

Graph Construction

Visualizzazione dell'Informazione Quantitativa



SoftEng
<http://softeng.polito.it>

<https://softeng.polito.it/courses/VIQ>

Version 3.1.0
© Marco Torchiano, 2020



This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License.

To view a copy of this license, visit
<http://creativecommons.org/licenses/by-nc-nd/4.0/>.

You are free: to copy, distribute, display, and perform the work

Under the following conditions:

 **Attribution.** You must attribute the work in the manner specified by the author or licensor.

 **Non-commercial.** You may not use this work for commercial purposes.

 **No Derivative Works.** You may not alter, transform, or build upon this work.

- For any reuse or distribution, you must make clear to others the license terms of this work.
- Any of these conditions can be waived if you get permission from the copyright holder.

Your fair use and other rights are in no way affected by the above.

Grammar of Graphics

- Theory behind graphics construction
 - ◆ Separation of data from aesthetic
 - ◆ Definition of common chart elements
 - ◆ Composition of such elements
- Building a graphic involves
 1. Specification
 2. Assembly
 3. Display

Leland Wilkinson, *The grammar of graphics*

3

Specification

- **DATA**: a set of data operations that create variables from datasets
 - ◆ Link variables (e.g., *by index* or *id*)
- **TRANS**: variable transformations (e.g., *rank*)
- **SCALE**: scale transformations (e.g., *log*)
- **COORD**: a coordinate system (e.g., *polar*)
- **ELEMENT**: visual objects (e.g., *points*) and their aesthetic attributes (e.g., *color*, *position*)
- **GUIDE**: guides (e.g., *axes*, *legends*)

4

Specification for a scatter plot

- DATA: $x = x$
 - DATA: $y = y$
 - TRANS: $x = x$
 - TRANS: $y = y$
 - SCALE: $\text{linear}(\text{dim}(1))$
 - SCALE: $\text{linear}(\text{dim}(2))$
 - COORD: $\text{rect}(\text{dim}(1, 2))$
 - ELEMENT: $\text{point}(\text{position}(x*y))$
 - GUIDE: $\text{axis}(\text{dim}(1))$
 - GUIDE: $\text{axis}(\text{dim}(2))$
-

5

Graph visual components

- Data components
 - ◆ Visual objects associated to measures
 - ◆ Visual attributes
- Layout
 - ◆ Positioning rules (e.g. cartesian coord)
- Support components
 - ◆ Axes
 - ◆ Labels
 - ◆ Legends

6

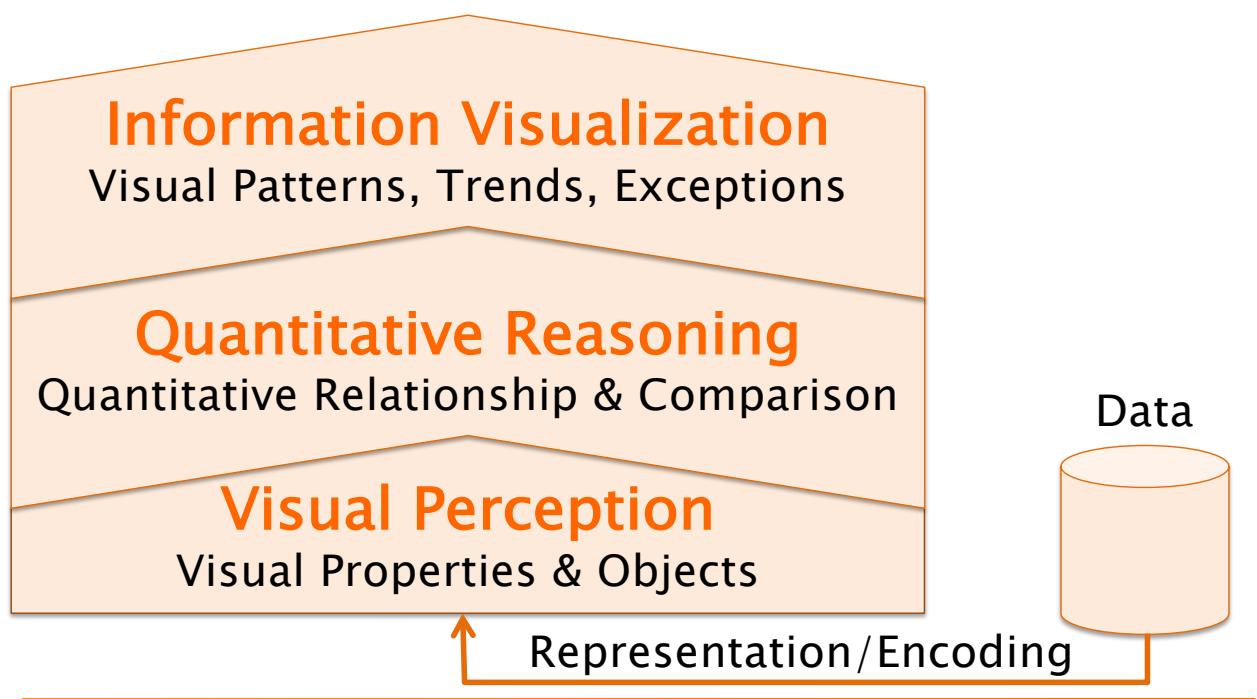
Visualizzazione dell'Informazione Quantitativa

VISUAL RELATIONSHIPS

7

Data Visualization

Understanding



Visual Encoding

- Given a variable (measure), identify:
 - ◆ Visual object
 - ◆ Visual attribute
- Main distinction
 - ◆ Quantitative (interval, ratio, absolute)
 - ◆ Categorical (nominal, ordinal)

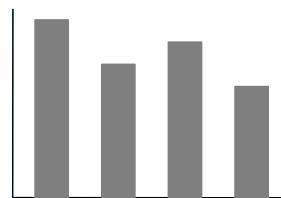
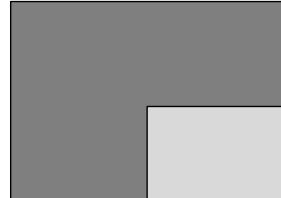
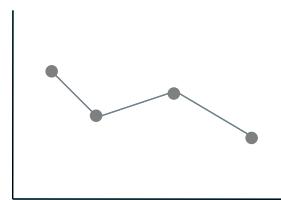
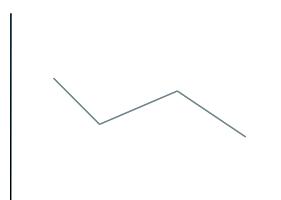
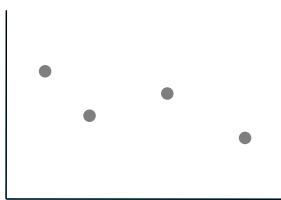
9

Relationships

- Within a category
 - ◆ Nominal comparison
 - ◆ Ranking
 - ◆ Part-to-whole
 - ◆ Distribution
- Between measures
 - ◆ Time series
 - ◆ Deviation
 - ◆ Correlation

10

Quantitative encoding



11

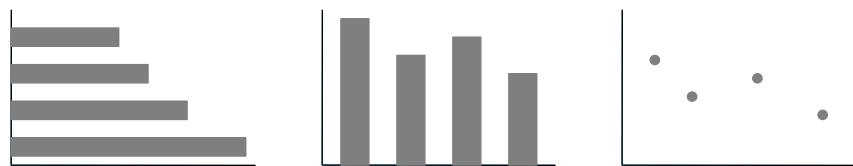
Sample data

Region	turnout (2018)	turnout (2013)	Region	turnout (2018)	turnout (2013)
ABRUZZO	75.3%	75.9%	MOLISE	71.6%	78.1%
BASILICATA	71.1%	69.5%	PIEMONTE	75.2%	77.3%
CALABRIA	63.6%	63.1%	PUGLIA	69.1%	69.9%
CAMPANIA	68.2%	67.9%	SARDEGNA	65.5%	68.5%
EMILIA-ROMAGNA	78.3%	82.1%	SICILIA	62.8%	64.6%
FRIULI-VENEZIA GIULIA	75.1%	77.2%	TOSCANA	77.5%	79.2%
LAZIO	72.6%	77.5%	TRENTINO-ALTO ADIGE	74.3%	81.0%
LIGURIA	72.0%	75.1%	UMBRIA	78.2%	79.5%
LOMBARDIA	76.8%	79.6%	VALLE D'AOSTA	72.3%	77.0%
MARCHE	77.3%	79.8%	VENETO	78.7%	81.8%

12

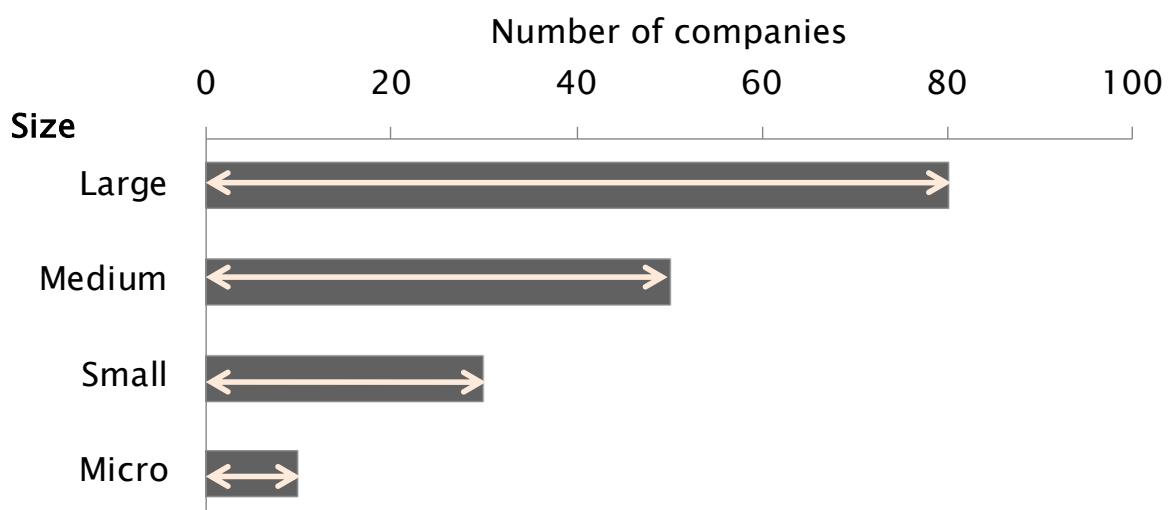
Nominal comparison

- Compare quantitative values corresponding to categorical levels
 - ◆ Small differences are difficult to see
 - Non zero-based scale can emphasize
 - ◆ Dot plots can be used for small differences
 - They do not require zero based scale



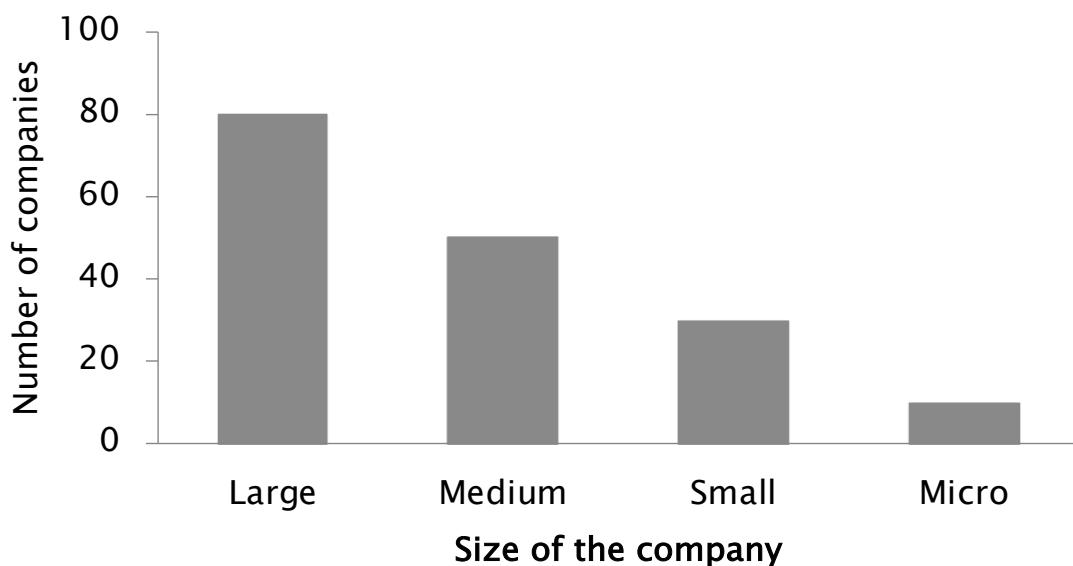
13

Line length – Bars chart



14

Vertical Bars (aka Columns)



15

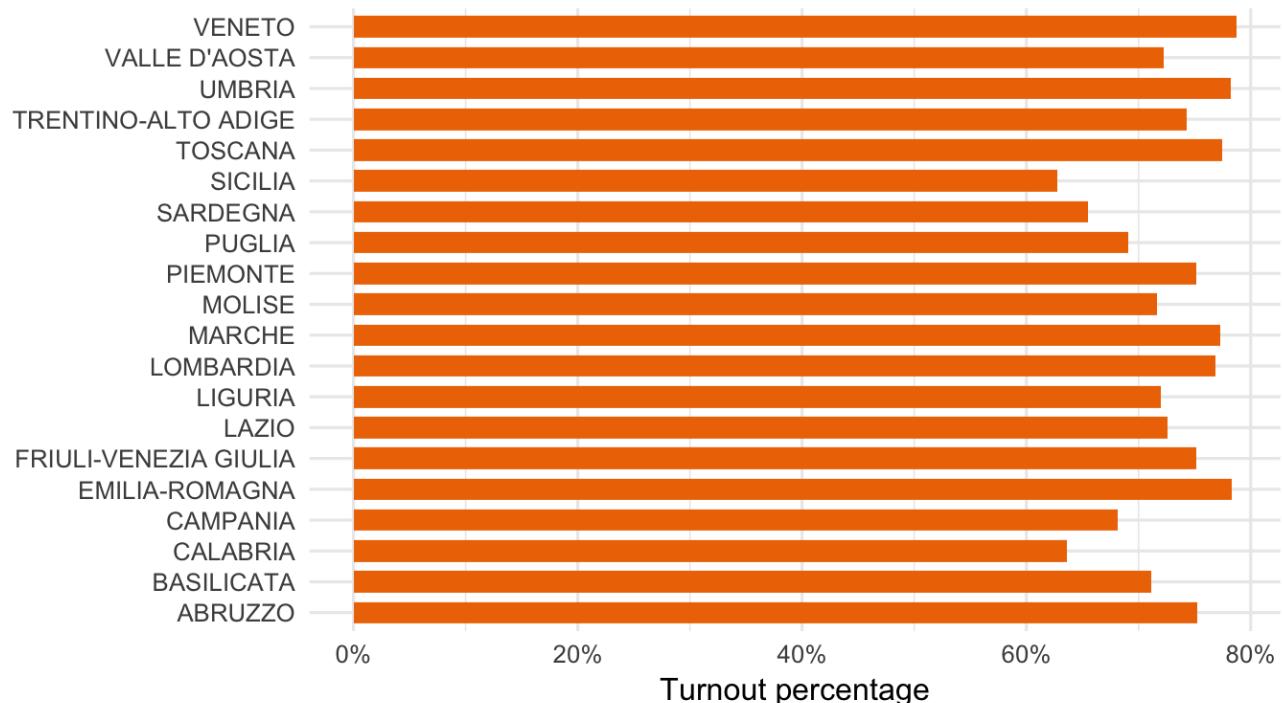
Bar charts

- Categorical values are encoded as position along an axis
- Quantitative values are encoded only as length of the bars
 - ◆ The axis is a supporting element
- Width of bars plays no role
 - ◆ Bars are just very thick lines
- Bars require a zero-based scale
 - ◆ See: Lie factor!

16

Comparison – Barplot

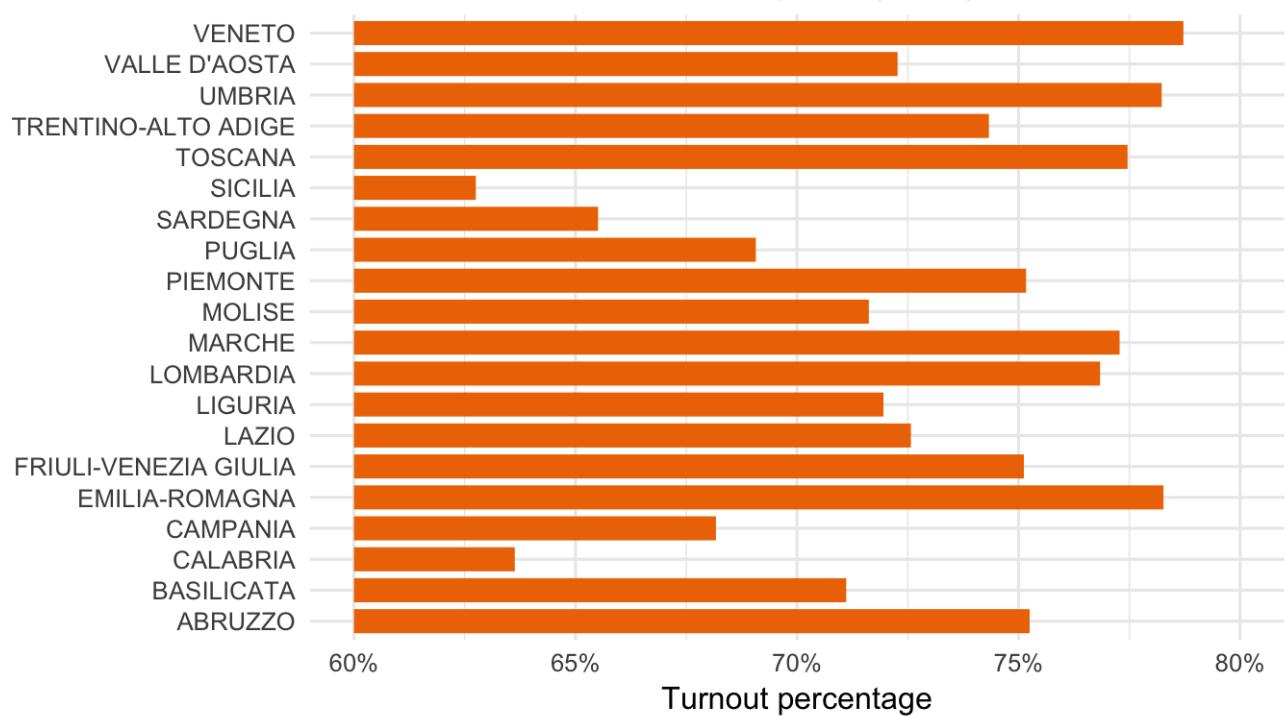
Electoral turnout in italian regions (2018)



17

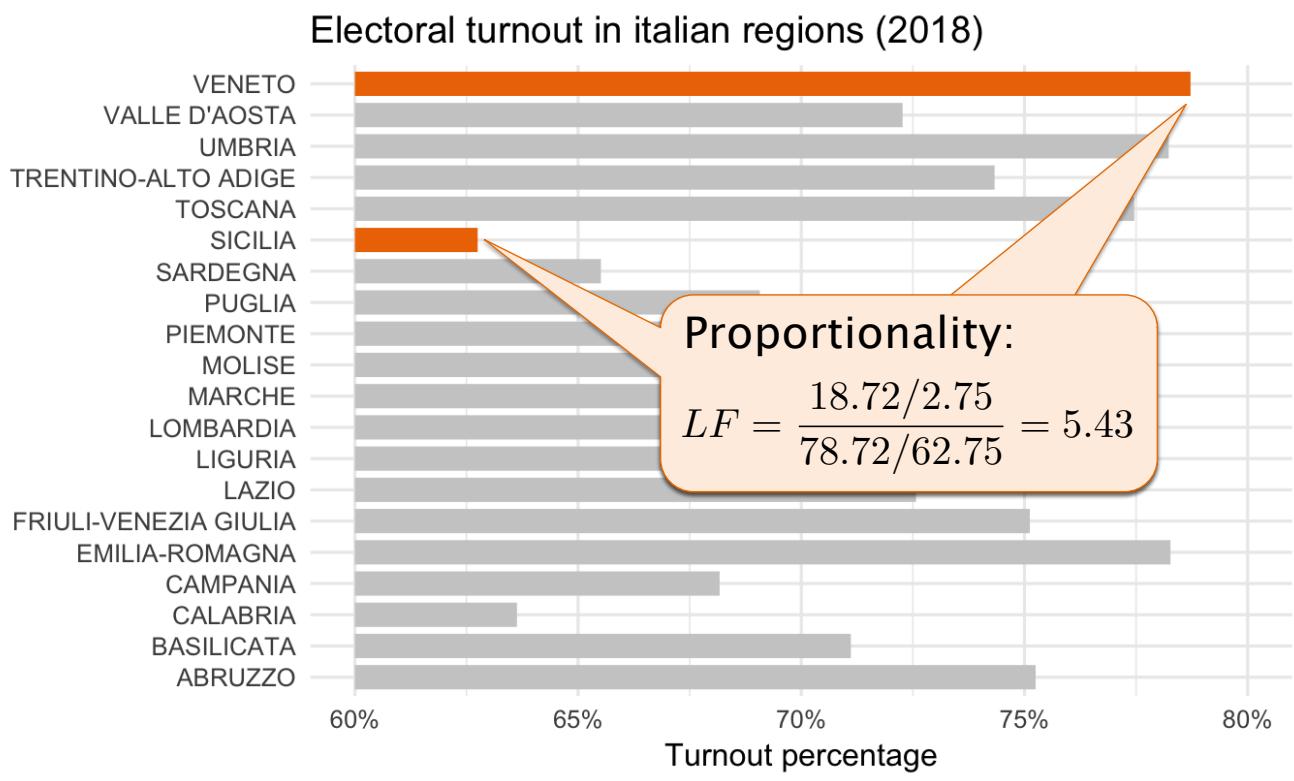
Barplot (non zero based scale)

Electoral turnout in italian regions (2018)



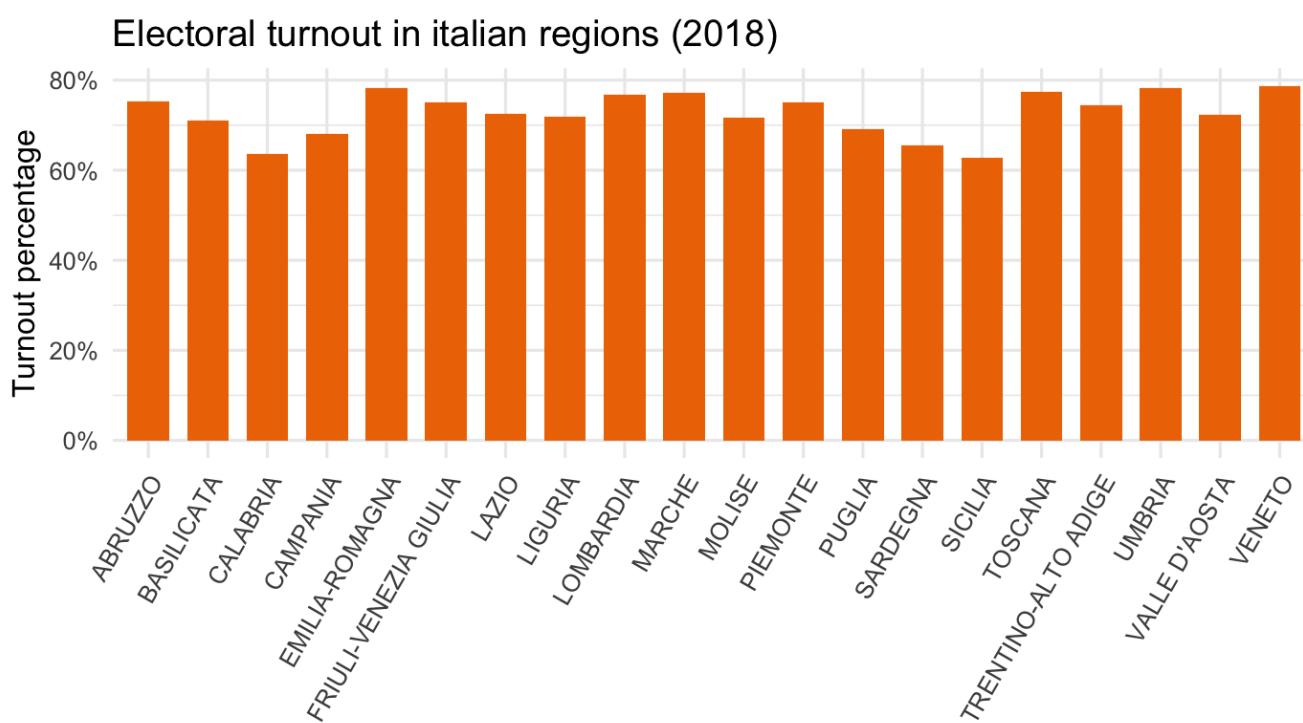
18

Barplot (non zero based scale)



19

Barplot vertical labels



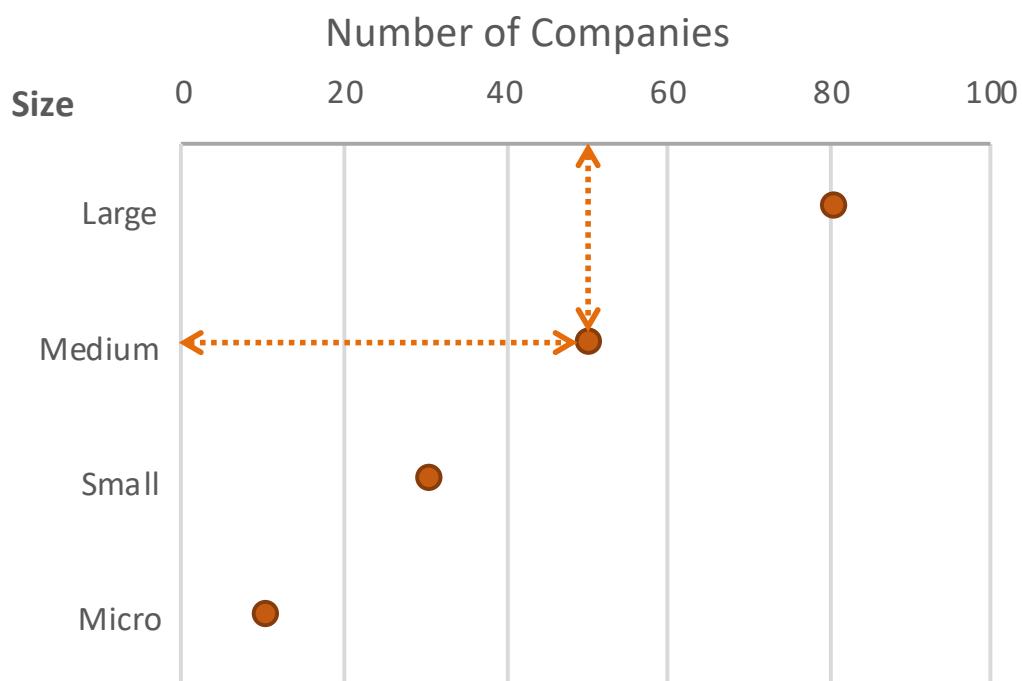
20

Bars Guidelines

- Use horizontal bars when
 - ◆ A descending order ranking
 - ◆ Categorical label don't fit
- Proximity
 - ◆ Use a 1:1 bar:spacing ratio $\pm 50\%$
 - ◆ No spacing between bars that are not labeled on the axis (legend categories)
 - ◆ No overlapping bars

21

Position – Dots plot



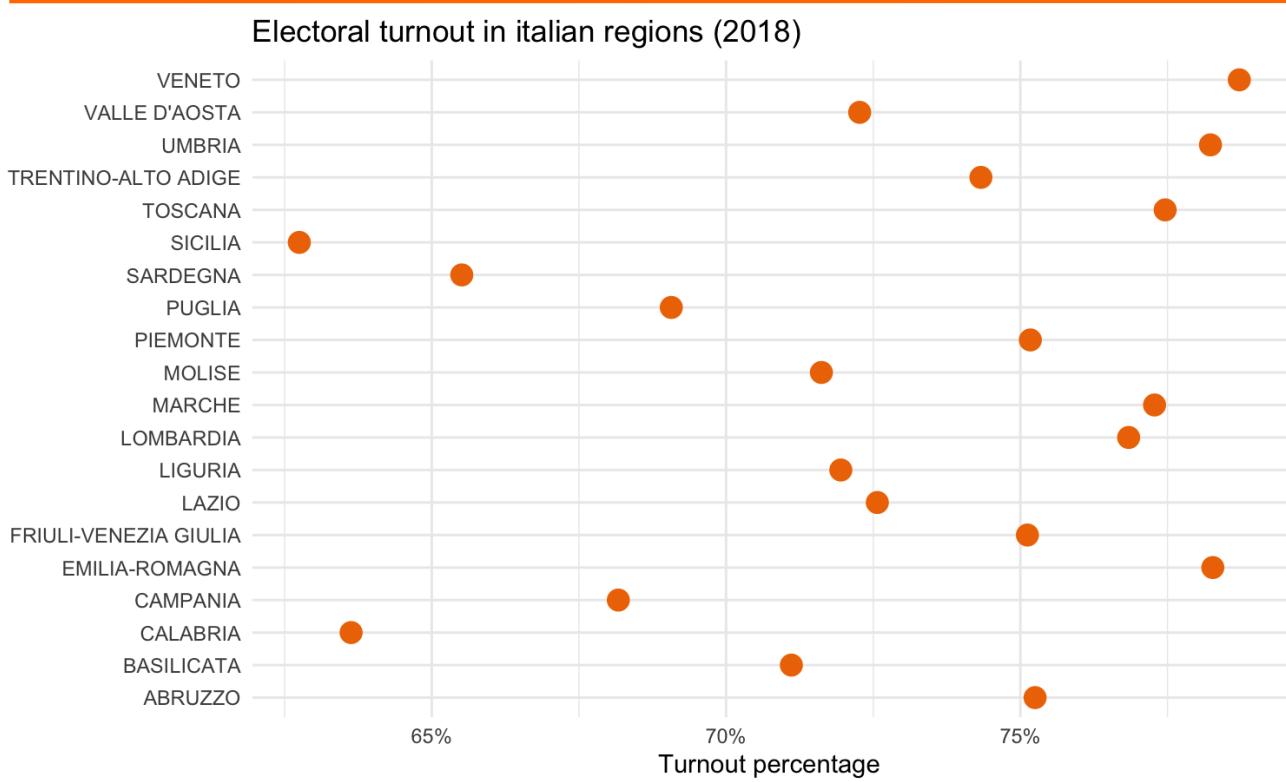
22

Dot plots

- Categorical values are encoded as position along an axis
- Quantitative values are encoded as position along an axis
 - ◆ There is no need to have a zero based axis range

23

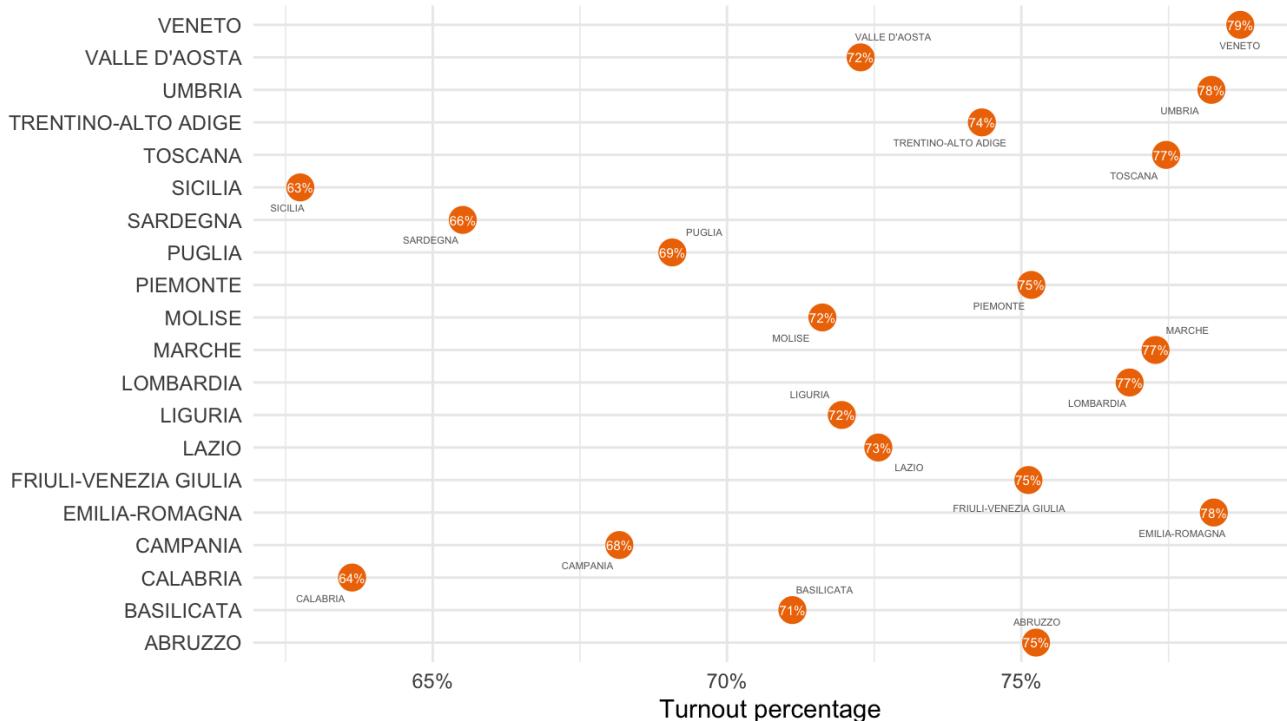
Comparison – Dot plot



24

Comparison – Dot plot

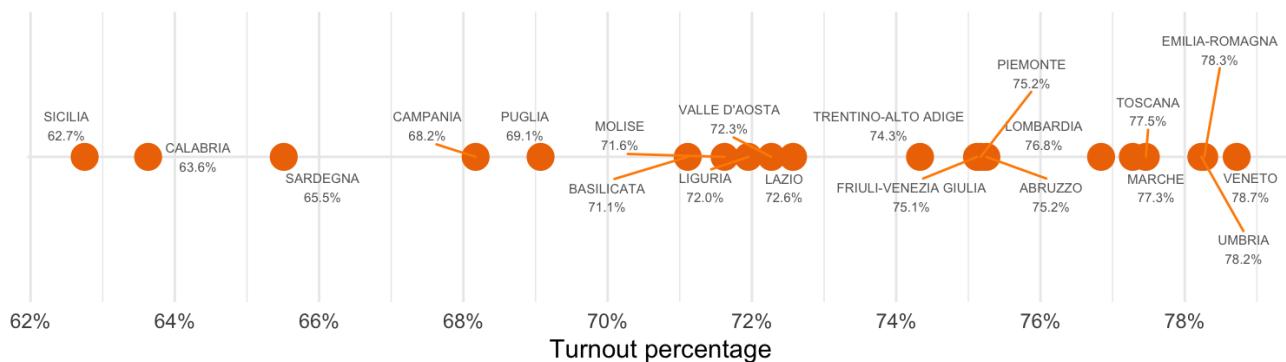
Electoral turnout in italian regions (2018)



25

Comparison – Strip plot

Electoral turnout in Italian regions (2018)

