# Visual Analytics

## Communicating data-driven insights through data visualization techniques and useful dashboards

Mikel Madina &

### Miren Berasategi

miren.berasategi@deusto.es

Deusto

# 0. Introduction

# 0.1 Key points

- **Data driven**: as seen in Onieva's and Lorenzo's lectures
- **Insights**: que usen las características gráficas
- **Data visualization techniques**: para obtener las los insights
- **Dashboards**: as *situattion awareness* tools

Tableau Desktop para prácticar

# 0.2 Why use visualization

- Sight is our most developed sense
- The visual system provides a very high-bandwidth channel to our brains
- A significant amount of visual information processing occurs in parallel at the preconscious level
- The human brain is *trained* to identify visual patterns
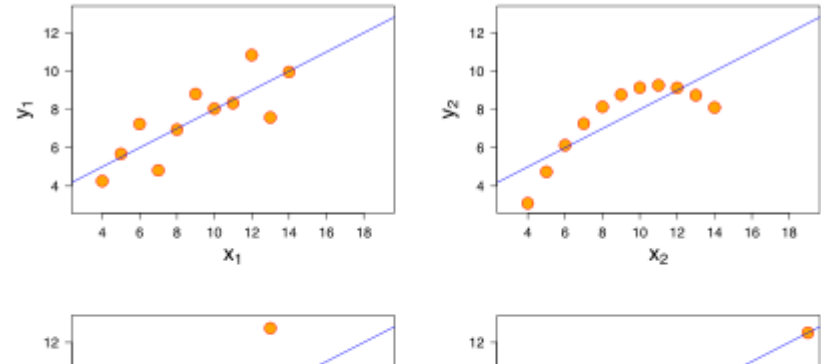- Summary statistics have the intrinsic limitation of data loss

Speaker notes

Implications of visual perception relevant to visualization design on next section

# 0.2 Why use visualization



Anscombe's quartet

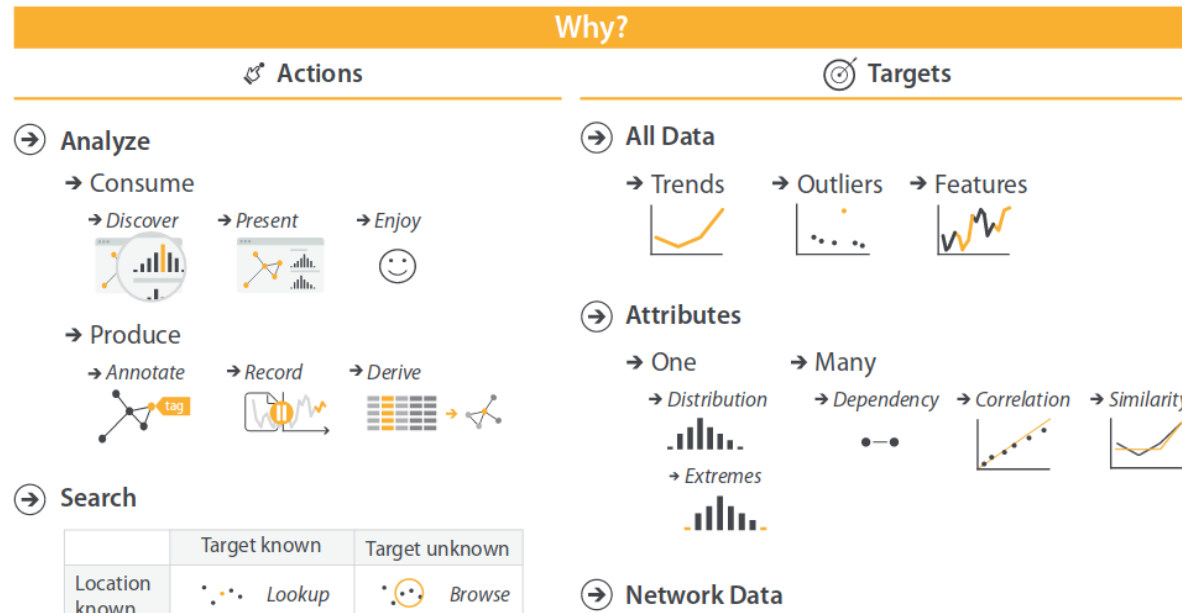| I | | II | | III | | IV | |
|---|---|---|---|---|---|---|---|
| x | y | x | y | x | y | x | y |
| 10.0 | 8.04 | 10.0 | 9.14 | 10.0 | 7.46 | 8.0 | 6.58 |
| 8.0 | 6.95 | 8.0 | 8.14 | 8.0 | 6.77 | 8.0 | 5.76 |
| 13.0 | 7.58 | 13.0 | 8.74 | 13.0 | 12.74 | 8.0 | 7.71 |
| 9.0 | 8.81 | 9.0 | 8.77 | 9.0 | 7.11 | 8.0 | 8.84 |
| 11.0 | 8.33 | 11.0 | 9.26 | 11.0 | 7.81 | 8.0 | 8.47 |
| 14.0 | 9.96 | 14.0 | 8.10 | 14.0 | 8.84 | 8.0 | 7.04 |
| 6.0 | 7.24 | 6.0 | 6.13 | 6.0 | 6.08 | 8.0 | 5.25 |

Speaker notes

Traditional summary statistics can be misleading. These datasets share almost identical mean, variance, correlation and linear regression lines.

1. Shows *normal* distribution
2. There is correlation, but it's not linear
3. There is correlation, it IS LINEAR, but different from what emerged from the data
4. No relationship whatsoever. Outlier is enough for *apparent* correlation

**Bottom line**: summary statistics, although very sensitive to outliers, are important, but plotting the data (visualization) is also a necessary step during the first stages of the data analytics process, before making any assumptions.

Munzner 2014 p.8

# 0.3 What to use visualization for



Why?

**Actions**

→ **Analyze**
  → Consume
    → *Discover*  → *Present*  → *Enjoy*
  → Produce
    → *Annotate*  → *Record*  → *Derive*

→ **Search**

|  | Target known | Target unknown |
|---|---|---|
| Location known | Lookup | Browse |

**Targets**

→ **All Data**
  → Trends  → Outliers  → Features

→ **Attributes**
  → One
    → *Distribution*
    → *Extremes*
  → Many
    → *Dependency*  → *Correlation*  → *Similarity*

→ **Network Data**

# 0.3 What to use visualization for



Speaker notes

- Very broadly relevant for all kinds of data: trends, outliers, features
- One attribute/variable: find an individual value, distribution of all values, find extremes
- Many attributes: dependency, correlation, similarity
- The rest (network/spatial) are specific to certain types of datasets

The abstract task of understanding **trends, outliers, distributions and correlations** are extremely common reasons to use data visualization.

# 0.3 What to use visualization for



| | Target known | Target unknown |
|---|---|---|
| Location known | Lookup | Browse |
| Location unknown | Locate | Explore |

**Search**

**Actions**

**Analyze**

→ Consume
→ Discover   → Present   → Enjoy

Speaker notes

- Search for an element of interest (based on knowledge of identity & location)
    - lookup: humans in tree vis of species of mammals
    - locate: where are rabbits? lagomorphs – not rodents
    - browse: when users don't know exactly what they're looking for but have an idea of characteristics
    - explore: not sure of either one
- Once a target or set of targets is identified, query those at one of 3 scopes (progression from one to many to all targets)
    - identify (US election map example)
    - compare, more difficult, requires more sophisticated idioms
    - summarize = overview (extremely common)
- Analyze data:
    - consume (most common use case): discover, present (known info, i.e. insights), enjoy (for the fun of it)
    - produce: annotate, record (capture steps), derive (produce new data elements)
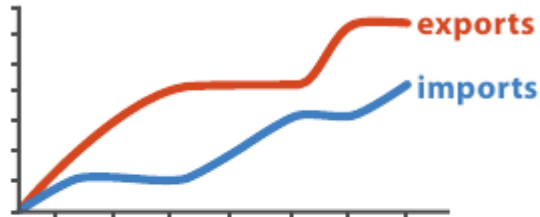
# 0.3 What to use visualization for

*There is a **strong relationship** between the form of the data (the attribute/variable and dataset types) and what kinds of vis[ualization] idioms are effective at displaying it. (…) Don't just draw what you are given; decide what the right thing to show is, create it with a series of **transformations** from the original database, and draw*
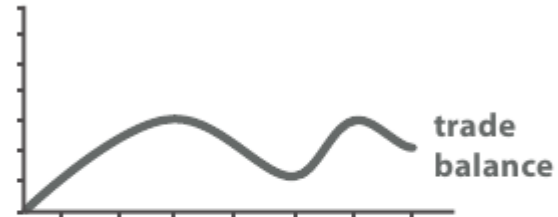
Original Data

Derived Data

trade balance = exports − imports

Derived attributes can be directly visually encoded. Munzner 2014, p.52

# Practice: meet our sample data

Download and open `data.xls` from Google Drive: fake data for online marketing goals and tools

