# Graph Construction

## Visualizzazione dell'Informazione Quantitativa

# 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*

# 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*)
- AESTHETIC: attributes (*e.g., color, position*)
- GUIDE: guides (e.g., *axes*, *legends*)

# Specification for a scatter plot

- DATA: x, y, group
- TRANS: identity
- SCALE: *linear*(*dim*(1)), *linear*(*dim*(2))
- COORD:*rect*(*dim*(1, 2))
- ELEMENT: *point*()
- AESTHETIC: *position*(x*y)
- GUIDE: *axis*(*dim*(1)), *axis*(*dim*(2))

# Graph visual components

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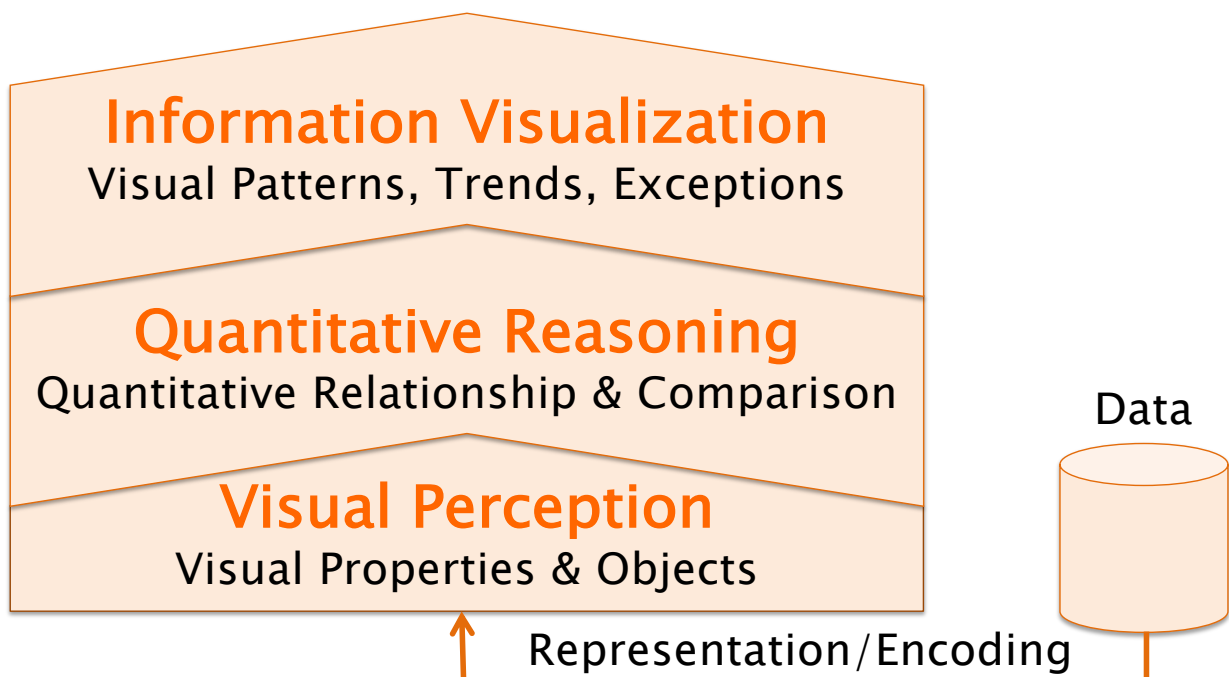  - ◆ Axes
  - ◆ Labels
  - ◆ Legends

Visualizzazione dell'Informazione Quantitativa

# VISUAL RELATIONSHIPS

# Data Visualization

## Understanding

# Visual Encoding

- Given a variable (measure), identify:
  - Visual object
  - Visual attribute

- Main distinction
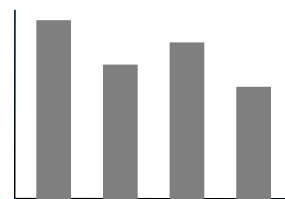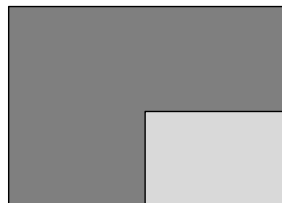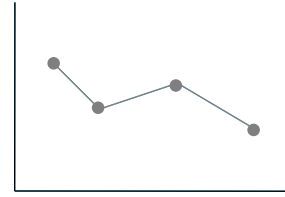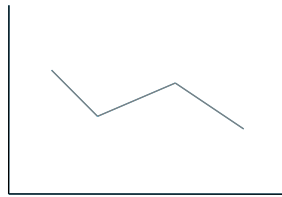  - Quantitative (interval, ratio, absolute)
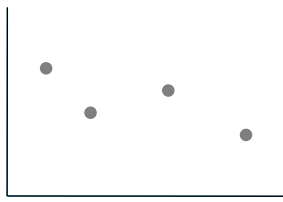  - Categorical (nominal, ordinal)