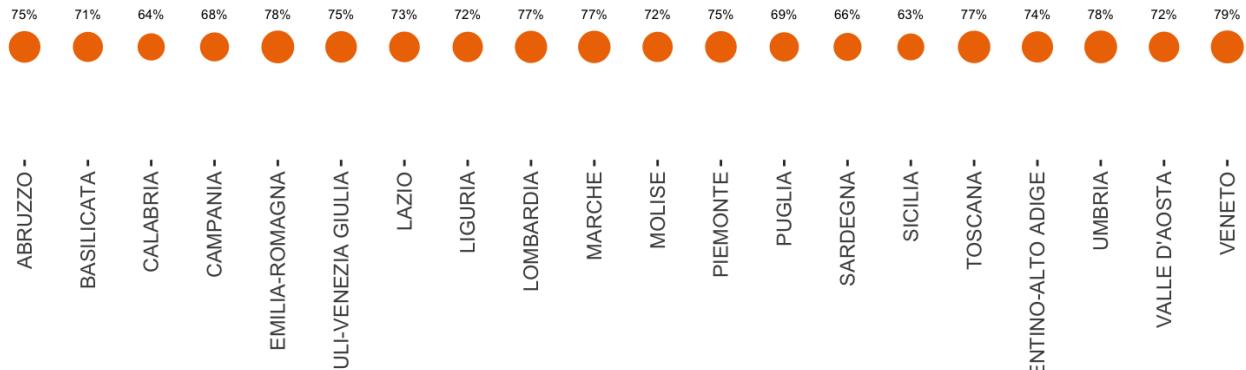


26

Comparison – Area – Bubbles

Electoral turnout in italian regions (2018)

Extremely difficult
to compare size



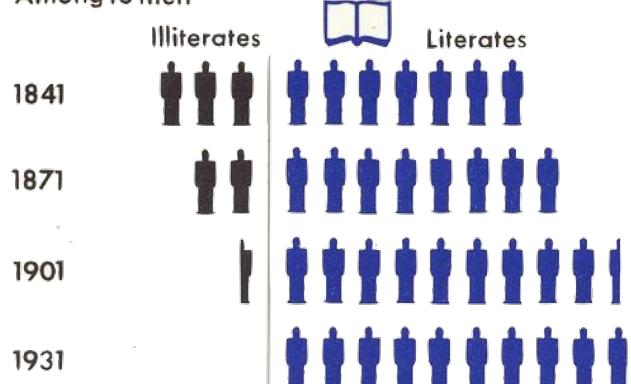
27

Count – Isotype

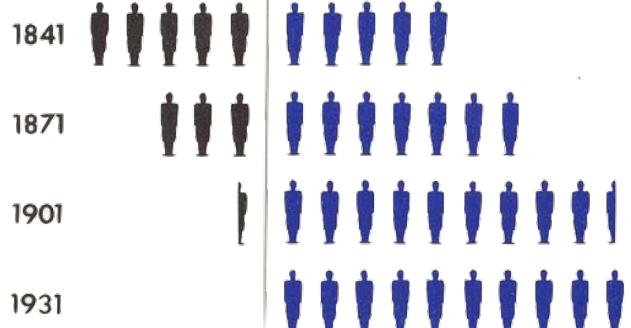
- Isotype
 - ◆ International System Of Typographic Picture Education
- Marie and Otto Neurath
 - ◆ Vienna, 1936

Literacy in England and Wales

Among 10 men

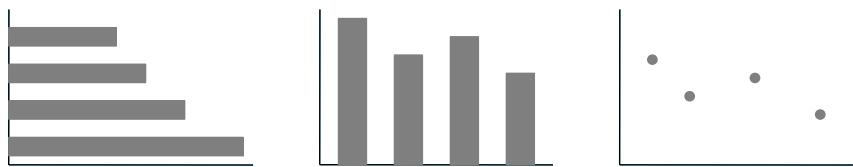


Among 10 women



Ranking

- Same type as nominal comparison
- Pay attention to order
 - ◆ Bar graphs
 - ◆ Dot plot
 - Allow non zero-based axes



29

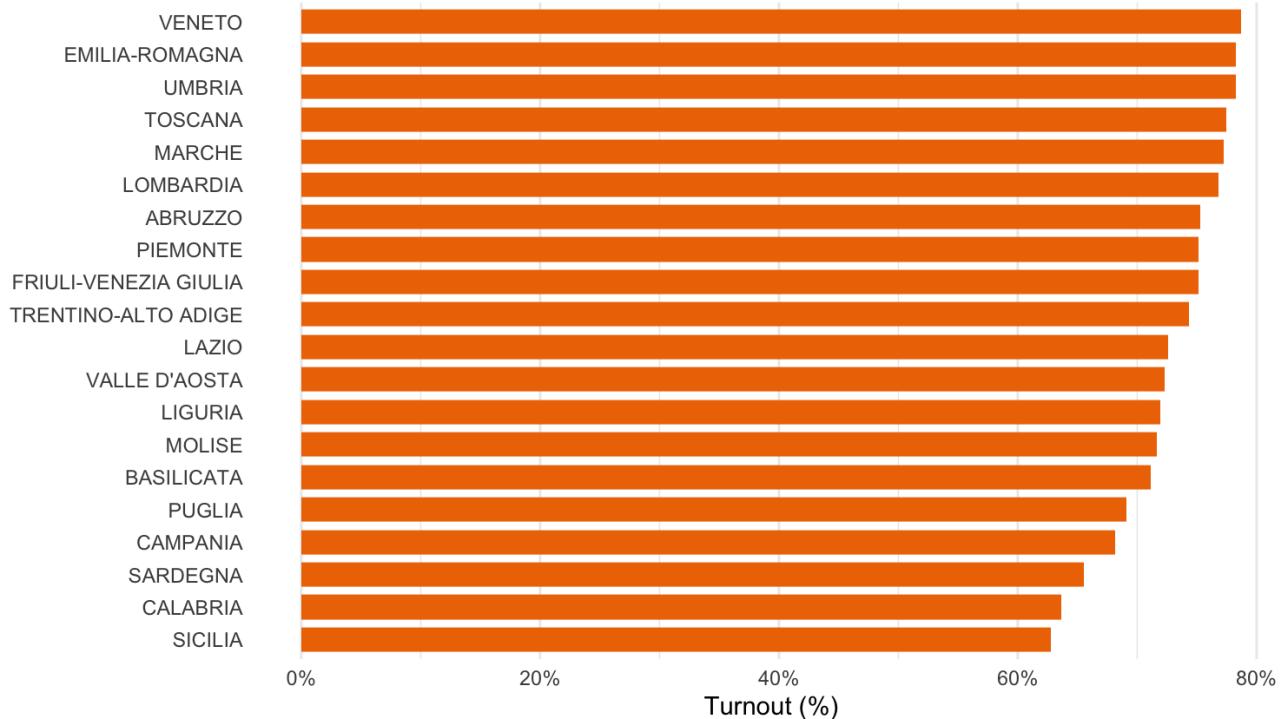
Ranking

Purpose	Sort order	Chart orientation
Highlight the highest value	Descending	H: highest on top V: highest on left
Highlight the lowest value	Ascending	H: lowest on top V: lowest on left

30

Ranking – Barplot

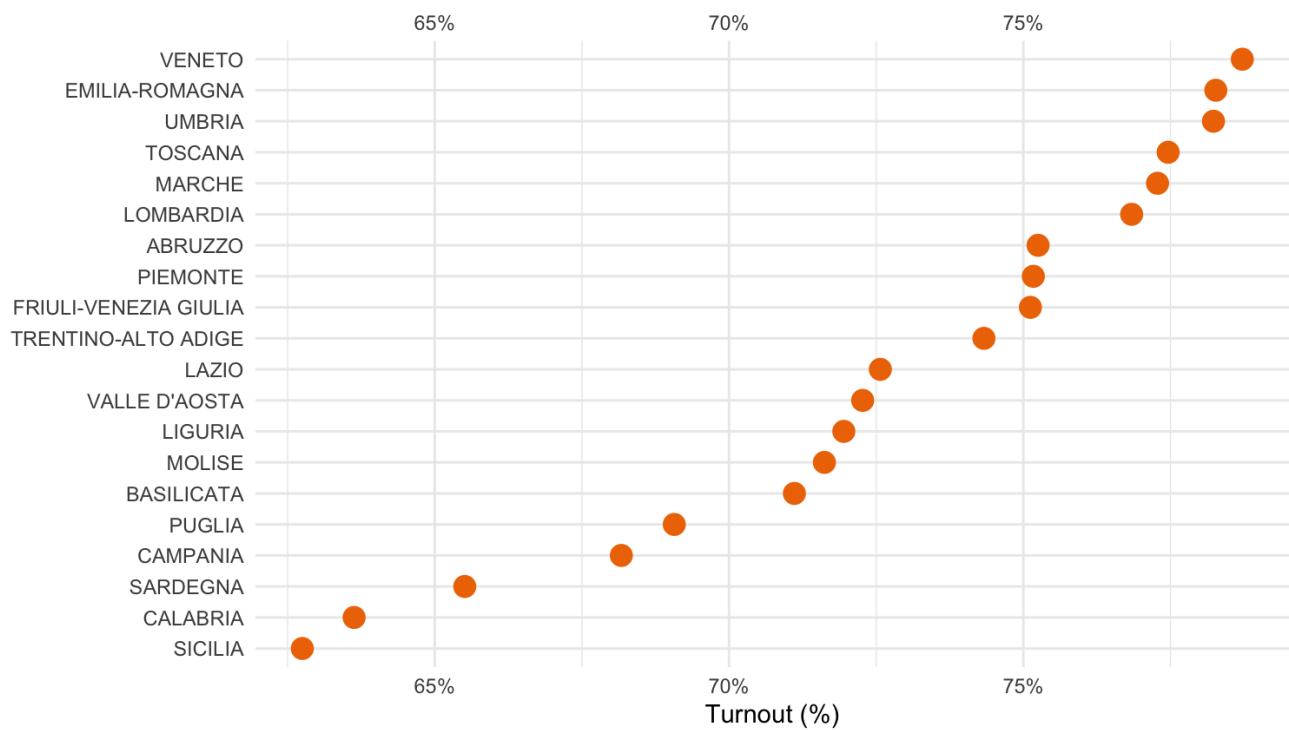
Electoral turnout in italian regions (2018)



31

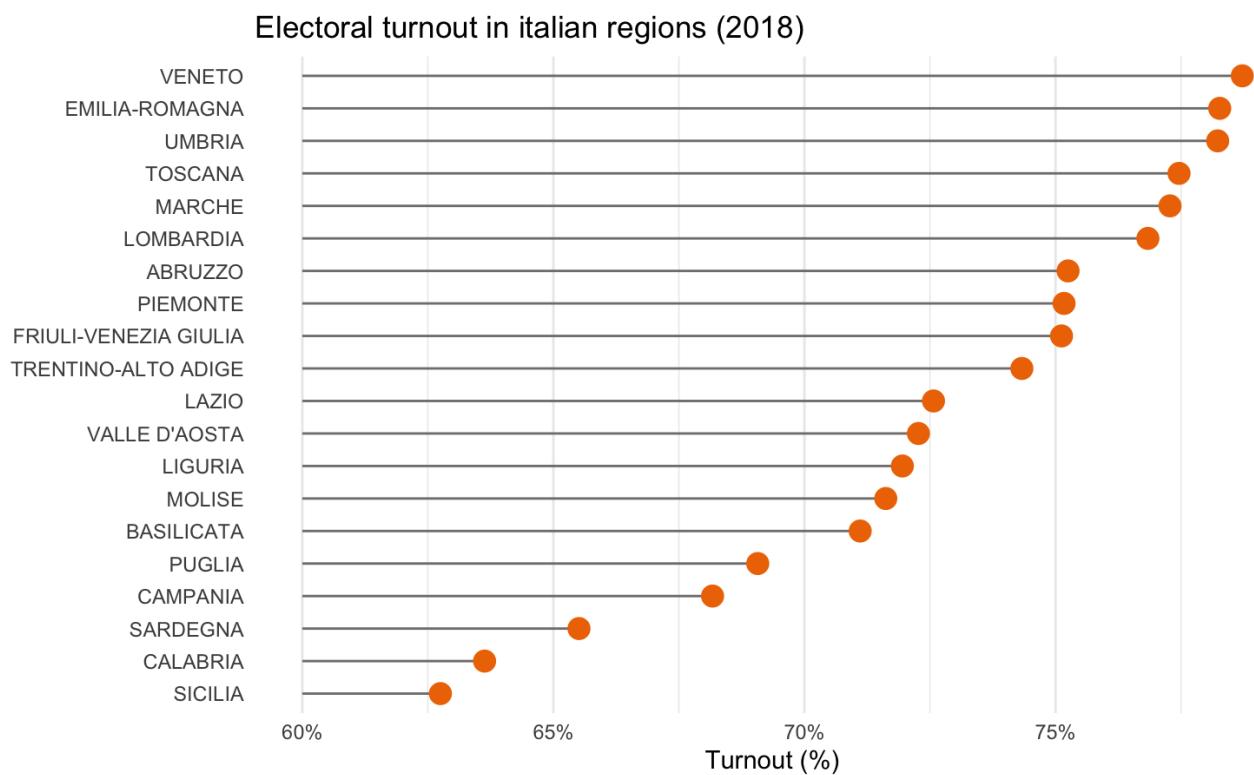
Ranking – Dot plot

Electoral turnout in italian regions (2018)



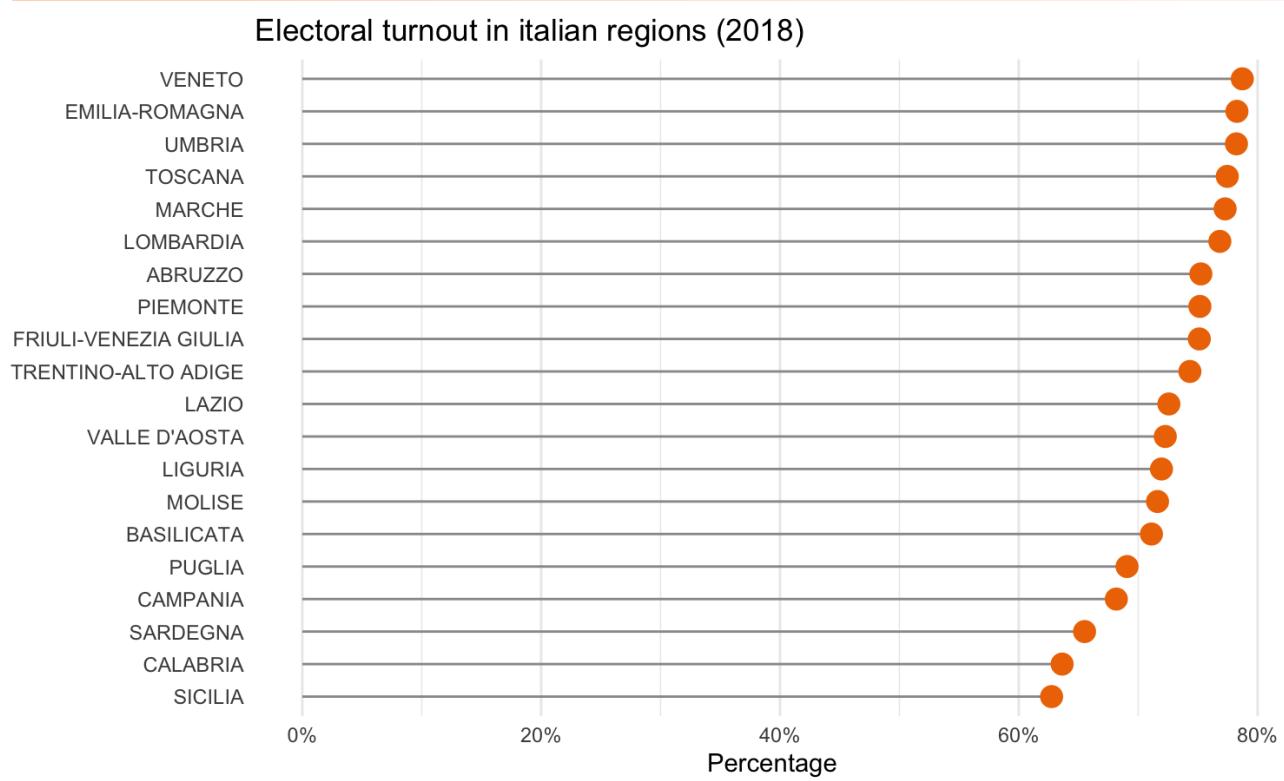
32

Lollypop (nonzero based scale)



33

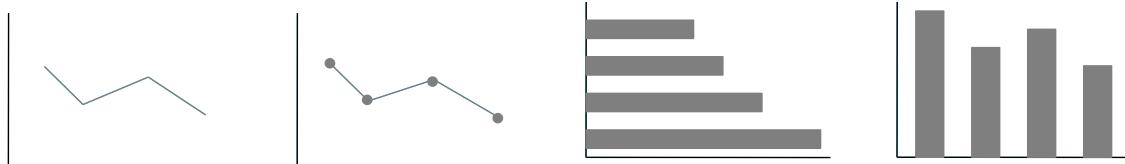
Lollypop (zero based scale)



34

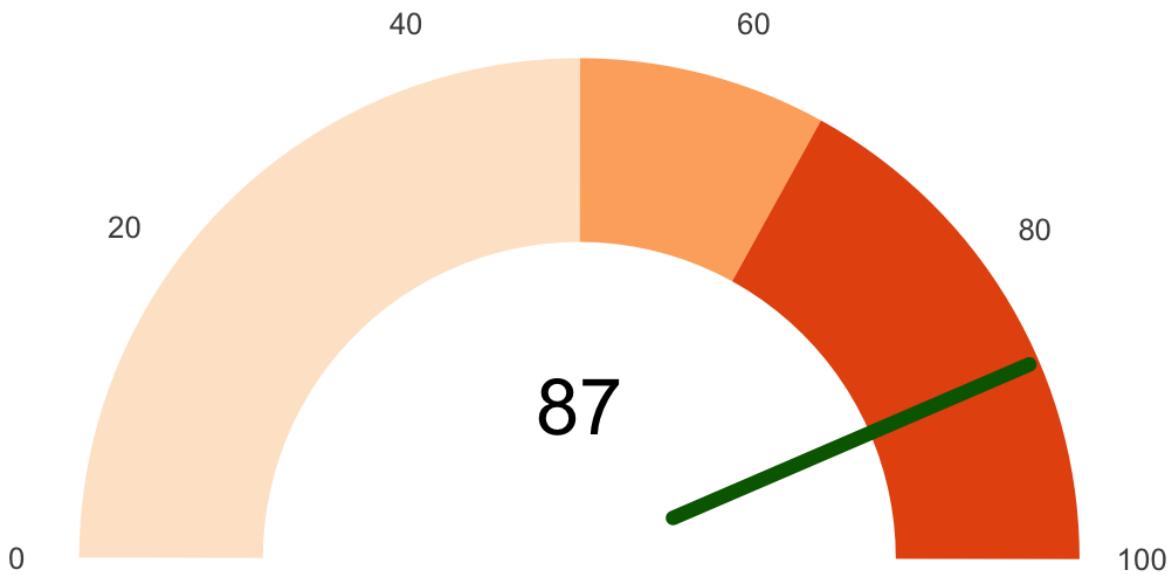
Deviation

- To what degree one or more sets of values differ in relation to primary values.
 - ◆ Points (dots)
 - ◆ Gauge
 - ◆ Bars
 - ◆ Bullet



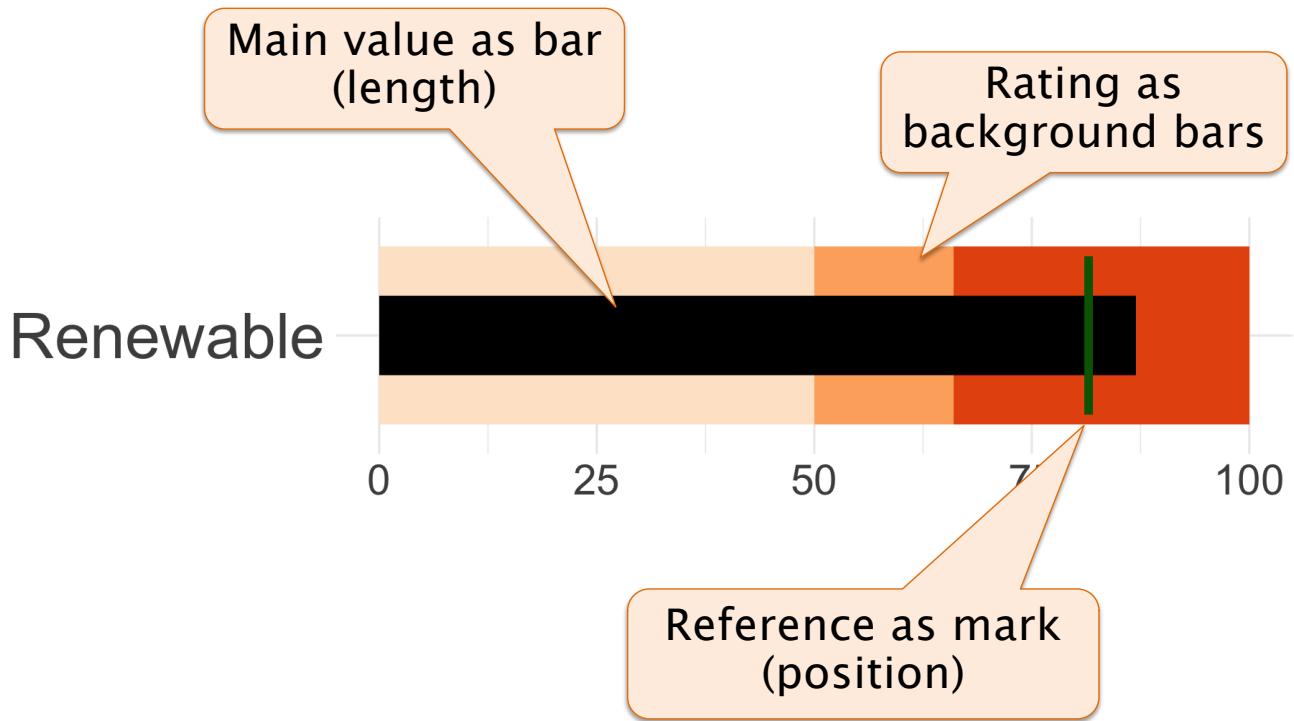
35

Angle + Position – Gauge



36

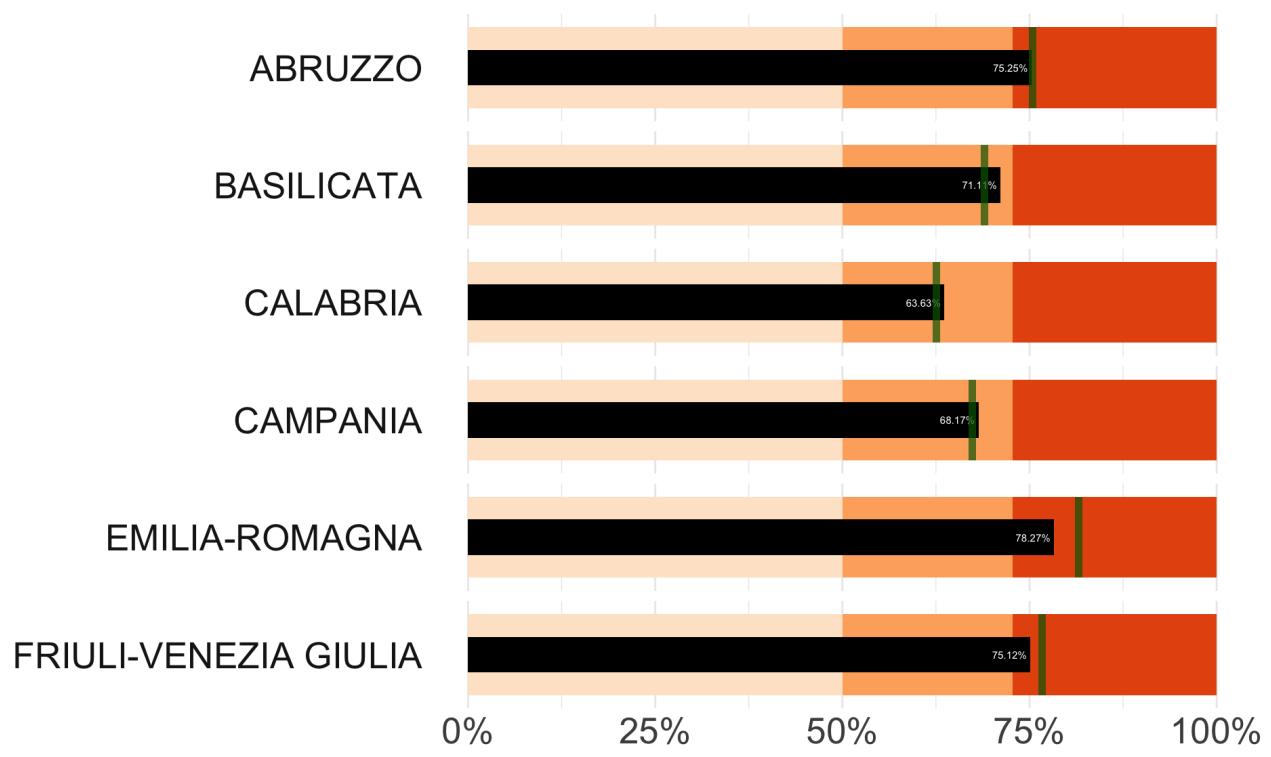
Length+Position– Bullet Graph



See: https://www.perceptualedge.com/articles/misc/Bullet_Graph_Design_Spec.pdf

37

Length+Position– Bullet Graph



38

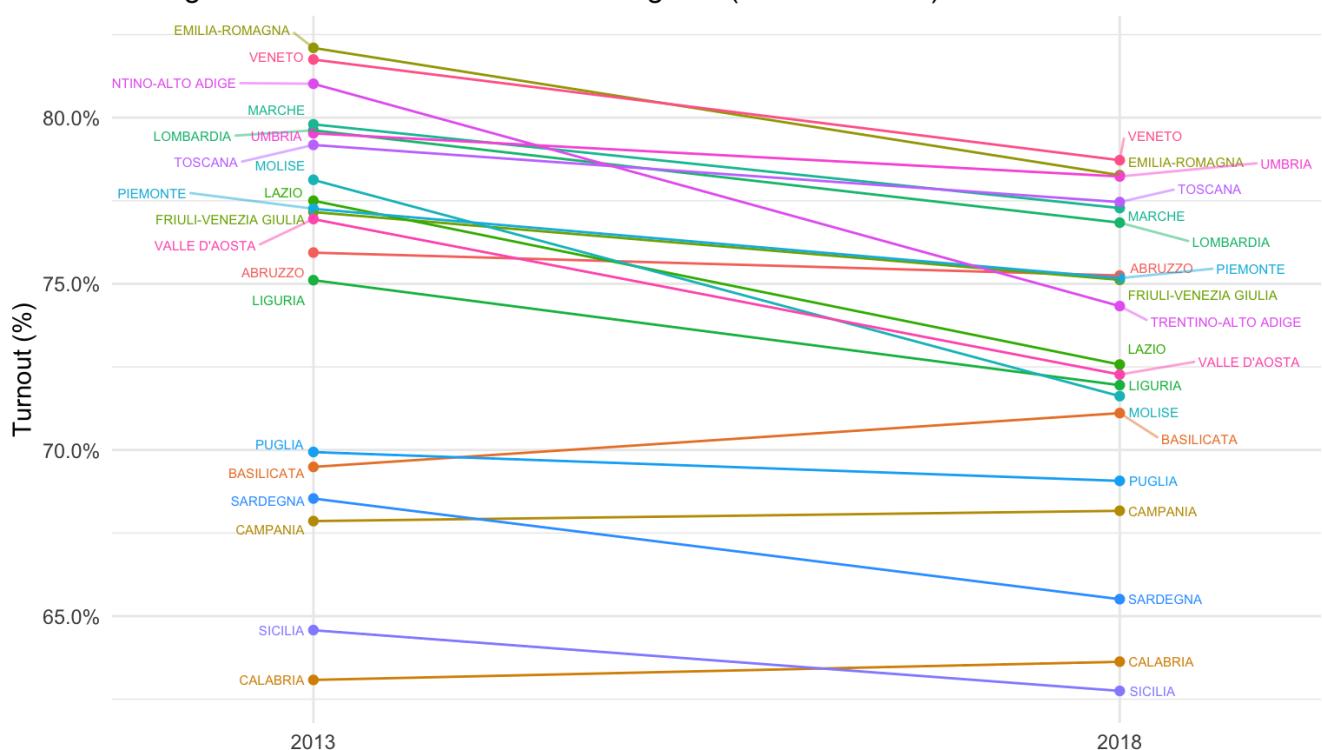
Pre–post variation

- Comparing several categorical values typically two conditions
 - ◆ Pre vs. post
 - ◆ With vs. without
 - ◆ ...

39

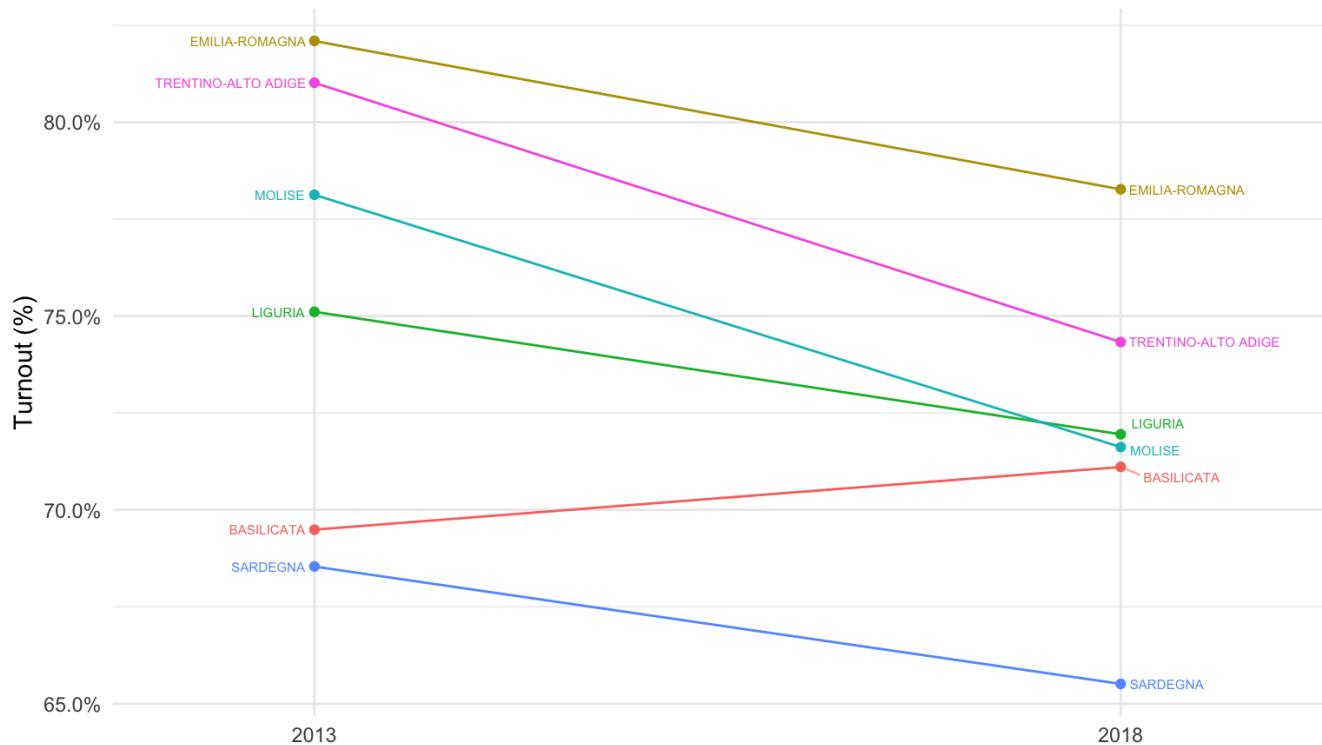
Slope chart

Change in electoral turnout for Italian regions (2018 vs. 2013)



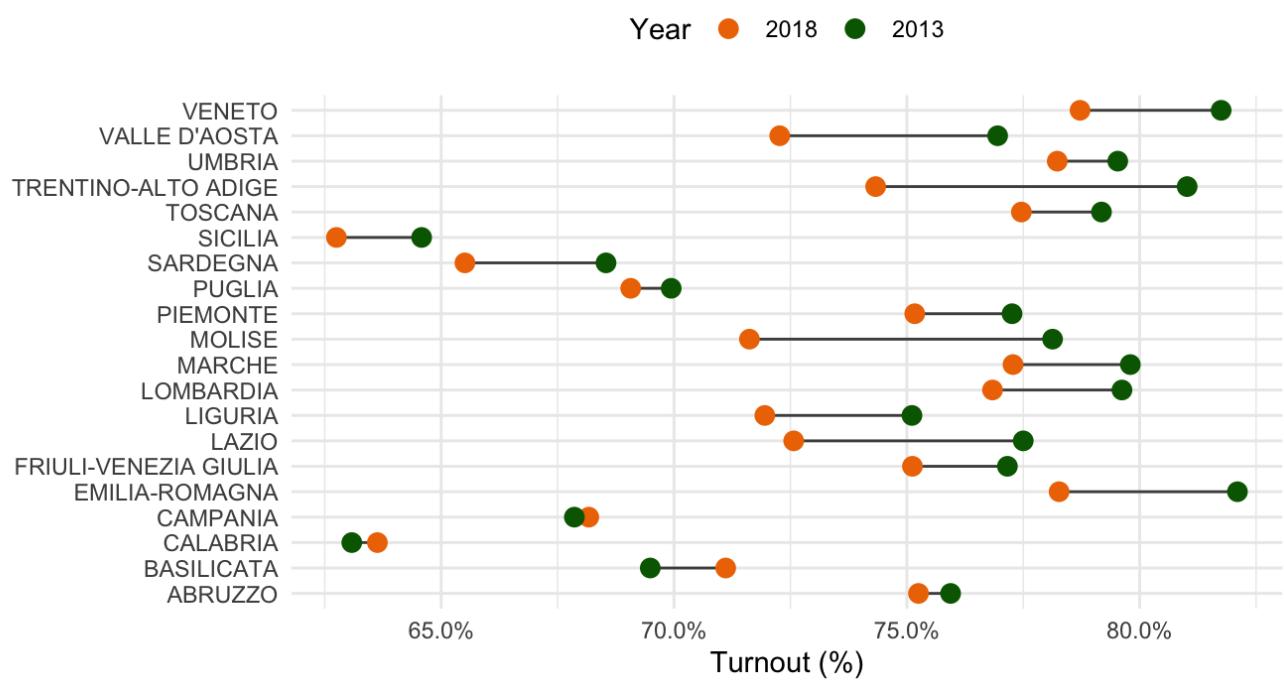
Slope chart

Change in electoral turnout for italian regions (2018 vs. 2013)