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | source | quarter | spent | visits | income | goal |
| 2 | Adwords | 20160101 | 1000 | 50000 | 900 | 1500 |
| 3 | Twitter | 20160101 | 200 | 8500 | 1300 | 1000 |
| 4 | Facebook | 20160101 | 500 | 20000 | 800 | 1500 |
| 5 | Adwords | 20160401 | 1000 | 48000 | 1200 | 1500 |
| 6 | Twitter | 20160401 | 300 | 9000 | 1400 | 1000 |
| 7 | Facebook | 20160401 | 750 | 21500 | 1400 | 1500 |
| 8 | Adwords | 20160701 | 1000 | 50000 | 1500 | 1500 |
| 9 | Twitter | 20160701 | 400 | 10000 | 1000 | 1000 |
| 10 | Facebook | 20160701 | 750 | 23000 | 200 | 1500 |
| 11 | Adwords | 20161001 | 1000 | 45000 | 1250 | 1500 |
| 12 | Twitter | 20161001 | 500 | 11000 | 1000 | 1000 |
| 13 | Facebook | 20161001 | 1000 | 25000 | 2000 | 1500 |
| 14 | Adwords | 20170101 | 1000 | 50000 | 1100 | 1500 |
| 15 | Twitter | 20170101 | 500 | 8500 | 1300 | 1000 |
| 16 | Facebook | 20170101 | 1000 | 20000 | 800 | 1500 |
| 17 | Adwords | 20170401 | 1000 | 48000 | 1500 | 1500 |
| 18 | Twitter | 20170401 | 500 | 9000 | 1400 | 1000 |
| 19 | Facebook | 20170401 | 1000 | 21500 | 1400 | 1500 |
| 20 | Adwords | 20170701 | 1000 | 50000 | 1500 | 1500 |
| 21 | Twitter | 20170701 | 500 | 10000 | 1000 | 1000 |
| 22 | Facebook | 20170701 | 1000 | 23000 | 200 | 1500 |
| 23 | Adwords | 20171001 | 1000 | 45000 | 1250 | 1500 |
| 24 | Twitter | 20171001 | 400 | 11000 | 1000 | 1000 |
| 25 | Facebook | 20171001 | 1000 | 25000 | 2000 | 1500 |

Speaker notes

Shows results of campaigns in three different sources, per quarters (trimester). Money spent in each source, number of visits got from that campaign, income generated, and goal income for each source.

Takes for granted many things, such as validity/reasonability of goals, of spent money per goal…
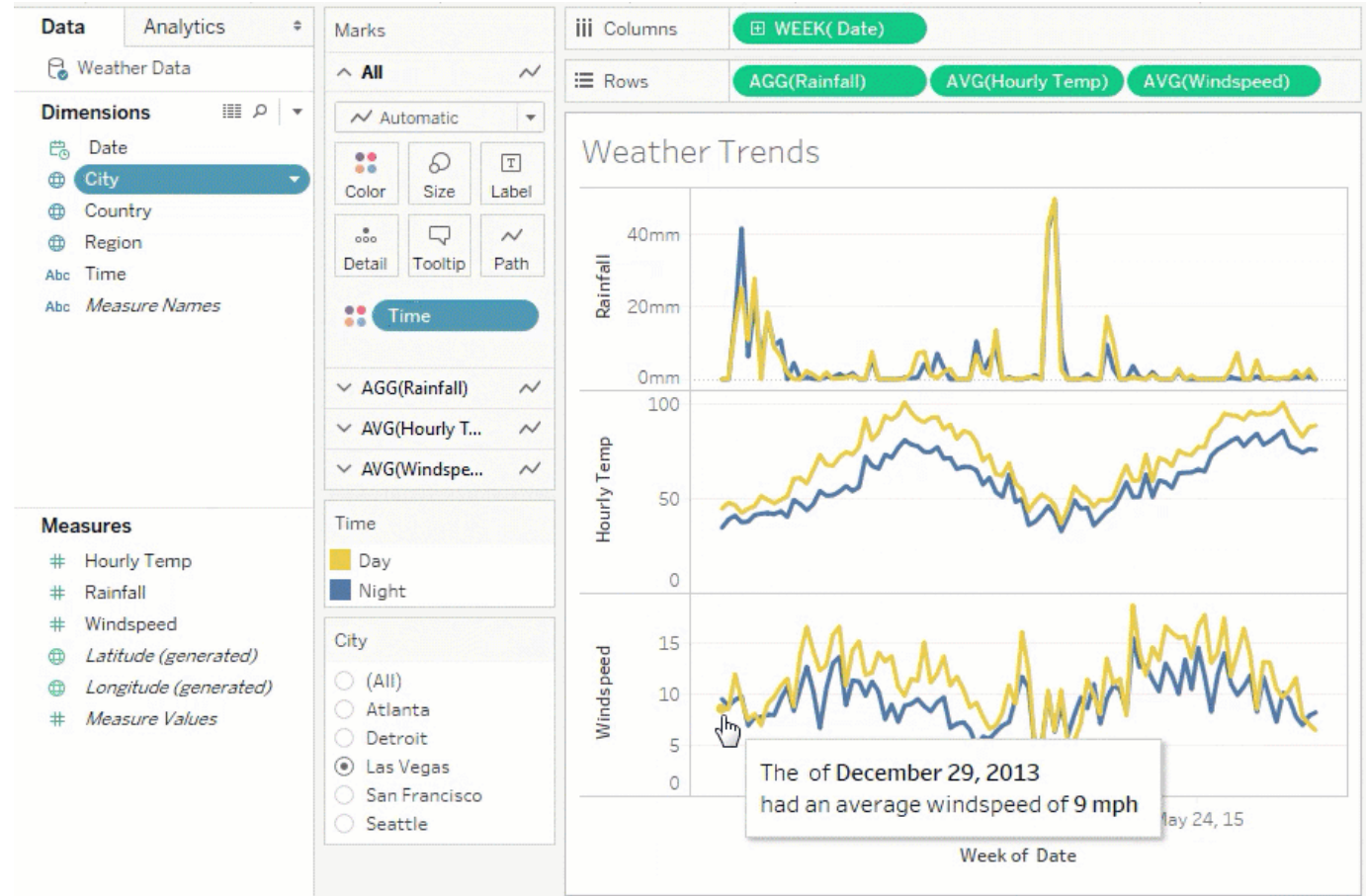
# Tableau Software



Tableau Software

# Tableau

## 1. Load data

- Possible to *merge* more sources (or **Connections**), only one for our practice – might require to define `New Union` to determine how merge takes place
- Shows the only sheet in the file, may be more
- `Connection > Live/Extract`, first one allows to modify source file and update work accordingly

Interprets the `source` field as Abc, rest as # (well done). Only, we would rather improve the `quarter` interpretation as `Date` instead of `Number`

Green & blue: **blue fields** are discrete, **green fields** are continuous. We'll see implications of this later on. For now, it's all good

Rename, hide… fields

The Excel > Tableau convert is quite straightforward (does a good job), more complex sources may require more manipulation at this point

This screen allows to preview how Tableau is interpreting our dataset: we are not allowed to change values, only variable interpretations

The source file (`data.xls`) is never modified

# 1. Graphics

ABERASTURI;1025880,88
ABETXUKO;1626,20
ABEZIA;183184,81
ABORNIKANO;54530,28
ACEBEDO;13519,09
ACOSTA;64930,00
ADANA;53139,42
AGIÑAGA;314344,94
ALAITZA;75534,95
ALBENIZ;61152,16
ALCEDO;21313,54
ALDA;27922,86
ALEGRIA-DULANTZI;142607
3,93



## Speaker notes

Dataset showing town name and amount of subsidies (funds) received in a given year. How do you think we could visualize this?

- visualization requires *mapping* or *translation* from *data* to *visual language* (idioms)
- this can be done in many ways
1. subsidy amount = bar height, color irrelevant
2. to show in map **transformation is required** (not possible directly)
- color IS relevant (subsidy amount = color saturation)
- comparison is more difficult: not *ordered* (althoug not impossible)
- allows for more direct insight (if context - knowledge of the area): Vitoria-Gasteiz not getting the higher amount
3. anything else? (i.e. pie chart, treemap…)

# 1.1 Reminder: variable types

- Quantitative
  - Continuous
  - Discrete
- Qualitative
  - Categorical
  - Ordinal

**A question of time**

Spatial and time/hour variables are special variable types. **Time variables** are specially complex:

- are there 365 days in every year? 30 days in every month? 24 hours in every day?
- *timezones* make it even more complex to use hours or time of day

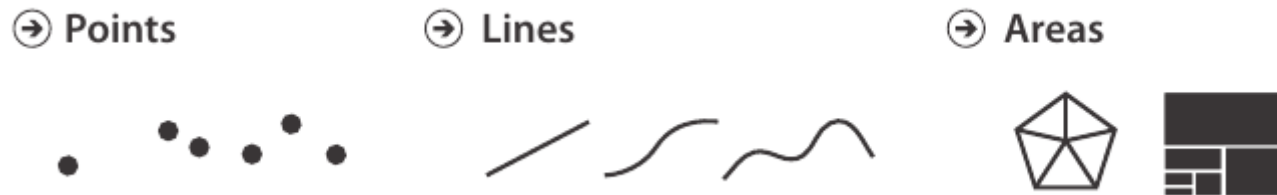Time may be used as a continuous or as a qualitative variable.

Speaker notes

- continuous variables can take on infinite/uncountable values (time, real numbers)

- discrete variables can take a finite/countable number of values (periods of time, integers)

- categorical/nominal: two or more categories, any order (gender, political party…). Cannot be quantified, do not allow arithmetic operations, cannot be assigned any order

- ordinal: allows for rank order (1st, 2nd, 3rd…). Can be dichotomous (yes/no, right/wrong, sick/healthy…) or non-dichotomous (spectrum of values: completely agree, mostly agree, mostly disagree, completely disagree)

# 1.2 Mapping variables to graphs

Understanding **marks and channels** provides the building blocks for analyzing visual encodings (Munzner 2014, p.95)

# 1.2.1 Marks

A **mark** is a basic graphical element in an image



Marks are geometric primitives (Munzner 2014, p.96)

Speaker notes

- a zero-dimensional (0D) mark is a point
- a one-dimensional (1D) mark is a line
- a two-dimensional (2D) mark is an area
- a three-dimensional (3D) mark is possible but not frequently used (will see why)

# 1.2.2 Channels

A visual **channel** is a way to control the appearance of marks

# 1.2.2 Channels

**One and only one** attribute/variable should be used per channel.

Multiple channels per attribute are possible (**redundant encoding**), but this approach has limitations.

