isn't adding value and **remove**

- think of cognitive load as mental effort required to process information (lower is better)

- Tufte refers to the **data to ink ratio** – "the larger the share of a graphic's ink devoted to data, the better"

- in Resonate, Duarte refers to this as "**maximizing the signal-to-noise ratio**" where the signal is the information or the story we want to communicate.
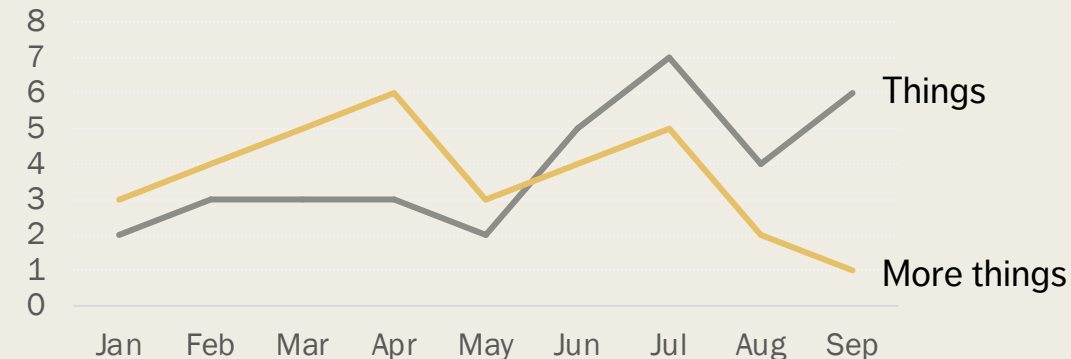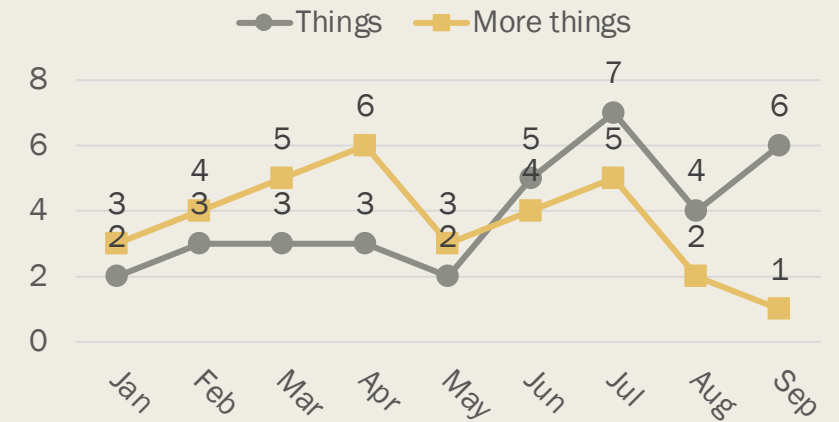
# DECLUTTERING (or, Less is More)

Use **Gestalt Principles** to organize/highlight data in a chart.

Align all the elements (graphs, text, lines, titles, etc.)

- DON'T rely on eye, use position boxes and values

**Charts:**

- remove border, gridlines, data markers
- clean up axis labels
- label data directly

# DECLUTTERING (or, Less is More)

Use **consistent** font, font size, colour and alignment.

Don't rotate text to anything other than 0 or 90 degrees.

Use **white space:**
- margins should remain free of text and visuals
- don't stretch visuals to edge of page or too close to other visuals
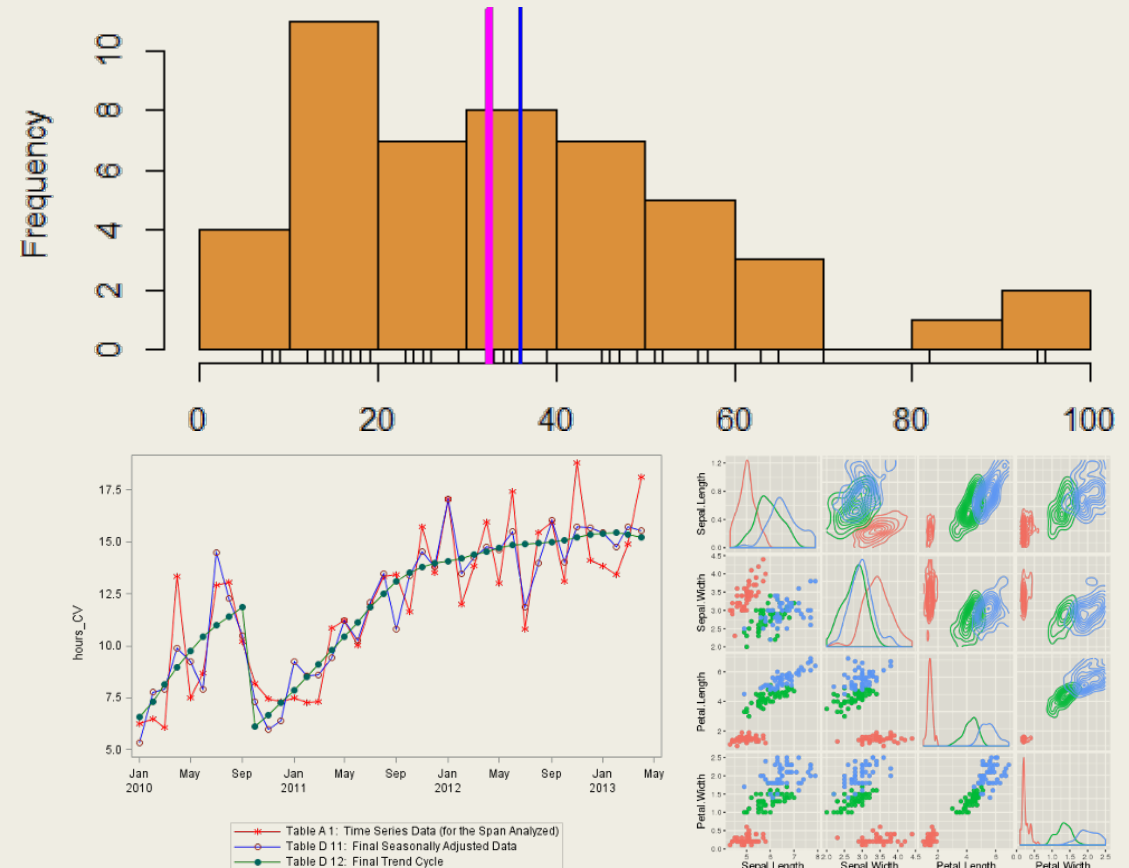- think of white space as a border