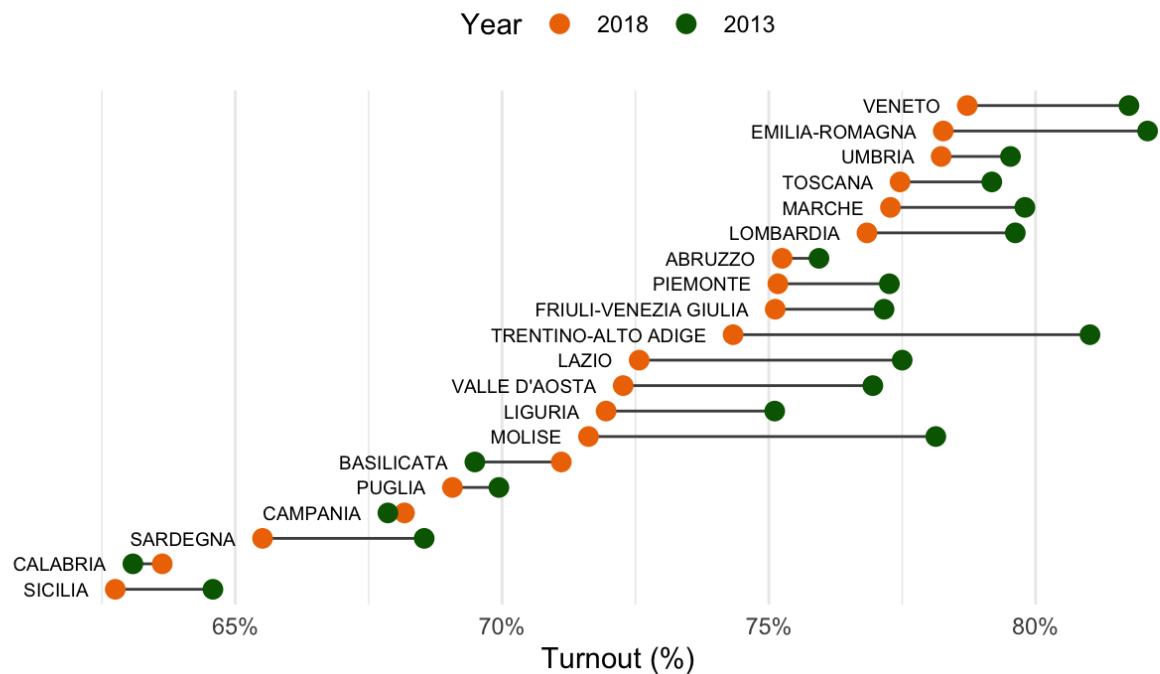
Dumbbell plot

Change in electoral turnout for italian regions (2018 vs. 2013)



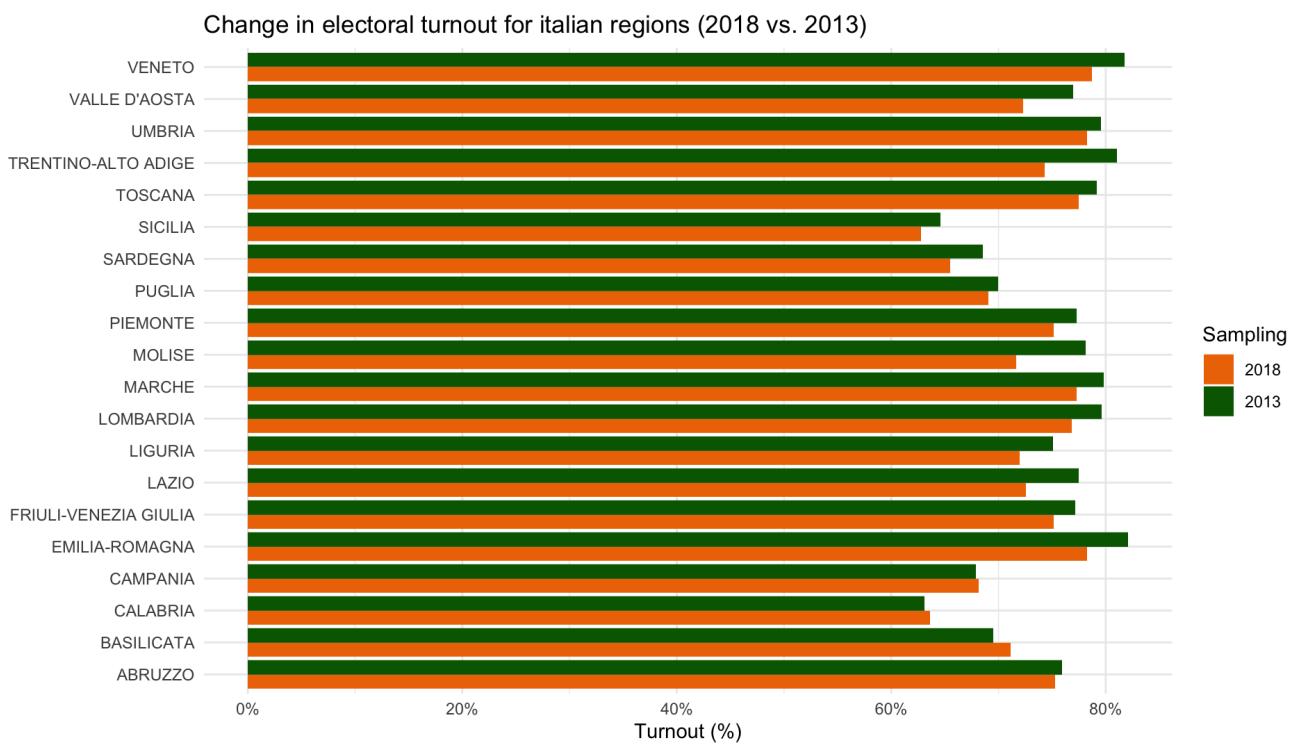
Dumbbell plot (sorted)

Change in electoral turnout for italian regions (2018 vs. 2013)



43

Clustered bars



44

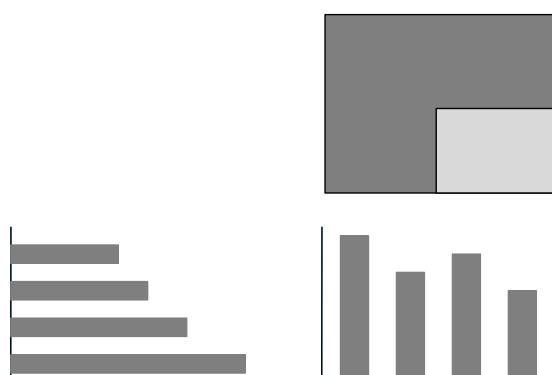
Proportion (Part-to-whole)

- Represent the frequency of different values within a given category
 - ◆ Be careful to use all values within the same category
- Can be used to compare frequency distribution across different categories sharing the same levels

45

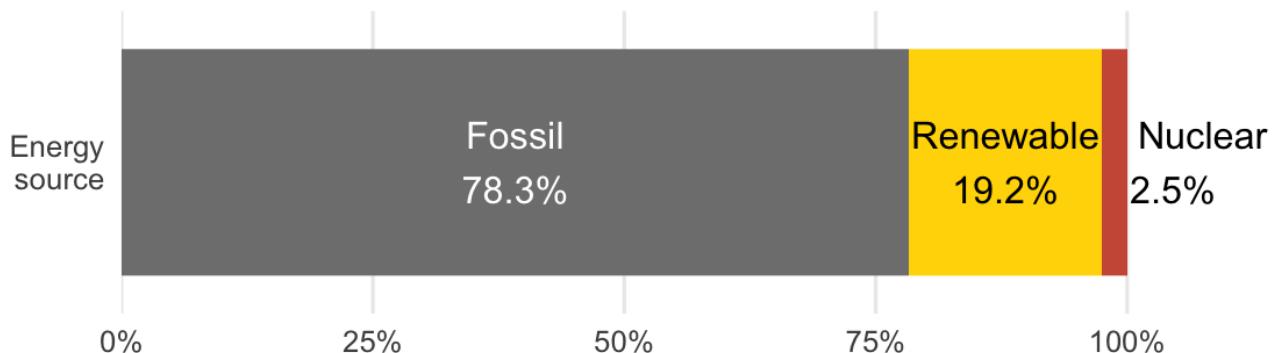
Proportion (Part-to-whole)

- Best unit: percentage
- Stacked bar graph
 - ◆ Difficult to read individual values
- Stacked area
- Treemap
- Gridplot
- Pie / Donut
- Marimekko



46

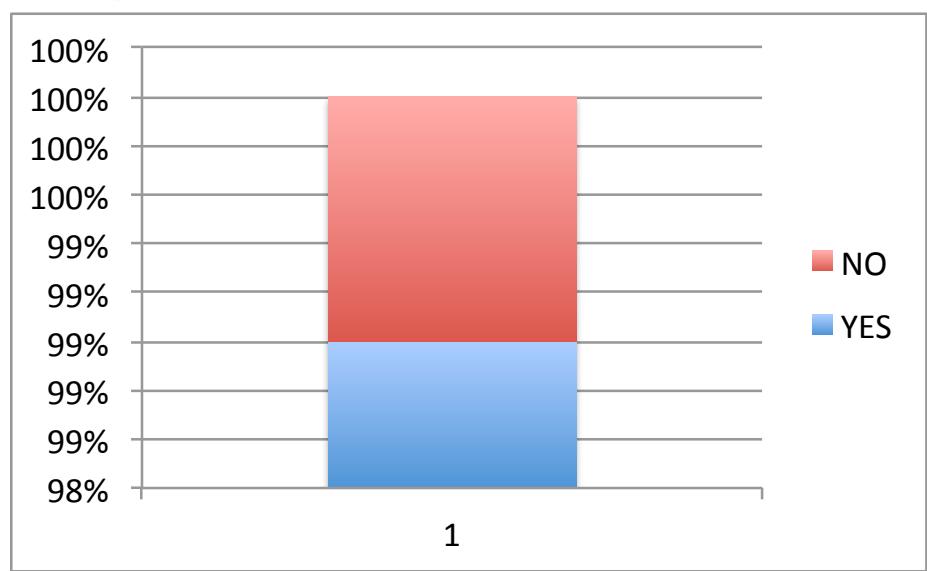
Length – Stacked Bar



47

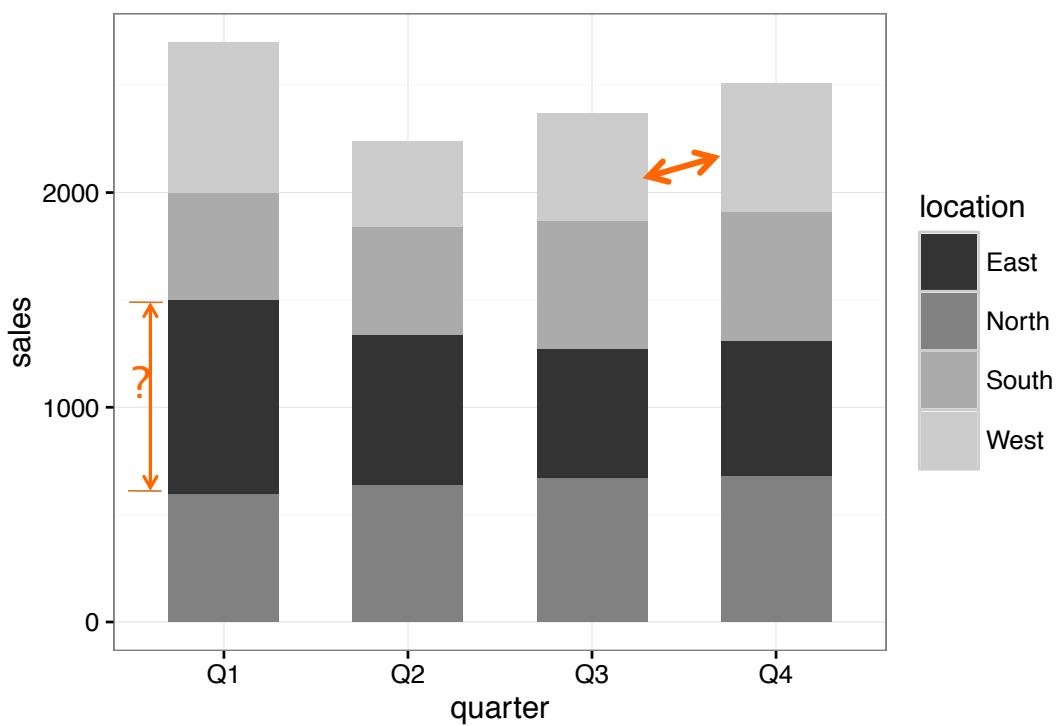
Beware of MS-Excel Defaults

	A	B
1	YES	99%
2	NO	1%



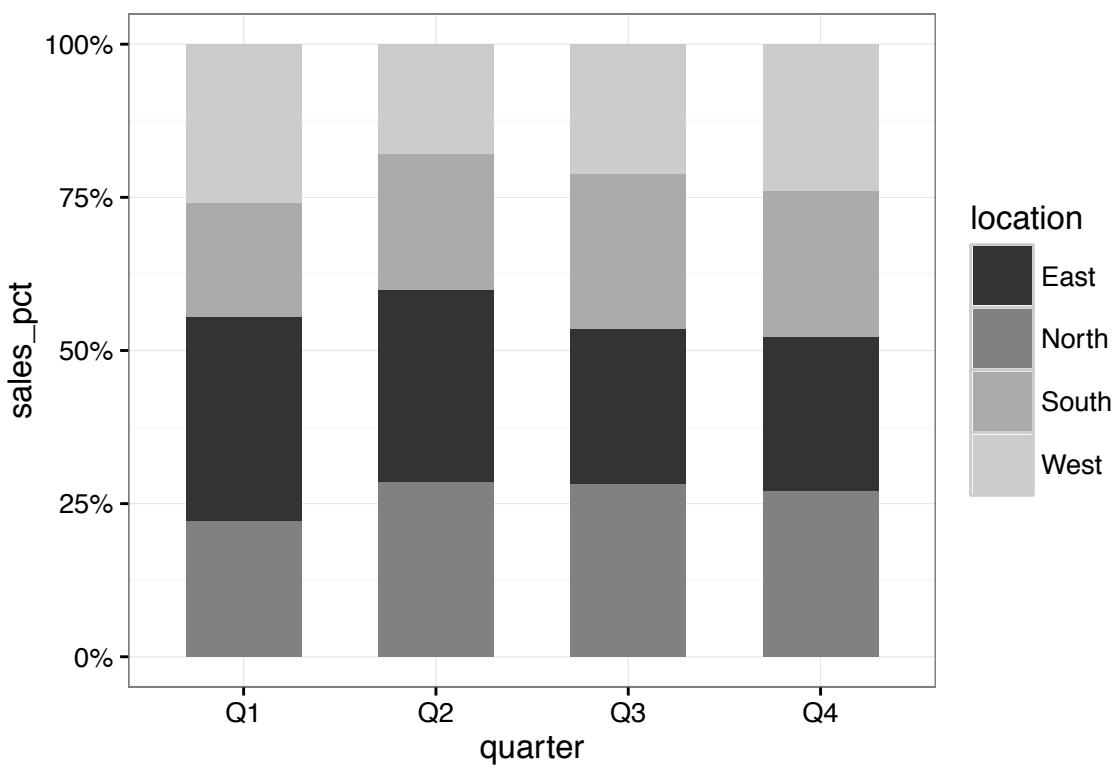
48

Stacked bar graph



49

Stacked bars w/percentage



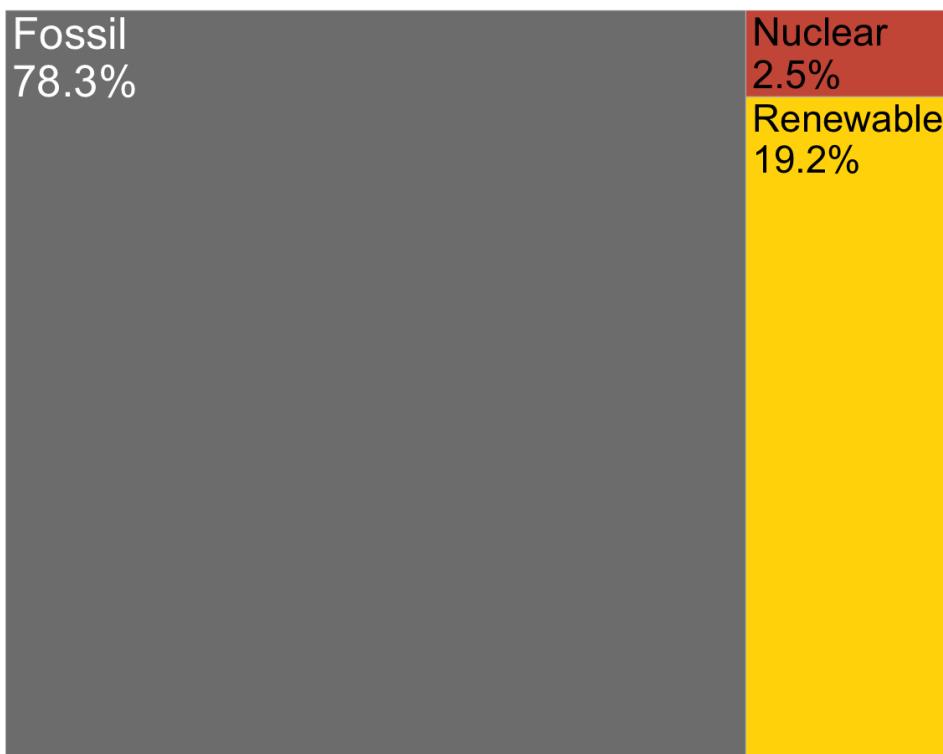
50

Area – Treemap



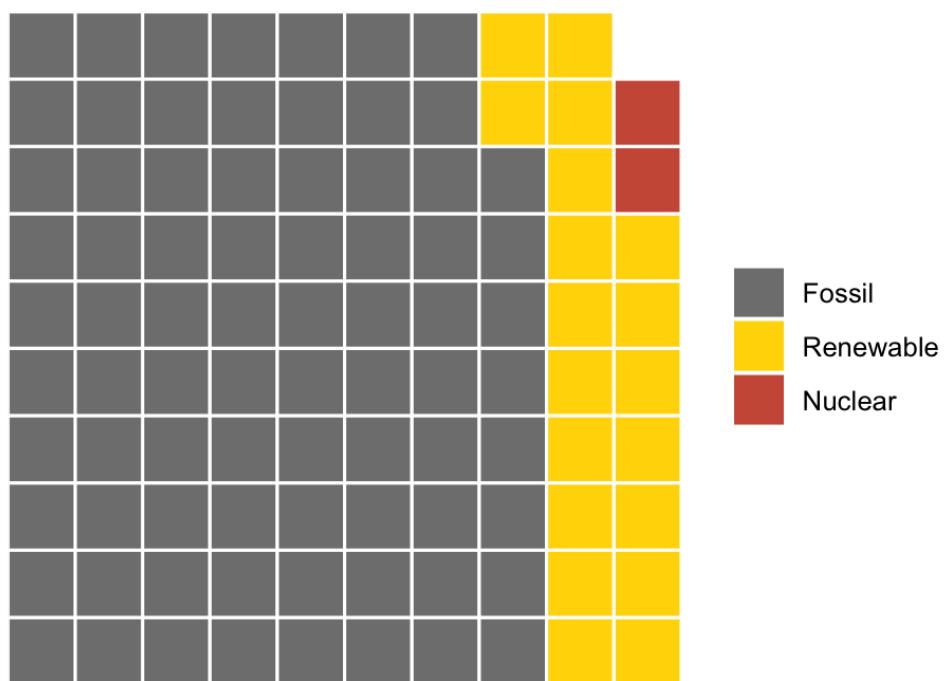
51

Area – Treemap



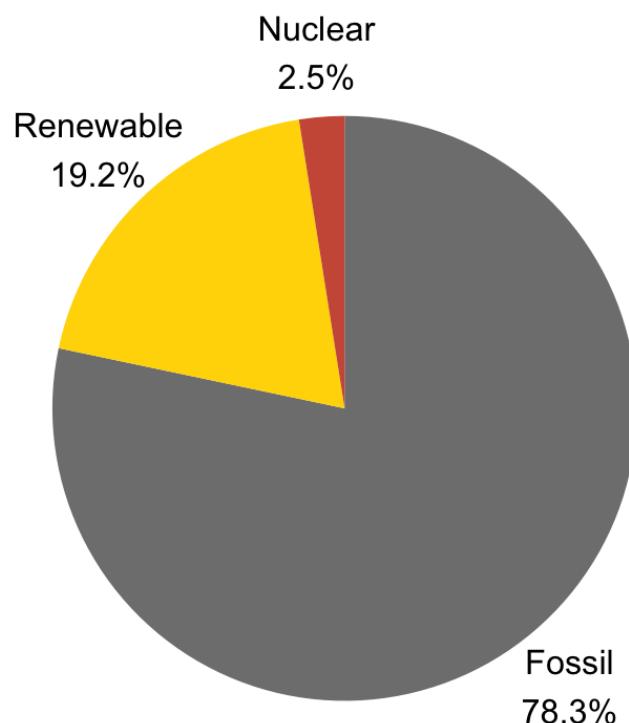
52

Area + Count – Waffle / Grid



53

Area + Angle – Pie Chart



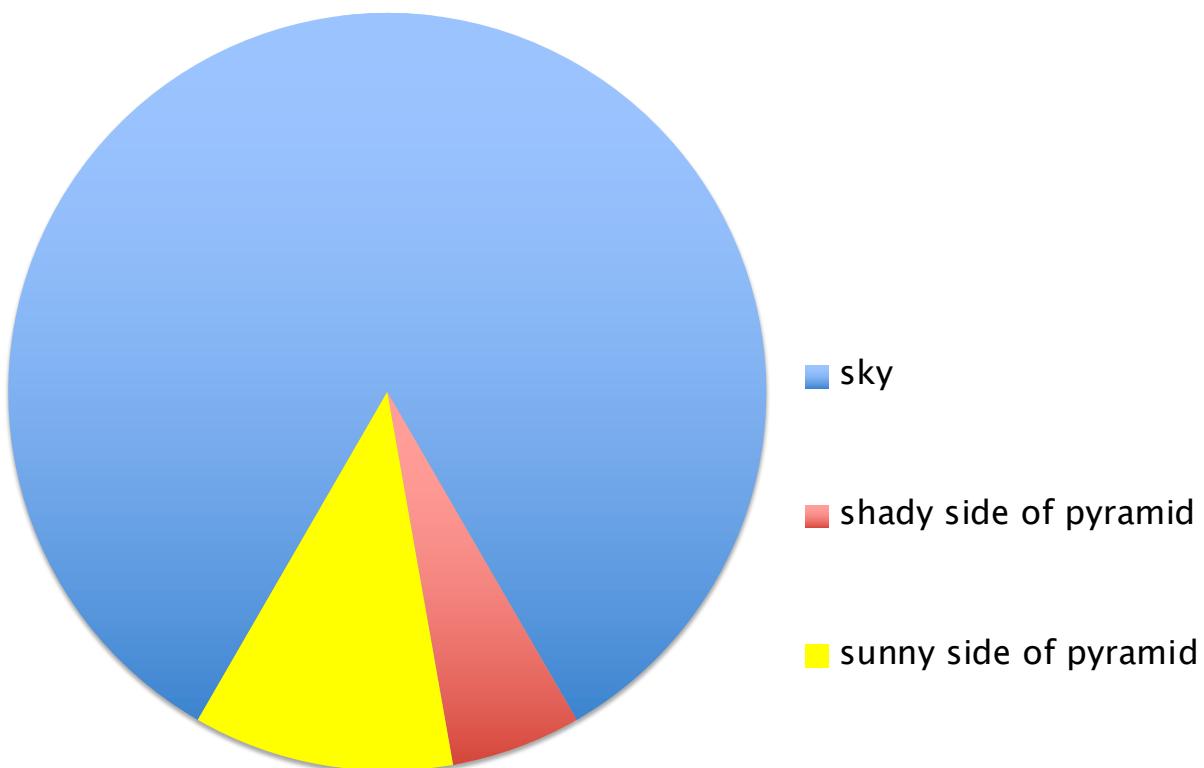
54

Pie charts



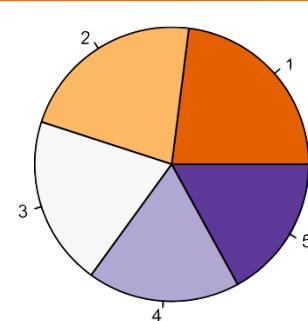
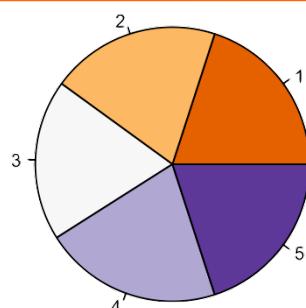
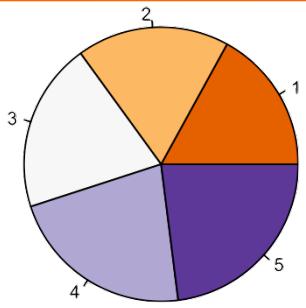
55

Pie charts



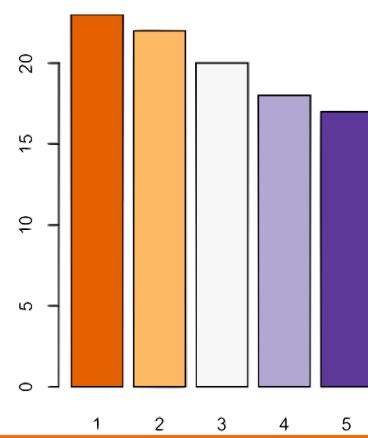
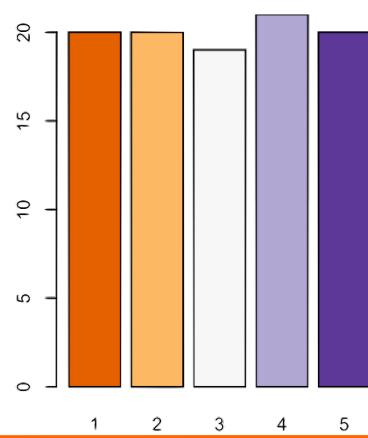
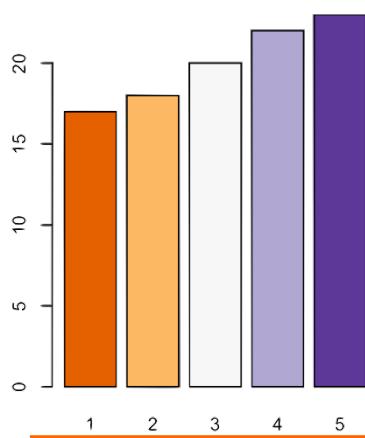
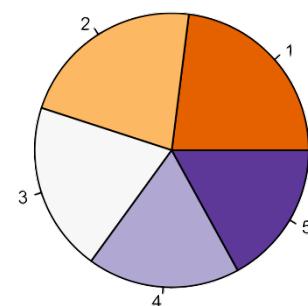
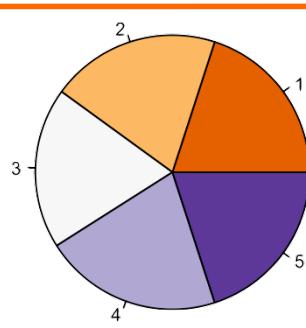
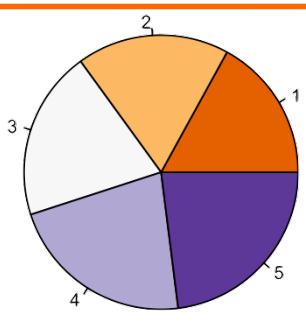
56

Pies



57

Pies vs. Bars



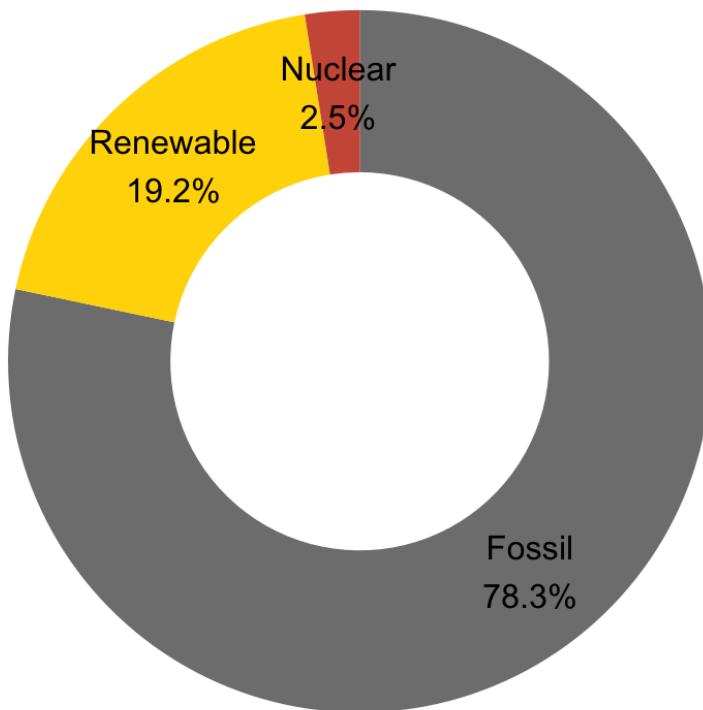
58

Pie Charts: guidelines

- Have serious limitations
 - ◆ To represent part–whole relationship
 - ◆ Only with a small number of categories
 - Up to four
 - Avoid rainbow pie
 - ◆ When proportions are distinct enough
- Remember to ease reading
 - ◆ Labels placed close to slices
 - ◆ Labels include values (percentages)

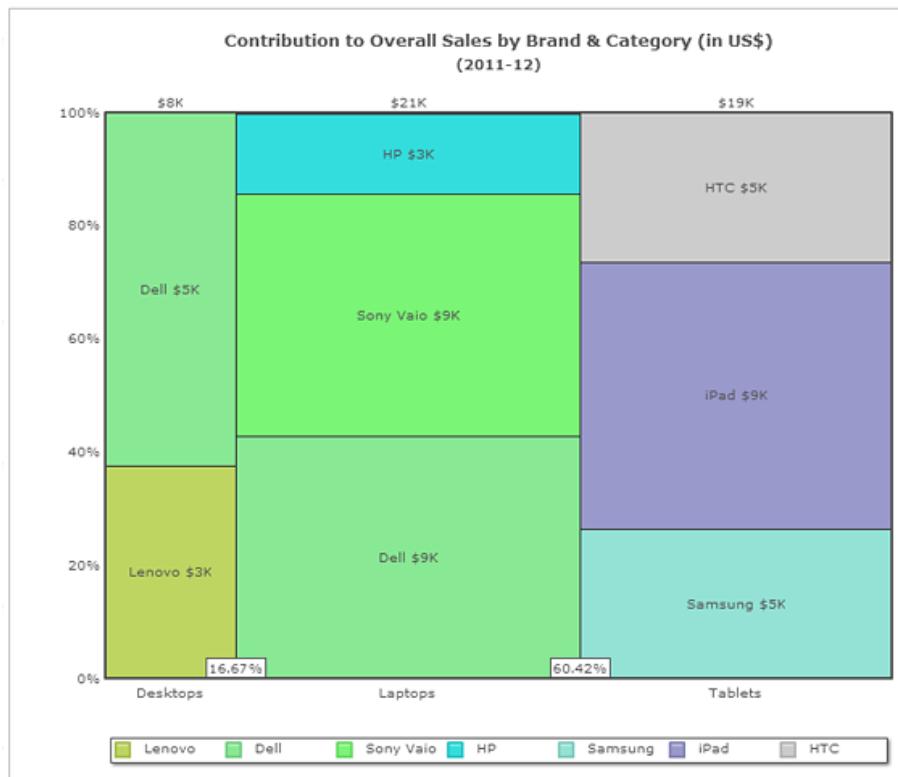
59

Area+Angle+Length – Donut



60

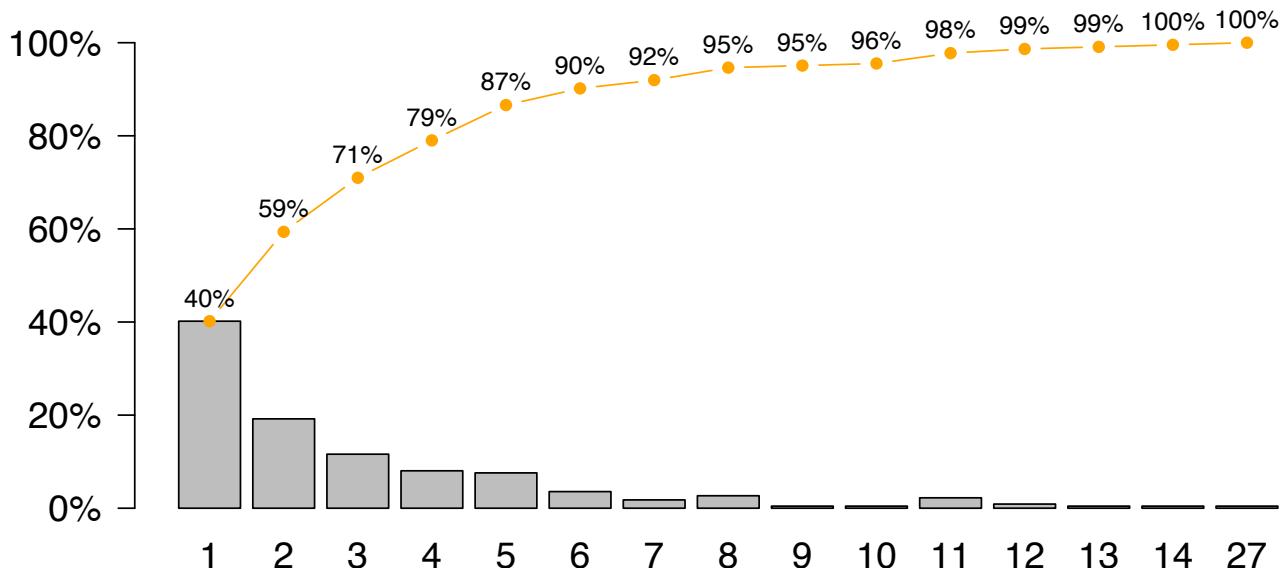
Marimekko Chart



<https://www.fusioncharts.com/chart-primer/marimekko-chart/>

61

Pareto chart



62

Distribution

- Continuous values
 - ◆ Show distribution of single set of values
 - ◆ Show and compare two or more distributions

63

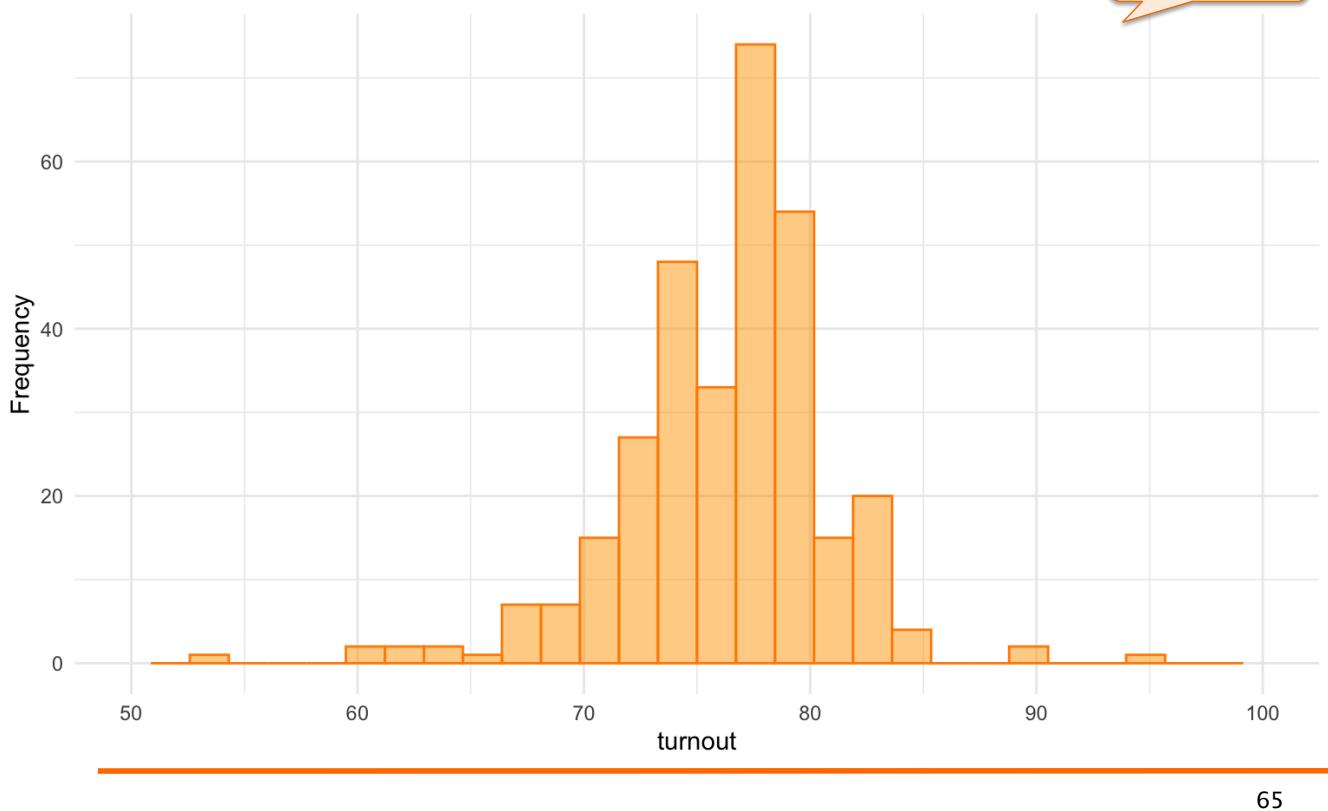
Single distribution

- Histogram
 - ◆ Vertical bar graph
 - ◆ Frequency for subdivision
 - Quantitative ranges
 - Categories
 - ◆ Emphasis on number of occurrences
- Frequency polygon
 - ◆ Line graphs
 - ◆ Frequency density function
 - ◆ Emphasis on the shape of the distribution
- Boxplot
 - ◆ Summary