# 1.2.2 Channels

The **size** and **shape** channels cannot be used on all types of marks, but most combinations are still possible:

- lines have two *size channels*: length + width
- points refer to location but can be *size* and *shape* coded

---

Speaker notes

Area marks cannot typically be size or shape coded: state or country already has a certain size and shape

- lines: if length is *taken* by a variable, it can't be used for another one, but width can be used to size code. They can be made wider on an individual basis to encode an additional attribute, or an entire set of bars can simply be made wider in a uniform way to be more visible
- points: intrinsically convey information only about position and are exactly the vehicle for conveying additional information through area and shape (and color!)

# 1.2.3 Channel types

Two kinds of sensory modalities:

1. **Identity**: what, where
2. **Magnitude**: how much

It does not make sense to ask magnitude questions for shape, color hue. We can ask about magnitudes with length, area or volume; color luminance or saturation;

Speaker notes

*The human perceptual system has two fundamentally different kinds of sensory modalities. The **identity** channels tell us information about what something is or where it is. In contrast, the **magnitude** channels tell us how much of something there is (Munzner 2014, p.99).*

# 1.2.4 Using marks and channels

All channels are not equal.

The selection of marks and channels should be guided by

- expresiveness: the visual encoding should express all of, and only, the information in the dataset attributes. Ordered data should be shown in a way that we perceive as ordered; unordered data SHOULD NOT be shown in a way that implies an ordering that does not exist.
- effectiveness: the importance of the channel (task abstraction, targets and actions) should match its salience, how noticeable it is. The most important attributes/variables should be encoded with the most effective channels, and less important attributes can be matched with less effective channels.

These can be combined to create a ranking of channels according to the type of data that is being visually encoded.

# 1.2.4 Using marks and channels



Channels: Expressiveness Types and Effectiveness Ranks

Speaker notes

Channels related to spatial position are at the top of both lists, and they are the only ones appearing on both lists (none of the others are effective for both data types). This primacy applies only to 2D positions, 3D depth is a much lower-ranked channel

# 1.2.4 Using marks and channels

The choice of **which attributes/variables to encode with position** is the most central choice in visual encoding.

# 1.2.4 Using marks and channels



Cleveland & McGill's Results

Crowdsourced Results

# 1.3 So, which graph?



Chart Suggestions—A Thought-Starter

Left column, `Data` tab:

- dimensions contain qualitative values
- measures contain numeric, quantitative values (you can apply calculations to them and aggregate them)

Reminder (blue *vs.* green, learn more):

- green = continuous. Its values are treated as an infinite range. Generally, continuous fields add axes to the view.
- blue = discrete. Its values are treated as finite. Generally, discrete fields add headers to the view.

Mapping variables to graphs is done, in Tableau, by dragging items from the left column to either: `Columns/Rows`, or `Marks`.

`Show me` tab displays the most common graphs and the minimum requirement of data for each graph. If graphs = recipes, lets us know which ingredients are required to cook a certain recipe.

Also works the other way round: if I select my ingredients from the left bar, the `Show me` tab highlights the recipes available for those ingredients.

(test some graphs) - Source only > table - add spent > table, pie chart etc

Fields can be used more than once (redundant coding) for easier identification, by dragging to more than one place. barchart income (rows) + source (columns), source to colour

Options in `Show me` are different from the `Marks` box. Show me only allows graphs that *make sense*, very basic and tested graphs. Marks allows to *force* graphs in case we need more complex displays. USE WITH CAUTION, possible to create graphs that make no sense or are misleading (change to "line" in Marks box)
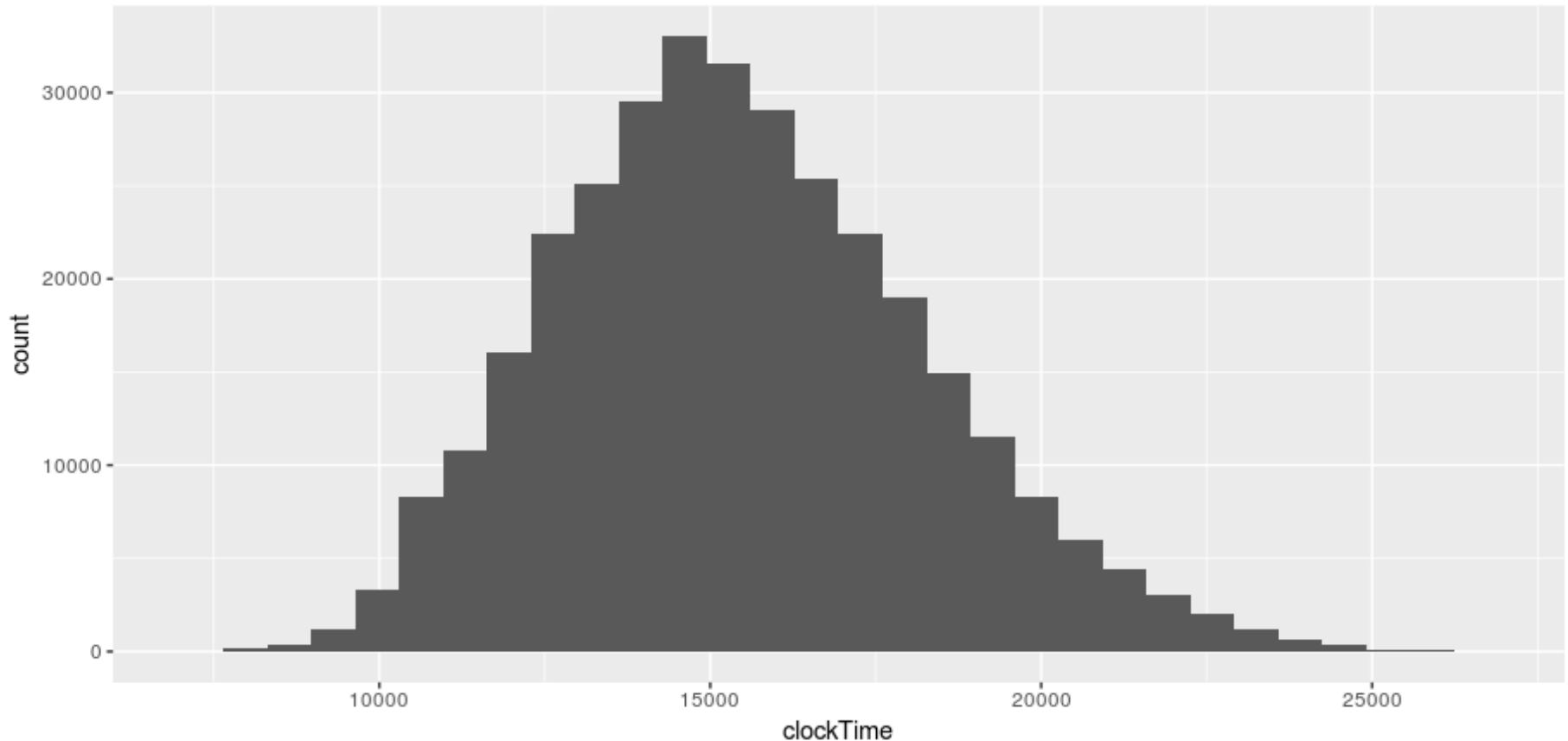
# 2. Provide easier analysis

# Section outline

How can we enable easier insight through data visualization?