# Relationships

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

# Quantitative encoding

# Sample data

| Region | turnout (2018) | turnout (2013) |
|---|---|---|
| ABRUZZO | 75.3% | 75.9% |
| BASILICATA | 71.1% | 69.5% |
| CALABRIA | 63.6% | 63.1% |
| CAMPANIA | 68.2% | 67.9% |
| EMILIA–ROMAGNA | 78.3% | 82.1% |
| FRIULI–VENEZIA GIULIA | 75.1% | 77.2% |
| LAZIO | 72.6% | 77.5% |
| LIGURIA | 72.0% | 75.1% |
| LOMBARDIA | 76.8% | 79.6% |
| MARCHE | 77.3% | 79.8% |

| Region | turnout (2018) | turnout (2013) |
|---|---|---|
| MOLISE | 71.6% | 78.1% |
| PIEMONTE | 75.2% | 77.3% |
| PUGLIA | 69.1% | 69.9% |
| SARDEGNA | 65.5% | 68.5% |
| SICILIA | 62.8% | 64.6% |
| TOSCANA | 77.5% | 79.2% |
| TRENTINO–ALTO ADIGE | 74.3% | 81.0% |
| UMBRIA | 78.2% | 79.5% |
| VALLE D'AOSTA | 72.3% | 77.0% |
| VENETO | 78.7% | 81.8% |

# Nominal comparison

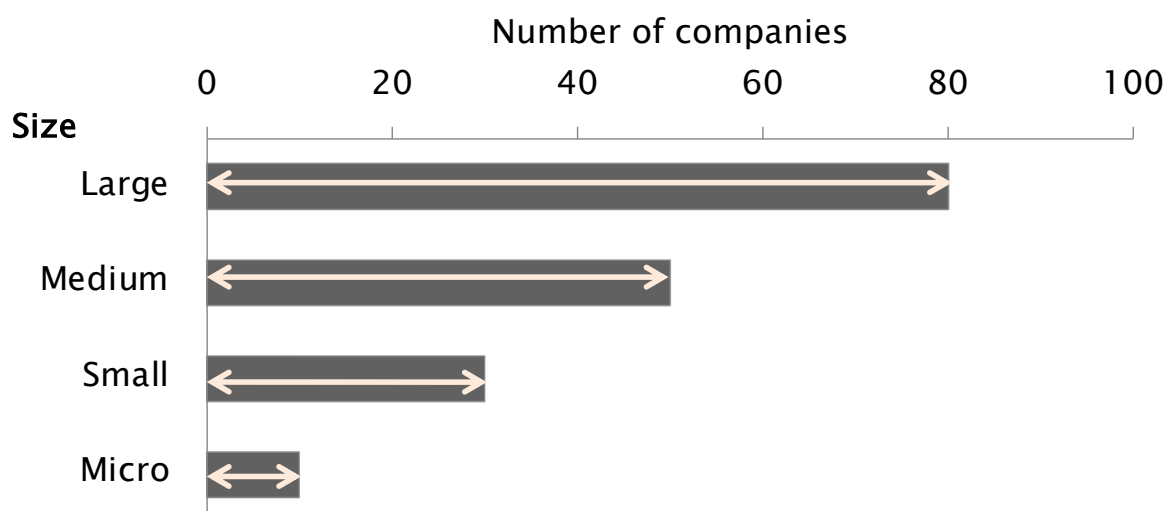- Compare quantitative values corresponding to categorical levels
  - Small differences are difficult to see
    - Non zero-based scale can emphasize
  - Dot plots can be used for small differences
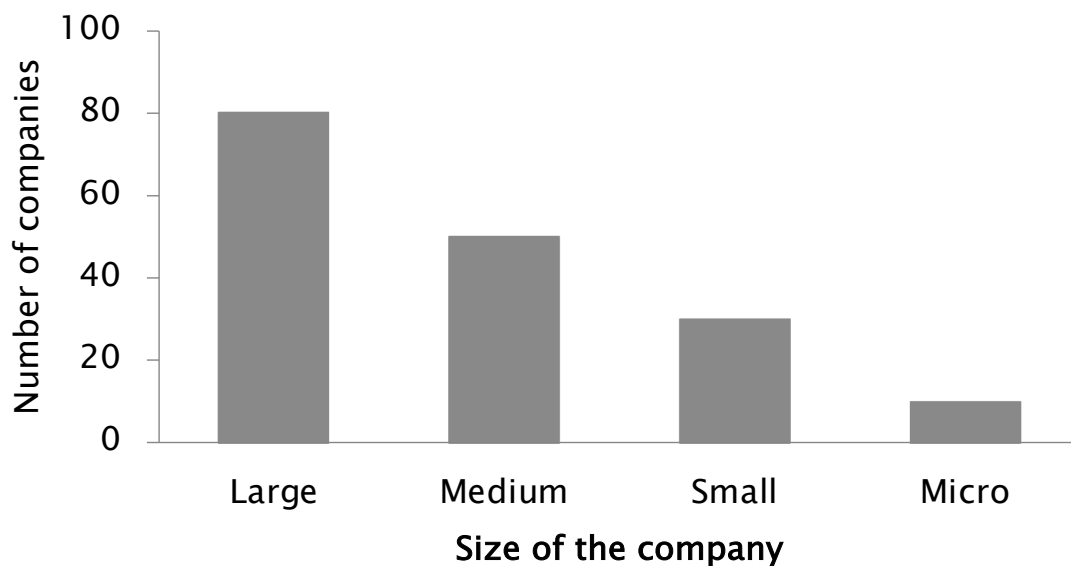    - They do not require zero based scale