64

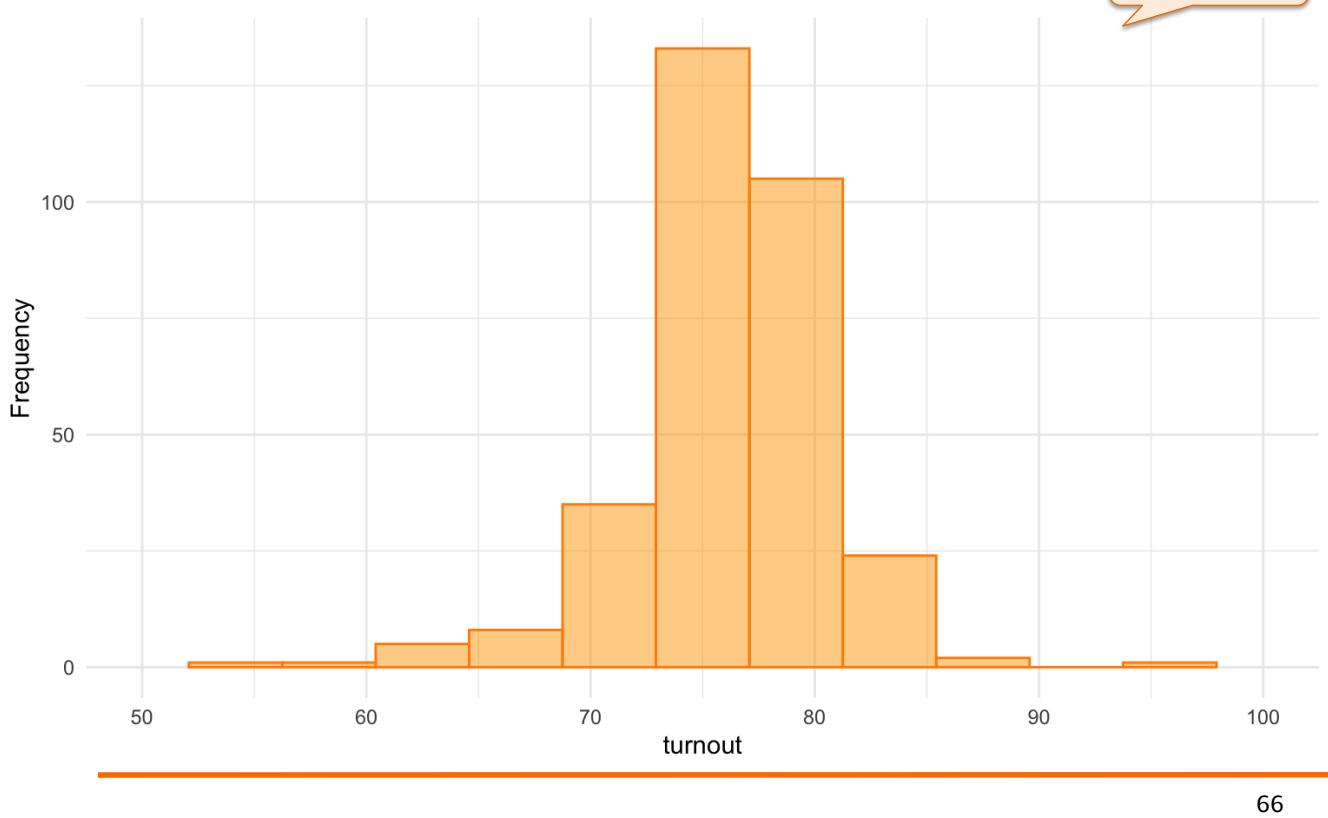
Histogram

30 bins



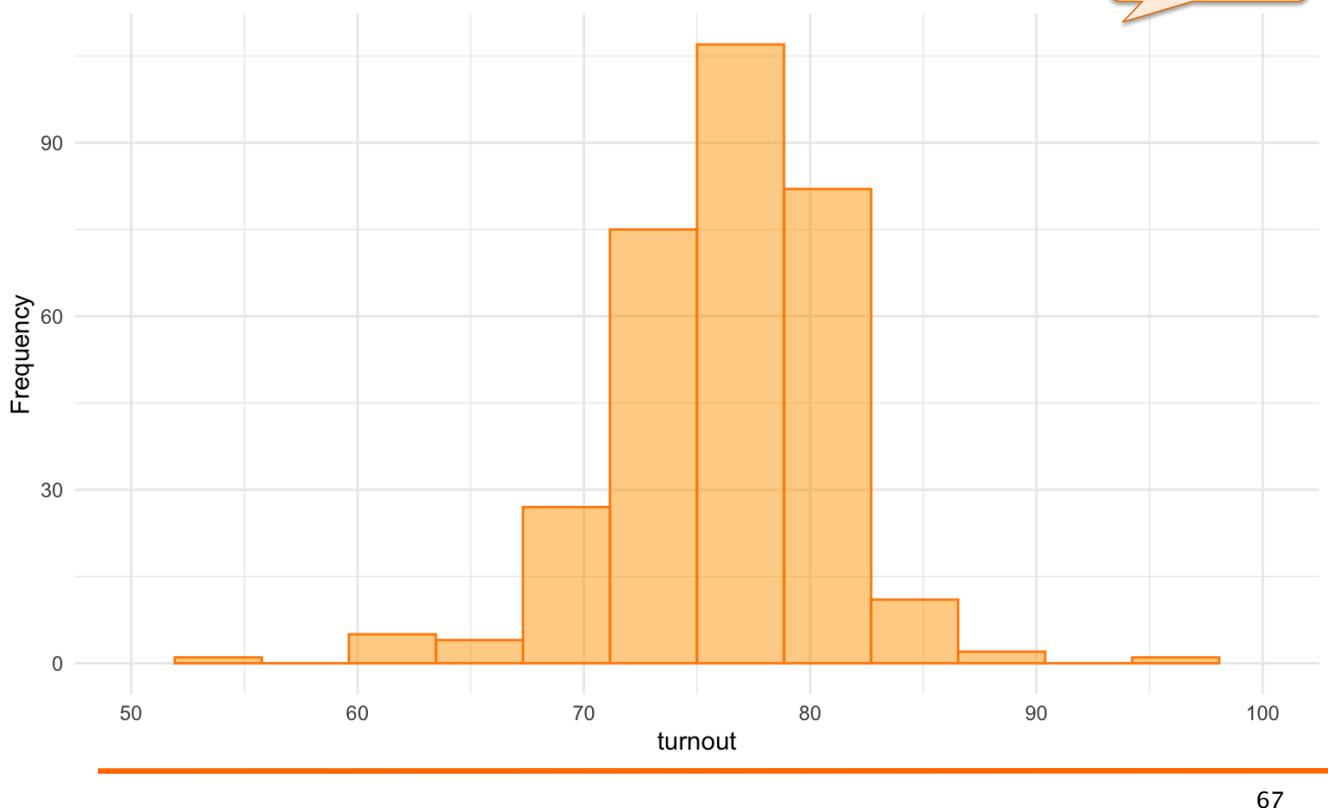
Histogram

13 bins

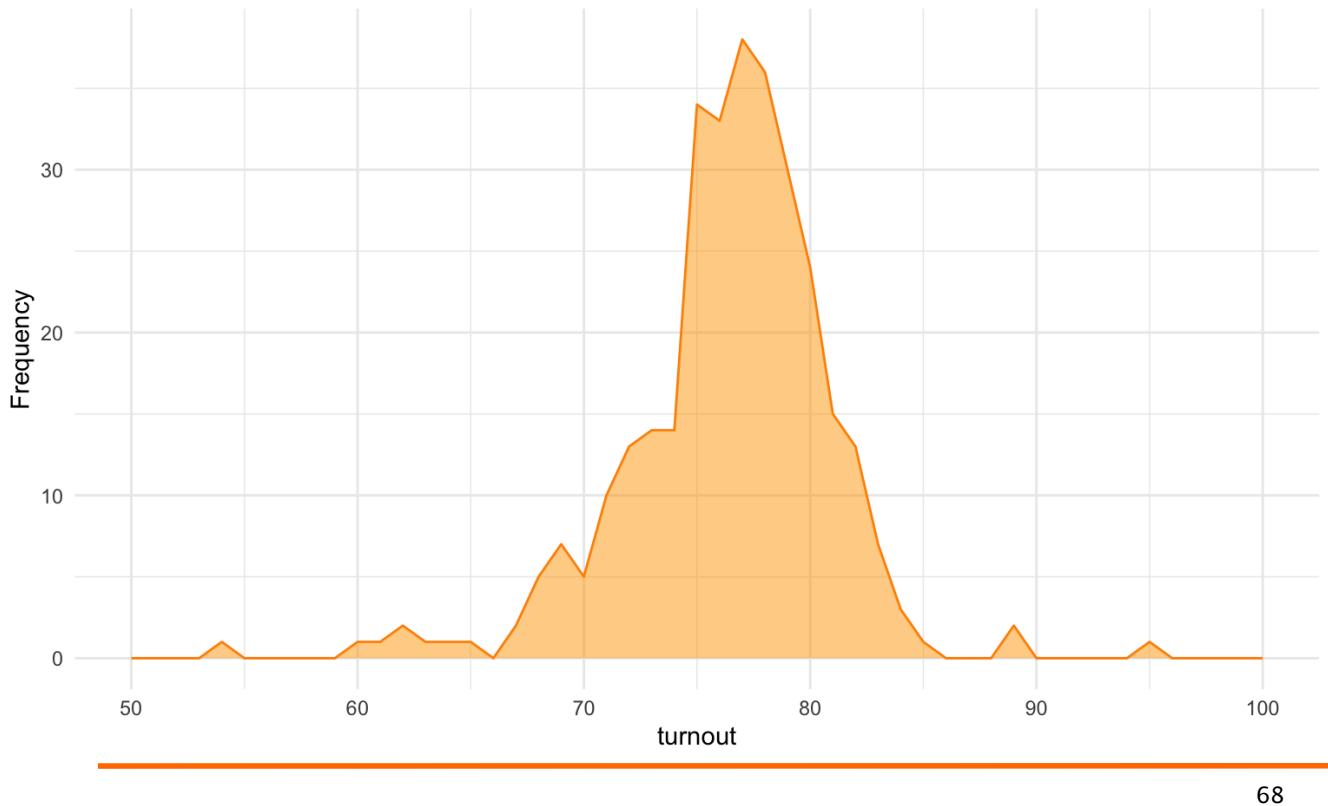


Histogram

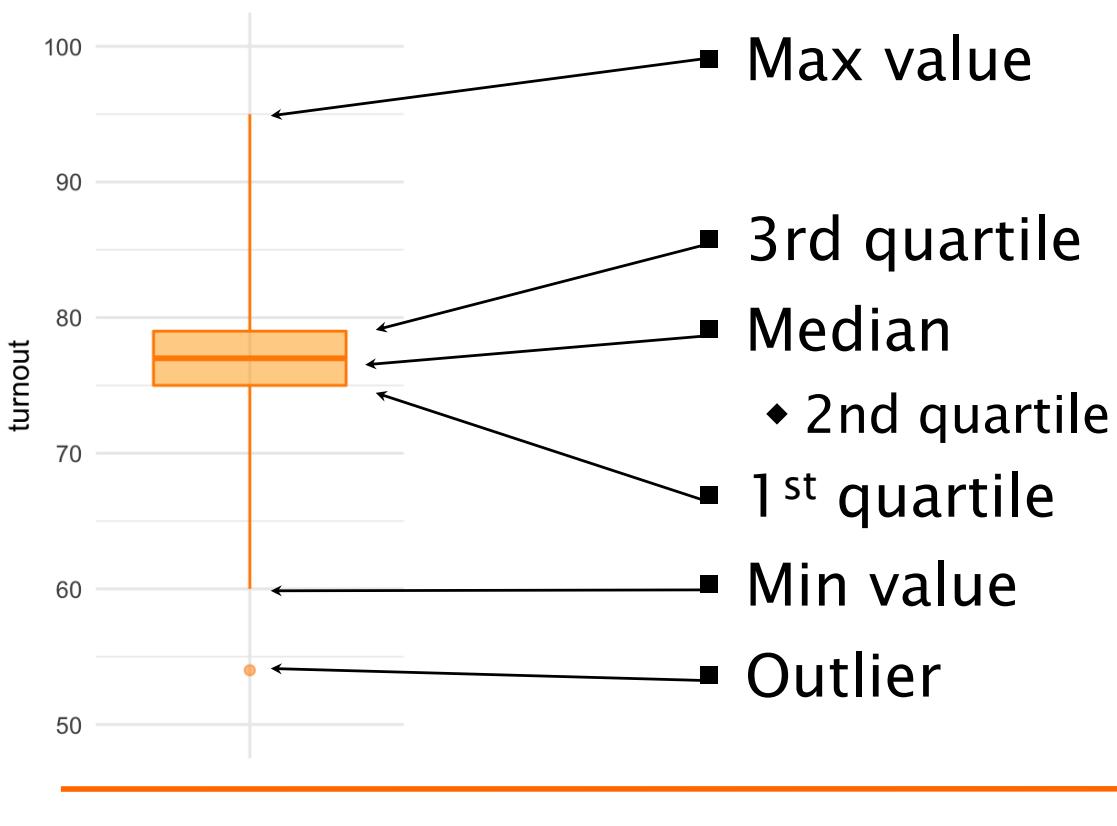
14 bins



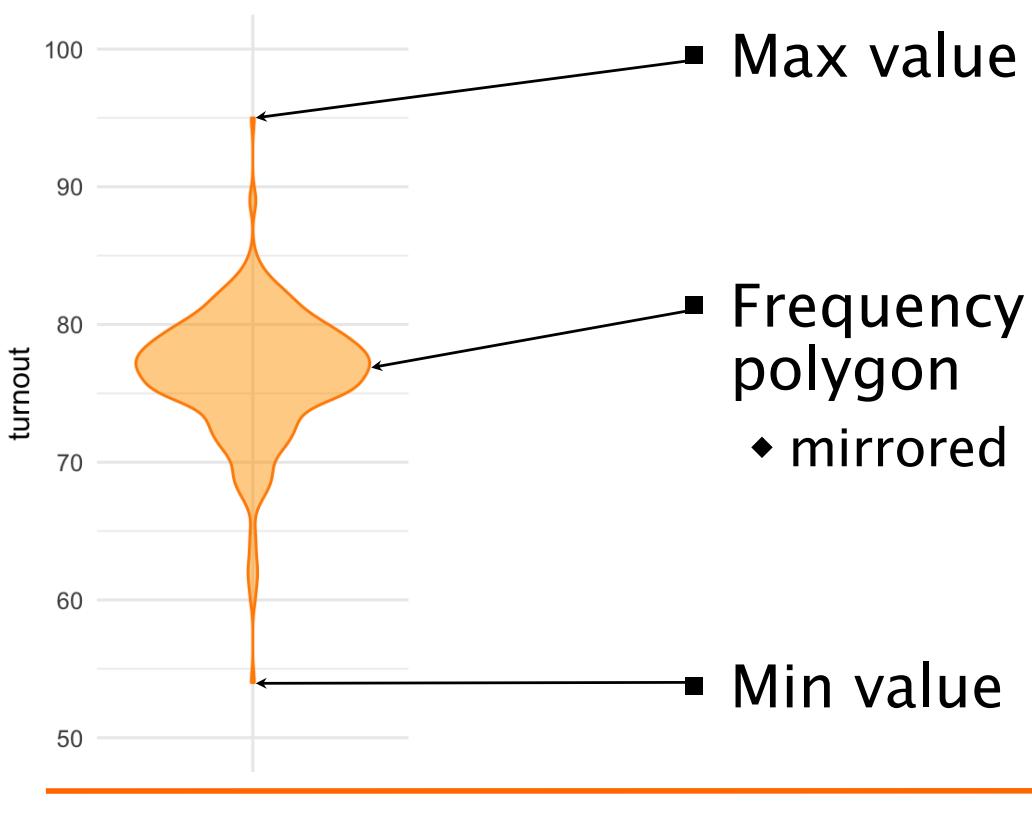
Frequency polygon



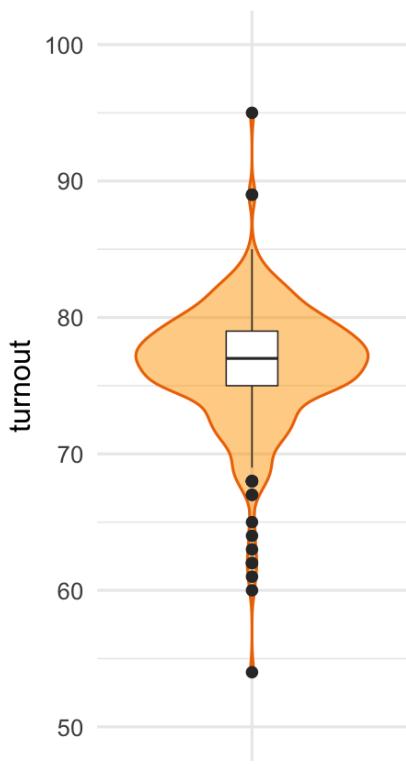
Boxplot



Violin plot



Violin + Boxplot



- Overlaying a box plot over the violin provides additional details

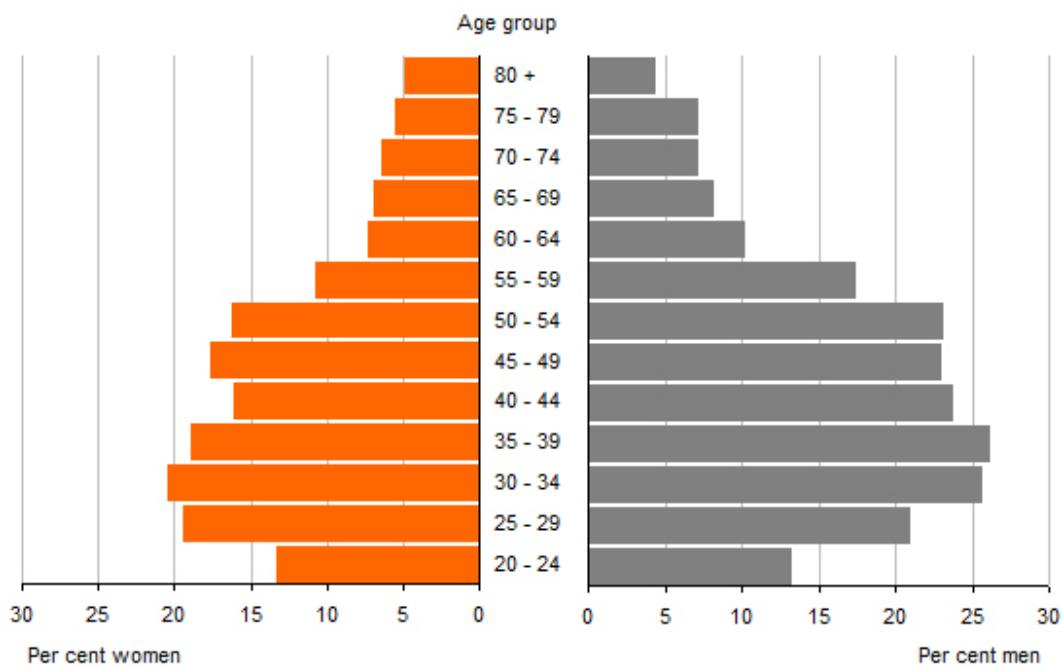
71

Multiple distribution

- Histogram is not suitable
- Frequency polygon
 - ◆ Line graphs
 - ◆ Frequency density function
- Boxplot
 - ◆ Summary
 - ◆ Less distracting with high number of categories

72

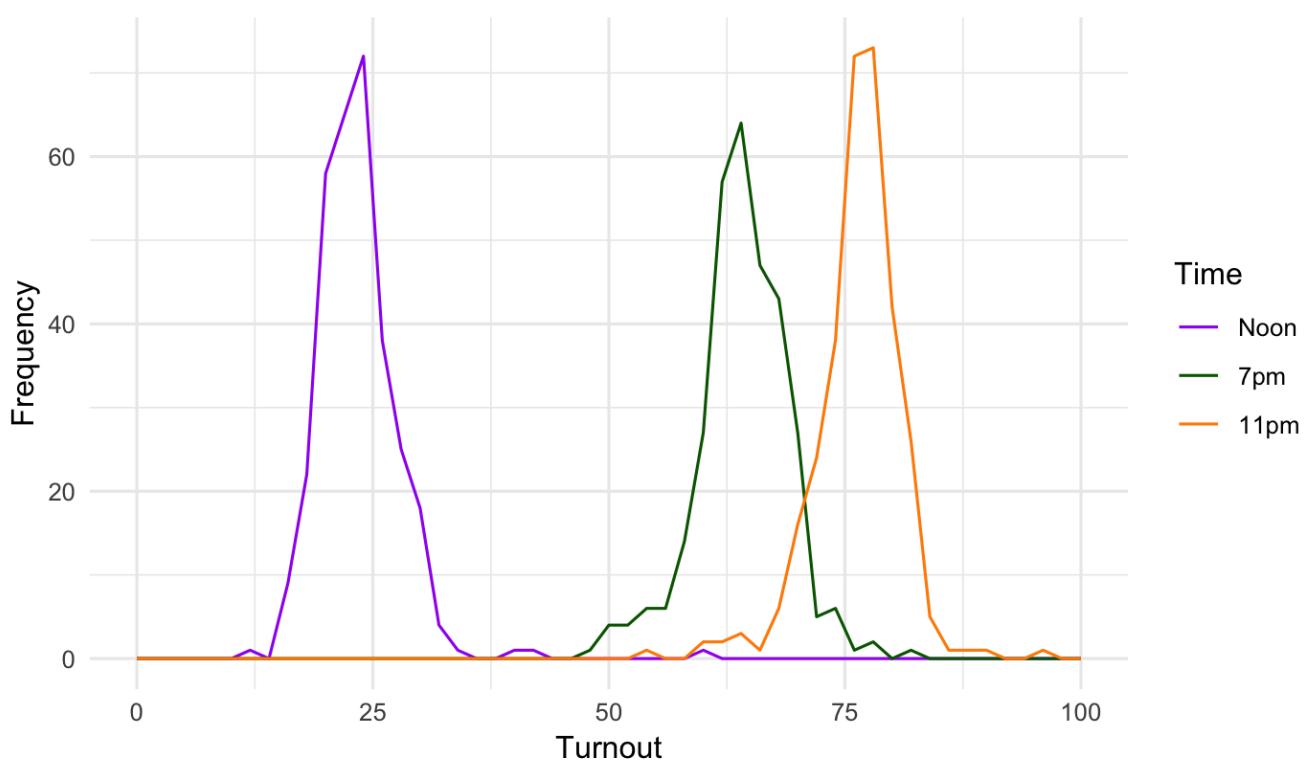
Paired diverging bargraph



<https://unstats.un.org/unsd/genderstatmanual/Print.aspx?Page=Presentation-of-gender-statistics-in-graphs>