1. Change default settings
2. Make simpler graphs
3. Highlight observations
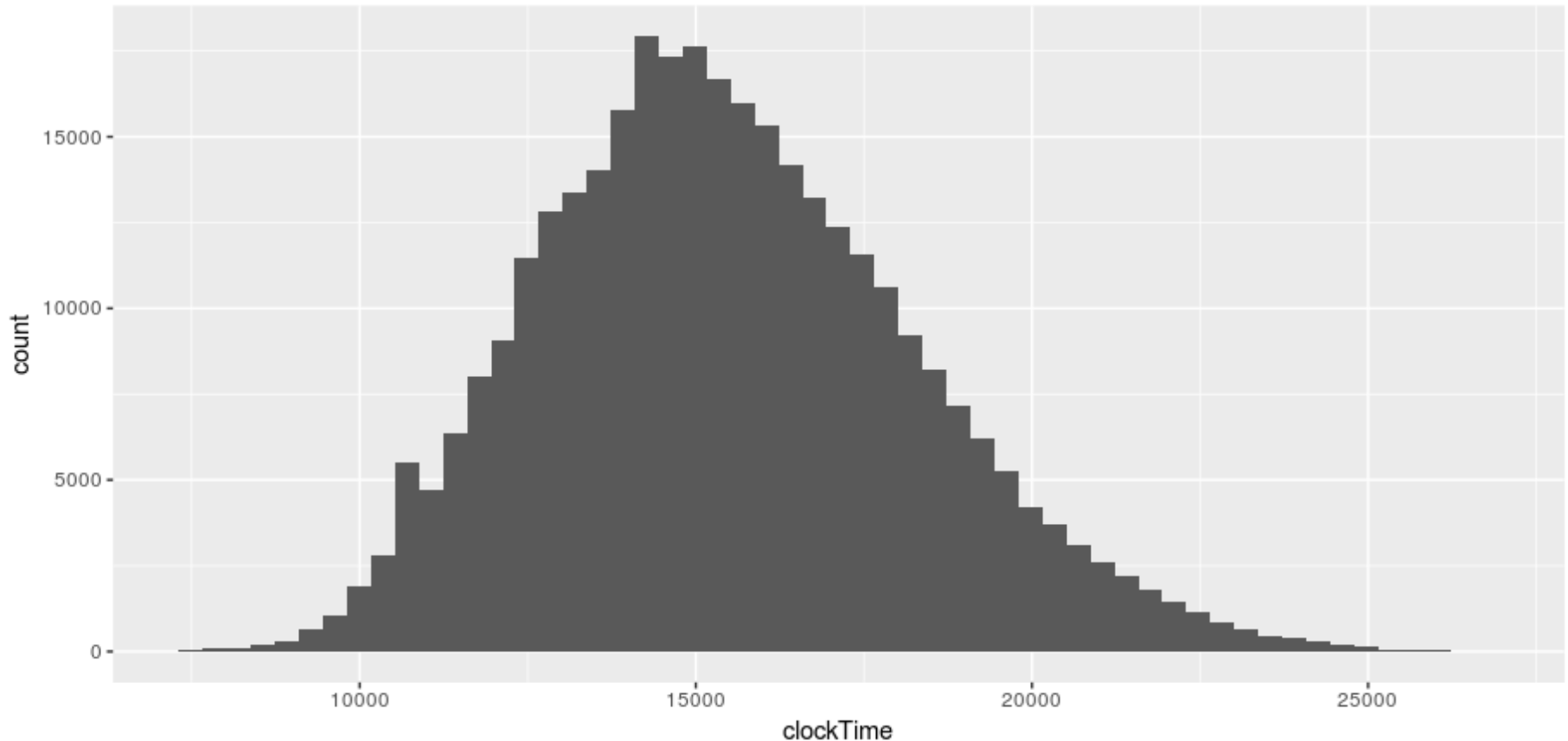4. Add attributes as context
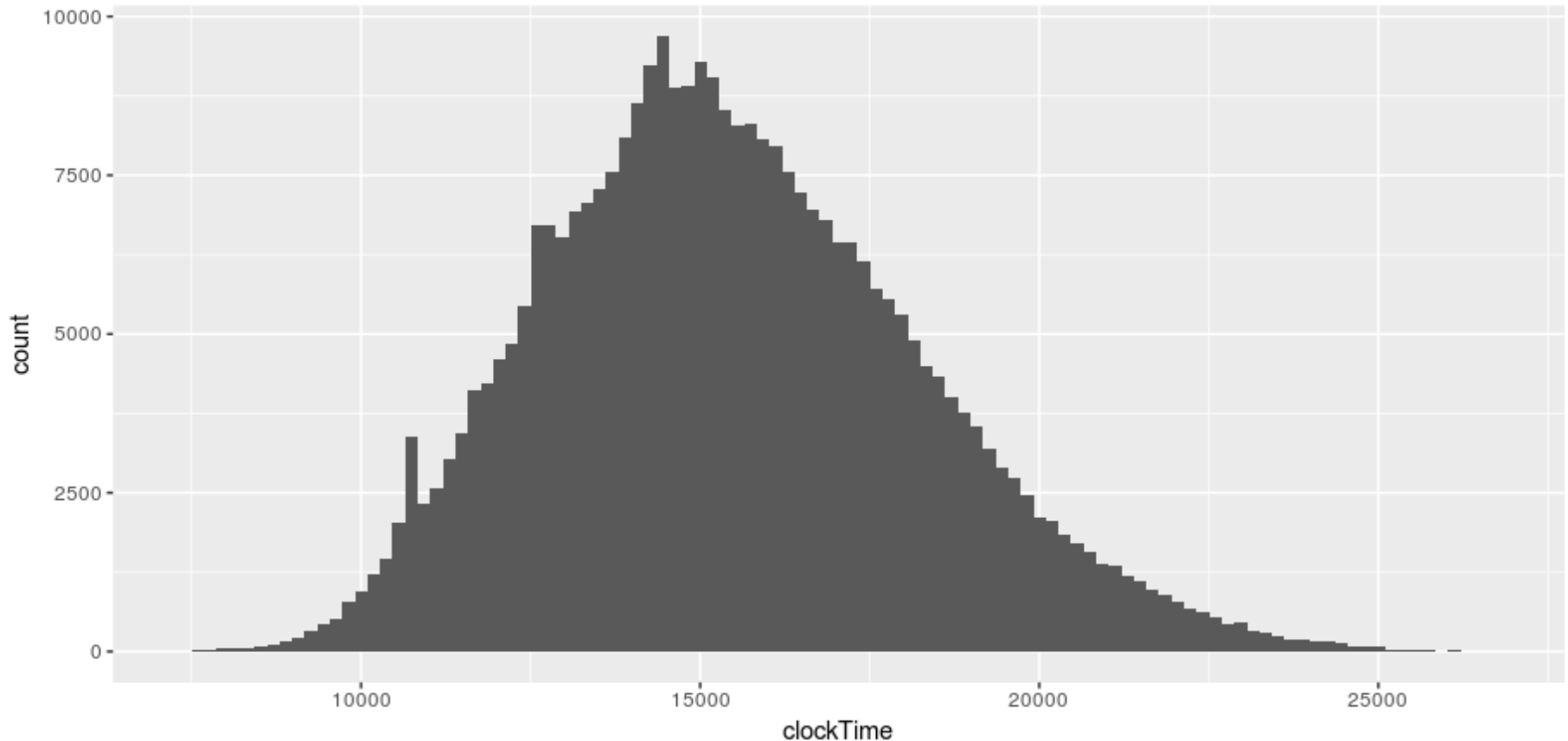5. Add statistical information

# 2.1 Change default settings

# 2.1 Change default settings



Speaker notes

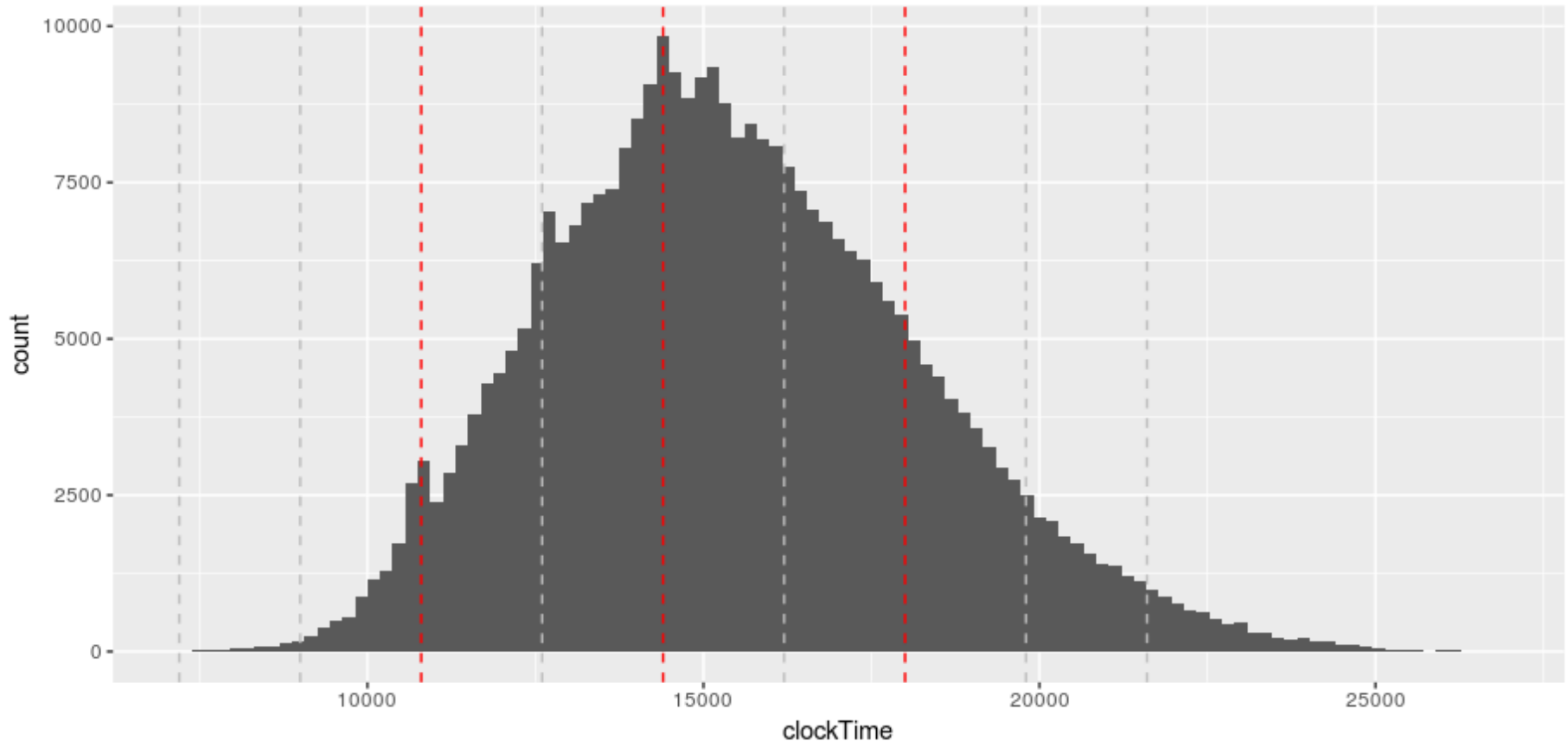More detail (less aggregation) allows to see some highlighted bars

# 2.1 Change default settings

# 2.1 Change default settings

# 2.2 Make simpler graphs

*Data-ink is the non-erasable core of the graphic, the non-redundant ink arranged in response to variation in the numbers represented.*

Speaker notes

we should remove all non-data-ink and redundant data-ink, within reason, to increase the data-ink-ratio and create a sound graphical design.

some redundancy is often more effective, however, most graphics don't struggle with understatement. In fact, most contain a stunning amount of excess ink (or pixels). Rather than dressing our data up we should be stripping it down.