# CHART SIZES

Assuming that the chart has been decluttered:

- things of equal importance size **similarly;**
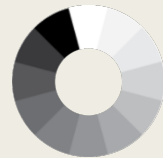
- other things scale to **importance**.

As one rarely puts more than 3-4 charts on a page, there are limited size options.

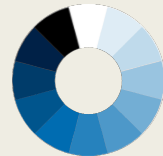Perennial exception: **geographical maps** may require more space.

# COLOUR SCHEMES

**Achromatic**

**Monochromatic**

Complementary

**Split complementary**

**Split-Left/Right Complementary**

**Analogous**

**Colour Diad**

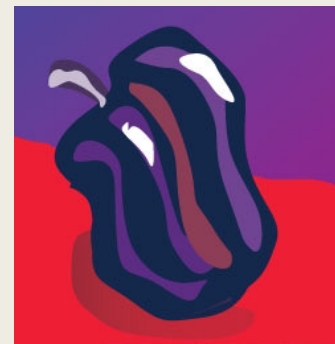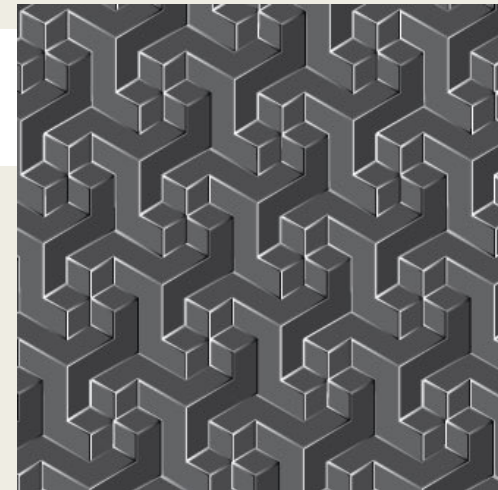**Colour Triad**

**Colour Tetrad**

Can you identify the colour schemes underlying each of these images?

Monochromatic (Blues)

Tetrad

Split Complementary (Green, Orange & Blue)

Achromatic

Diad (Blue & Green)

Triad (Primary Colors)

Diad (Green & Orange)

Analogous (Green & Yellow)

Complementary

Diad (Red & Violet)

Can you identify the colour schemes underlying each of these images?

Zeileis, Hornick & Murrell 24 Distinct Colors (1st group)

Bike ■ Car ■ Sled ■ Truck ■ Scooter ■ Skateboard

Kelly's 22 Colors of Maximum Contrast (1st 6)

Bike ■ Car ■ Sled ■ Truck ■ Scooter ■ Skateboard

Paul Tol 14 Rainbow Scheme (middle green – orange)