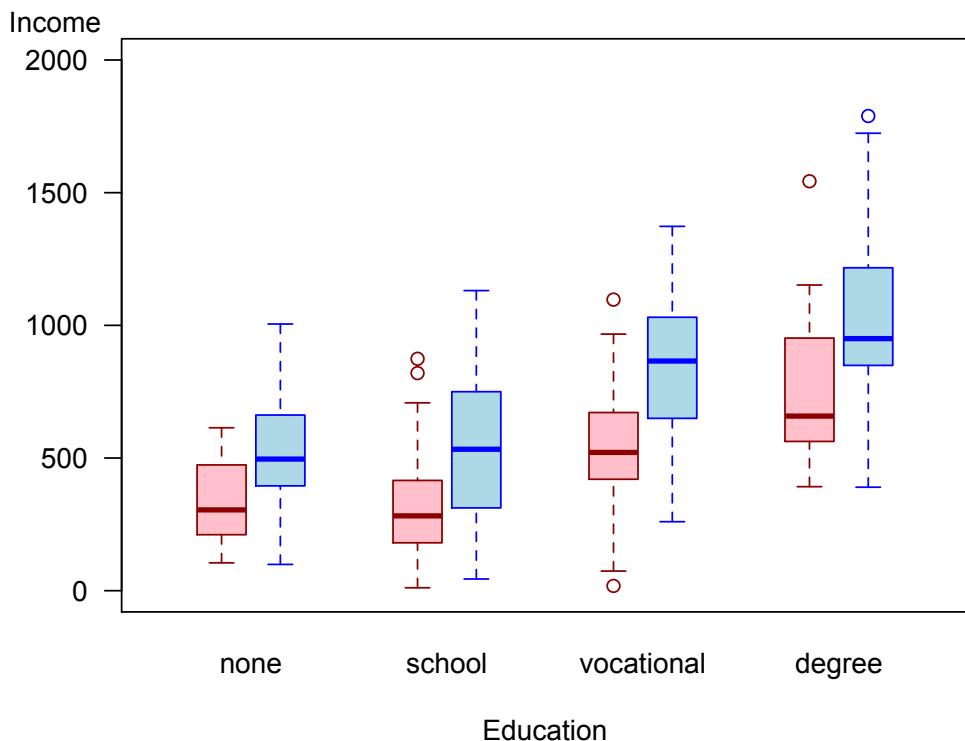
73

Multiple Frequency polygons



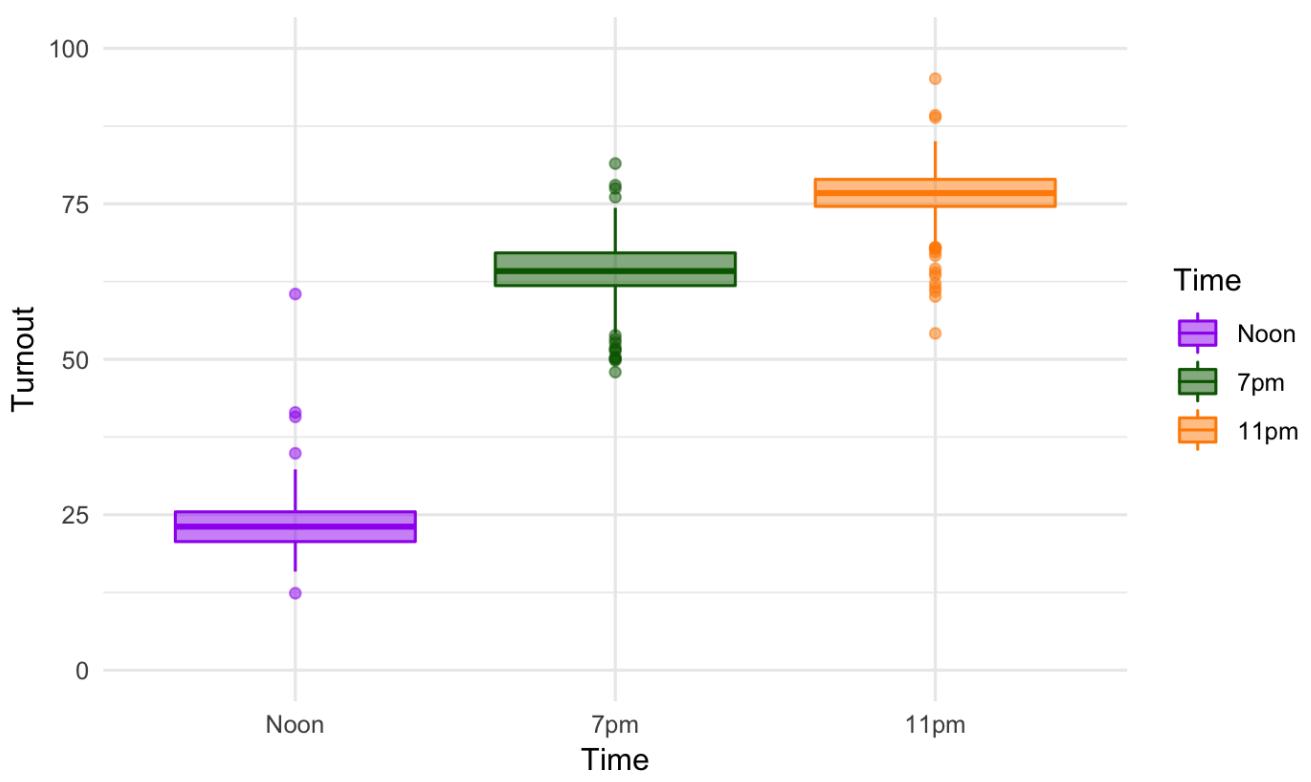
74

Box plot



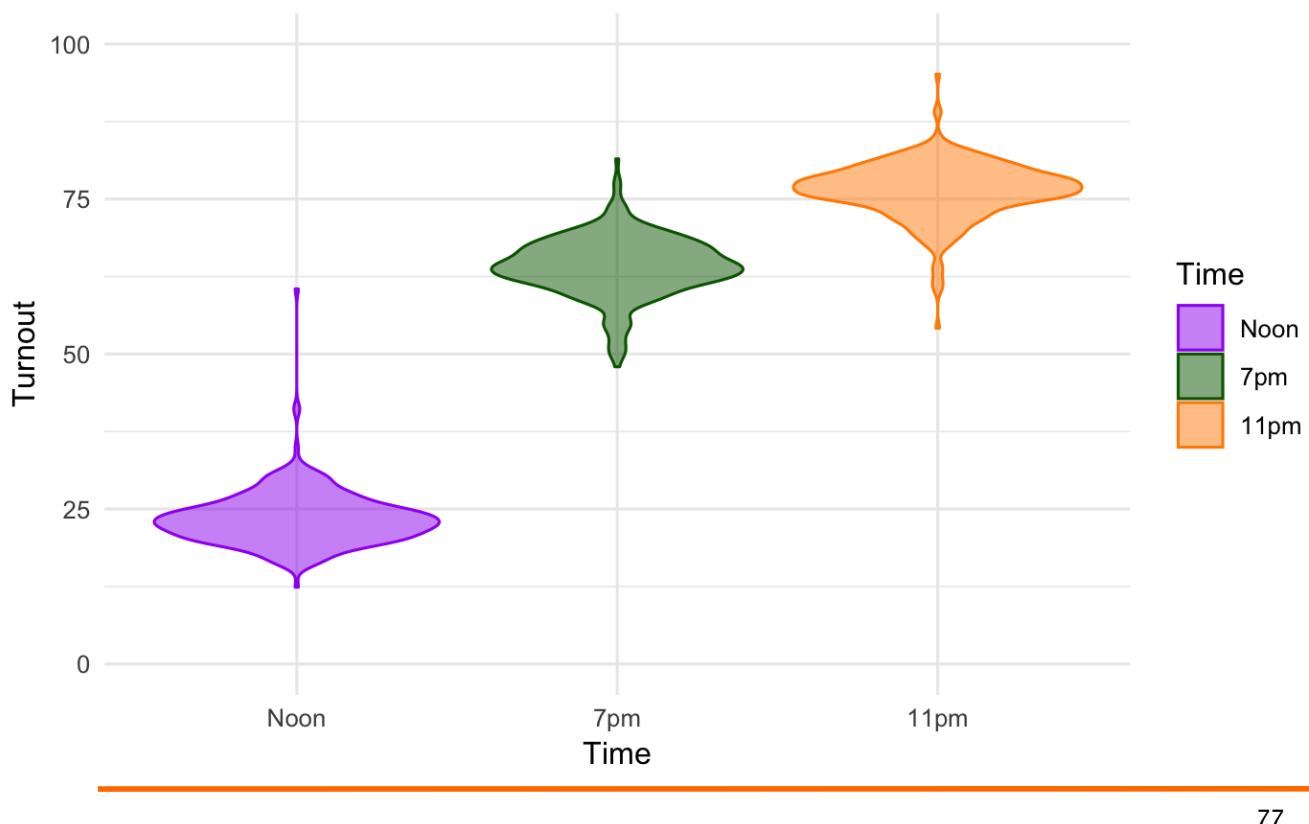
75

Multiple Box plot

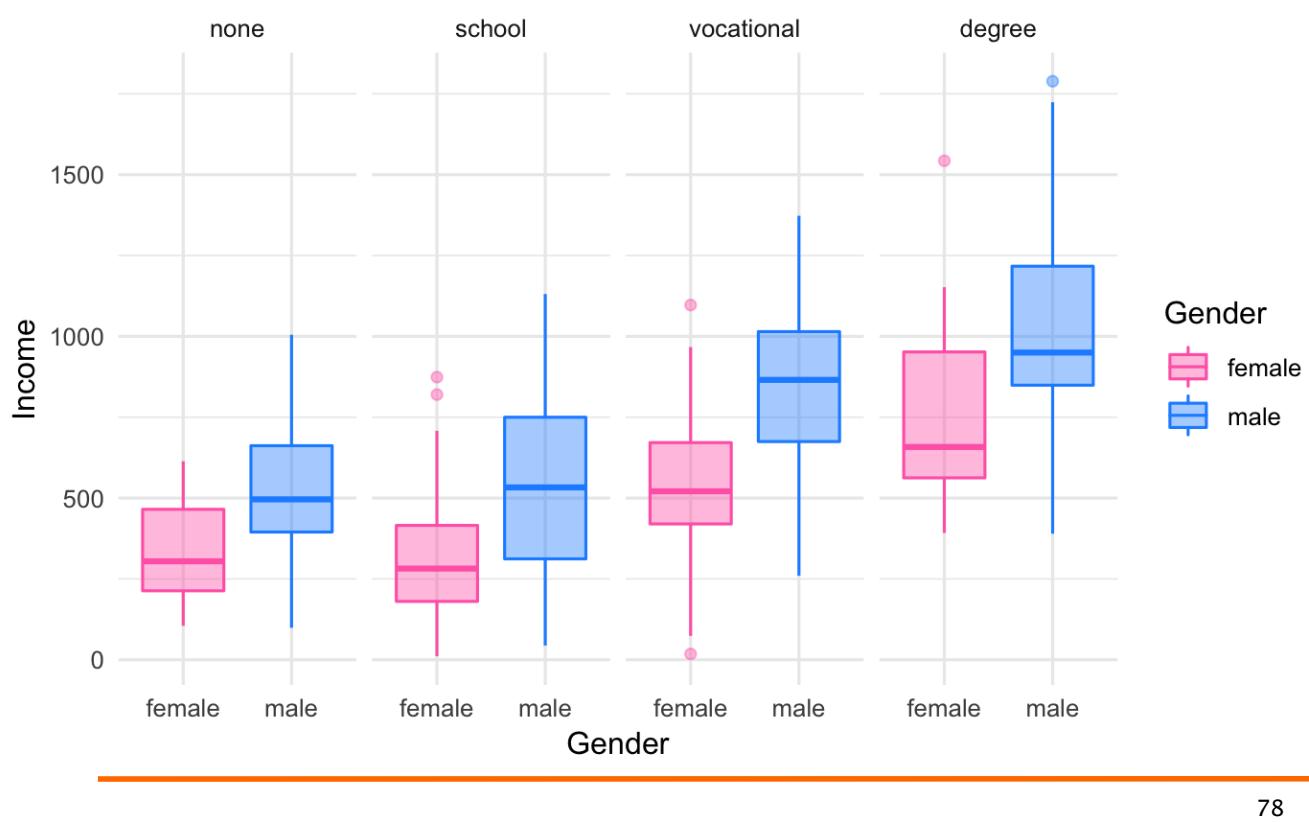


76

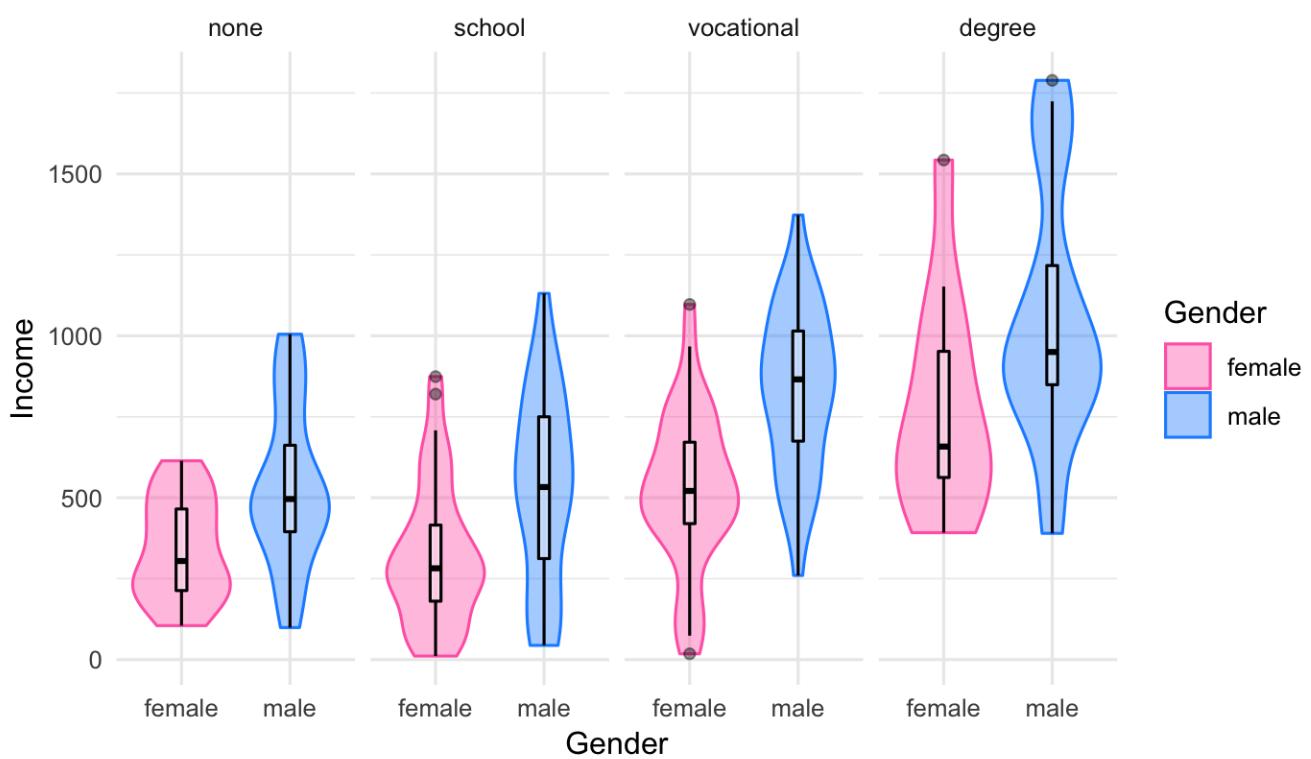
Violin plot



Multiple box plots

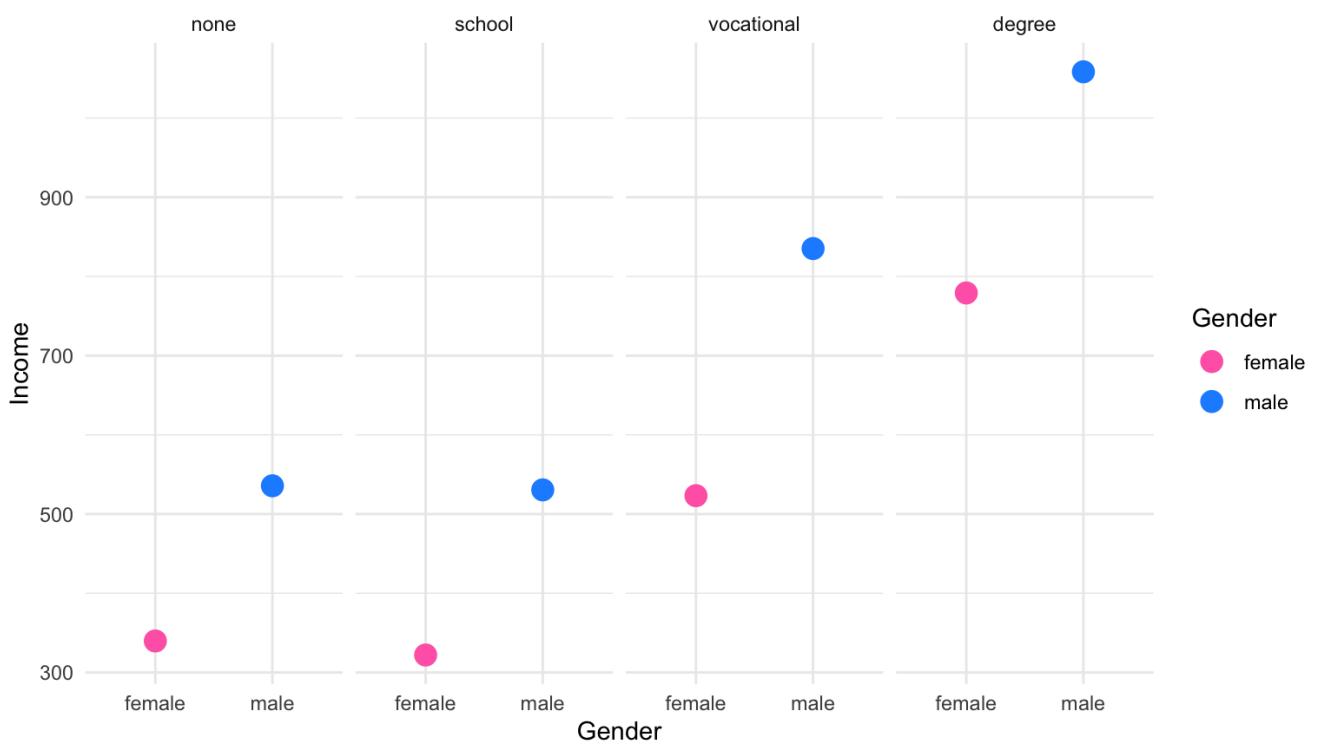


Multiple violin plots



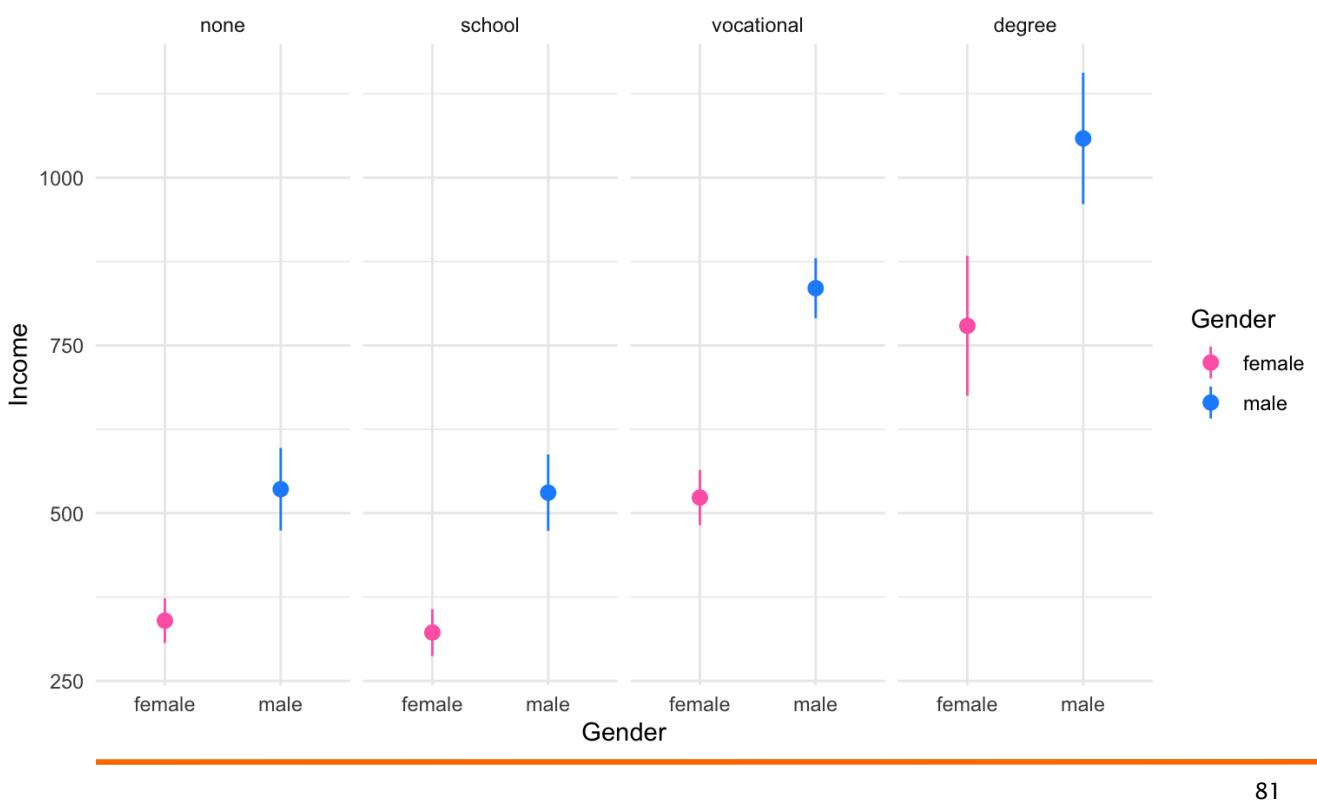
79

Just dots for mean values



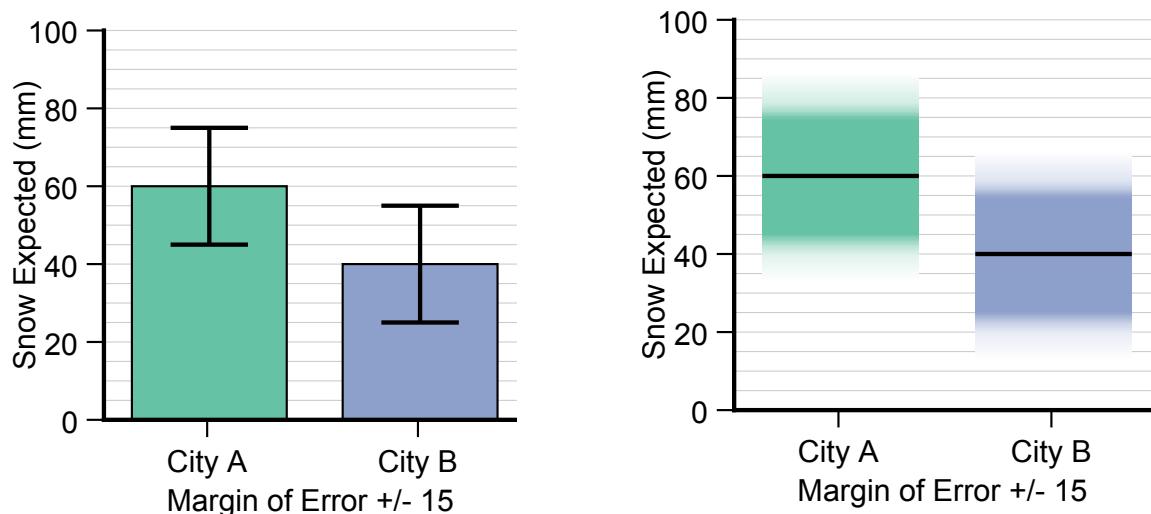
80

Confidence intervals



81

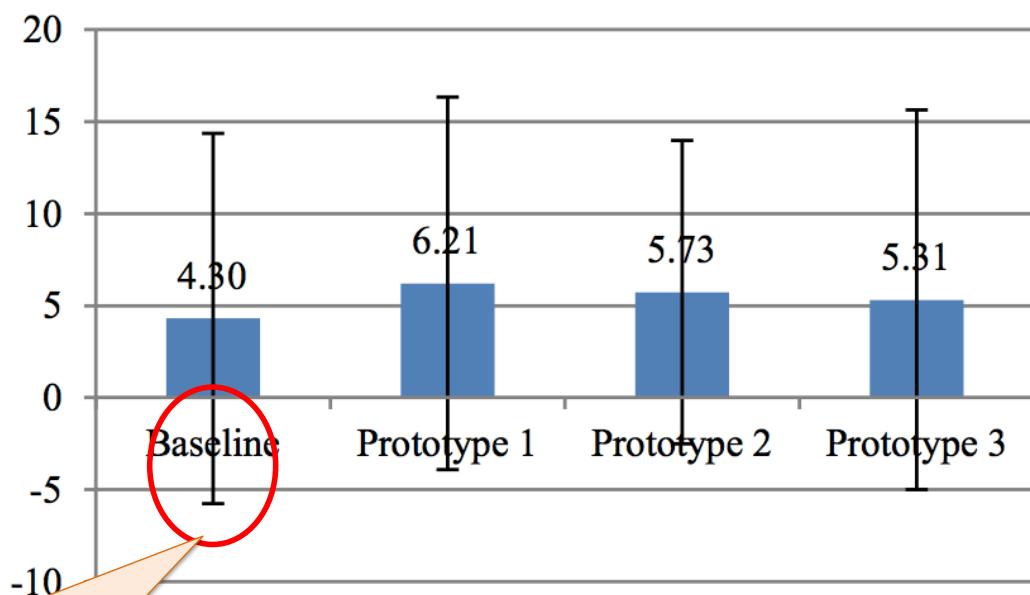
Confidence Intervals



Error Bars Considered Harmful: Exploring Alternate Encodings for Mean and Error
Michael Correll, and Michael Gleicher
IEEE Transactions on Visualization and Computer Graphics, Dec. 2014

82

Interval may be Asymmetric



It is physically impossible to modify -6 files

Figure 5. Mean files per changeset.

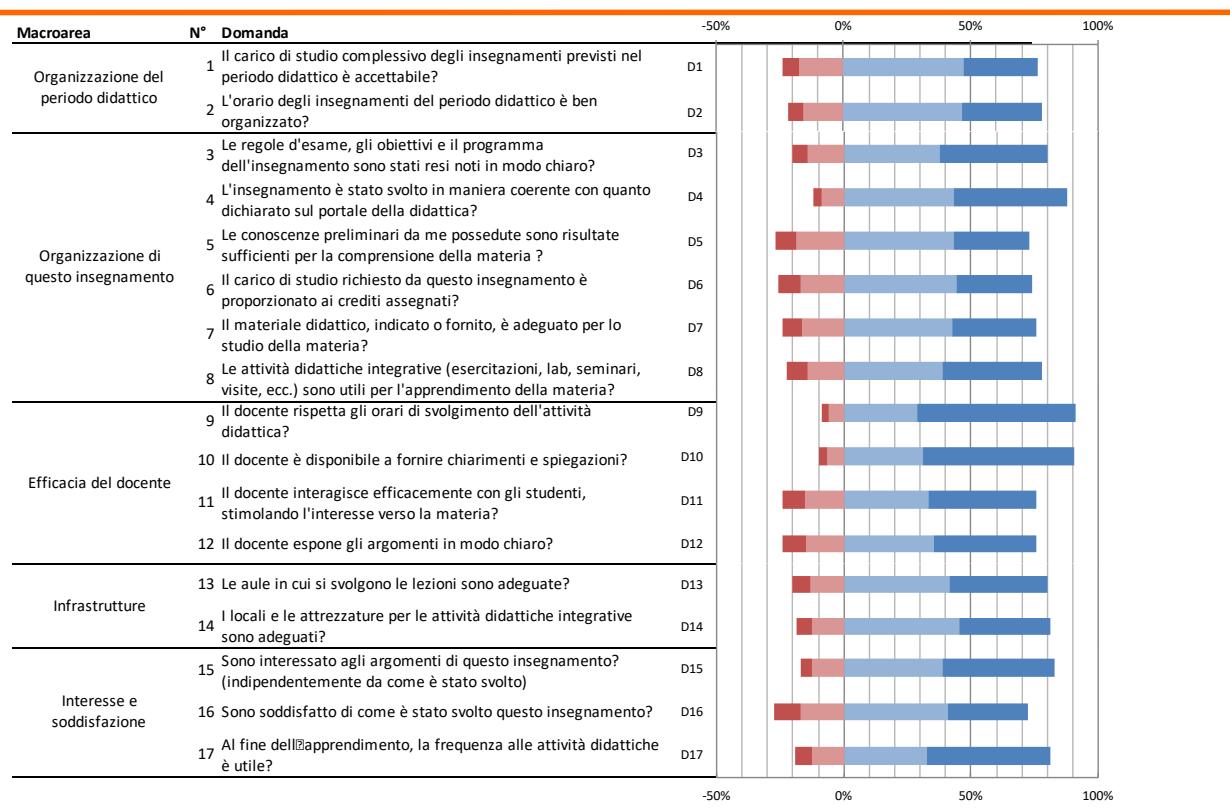
83

Likert / Agreement

- Likert scale:
 - ◆ Measures agreement / disagreement with a given statement
 - ◆ Response on an ordinal scale, e.g.
 - Definitely No
 - Mostly No
 - Undecided
 - Mostly Yes
 - Definitely Yes
- Often used to measure positive vs. negative perception

84

Diverging stacked bars



85

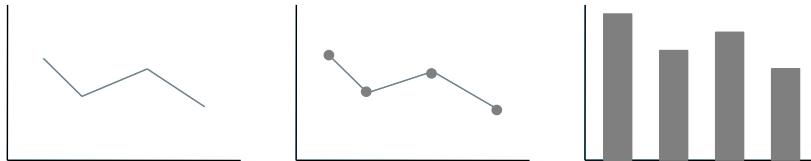
Time series

- Series of relationships between quantitative values that are associated with categorical subdivisions of time
- Communicate
 - ◆ Change
 - ◆ Rise
 - ◆ Increase
 - ◆ Fluctuate
 - ◆ Grow
 - ◆ Decline
 - ◆ Decrease
 - ◆ Trend

86

Time series

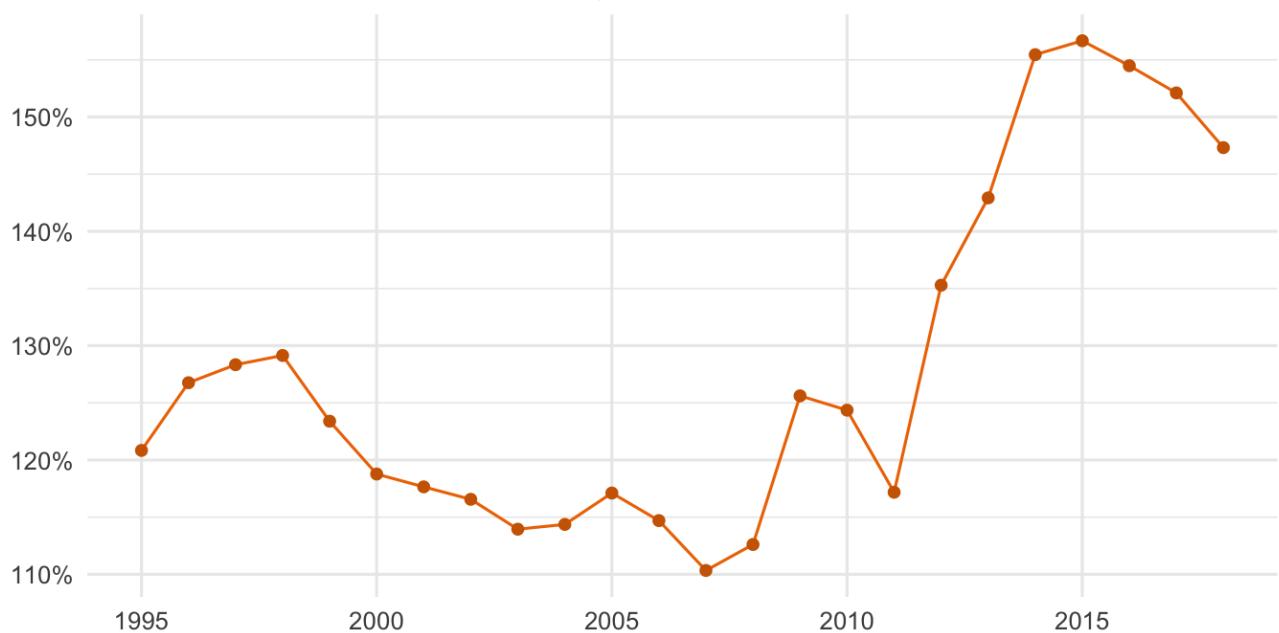
- Time grows from left to right
 - ◆ Cultural convention
- Vertical bars
 - ◆ highlight individual points in time
 - ◆ hide overall trend



87

Line plot

Italian Public Debt as percentage of GDP

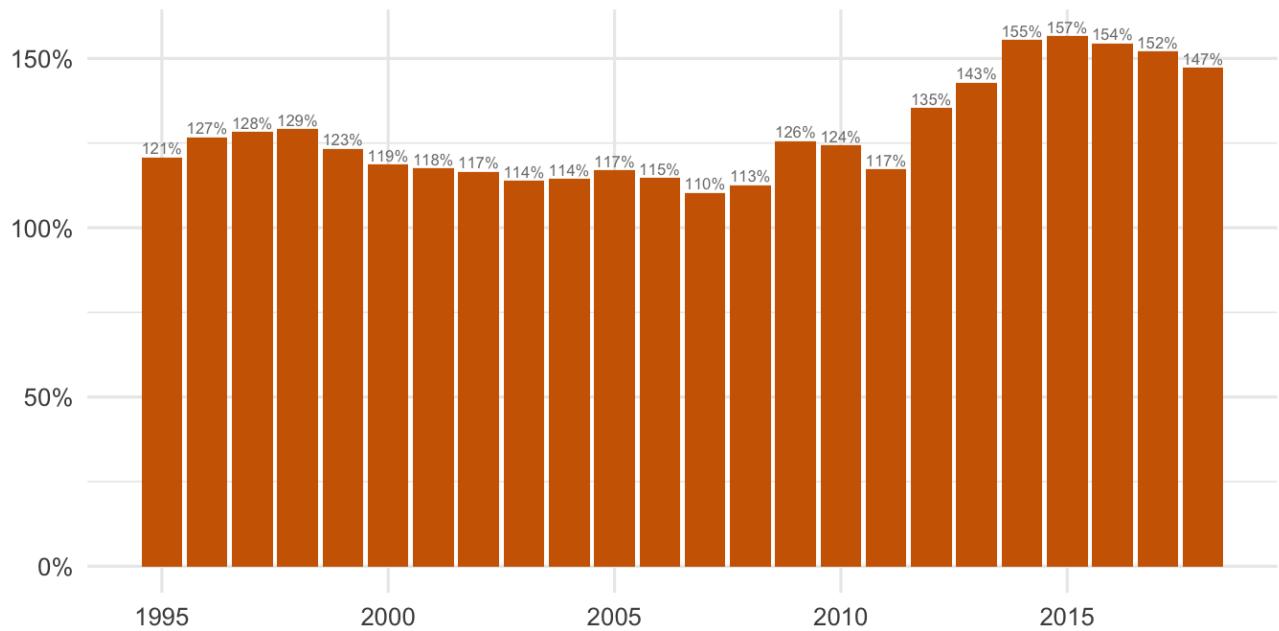


Source: OECD - <https://data.oecd.org/chart/5M2J>

88

Bars

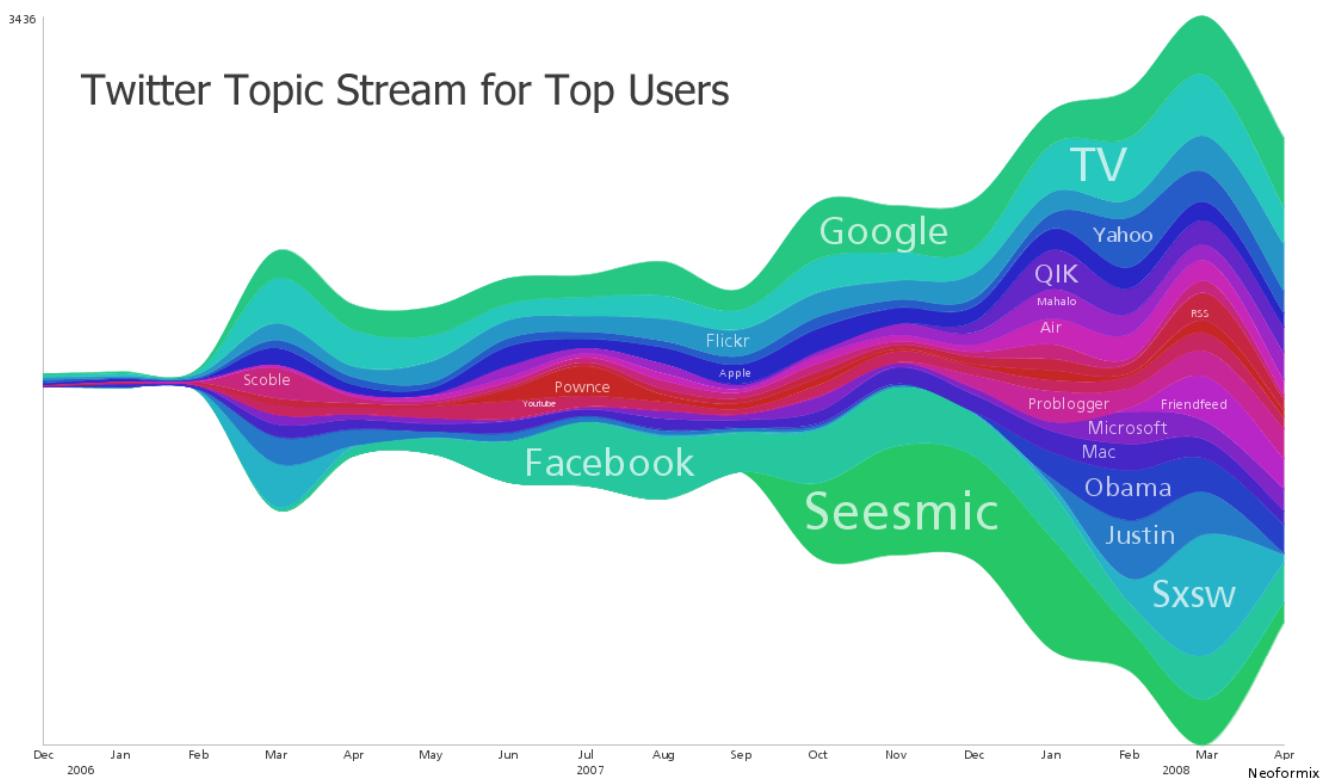
Italian Public Debt as percentage of GDP



Source: OECD - <https://data.oecd.org/chart/5M2J>

89

Streamgraph

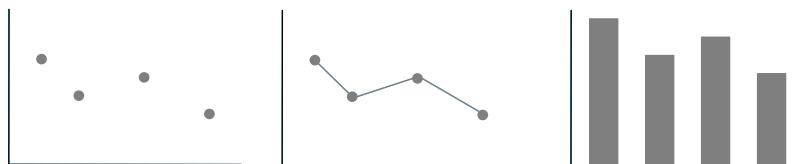


<http://www.neoformix.com/2008/TwitterTopicStream.html>

90

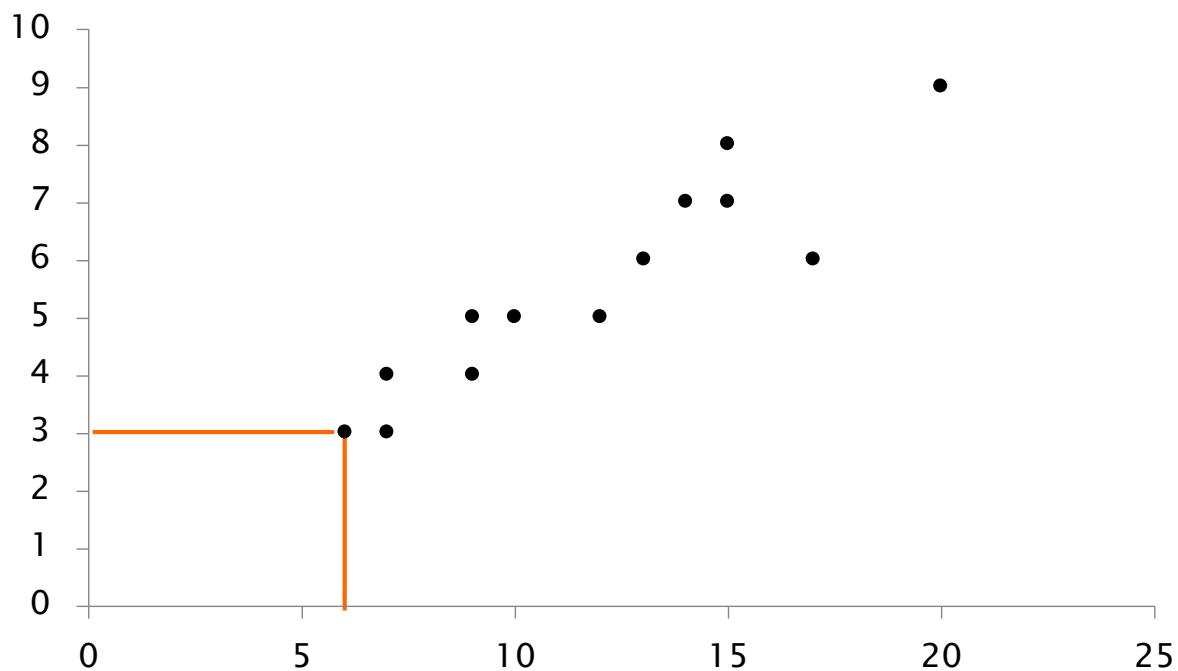
Correlation

- Relationships between two paired sets of quantitative values
 - ◆ Scatter plot w/possible trend line
 - Ok for educated audience
 - ◆ Paired bar graph



91

Points



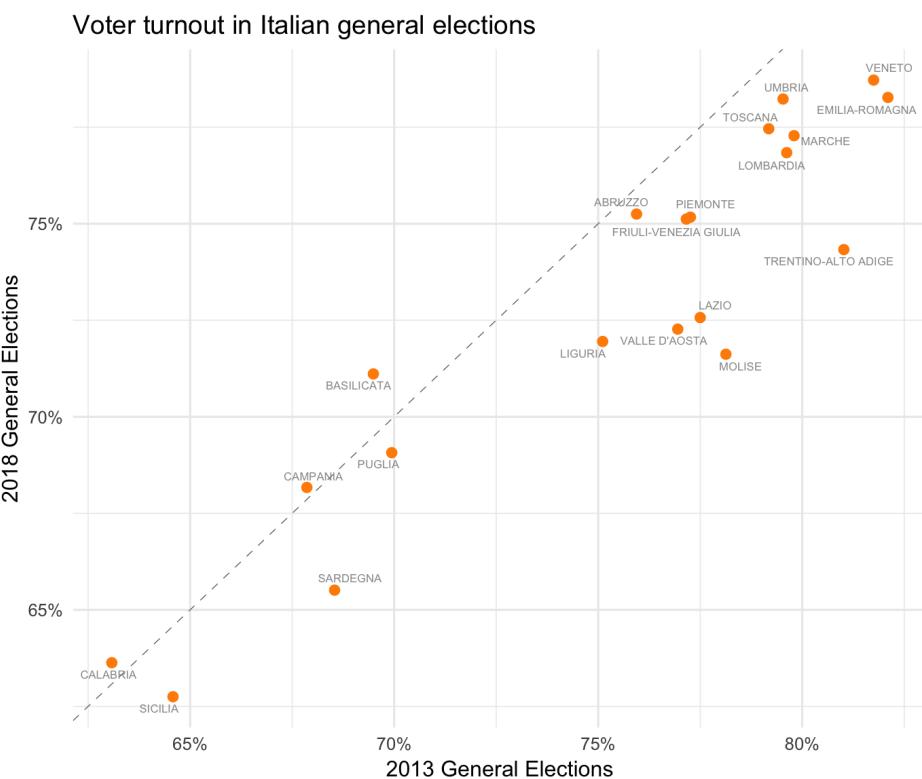
92

Points Guidelines

- Points must be clearly distinguished
 - ◆ Enlarge points
 - ◆ Select radically distinct shapes (+ ○)
 - ◆ Balance size of points and graph
 - ◆ Use outlined shapes
- Lines must not obscure points

93

Scatter plot



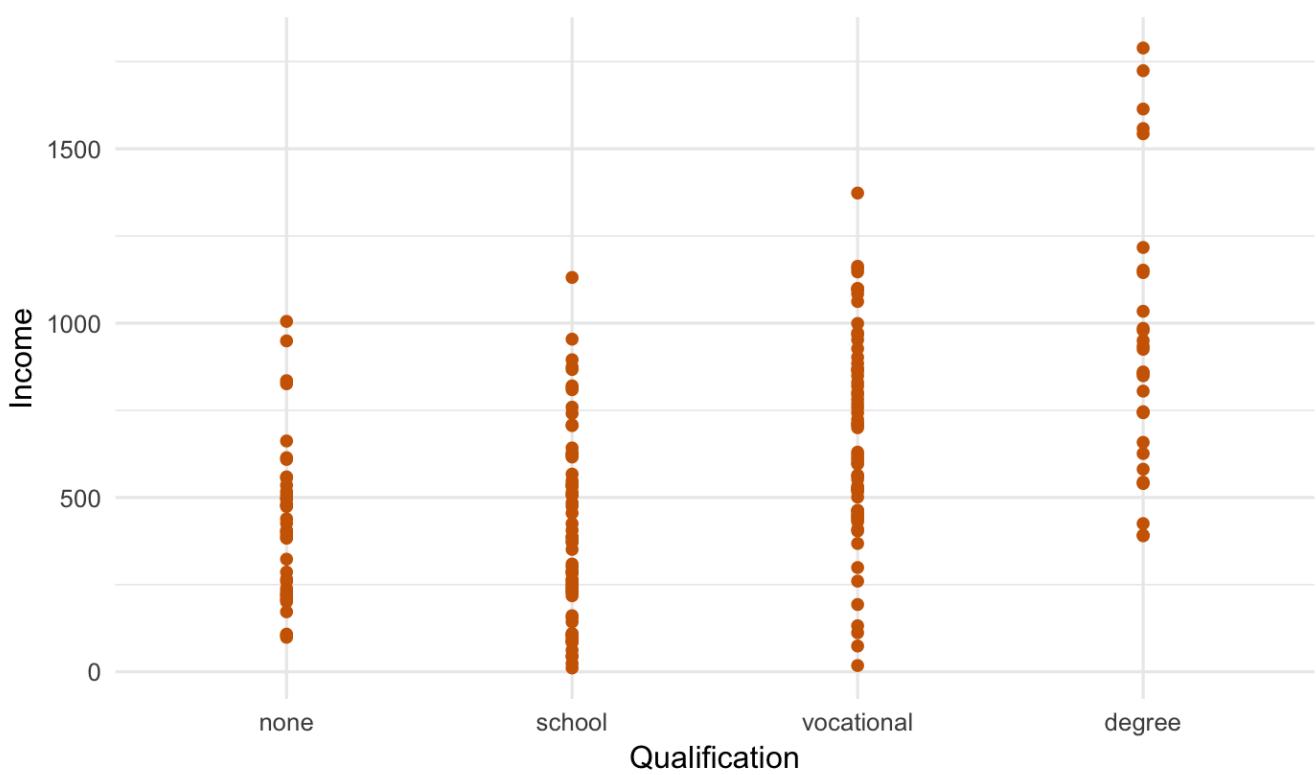
94

Overplotting

- Phenomenon related to multiple points (or shapes) overlapping
 - ◆ Discrete (integer) measure
 - ◆ Very large dataset
- Solutions
 - ◆ Small shapes
 - ◆ Outlined shapes
 - ◆ Transparent shapes (alpha)
 - ◆ Jittering