# 2.2 Make simpler graphs

Speaker Deck
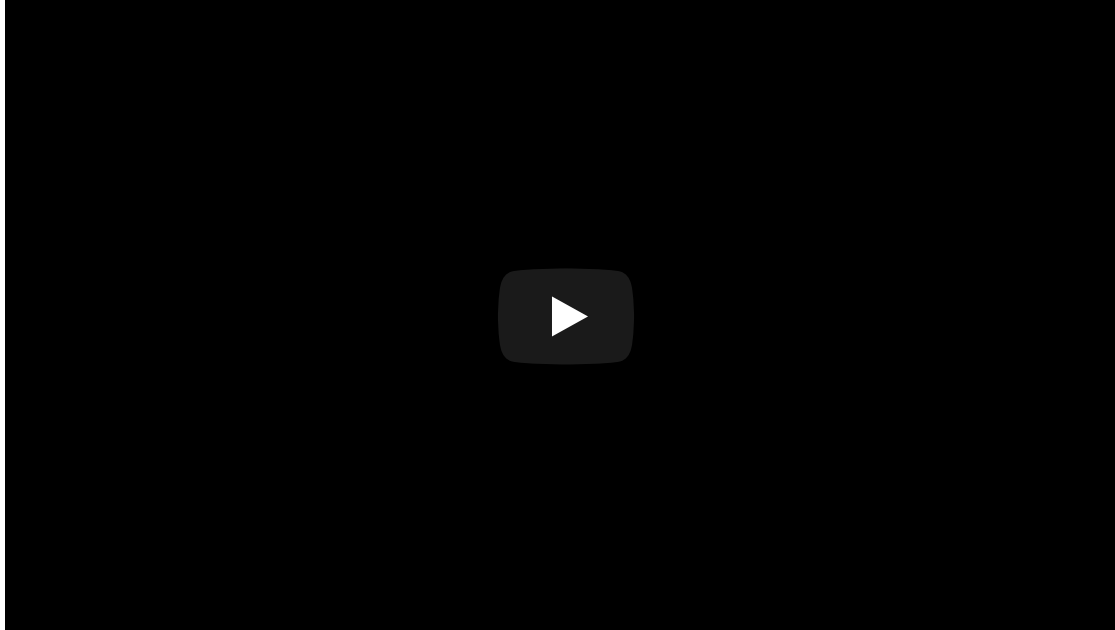
Talk by Joey Cherdarchuk

Full Screen

Speaker notes

In it we start with a chart, similar to what we've seen in many presentations, and vastly improve it with progressive deletions and no additions.

The next time you are trying to improve a chart, consider what you can take away rather than what you can add.

# 2.2 Make simpler graphs

More on decluttering:



Nussbaumer, Declutter Your Data Visualizations

# 2.3 Highlight observations

Through preattentive attributes:

- they are processed in spatial memory without our conscious action
- make it easier to understand what is represented through a design: saves from consciously processing data

Speaker notes

A preattentive visual property is one which is processed in spatial memory without our conscious action. In essence it takes less than 500 milliseconds for the eye and the brain to process a preattentive property of any image. These properties can be harnessed to make it easier for a user to understand what is presented through the design and save them from consciously processing all the data presented in short-term memory which requires more effort.

# 2.3 Highlight observations

7563950 8473
65866303 7576
86037265 8602
84658910 7830

FIGURE 4.2  Count the 3s example

Nussbaumer 2015, p.103

# 2.3 Highlight observations

7563950847**3**
6586**3**0**3**7576
860**3**72658602
846589107830

FIGURE 4.3 Count the 3s example with preattentive attributes

Nussbaumer 2015, p.104

# 2.3 Highlight observations



Orientation · Shape · Line length · Line width · Size · Curvature · Added marks · Enclosure
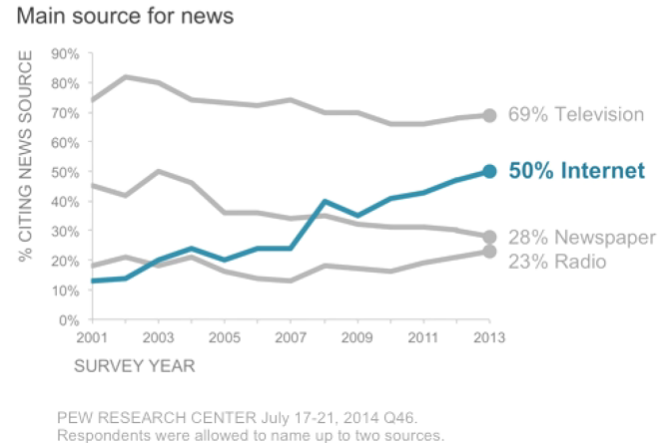
# 2.3 Highlight observations



**1. More Americans get news online...** 50% of the public now cite source for national and international news 🐦, still below television, newspapers and radio. (Report)

**Main Source for News**

Television
74    66    69
Newspaper
45
Internet    50
43
31    28
18
13    Radio    19    23
2001  2003  2005  2007  2009  2011  2013

PEW RESEARCH CENTER July 17-21, 2013. Q46.
Respondents were allowed to name up to two sources.

More Americans get news online

50% of the public cite the **internet** as a main source for national & international news. This remains below television, but is far above newspapers and radio.

Main source for news

% CITING NEWS SOURCE
90%
80%
70%    69% Television
60%
50%    **50% Internet**
40%
30%    28% Newspaper
20%    23% Radio
10%
0%
2001  2003  2005  2007  2009  2011  2013
SURVEY YEAR

PEW RESEARCH CENTER July 17-21, 2014 Q46.
Respondents were allowed to name up to two sources.

Source: http://www.pewresearch.org/fact-tank/2013/10/16/12-trends-shaping-digital-news/

storytelling·data

Speaker notes

After decluttering, more *channels* are available for highlight. This allows to

1. draw your audience's attention quickly to where you want them to look, and
2. create a visual hierarchy of information

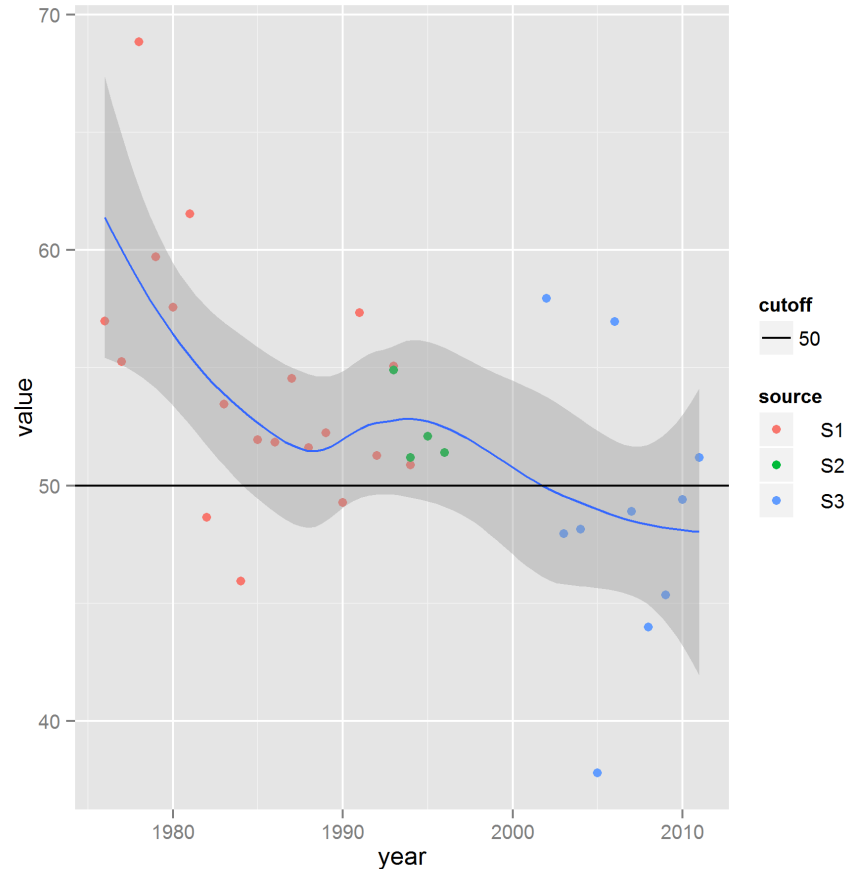# 2.4 Add variables (as context)

- Adding preexisting variables (in moderation)
- Creating conditional variables from preexisting variables
  - binaries or with few levels are best
  - example of calculated field or variable: weekend date

Speaker notes

useful only if addition conveys meaning / enables insight. In the graph, the colour coding is explanatory to the left, but it is not to the right.