Bike ■ Car ■ Sled ■ Truck ■ Scooter ■ Skateboard

# COLOUR SCHEMES

When it comes to colour, **less is more**: use it sparingly (graphic designers are taught to "get it right, in black and white").
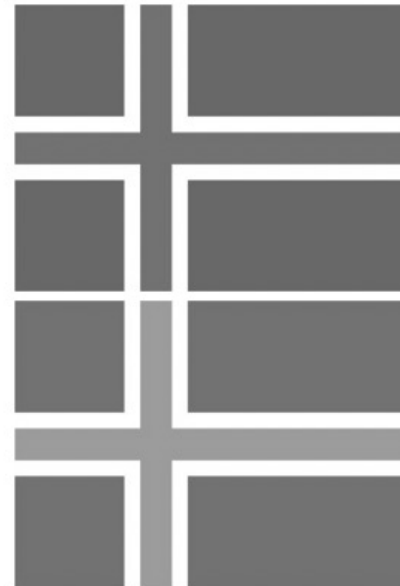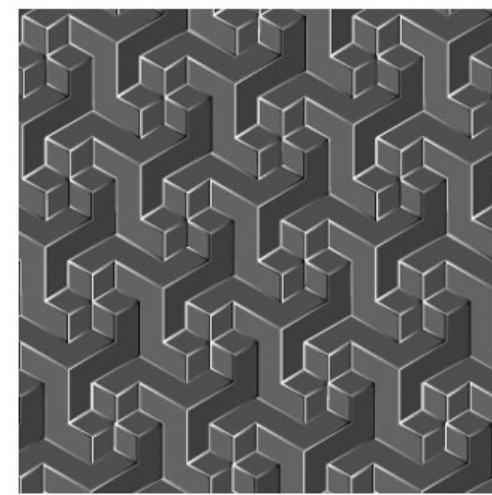
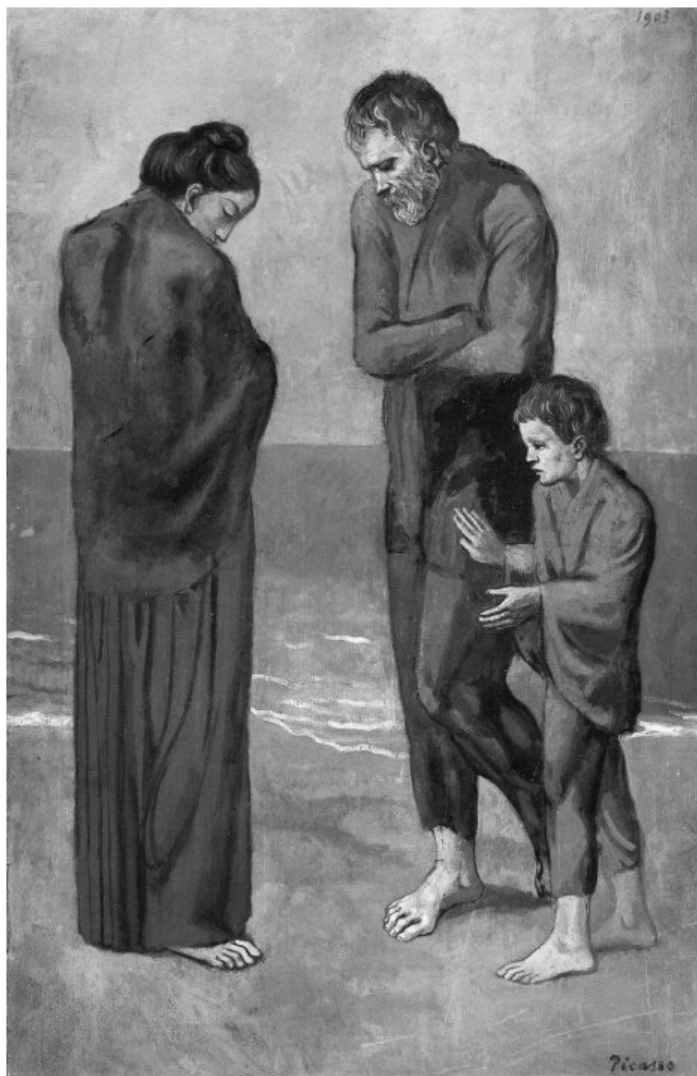Based on the Gestalt Principles, **monochrome** schemes can be particularly effective.

When appropriate, pick scheme based on corporate identity (this maximizes buy in).

Create a template (and stick to it).

Upload images to see what charts look like in various flavours of colour-blindness:
- ▪ https://www.color-blindness.com/coblis-color-blindness-simulator (there are other tools)

Monochromatic
achromatopsia
(greys)

# POSITION

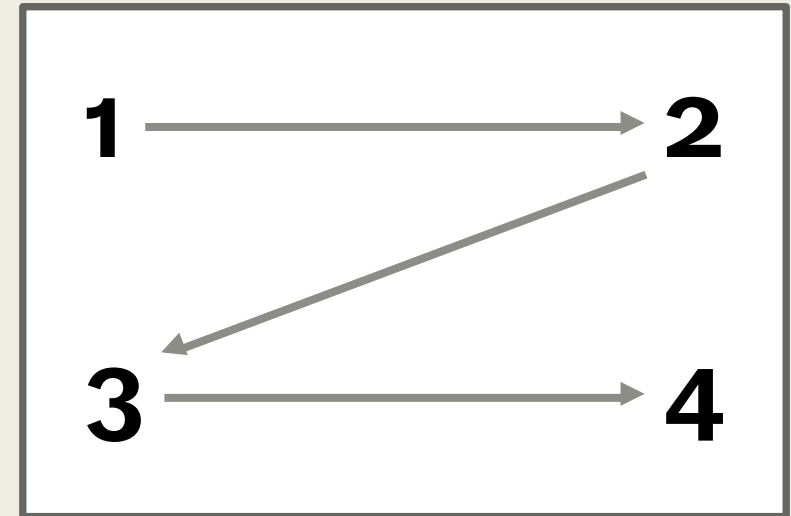How should the elements be placed in a chart or a dashboard?

In the West, most people start at the **top left** and zig- zag all the way to the **bottom right**.