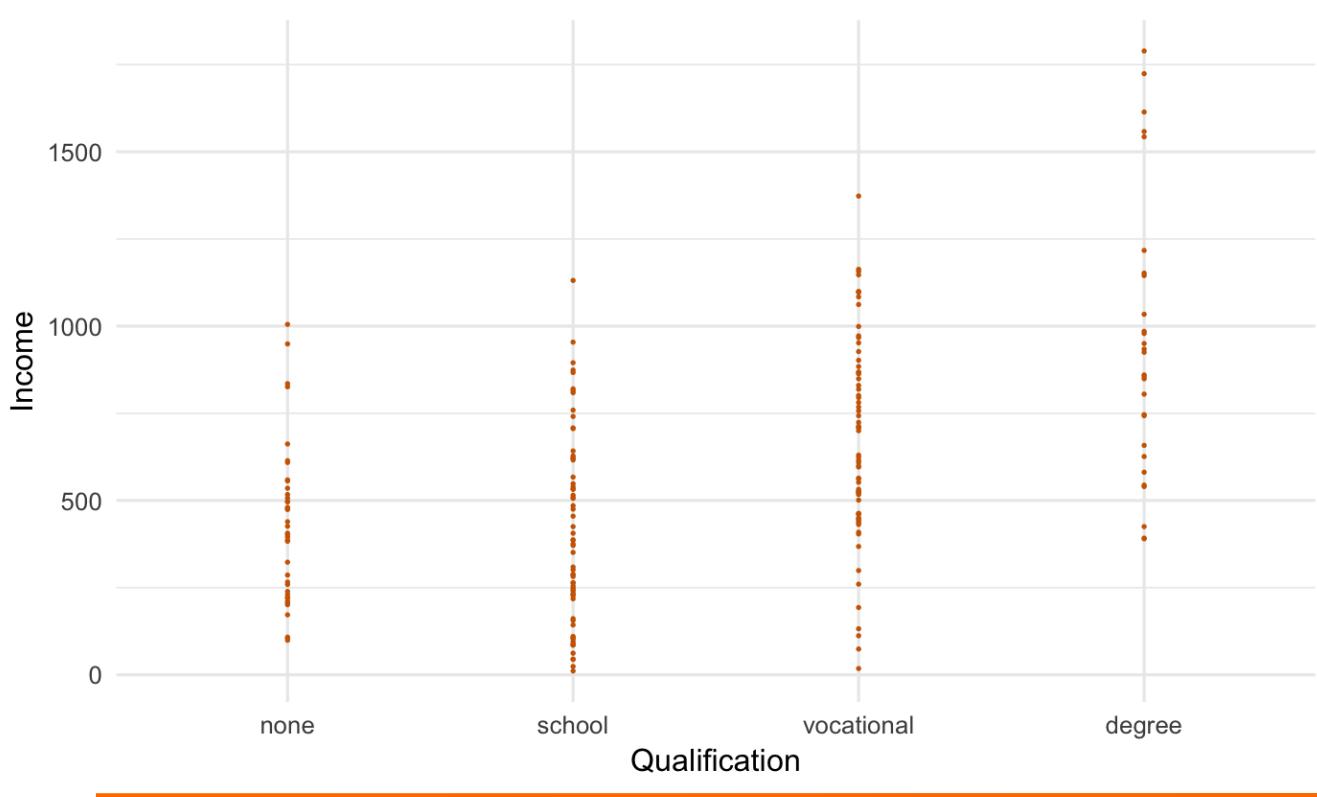
95

Overplotting example



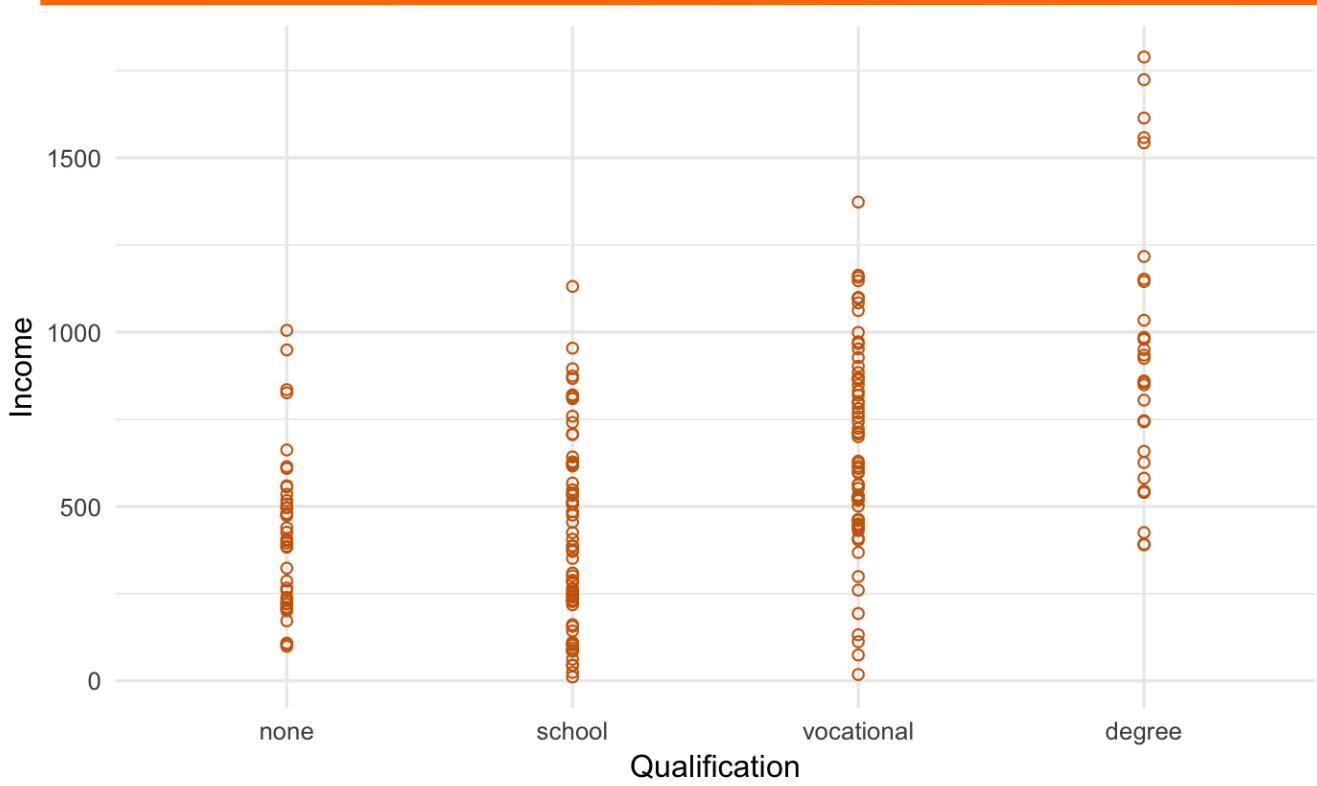
96

Overplotting – Small



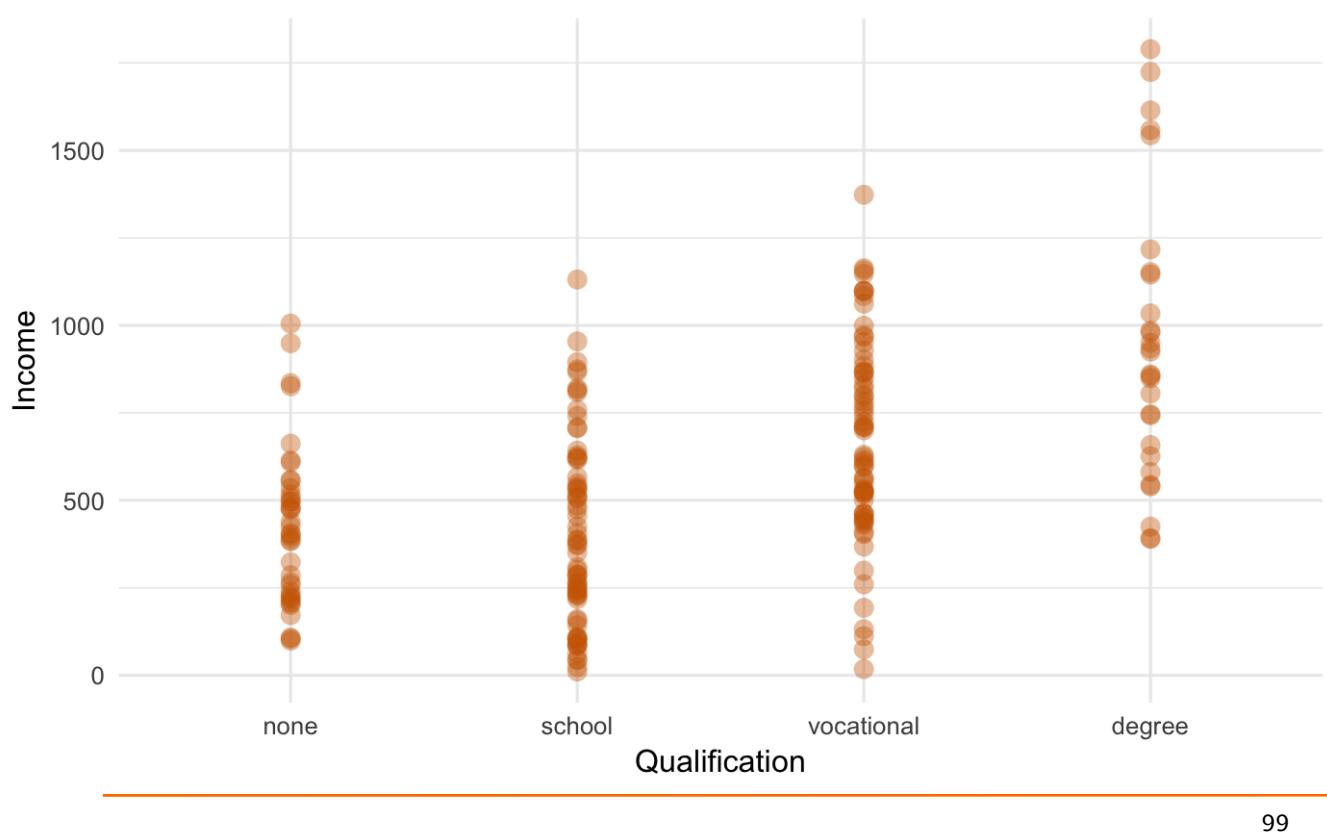
97

Overplotting – Outlined



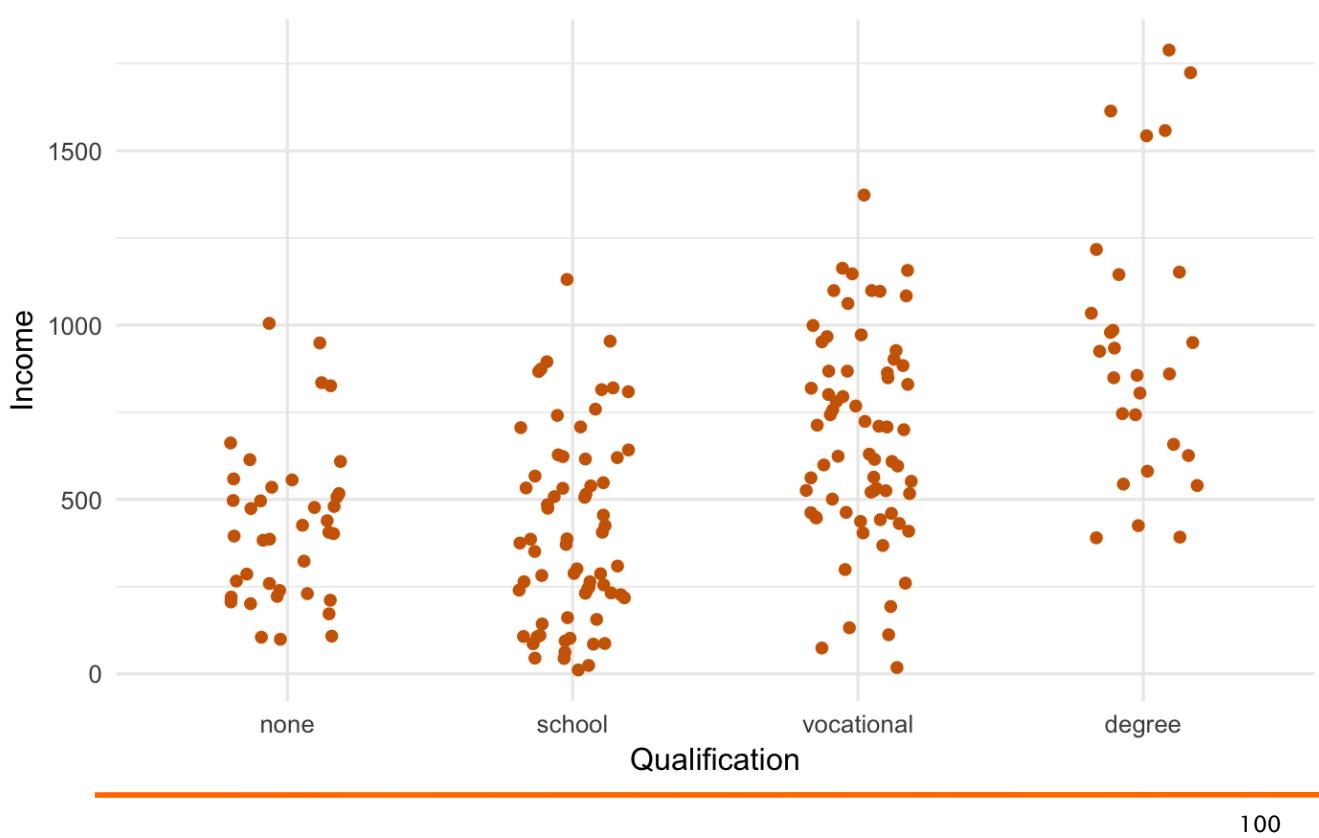
98

Overplotting – Transparent



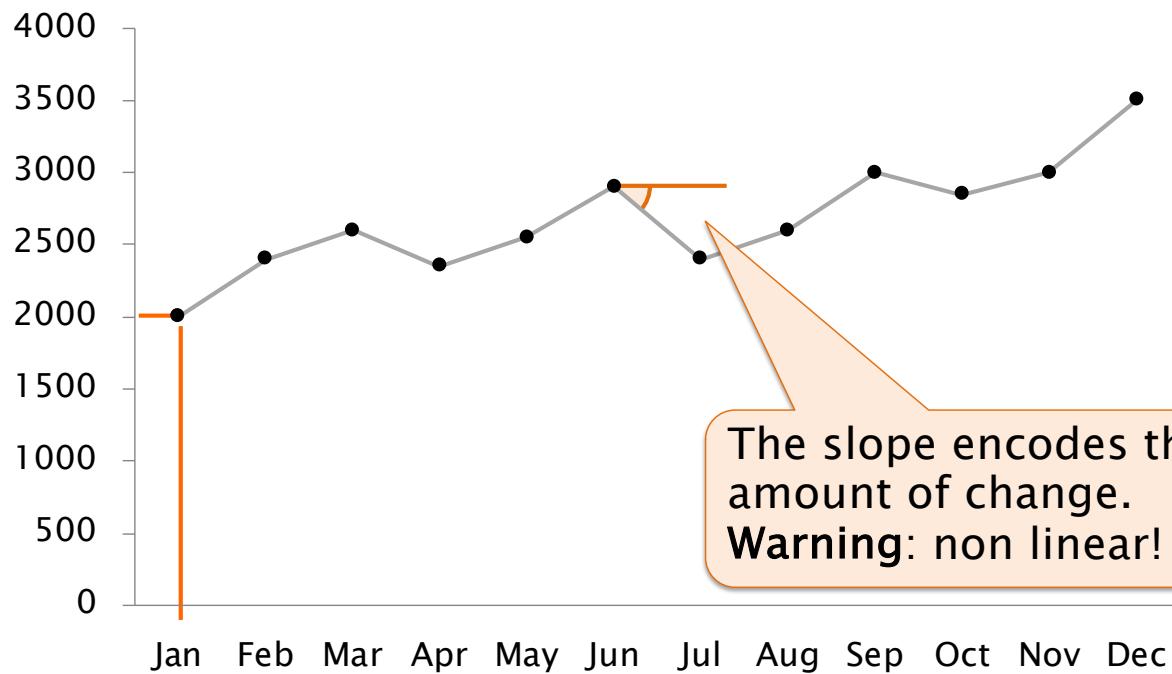
99

Overplotting – Jittering



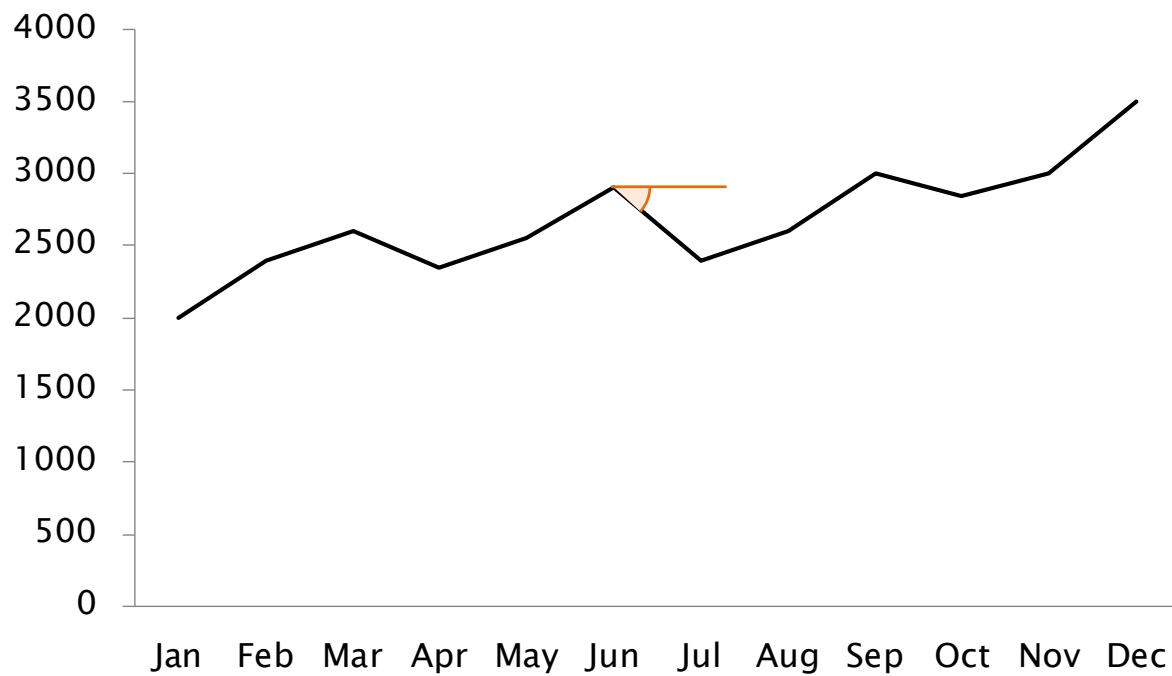
100

Points and Lines



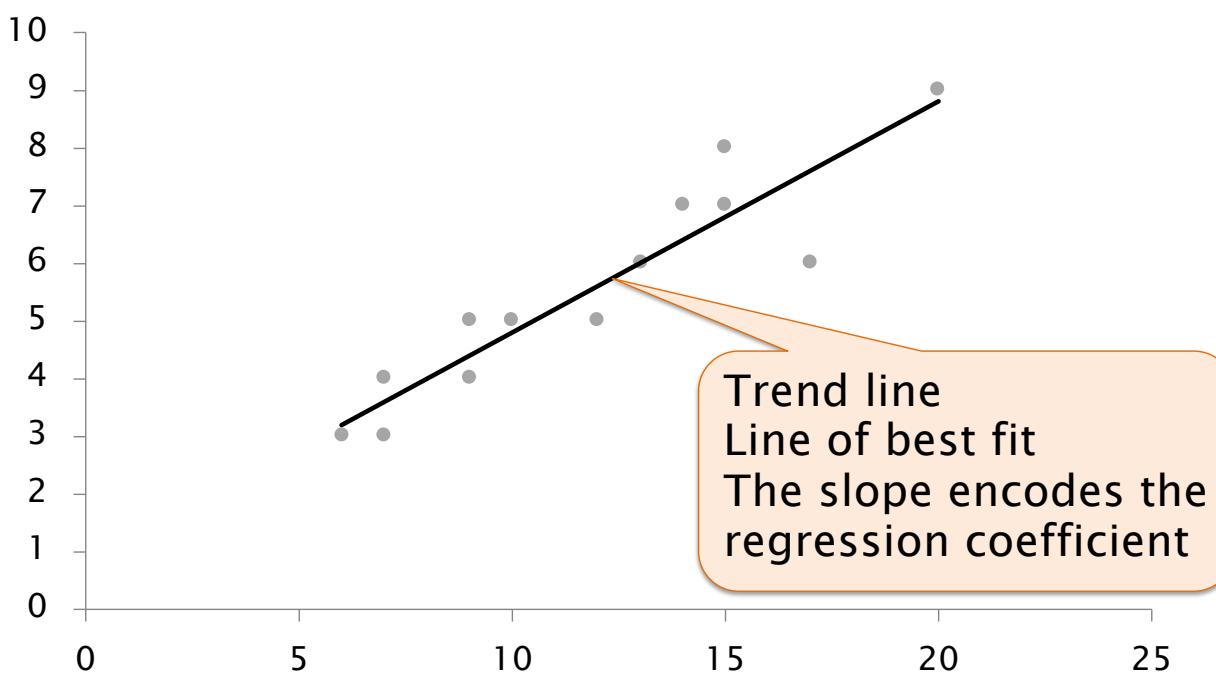
101

Slope of lines



102

Slope of lines



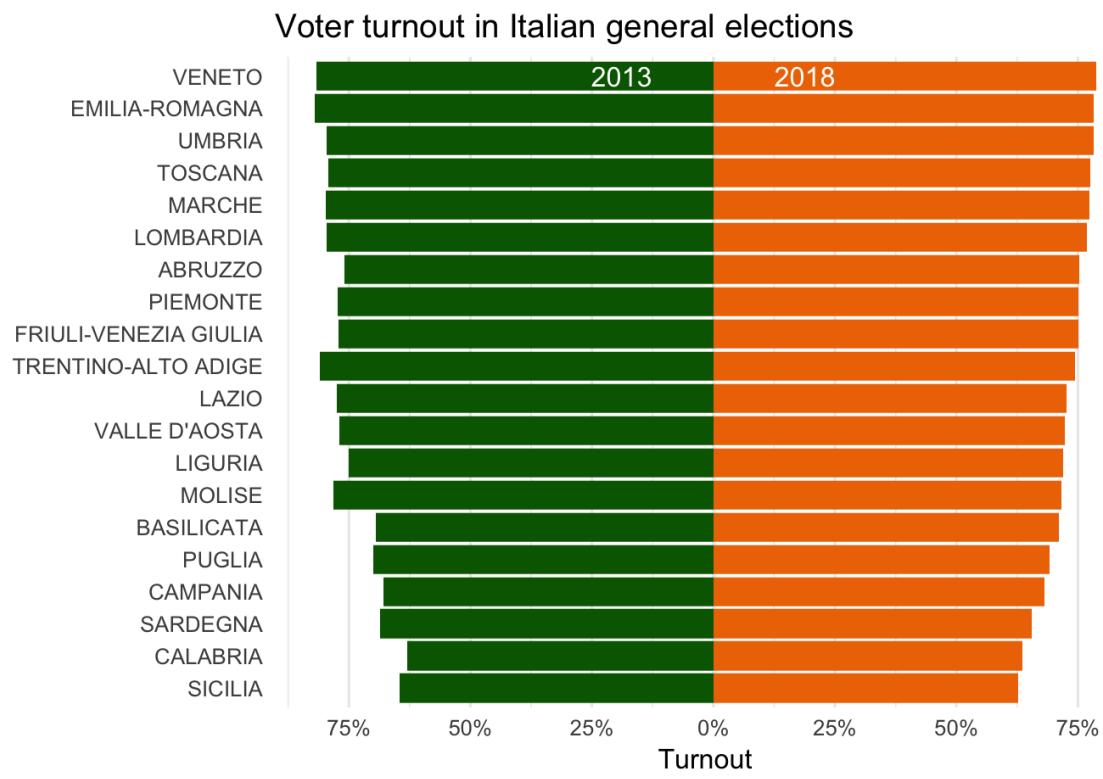
103

Lines

- Easy perception of trends and overall shape of data
- Best suited for time series
- Variation encoded as slope
 - ◆ Clear direction
 - ◆ Approximate magnitude

104

Paired diverging bars



105

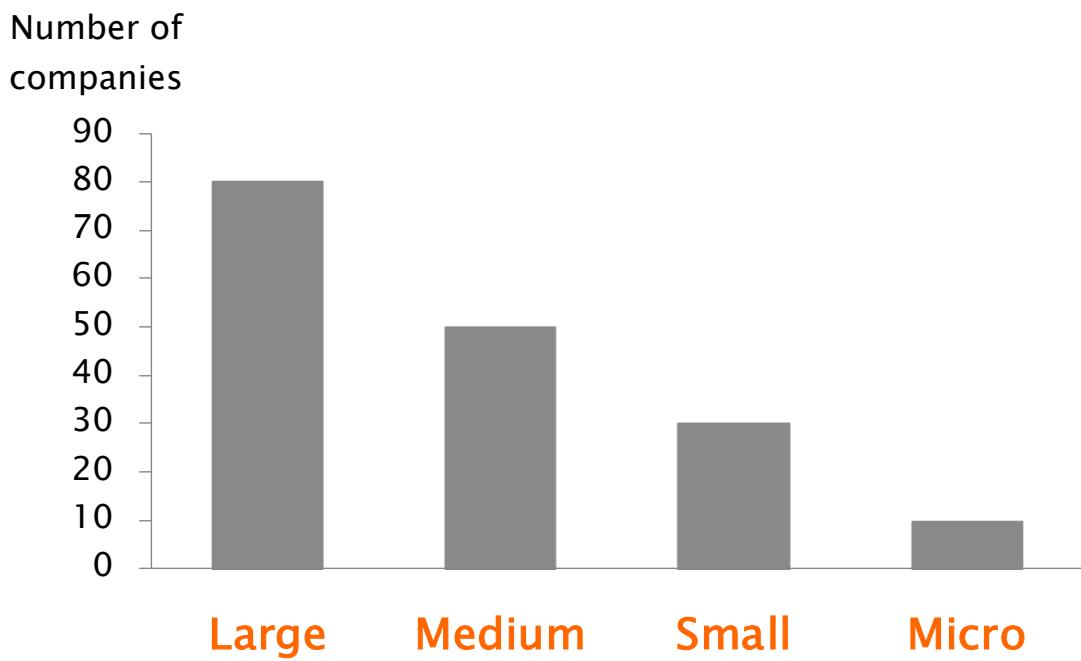
Categorical encoding attributes

- Encoding of categorical levels
 - ◆ Position (along an axis)
 - ◆ Size
 - ◆ Color
 - Intensity
 - Saturation
 - Hue
 - ◆ Shape
 - ◆ Fill pattern
 - ◆ Line style

Ordinal

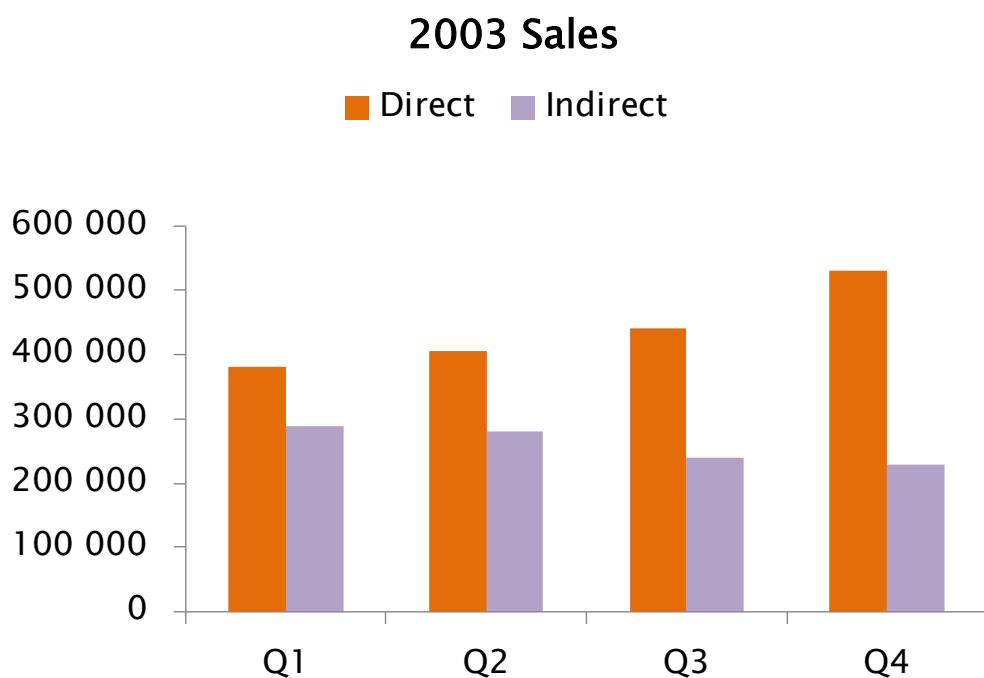
106

Position



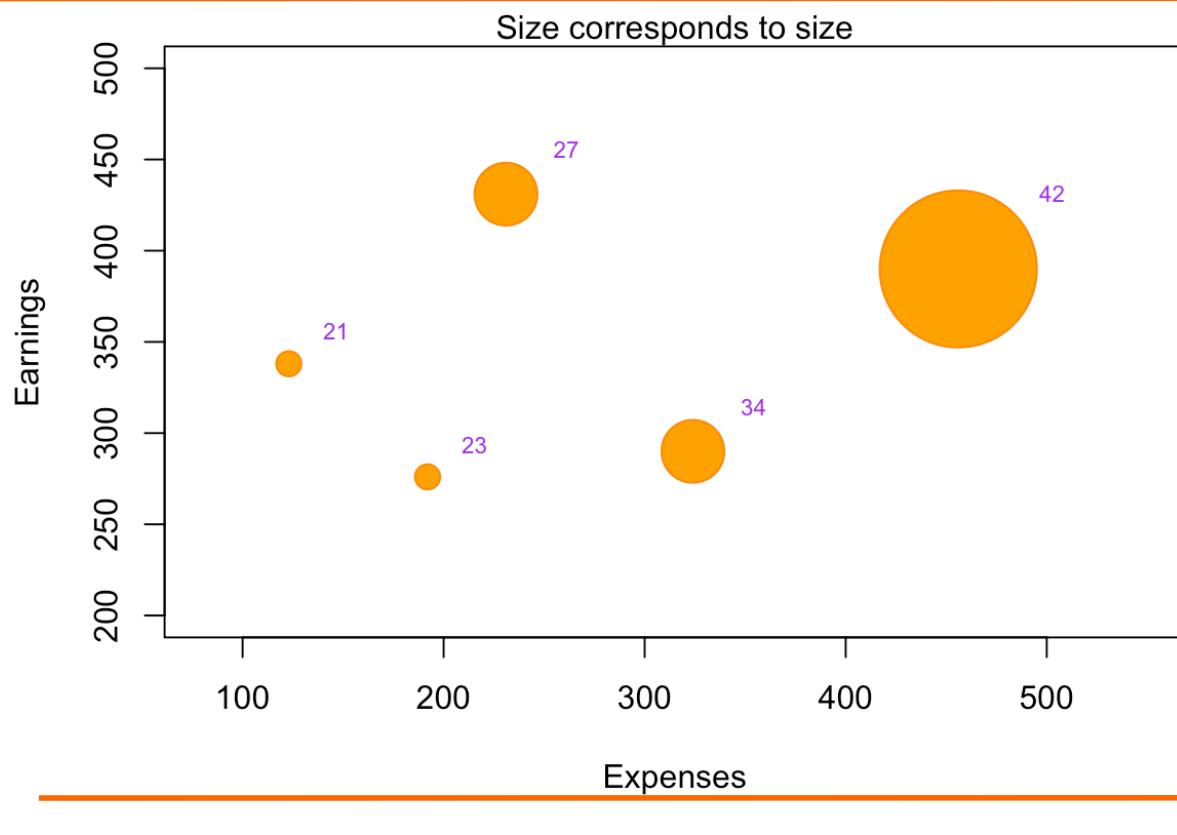
107

Position × Color (hue)