# 2.5 Add statistical information

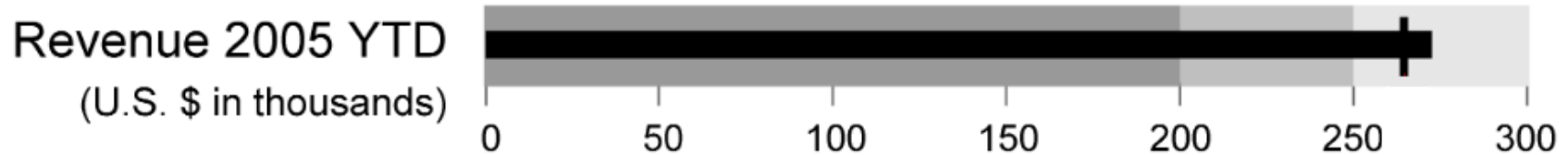- statistical summaries (mean, variance)
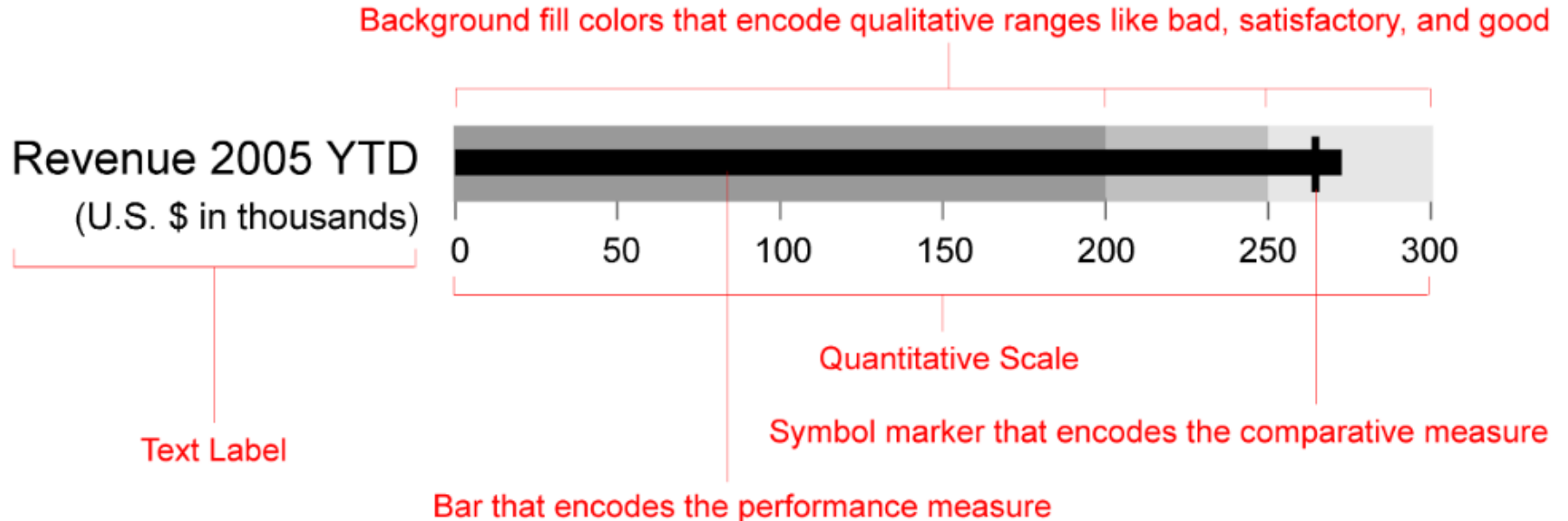- models



source

# Tableau: (not so) basic graphs



**93.39**

Sparklines (Tufte 2006)

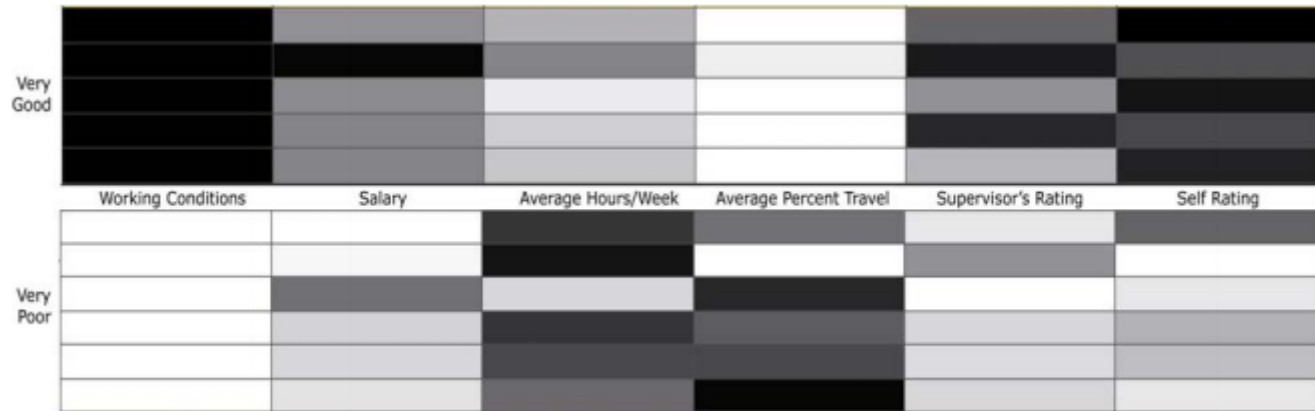# Tableau: (not so) basic graphs



Bulletgraphs (Few 2007)

# Tableau: (not so) basic graphs



Background fill colors that encode qualitative ranges like bad, satisfactory, and good

Revenue 2005 YTD
(U.S. $ in thousands)

0    50    100    150    200    250    300

Quantitative Scale

Symbol marker that encodes the comparative measure

Text Label

Bar that encodes the performance measure

A standard bullet graph with each of its parts labeled.

Bulletgraphs (Few 2007)

# Tableau: (not so) basic graphs



Heatmaps (Few 2006)

# 3. Dashboards

# 3.1 Dashboards for *situation awareness*

> *The term "dashboard" refers to a single screen information display that is used to monitor what's going on in some aspect of the business.*

Few (2007), Dashboard Design

# 3.1 Dashboards for *situation awareness*

- Perception of own's environment
- Comprehension of it's meaning
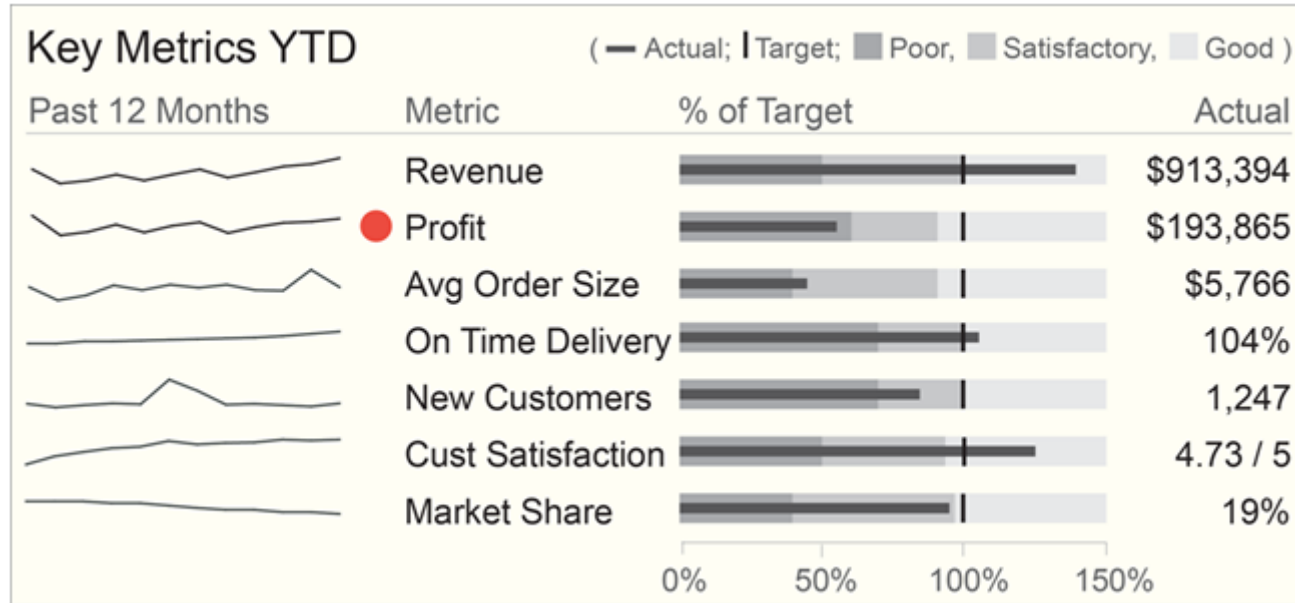- Projection of that understanding into the future

# 3.2 Dos: Principles you should follow

- Use flicker and sound to grab attention
- Encourage active thinking about the data, not just passive reaction to alarms
- Don't over-automate actions to the point where people become disengaged
- Provide smooth and simple means to respond
- Provide a common picture for the whole team
- Support projections for proactive responses
- Match the mental model