**Simple rule:** don't make people work too hard

- main message: top left/top right
- info in order of preference
- people concentrate less as they scan so get less complex as you move to bottom corner

# DASHBOARDS

DATA VISUALIZATION BASICS

# DASHBOARDS

A **dashboard** is any visual display of data used to monitor conditions and/or facilitate understanding.
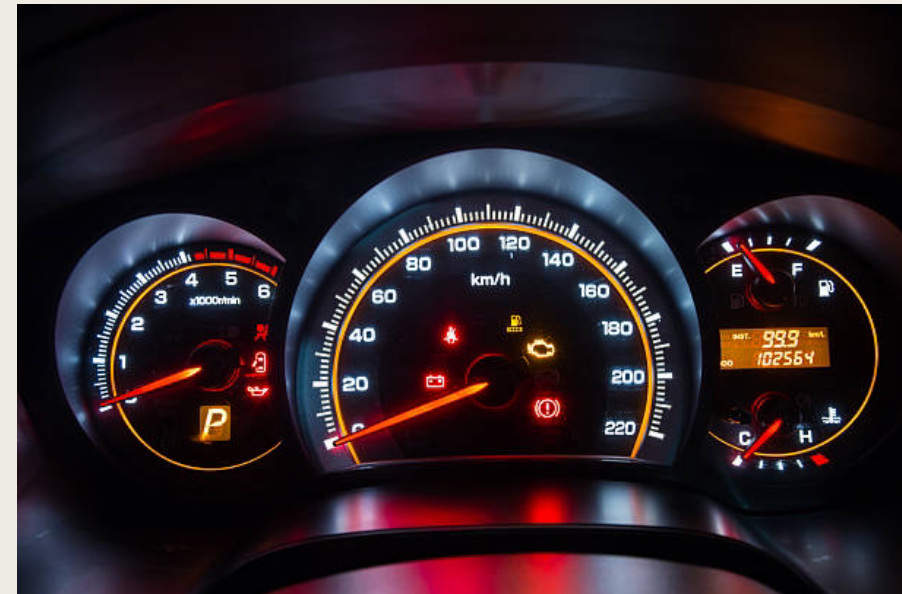
**Examples:**

- interactive display that allows people to explore motor insurance claims by city, province, driver age, etc.
- PDF showing key audit metrics that gets e-mailed to a Department's DG on a weekly basis.
- wall-mounted screen that shows call centre statistics in real-time.
- mobile app that allow hospital administrators to review wait times on an hourly- and daily-basis for the current year and the previous year.

# SOME QUESTIONS TO PONDER

In a car's dashboard, a small number of **key indicators** (speed, gasoline level, lights, etc.) need to be understood **at a glance**. A dashboard design that does not take these two characteristics under consideration can have catastrophic consequences.

The following questions need to be answered prior to the dashboard being designed:

- Who is the dashboard's **consumer**?
- What **story** does the dashboard tell?
- What data (categories) will be used?
- What will **appear** on the dashboard?
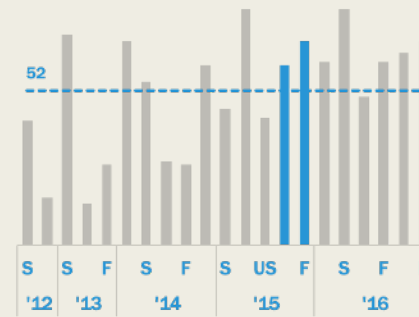- How can the dashboard **help** the consumer?

# Course Metrics

**Strengths:**

- Easy-to-see key metrics

- Simple color scheme

- Potential to be static or interactive
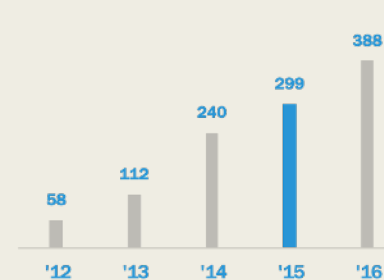
- Both overview and details are clear



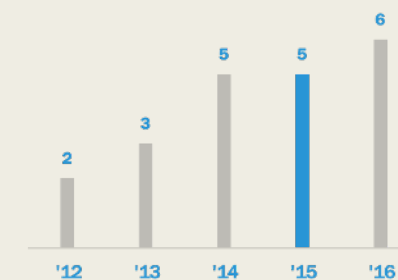Course Metrics Dashboard created by Jeffrey A. Shaffer. Data from University of Cincinnati Course Evaluations. Blue indicates the 2 most recent rating periods.

# DASHBOARD EVALUATION

There are no perfect dashboards – no collection of charts will ever suit everyone who encounters it.

All dashboards should be **truthful** and **functional**, but dashboards that are also **elegant** (delightful, enjoyable) will take you further.

All dashboards are **incomplete**. Good dashboards will still lead to dead ends, but they should allow users to ask: "Why? What is the root cause of a problem?"