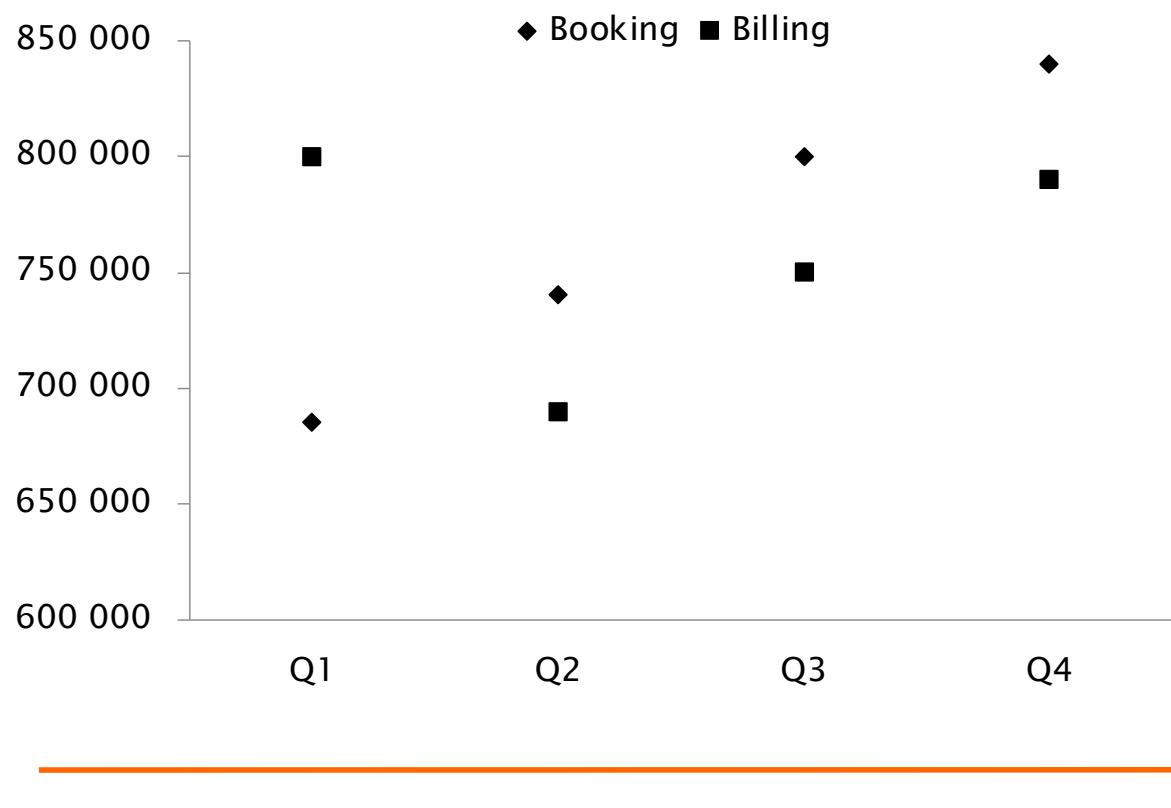


108

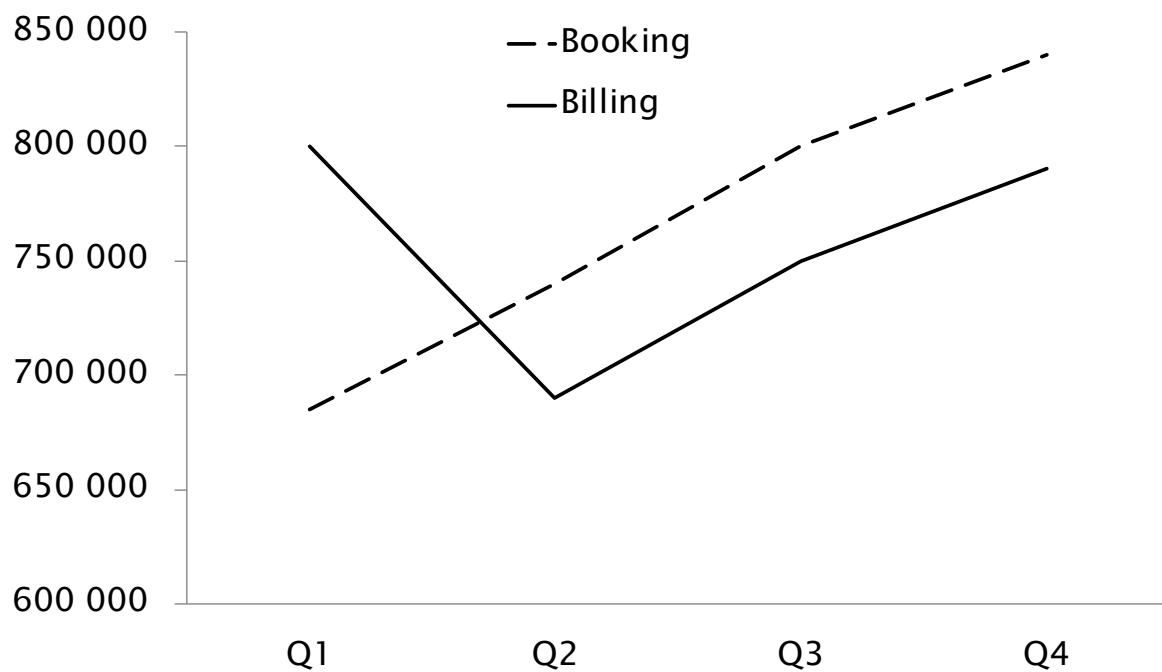
Size



Point shape

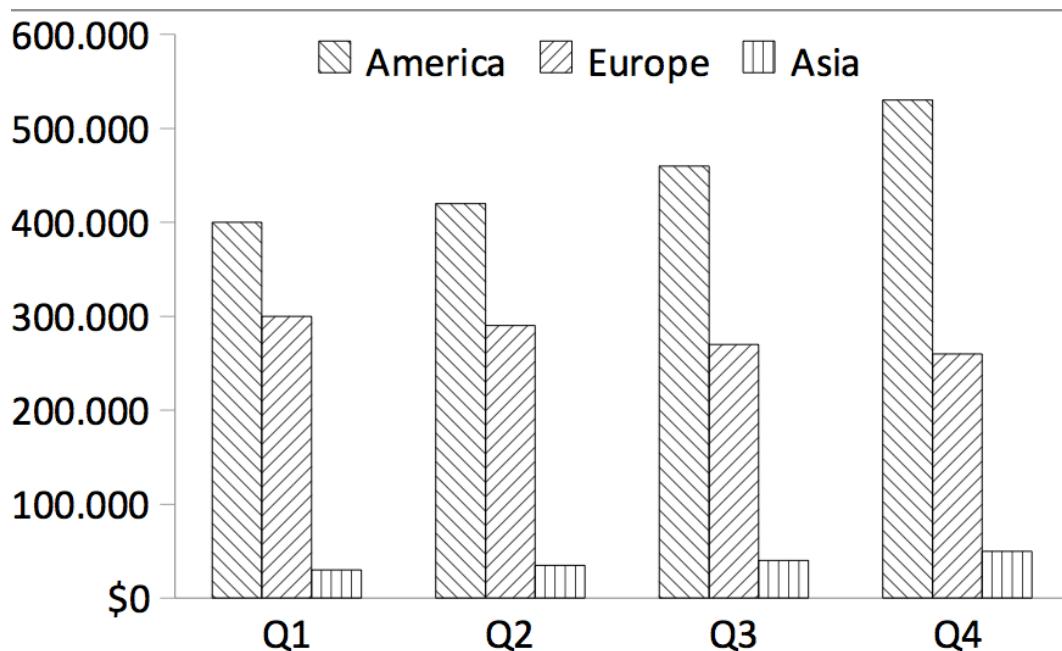


Line style



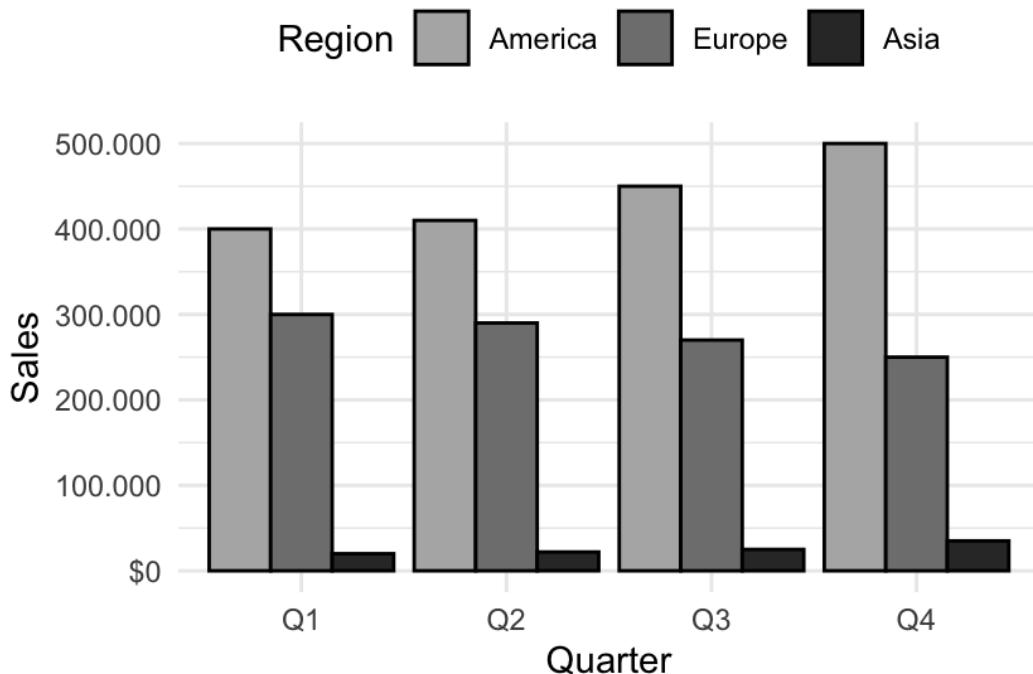
111

Fill Texture



112

Fill Texture



113

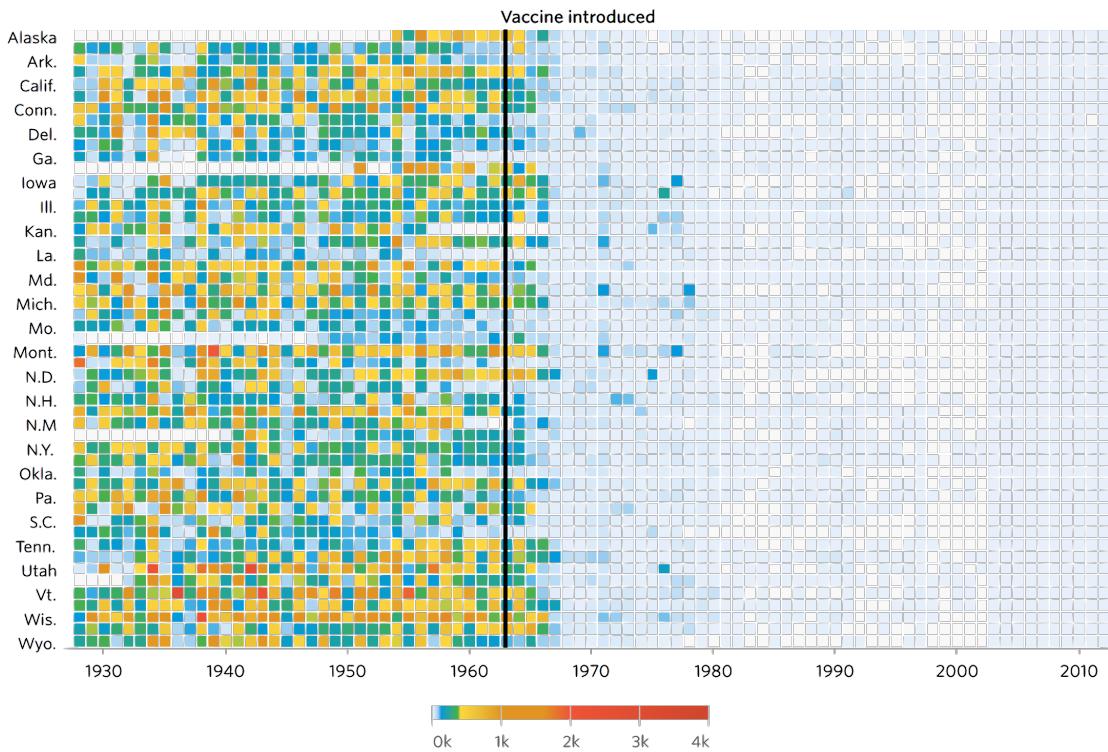
Discretization / Quantization

- A data transformation that maps a quantitative measure into an ordinal one
 - ◆ Based on the definition of intervals
- Discretized measures can be encoded using an ordinal-friendly visual attribute
 - ◆ Size
 - ◆ Color
- Warning: details are lost in the process

114

Heatmaps

Measles



<http://graphics.wsj.com/infectious-diseases-and-vaccines/>

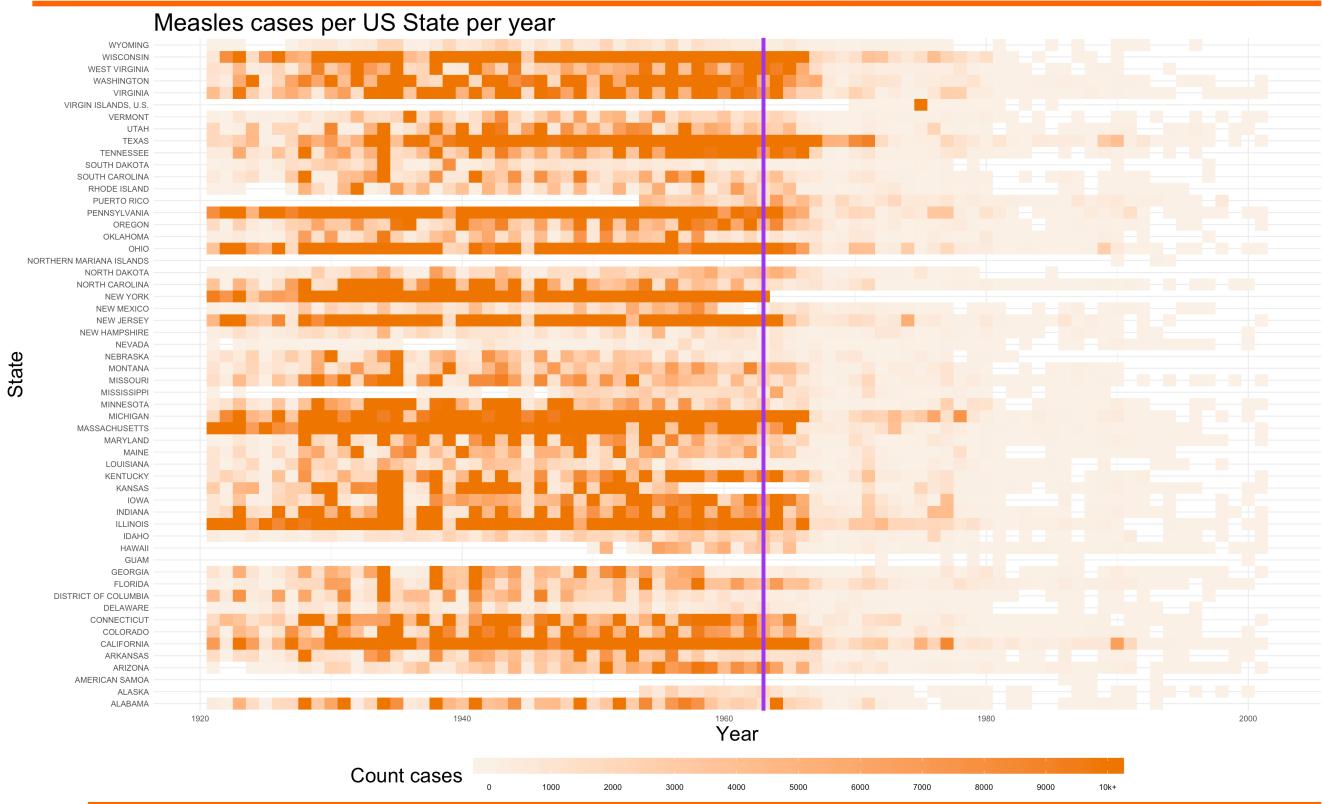
115

Heatmaps

- Hues have no unique order semantics
 - ◆ Only intensity has one
- Rainbow palette have serious problems for color blinds
 - ◆ Roughly 5% of the population

116

Heatmaps



117

SUPPORT ELEMENTS

118

Support elements

- Axes
 - ◆ Ticks
 - Graph area
 - ◆ Grids
 - Labels
 - Legends
 - References
 - Trellies
-

119

Axes

- Allow positioning of elements
 - ◆ Points
 - ◆ Extremes of bars and lines
- Labeled
 - ◆ What is the measure?
- Number of axis should be 2
 - ◆ 1 is fine for bars
 - continuity gestalt principle

120

Tick marks

- Must not obscure data objects
- Outside the data region
- Avoid for categorical scales
- Balanced number
 - ◆ Too many clutter the graph
 - ◆ Too few make difficult to discern reference for data objects
 - ◆ Intervals must be equally spaced

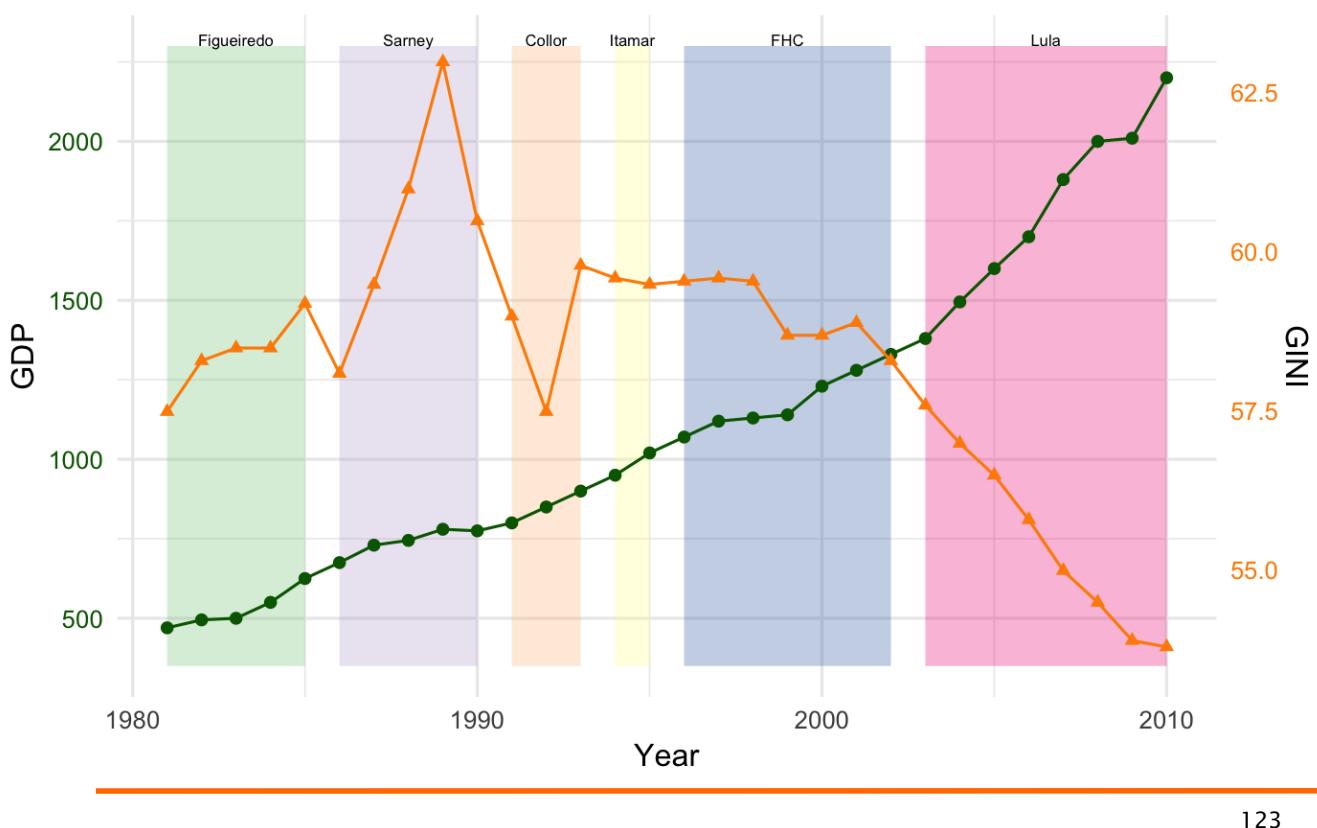
121

Multiple variables

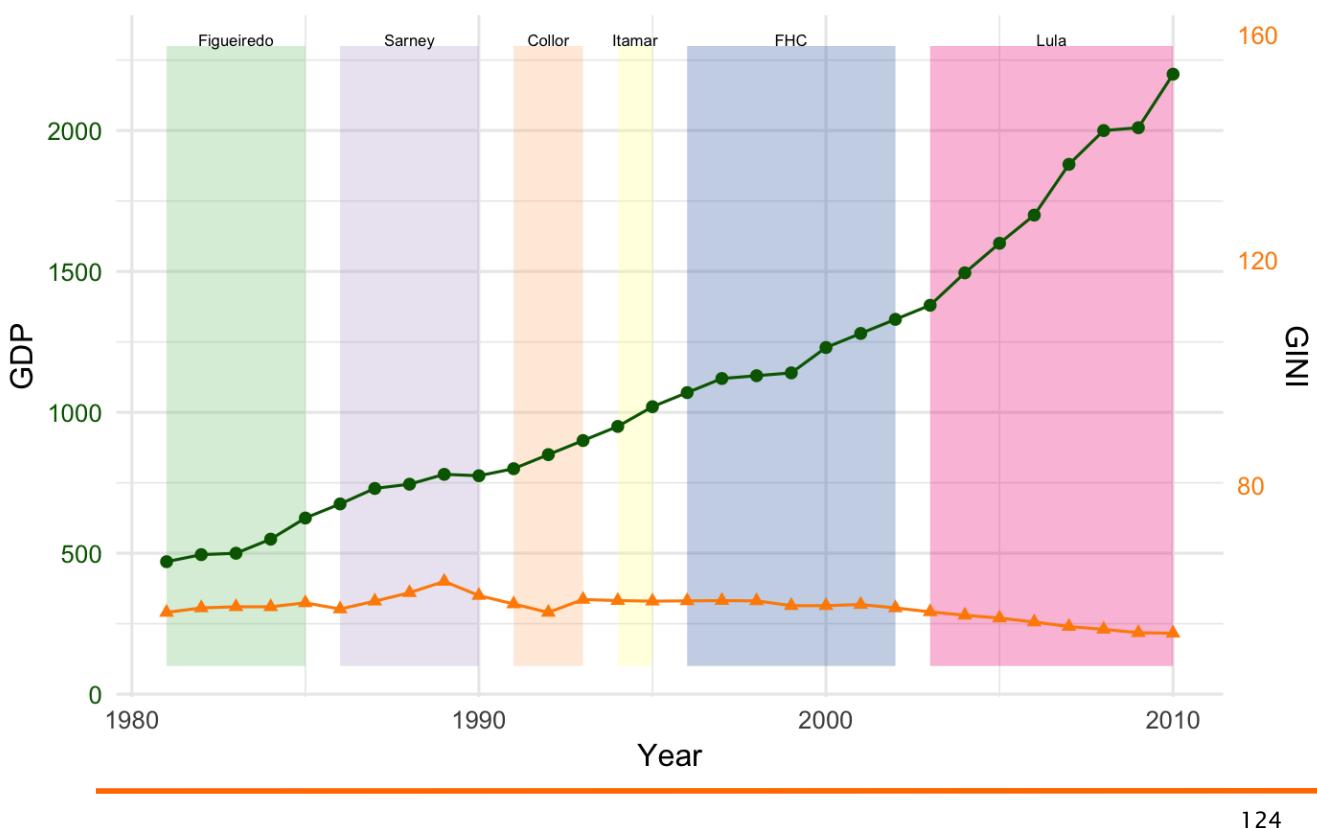
- Correlation between 3+ variables
 - ◆ E.g. two measures in time series
- Multiple units of measure
 - ◆ Double quantitative (y) axis
 - ◆ Multiple graphs
 - ◆ One variable not encoded explicitly

122

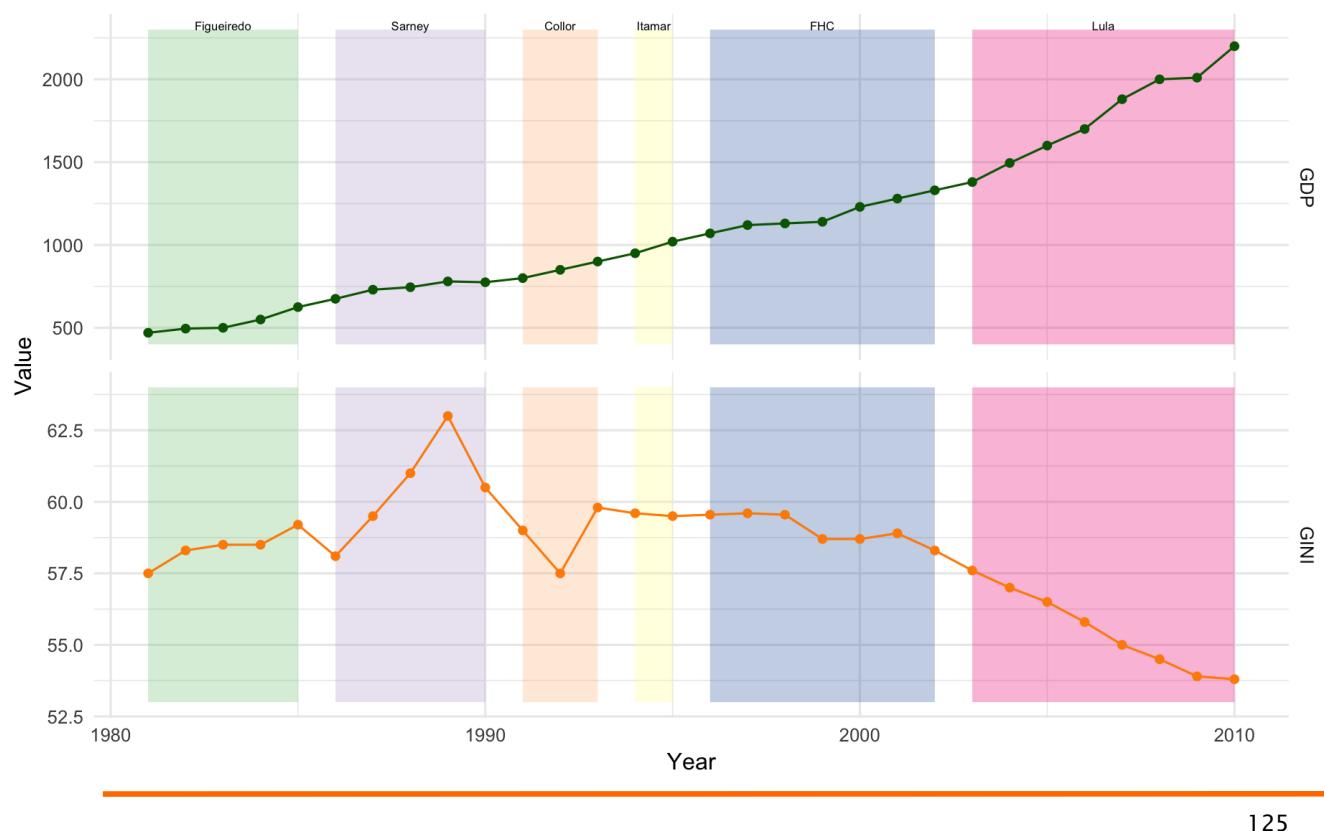
Double scale



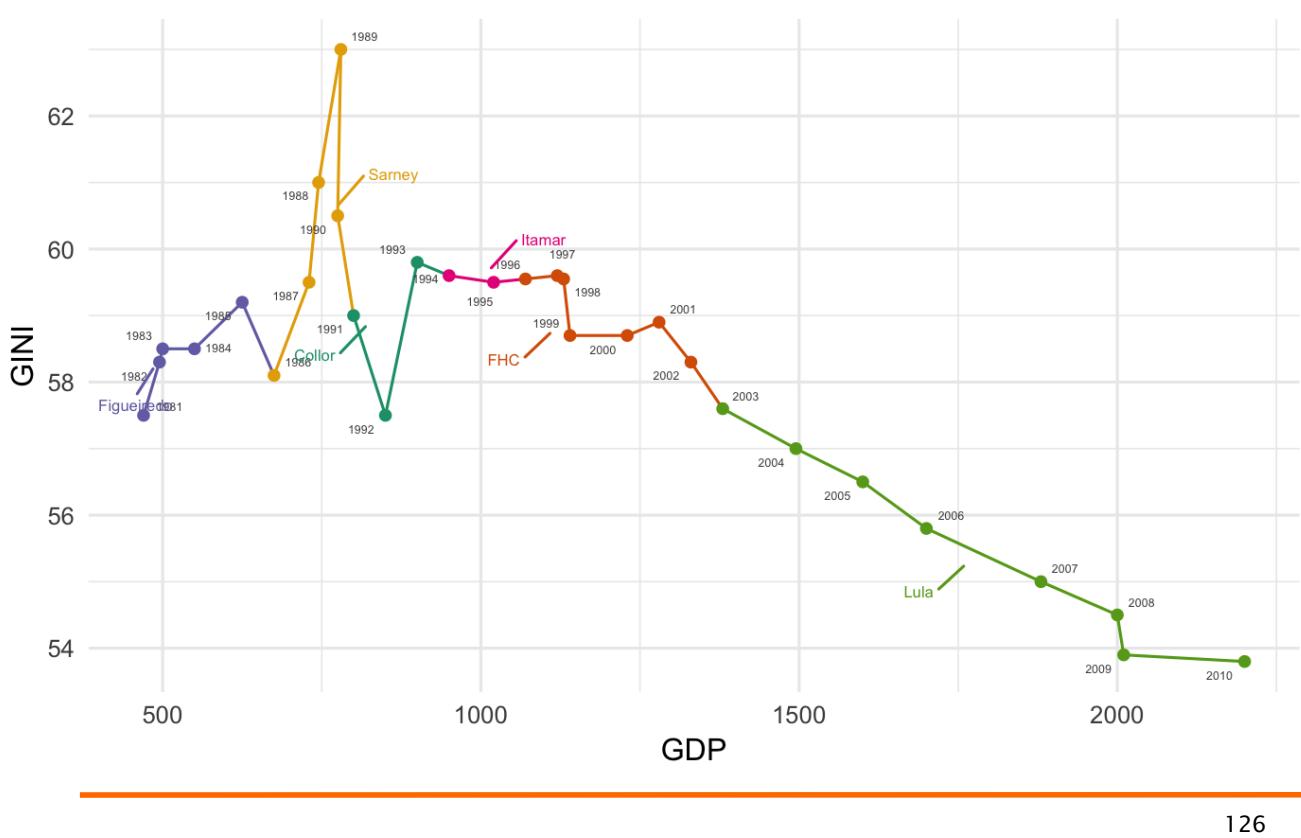
Double scale (alternative)



Multiple graphs



Path

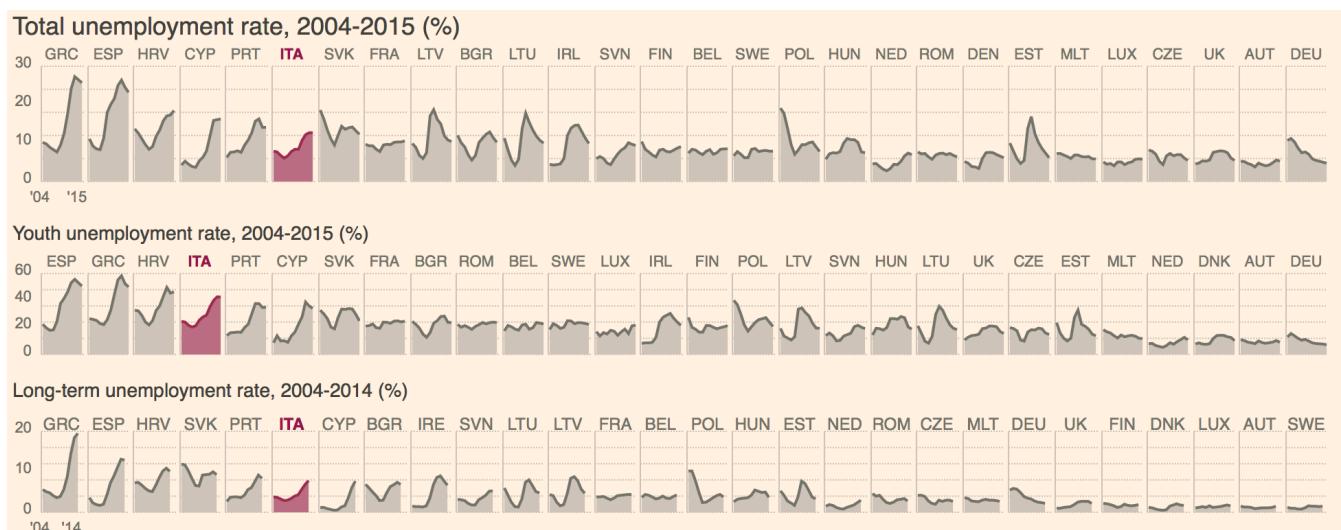


Small multiples

- A.k.a.
 - ◆ Trellis
 - ◆ Lattice
 - ◆ Grid
- Set of aligned graphs sharing (at least one) scale and axis
 - ◆ Enable ease of comparison among different measures

127

Small multiples



FT EU unemployment tracker
<http://blogs.ft.com/ftdata/2015/04/17/eu-unemployment-tracker/>

128

Small multiples

- Consistency
 - ◆ Same scale
 - ◆ Same categorical levels
 - ◆ Same ordering of categorical levels
- Arrangement
 - ◆ Align axis that involve comparison
 - Possibly along a matrix

129

Trellis

- Sequence
 - ◆ Intrinsic order
 - ◆ Order of relevance
 - ◆ Order by some quantitative attribute
- Rules and grids
 - ◆ Use when spacing is not enough
 - ◆ Can direct the reader to scan graphs horizontally or vertically

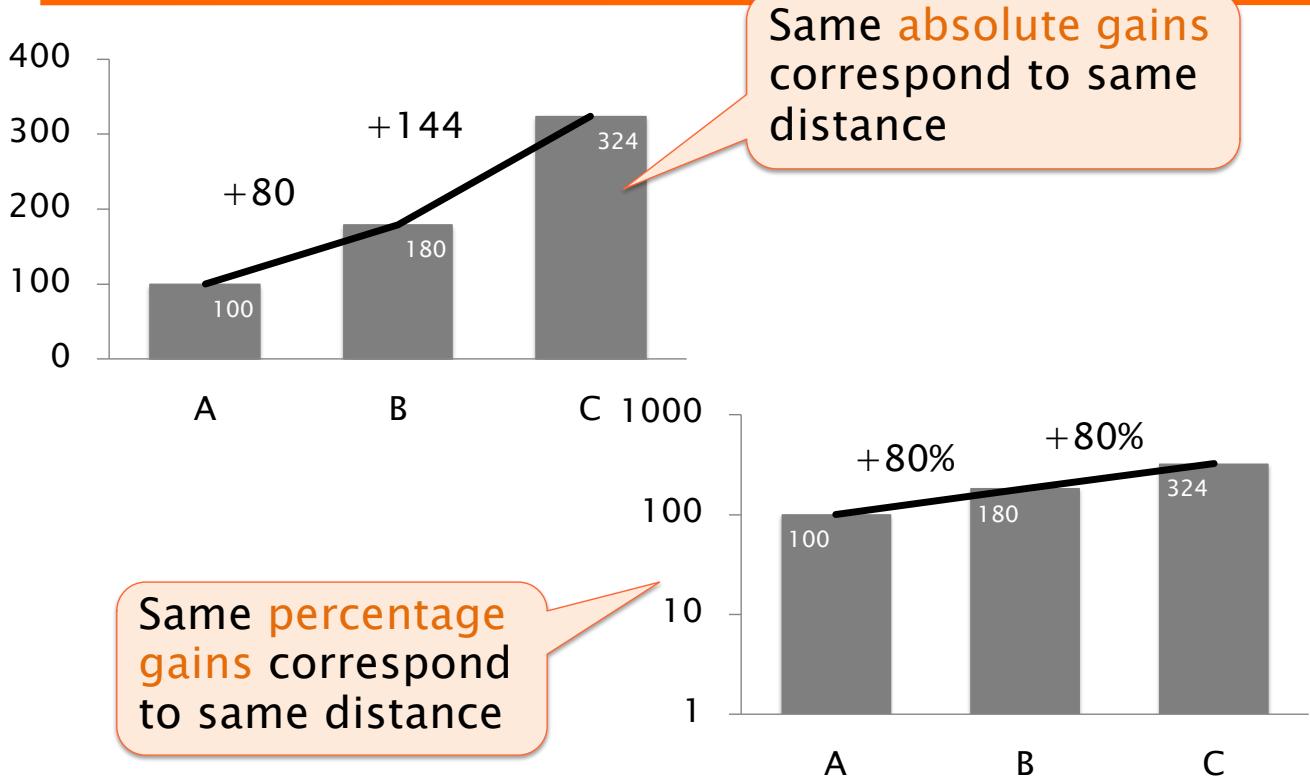
130

Log scale

- Reduce visual difference between quantitative data sets with significantly wide ranges
- Differences are proportional to percentages