# 3.3 ...and Don'ts: Design problems you should avoid

- Too much complexity
- Too many alert conditions
- Alerts that cannot be diff erentiated
- Overwhelming visuals
- Distracting visuals
- Inappropriate visual salience
- Mismatch between information and its visual representation
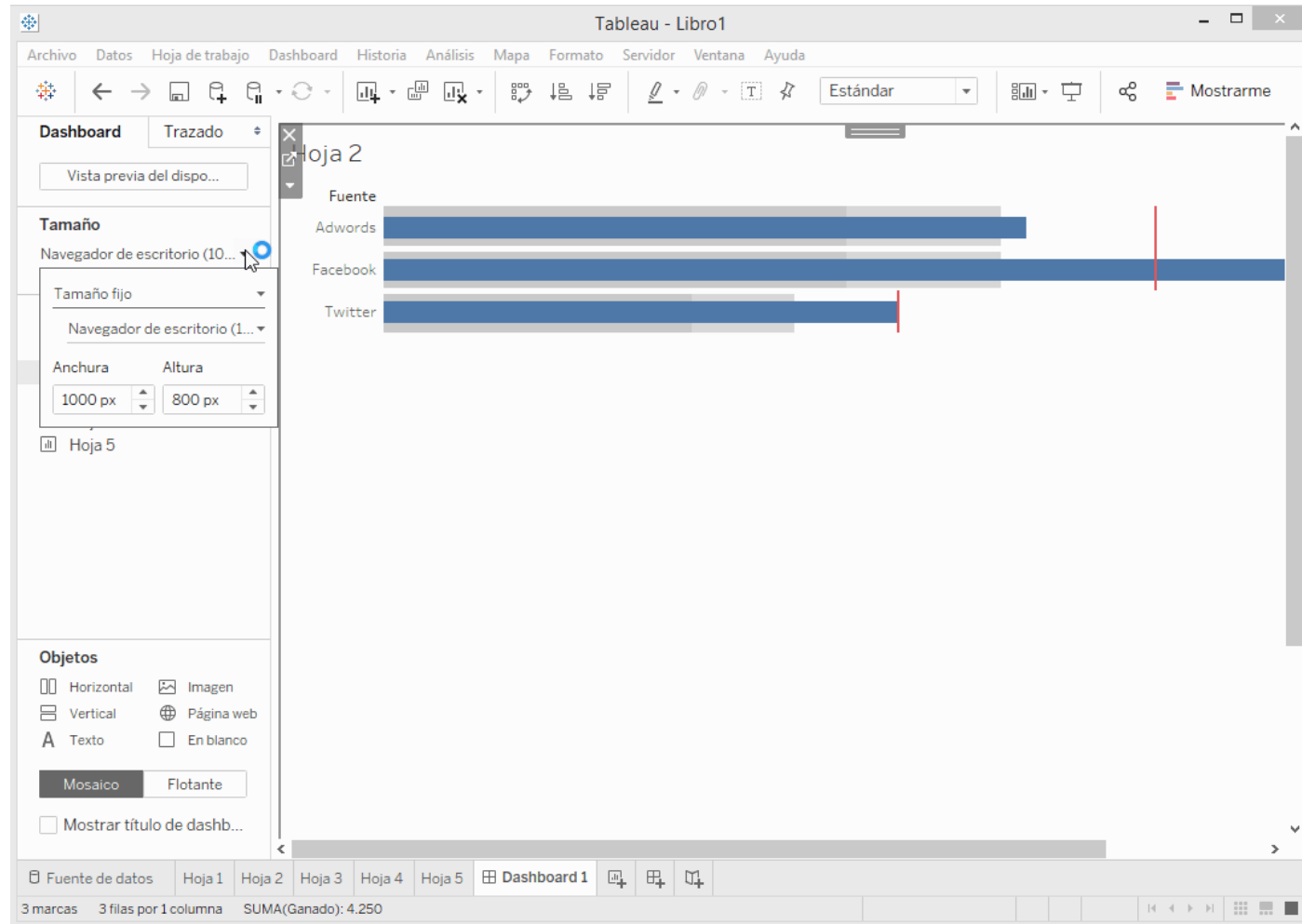- Indirect expression of measures
- Not enough context