**Tools:** Excel, Power BI, Tableau, R + Shiny, Geckoboard, Matillion, etc.
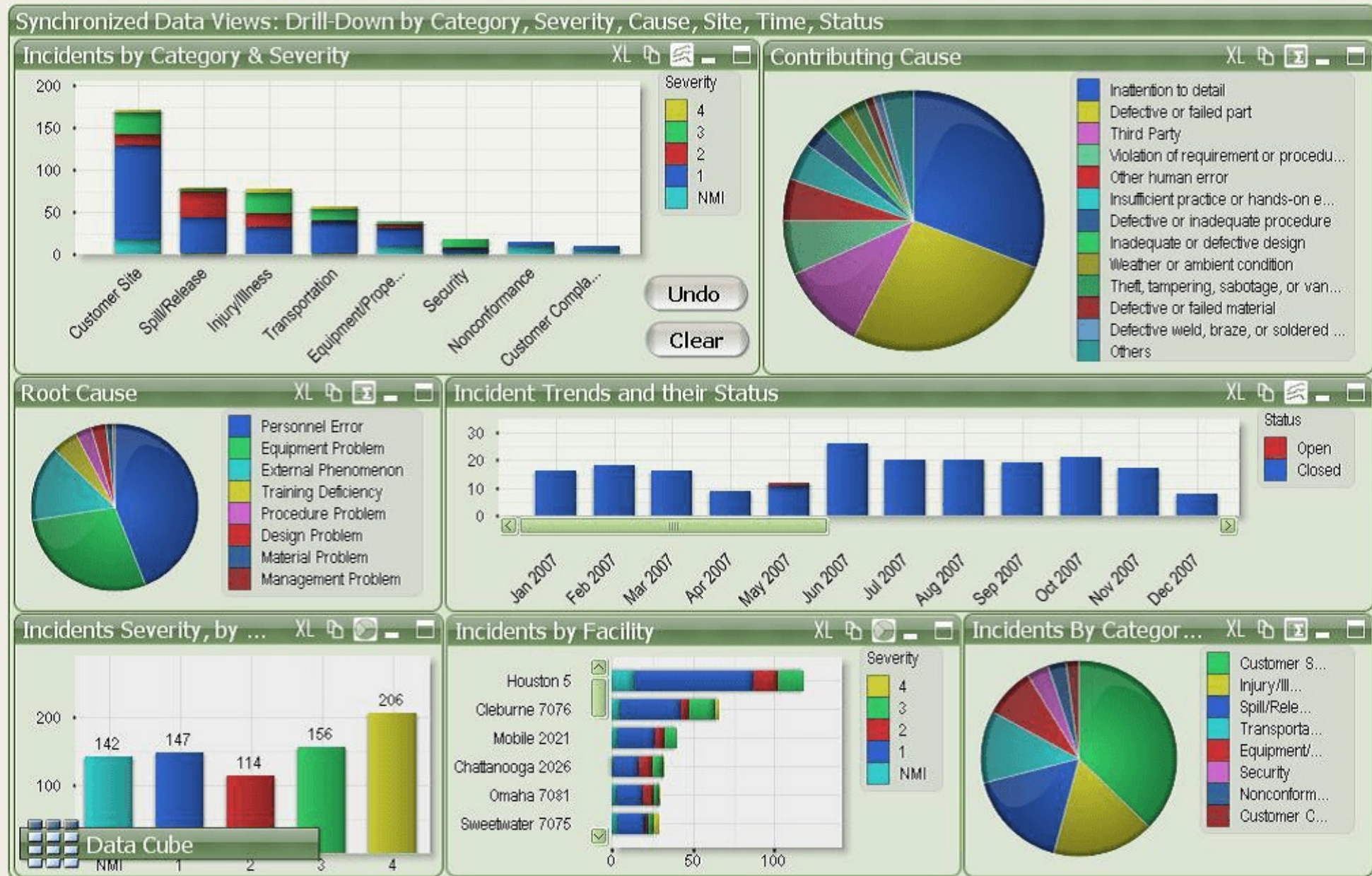
# EXERCISE

Consider the following dashboards.

Can you figure out, at a glance, who their audience is?

What are their strengths?
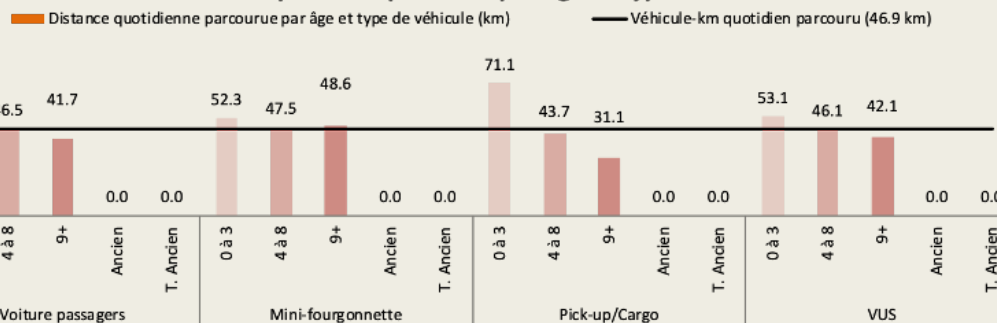
What are their limitations?
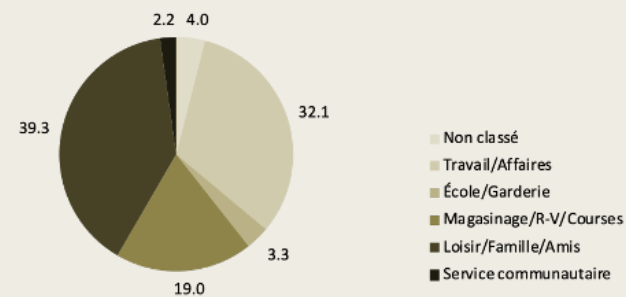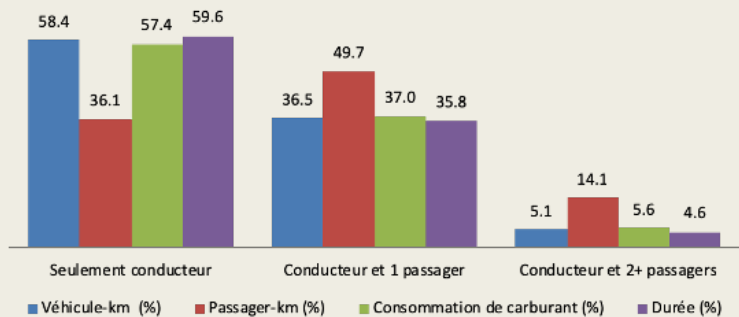
How would you improve them?