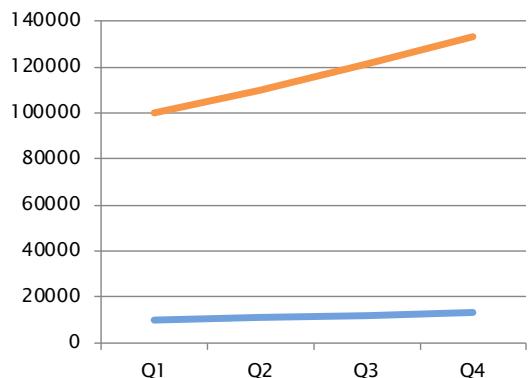
131

Log scale

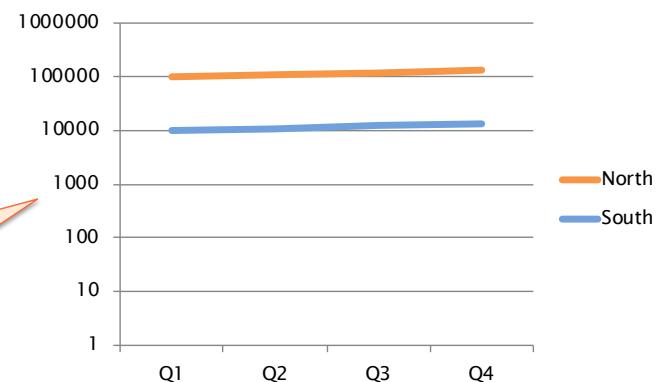


132

Log scale



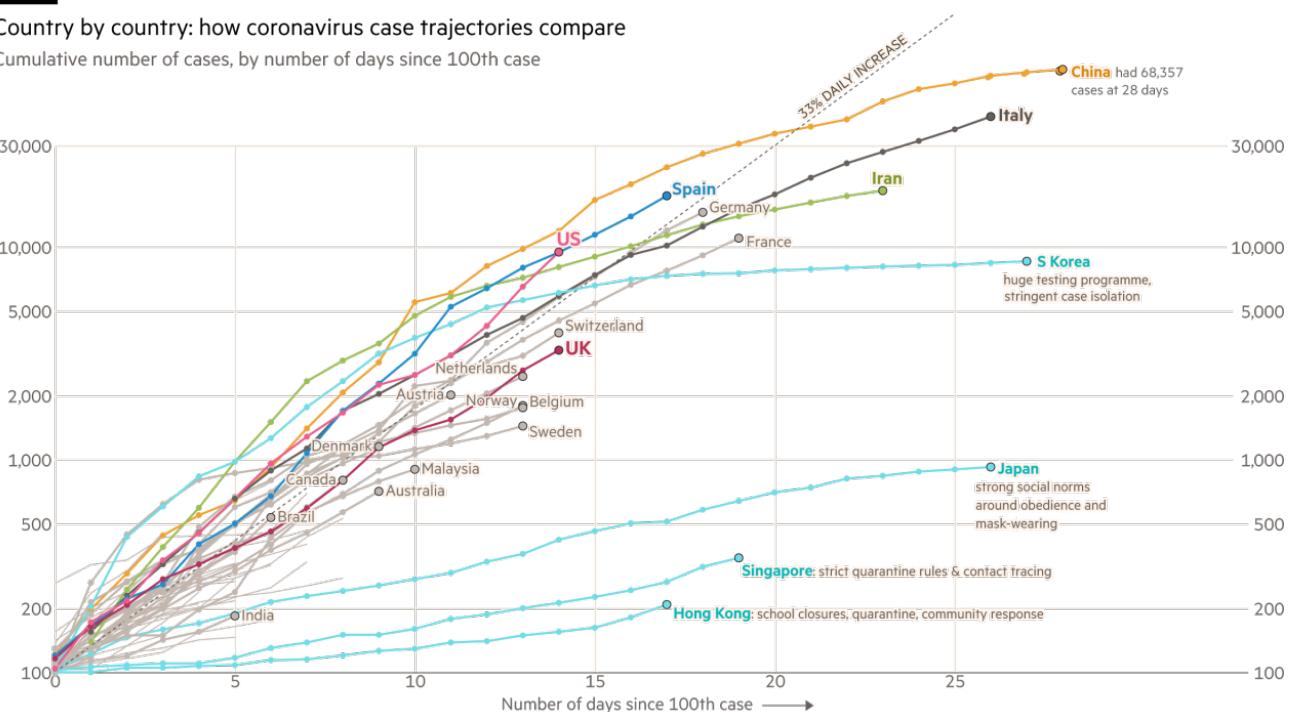
Parallel lines for same absolute gains



133

Log scale

Country by country: how coronavirus case trajectories compare
Cumulative number of cases, by number of days since 100th case



FT graphic: John Burn-Murdoch / @burnmurdoch
Source: FT analysis of Johns Hopkins University, CSSE; Worldometers. Data updated March 19, 19:00 GMT
© FT

Graph area

- Aspect ratio should not distort perception
 - ◆ Typically wider than taller
 - ◆ Scatter plots may be squared
- Grid lines must be thin and light
 - ◆ Useful to look-up values
 - ◆ Enhance comparison of values
 - ◆ Enhance perception of localized patterns

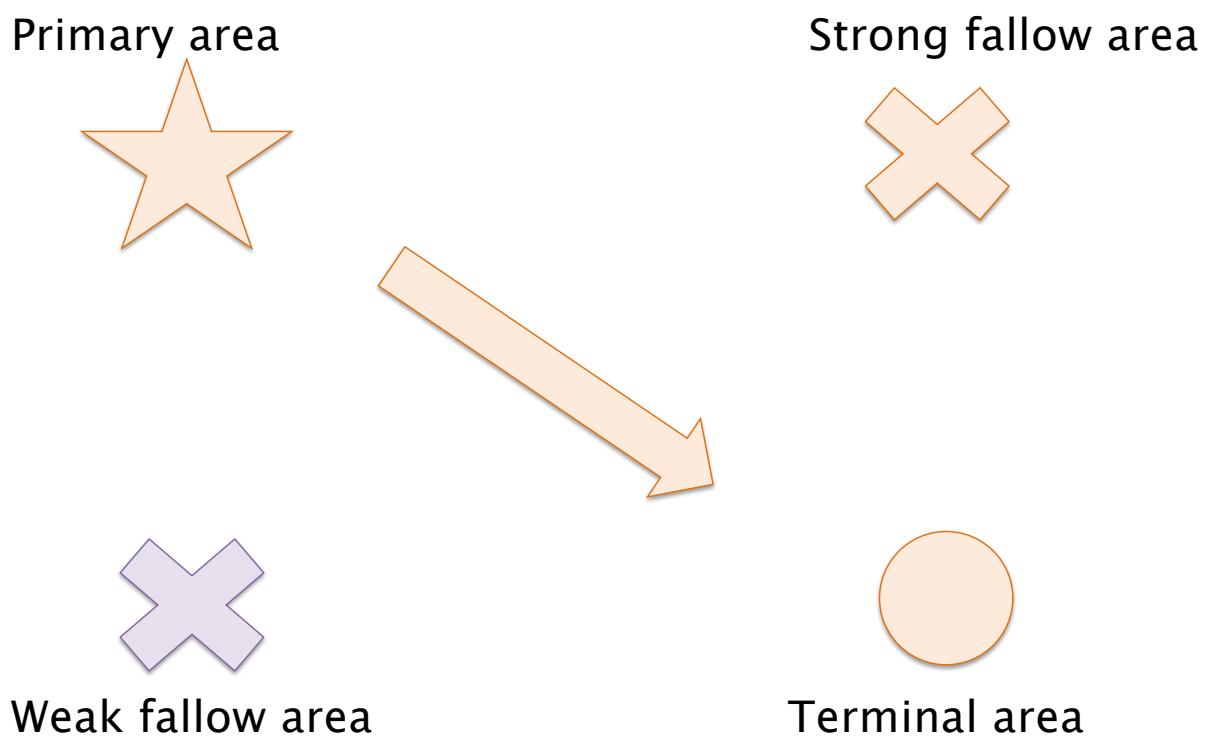
135

Labels

- Important elements (e.g. titles) should be prominent
 - ◆ Top
 - ◆ Larger

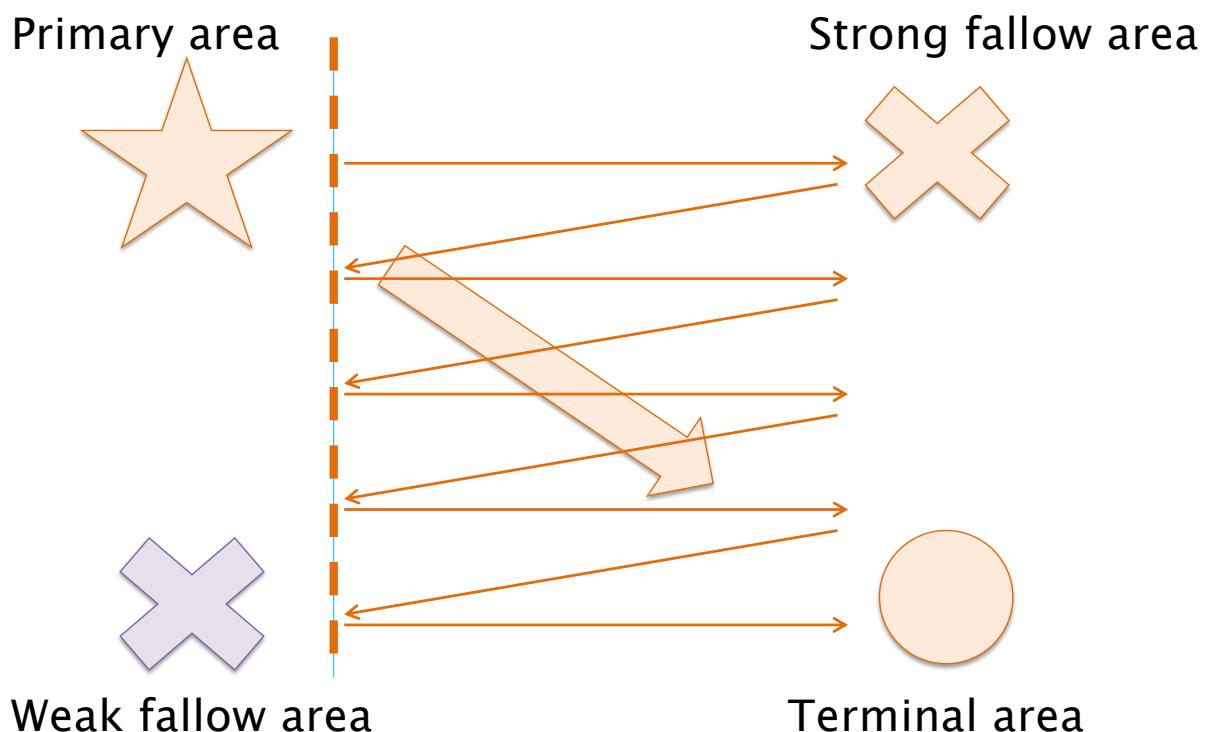
136

Guthenberg Diagram



137

Guthenberg Diagram



138

Legends

- Used for categorical attributes not associated to any axis
- As close as possible to the objects
- Less prominent than data objects
- Borders are used only when necessary to separate from other elements

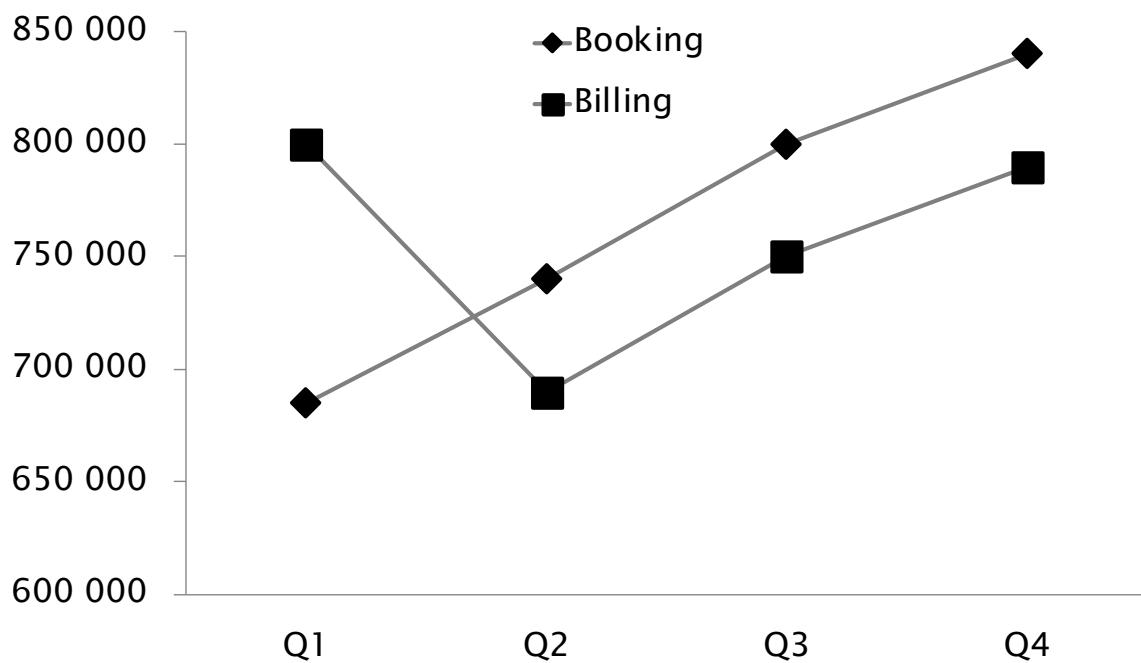
139

Legends

- Text should be as close as possible to the object it complements
 - ◆ Prefer direct labeling to separate legends
- Number of categorical subdivisions
 - ◆ Perceptual limit is between 5 and 8
 - ◆ Limit is independent of the visual attribute used to encode it
 - ◆ Joint use of attributes ease discrimination

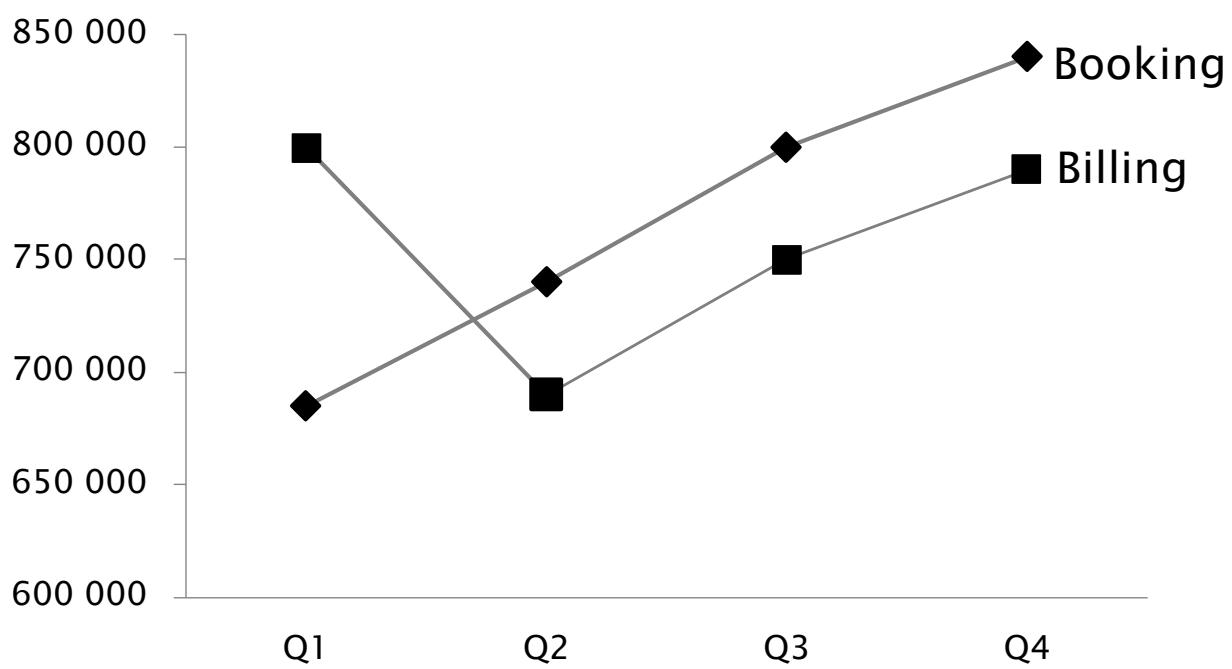
140

Legend



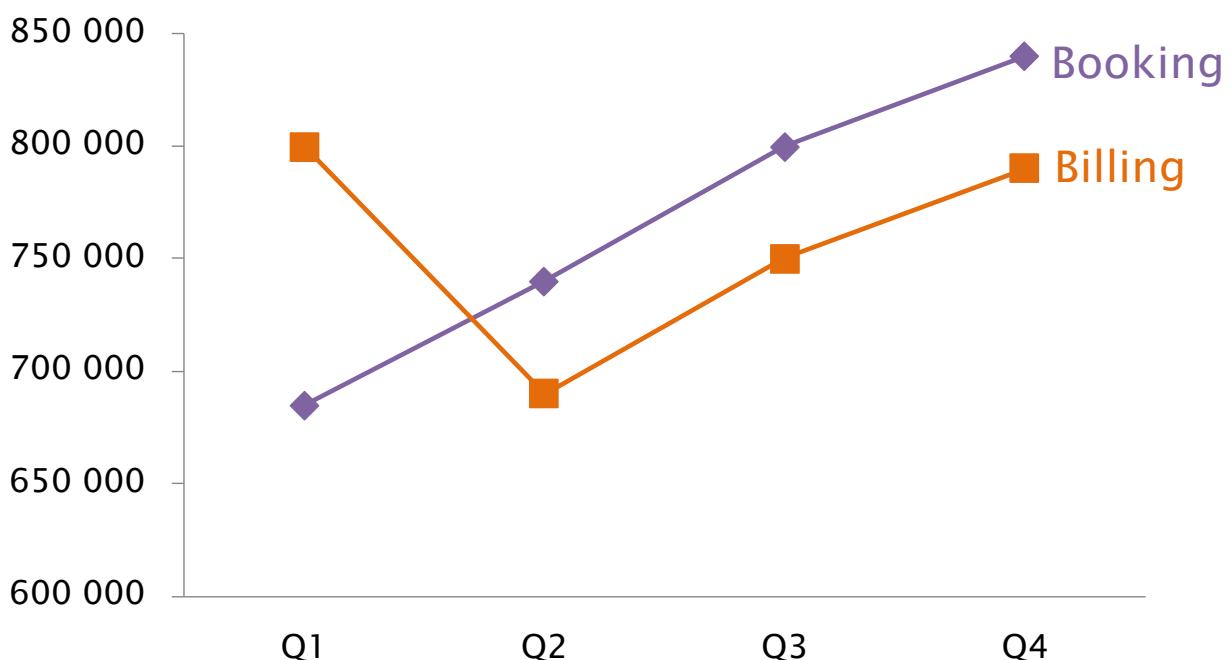
141

Direct labeling



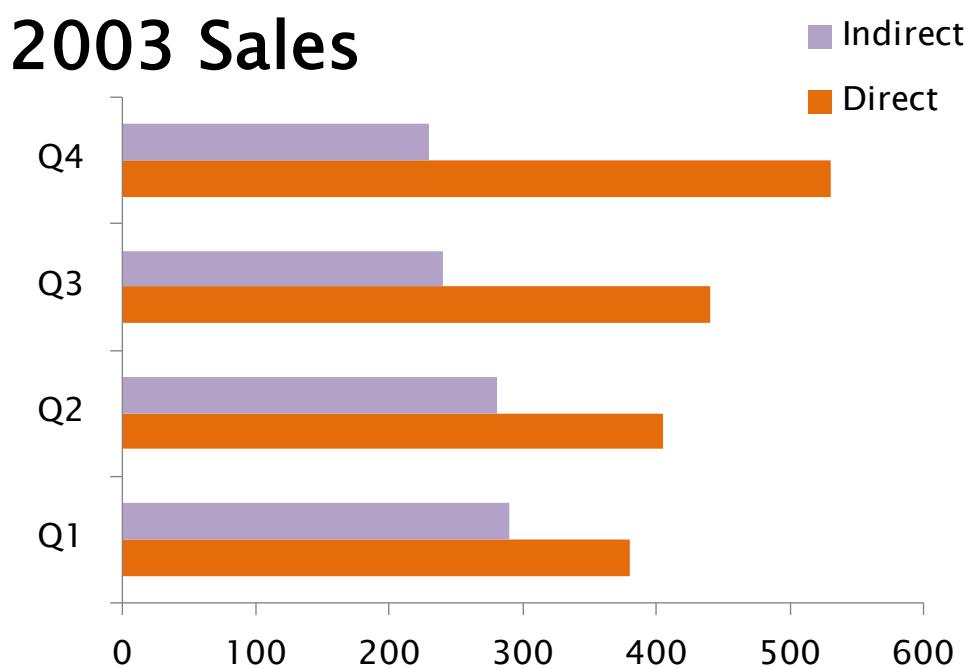
142

Direct labeling and color



143

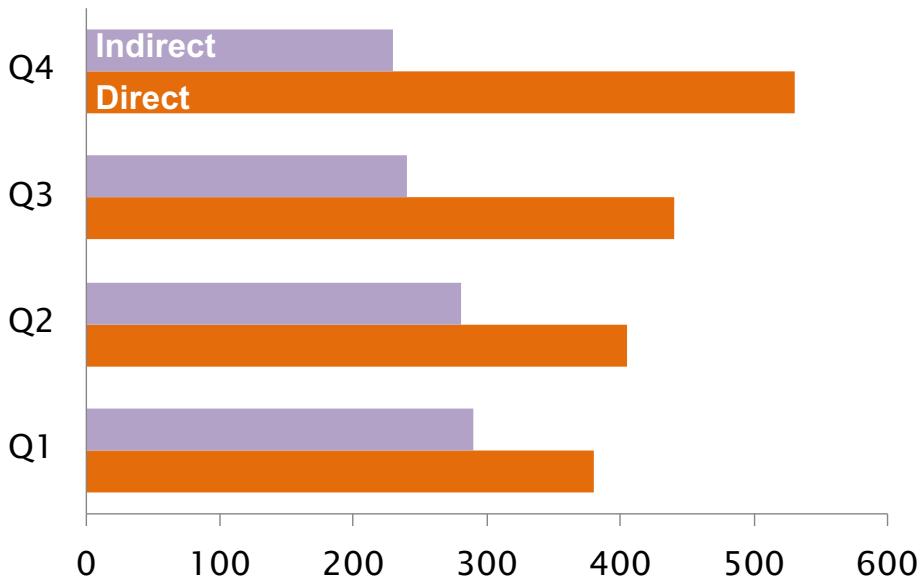
Legend



144

Direct labeling

2003 Sales



145

Reference lines and regions

- Reference lines support an easy comparison to a given value
 - ◆ Mean
 - ◆ Threshold
- Reference regions allow comparison with several values
 - ◆ Use background color

146

DASHBOARD

147

Dashboard

Visualization of the most relevant information needed to achieve one or more goals which fits entirely on a single screen so it can be monitored at a glance

148

Dashboard

- Dashboards display mechanisms are
 - ◆ small
 - ◆ concise
 - ◆ clear
 - ◆ intuitive
- Dashboards are customized
 - ◆ To suit the goals of person, group, function

149

Provide context for data

PUC

-
- References allow judging the data

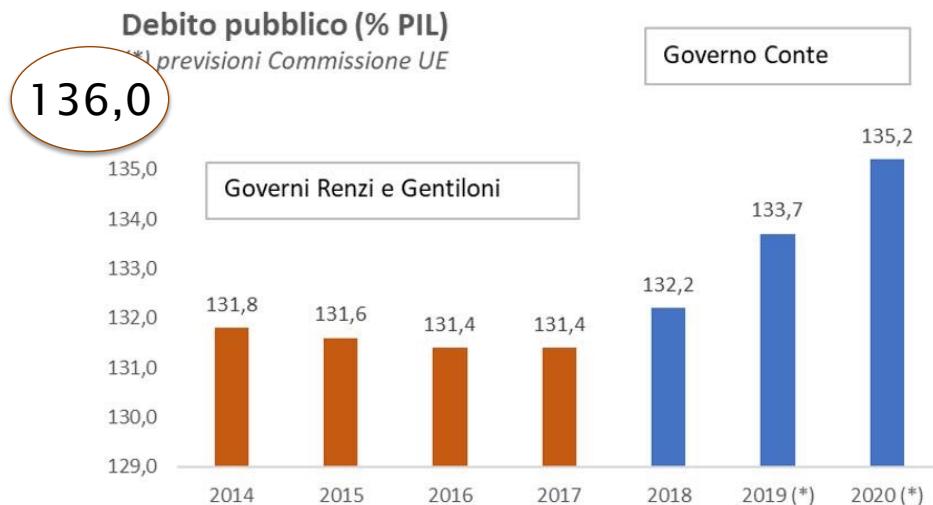


150

Use appropriate detail

PUC

- Typical counterexamples
 - ◆ Dates with seconds detail
 - ◆ Decimals



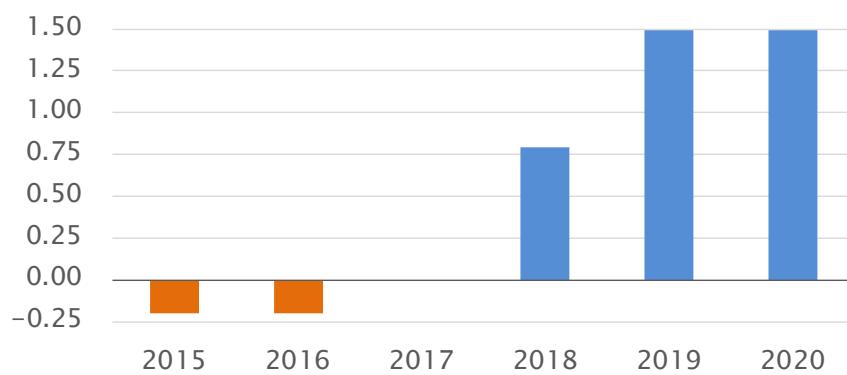
151

Use the right measures

- If you are interested in e.g. the difference, ratio, variation show such derived measure

Variazione Debito Pubblico (% PIL)

PD M5S-L



152

Use appropriate visualization

- Typical errors:
 - ◆ Any chart when a table would be better
 - ◆ Pie-charts not representing part-whole
 - ◆ Bubble charts

153

Visualization instruments

- Tables
 - ◆ Textual information
- Graphs
 - ◆ Visual information

154

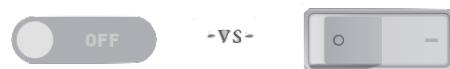
Avoid decorations

- Skeumorphic design
- Backgrounds motives
- Color gradients
- Variations not encoding any measure
 - ◆ Typically color

155

Avoid decorations

- Skeumorphic design
- Backgrounds motives
- Color gradients
- Variations not encoding any measure
 - ◆ Typically color



156

Avoid decorations

- Skeumorphic design
- Backgrounds motives
- **Color gradients**
- Variations not encoding any measure
 - ◆ Typically color

A

B

157

Avoid decorations

- Skeumorphic design
- Backgrounds motives
- Color gradients
- Variations not encoding any measure
 - ◆ Typically color

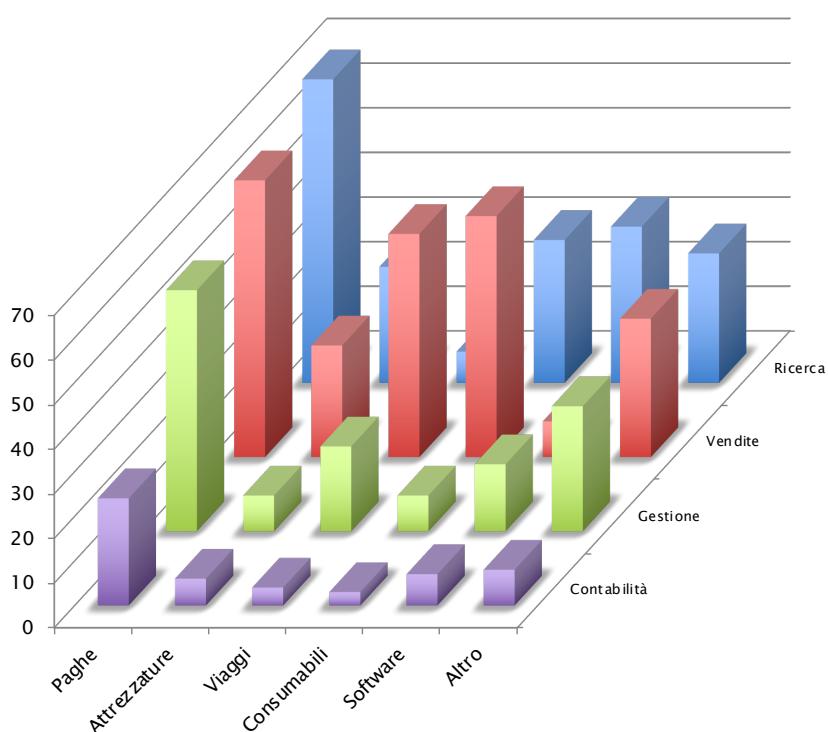
158

3D diagrams

- Encoding
 - ◆ Axonometry typically hides some data and makes comparison hard
- Not encoding
 - ◆ Perspective deform dimensions
 - ◆ Depth or height distract and make comparison more difficult

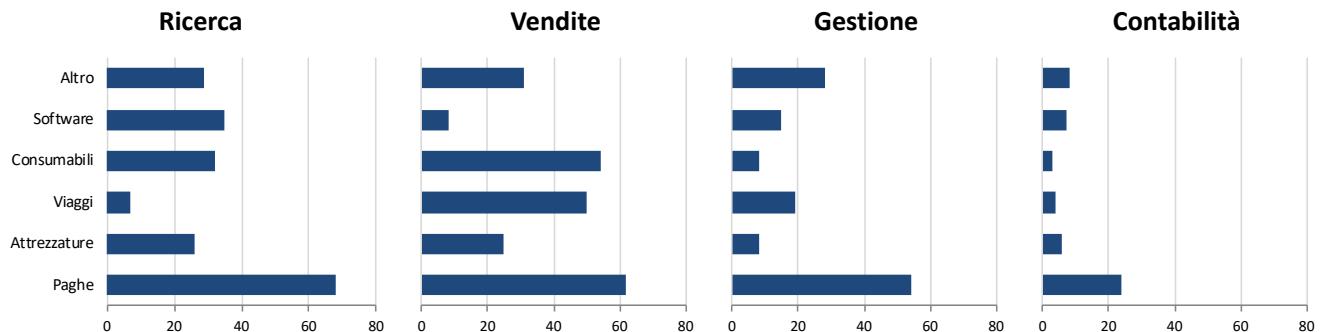
159

Encoding 3D



160

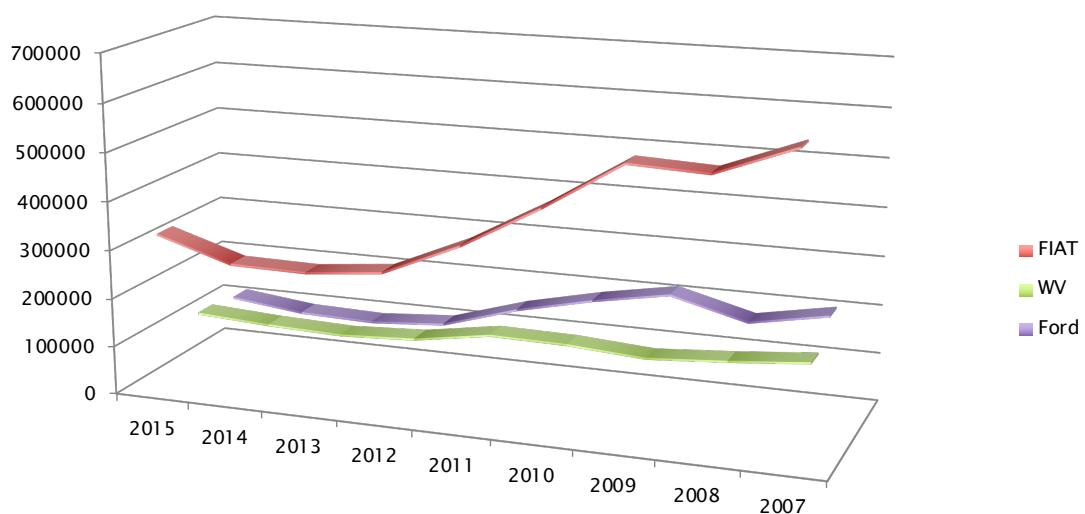
Encoding 3D → 2D



161

Decorative 3D

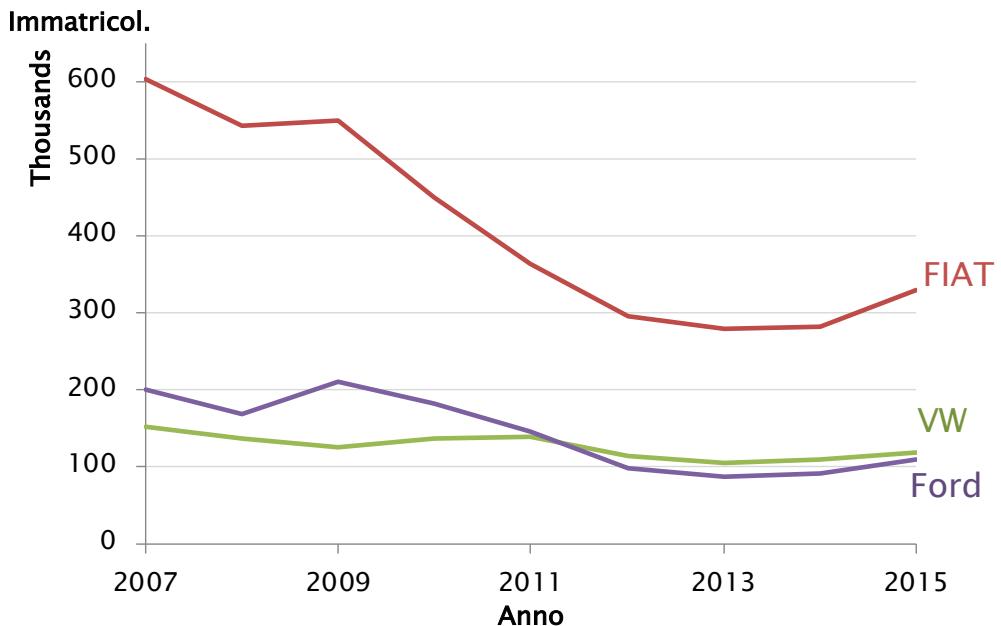
Immatricol.



162

Decorative 3D → 2D

Immatricolazioni auto per marchio
sul mercato italiano



163

References

- Stephen Few, 2004. Show me the numbers. Analytics Press.
 - ◆ <http://www.perceptualedge.com/blog/>
- Edward R. Tufte, 1983. The Visual Display of Quantitative Information. Graphics Press.

164

References

- Wilkinson, L. (2006). *The grammar of graphics*. Springer Science & Business Media.
- Wickham, H. (2010). A layered grammar of graphics. *Journal of Computational and Graphical Statistics*, 19(1), 3–28.
- Visual Vocabulary
<http://ft.com/vocabulary>

165

References

- R.Olson. Revisiting the vaccine visualization
 - ◆ <http://www.randalolson.com/2016/03/04/revisiting-the-vaccine-visualizations/>
- Nathan Yau. 9 Ways to Visualize Proportions – A Guide
 - ◆ <http://flowingdata.com/2009/11/25/9-ways-to-visualize-proportions-a-guide/>
- M.Correr, and M.Gleicher. Error Bars Considered Harmful: Exploring Alternate Encodings for Mean and Error *IEEE Transactions on Visualization and Computer Graphics*, Dec. 2014
 - ◆ <http://graphics.cs.wisc.edu/Papers/2014/CG14/Preprint.pdf>

166