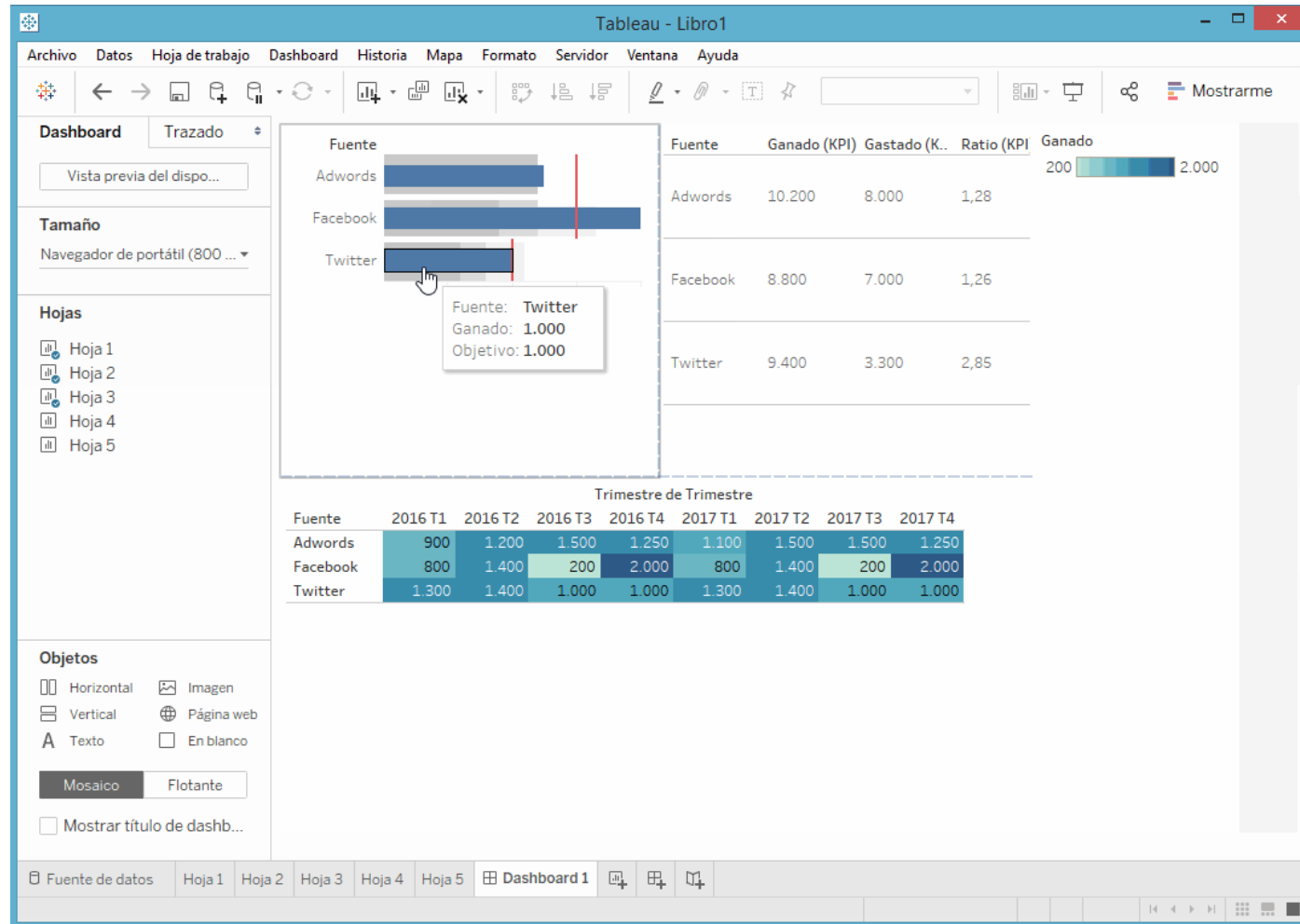# 3.4 Few's few examples



Few 2007

# Dashboards in Tableau

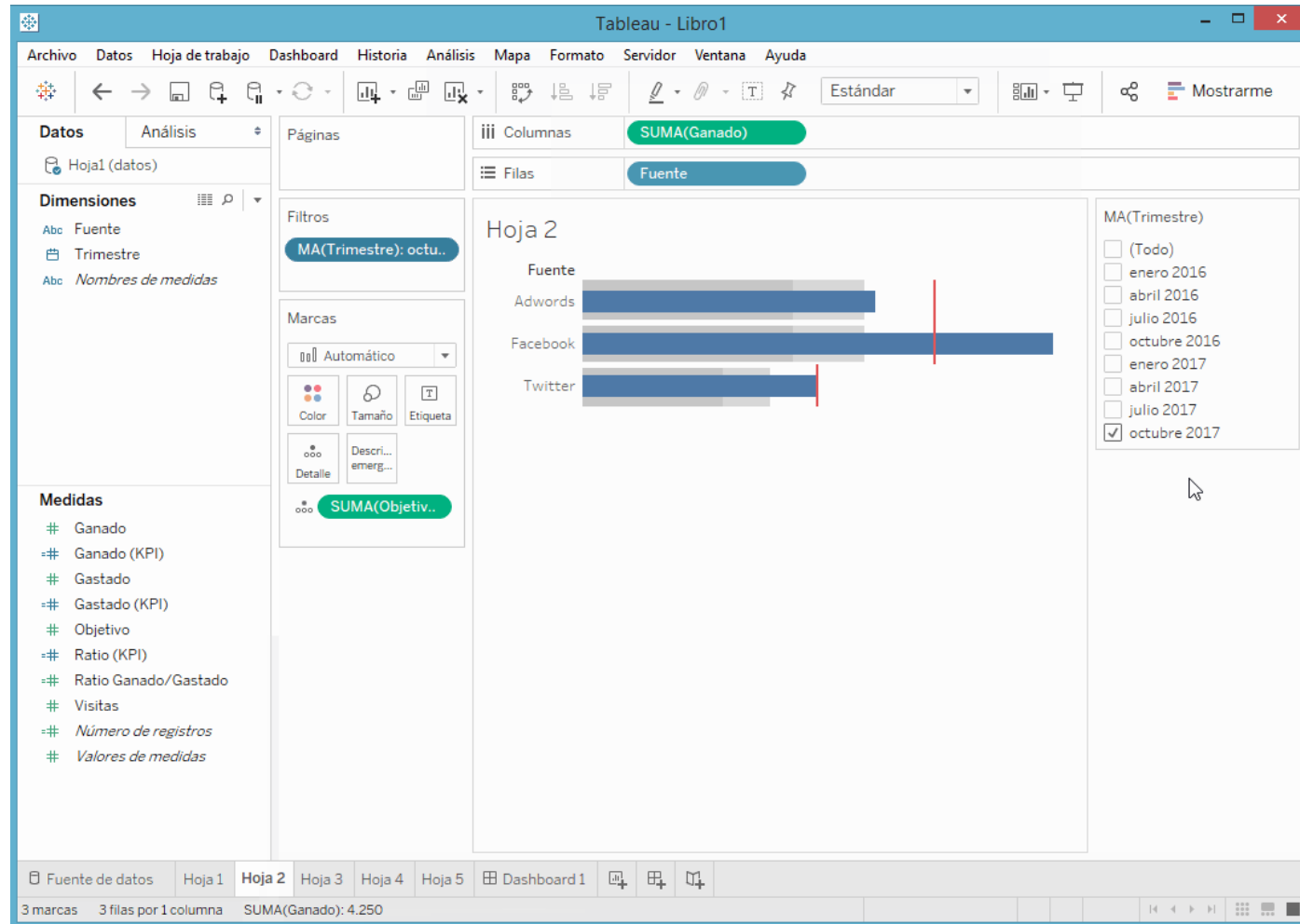# Tableau 2.1: basic dashboard



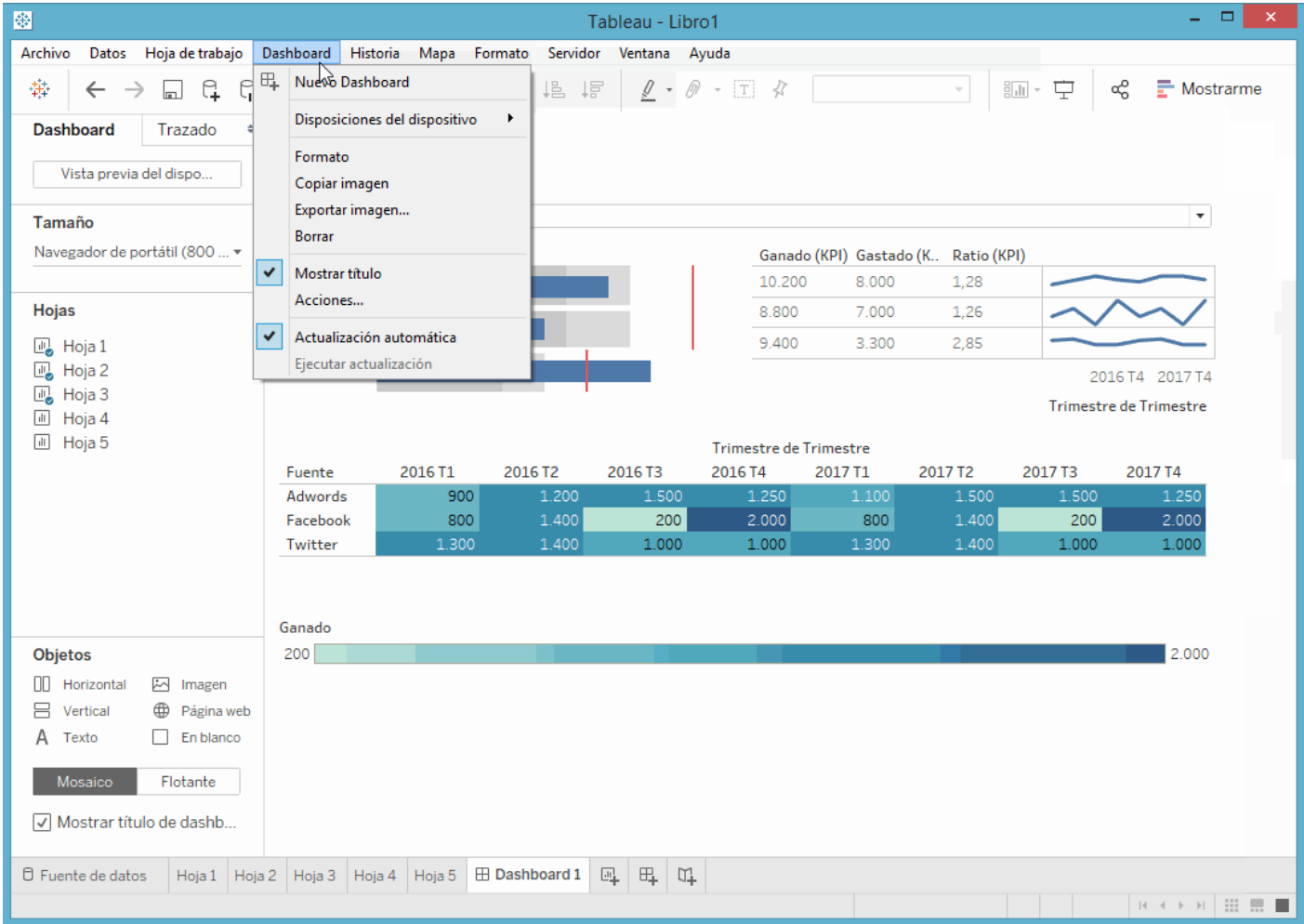Basic dashboard

# Tableau 2.2: basic formating



Basic formating

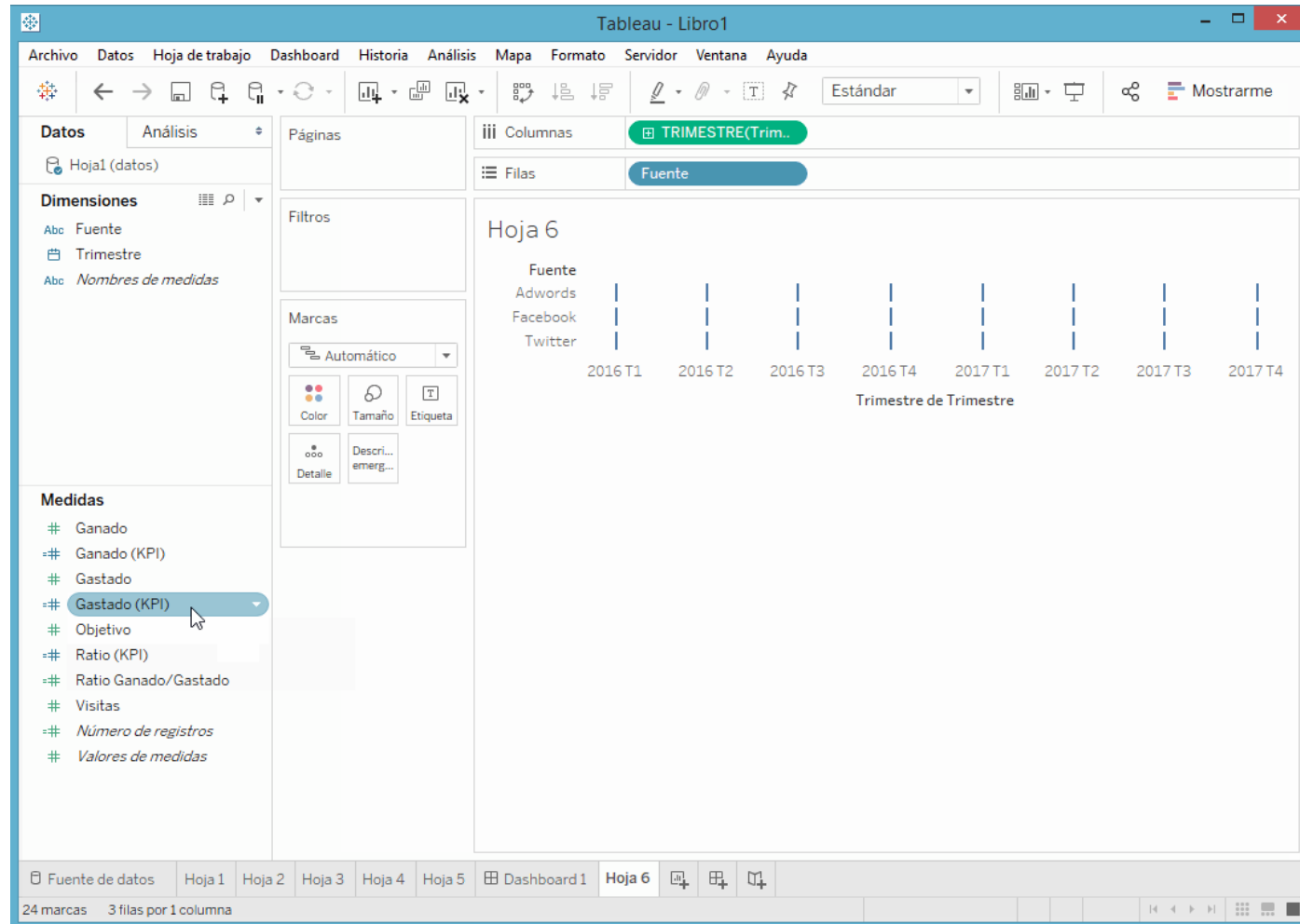# Tableau 2.3: show filters



Show filters

# Tableau 2.4: highlight action



Highlight action

# Tableau 2.5: filter action



Filter action

# Epilogue

# References

Abela, Andrew (2006). Choosing a good chart.

Kirk, Andy (2016). *Data Visualisation: A Handbook for Data Driven Design*. SAGE: London 316.763 K 63 a

Munzner, Tamara (2015). *Visualization Analysis and Design*. CRC Press: Boca Raton, Florida 316.763 M 92 t

Tufte, Edward R. (1983). *The Visual Display of Quantitative Information*. Graphics Press: California 316.763 T 87 e

# Thank you!

Miren Berasategi

miren.berasategi@deusto.es

# License