# Ontario – 1er trimestre 2012

**Caractéristiques des déplacements**

Véhicule-km quotidien parcouru par âge et type de véhicule
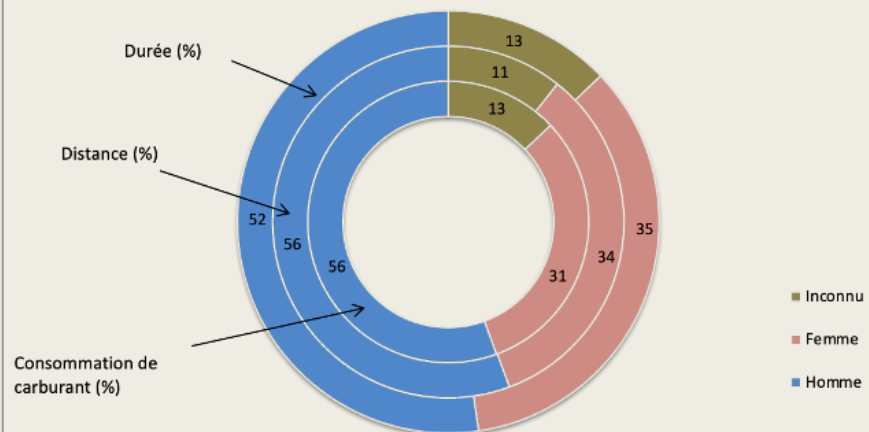
■ Distance quotidienne parcourue par âge et type de véhicule (km) — Véhicule-km quotidien parcouru (46.9 km)

Voiture passagers: 0 à 3 = 54.4, 4 à 8 = 46.5, 9+ = 41.7, Ancien = 0.0, T. Ancien = 0.0
Mini-fourgonnette: 0 à 3 = 52.3, 4 à 8 = 47.5, 9+ = 48.6, Ancien = 0.0, T. Ancien = 0.0
Pick-up/Cargo: 0 à 3 = 71.1, 4 à 8 = 43.7, 9+ = 31.1, Ancien = 0.0, T. Ancien = 0.0
VUS: 0 à 3 = 53.1, 4 à 8 = 46.1, 9+ = 42.1, Ancien = 0.0, T. Ancien = 0.0



Passager-km quotidien parcouru par but des déplacements (%)

- Non classé: 4.0
- Travail/Affaires: 32.1
- École/Garderie: 3.3
- Magasinage/R-V/Courses: 19.0
- Loisir/Famille/Amis: 39.3
- Service communautaire: 2.2



Distance, passager-km parcouru, durée et consommation de carburant par occupation

Seulement conducteur: 58.4, 36.1, 57.4, 59.6
Conducteur et 1 passager: 36.5, 49.7, 37.0, 35.8
Conducteur et 2+ passagers: 5.1, 14.1, 5.6, 4.6

■ Véhicule-km (%) ■ Passager-km (%) ■ Consommation de carburant (%) ■ Durée (%)



Proportion de déplacements par segments de distance

- 100+ km: 0.0
- 51 km à 100 km: 2.1
- 31 km à 50 km: 4.3
- 21 km à 30 km: 6.4
- 16 km à 20 km: 4.3
- 11 km à 15 km: 8.5
- 6 km à 10 km: 17.0
- 1 km à 5 km: 55.3
- 0 km: 2.1

Pourcentage du total



Durée, distance et consommation de carburant par sexe

Durée (%), Distance (%), Consommation de carburant (%)

52, 56, 56 / 35, 34, 31 / 13, 11, 13

■ Inconnu ■ Femme ■ Homme



Consommation de carburant, distance et durée par âge des conducteurs

Consommation de carburant (%), Distance (%), Durée (%)

13, 13, 13 / 11, 10, 12 / 4, 4, 4 / 31, 33, 31 / 41, 41, 40

■ Inconnu ■ 16-24 ■ 25-44 ■ 45-64 ■ 65+
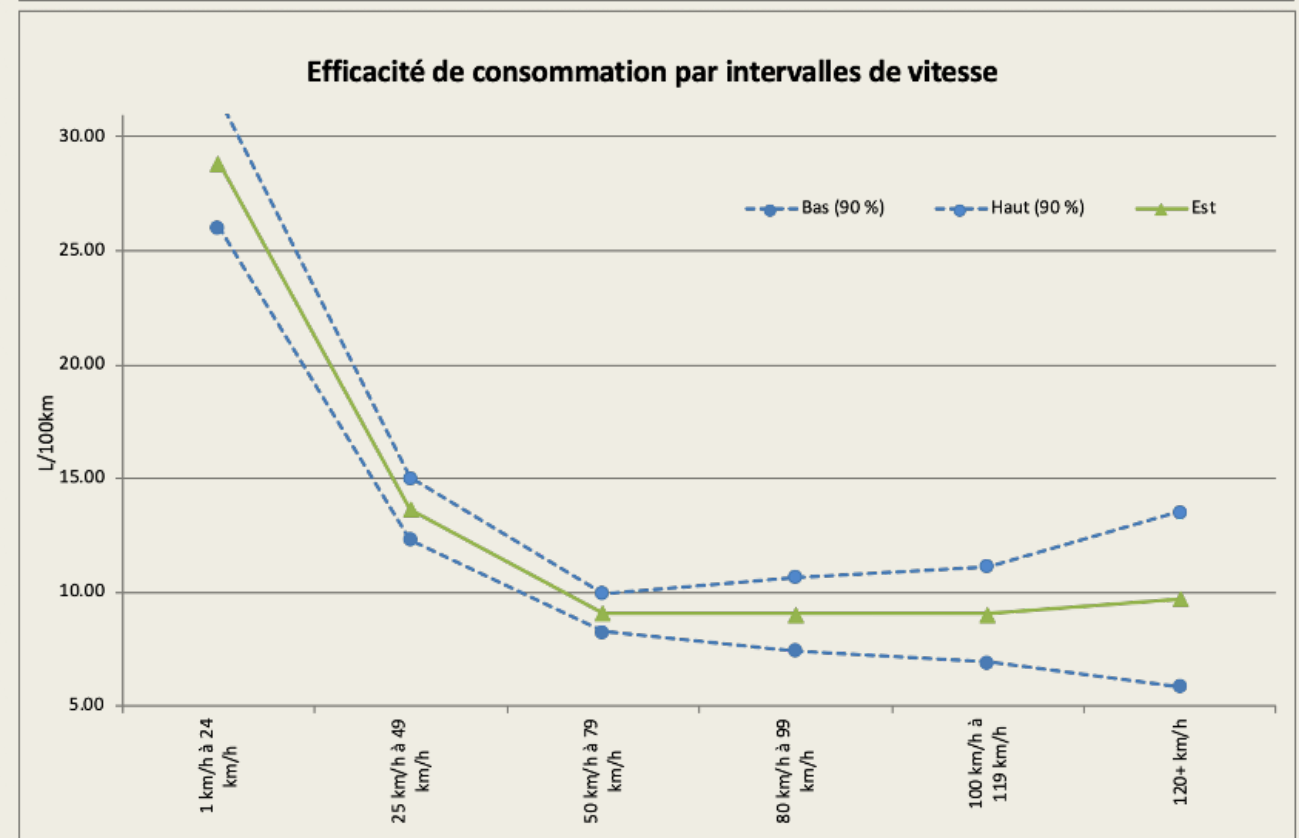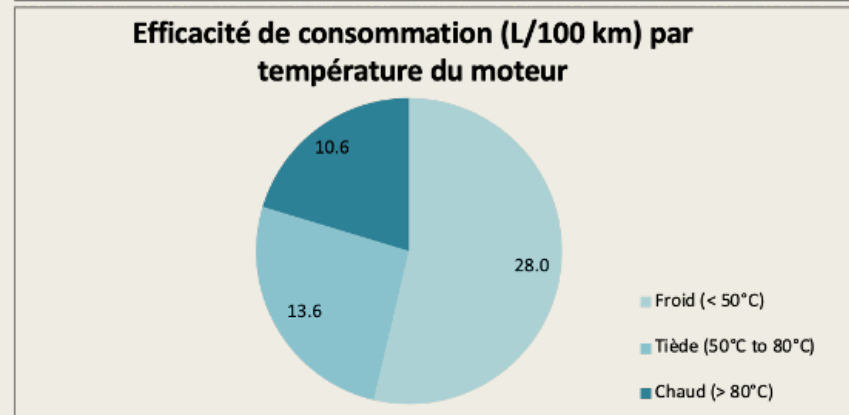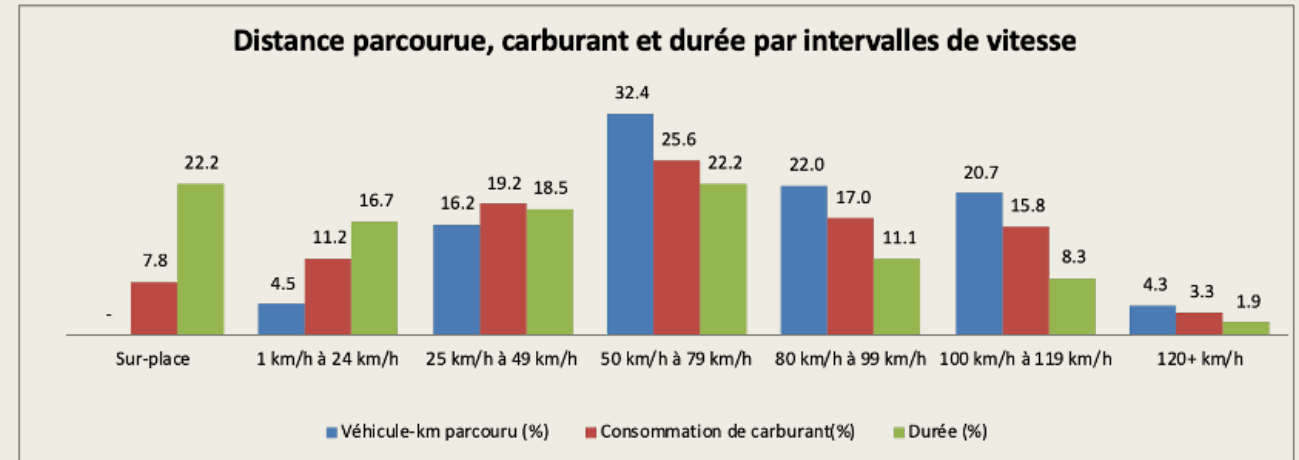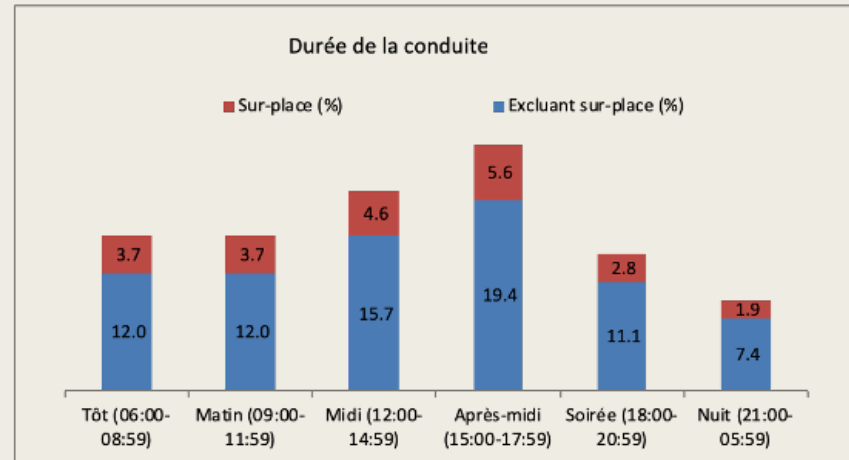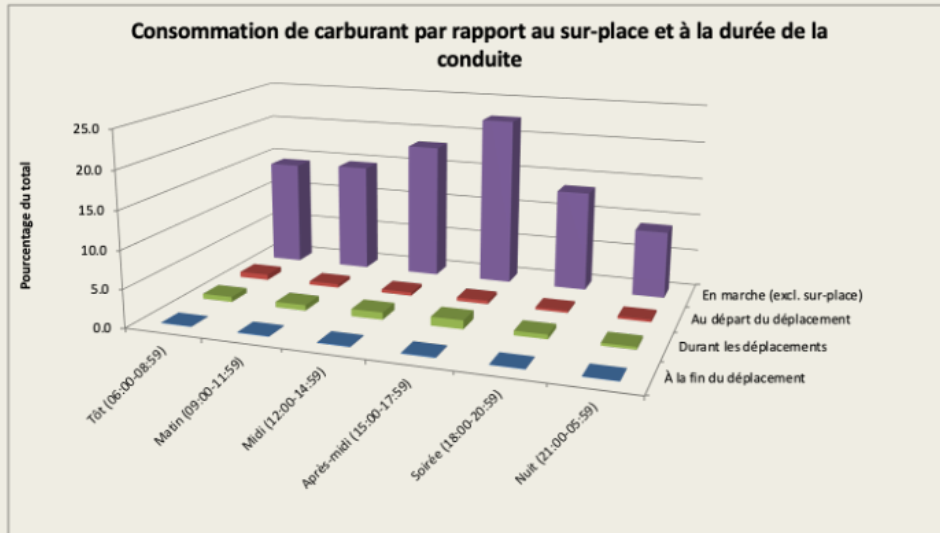
# Ontario – 1er trimestre 2012

**Sous-caractéristiques des déplacements**



Durée de la conduite

- Sur-place (%)
- Excluant sur-place (%)



Distance parcourue, carburant et durée par intervalles de vitesse

- Véhicule-km parcouru (%)
- Consommation de carburant(%)
- Durée (%)



## Efficacité de consommation (L/100 km) par température du moteur

- Froid (< 50°C)
- Tiède (50°C to 80°C)
- Chaud (> 80°C)



## Efficacité de consommation par intervalles de vitesse

- Bas (90 %)
- Haut (90 %)
- Est



## Durée et consommation de carburant par type de sur-place

- Durée de conduite quotidienne (%)
- Consommation de carburant (%)

# Ontario – 1er trimestre 2012

## Caractéristiques mixtes sur les déplacements



Durée de la conduite (min) par jour-type et occupation



Durée de la conduite (min) par but et occupation



Consommation de carburant par vitesse et température du moteur



Distance par occupation et durée de la conduite



Consommation de carburant par rapport au sur-place et à la durée de la conduite

# WHAT IS WRONG WITH THEM?

**Dashboard #1:** not glanceable, overuse of colour, pie charts, ...

**Dashboard #2:** 3D visualizations, distracting borders and background, lack of filtered data, insufficient labels and context, ...

**Dashboards #3:** where to begin ...

# EXERCISE

In teams or individually, identify a scenario for which a dashboard could prove useful.

Determine specific questions that the dashboard could help answer or insights that it could provide.

Identify data sources and data elements that could be fed into your dashboard.

Design a display (with pen and paper) with mock charts.

What are the strengths and limitations of your dashboard? Is it functional? Elegant?