# Line length – Bars chart

# Vertical Bars (aka Columns)



Bar chart with y-axis "Number of companies" (0 to 100) and x-axis "Size of the company" showing: Large 80, Medium 50, Small 30, Micro 10.
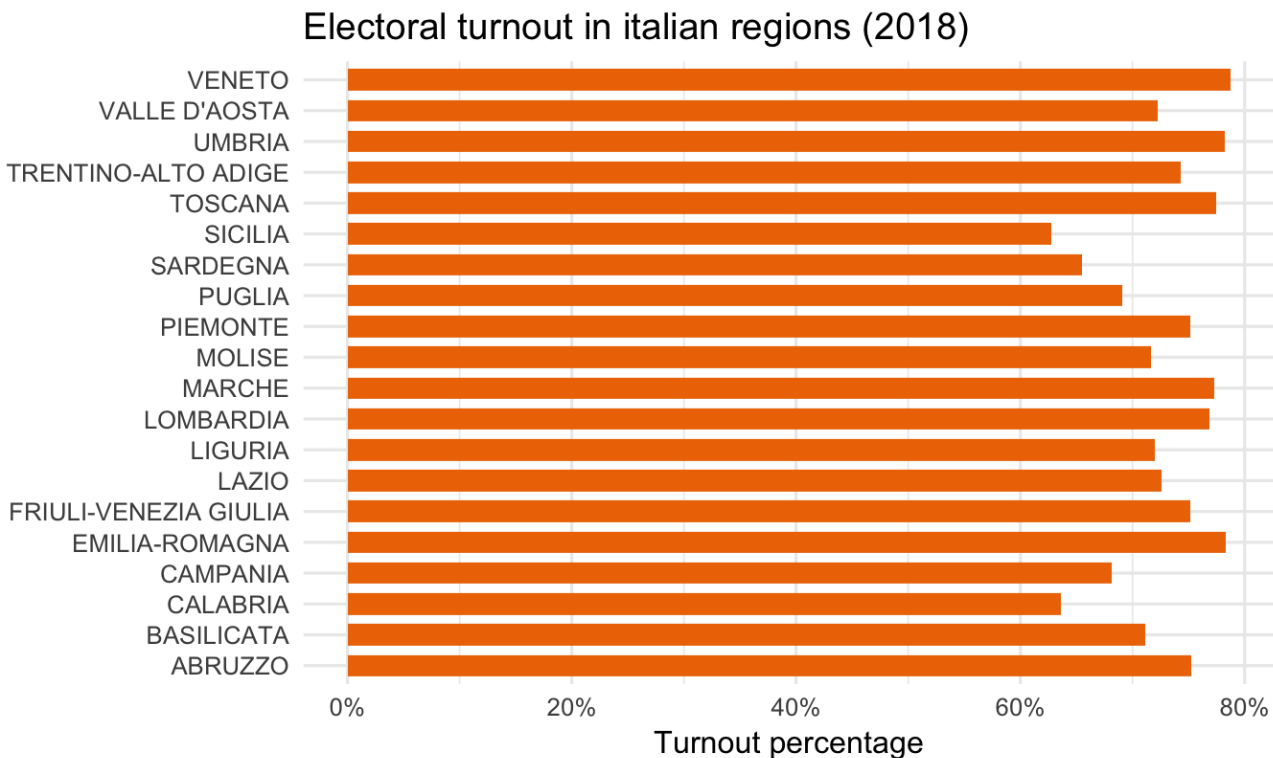
# Bar charts

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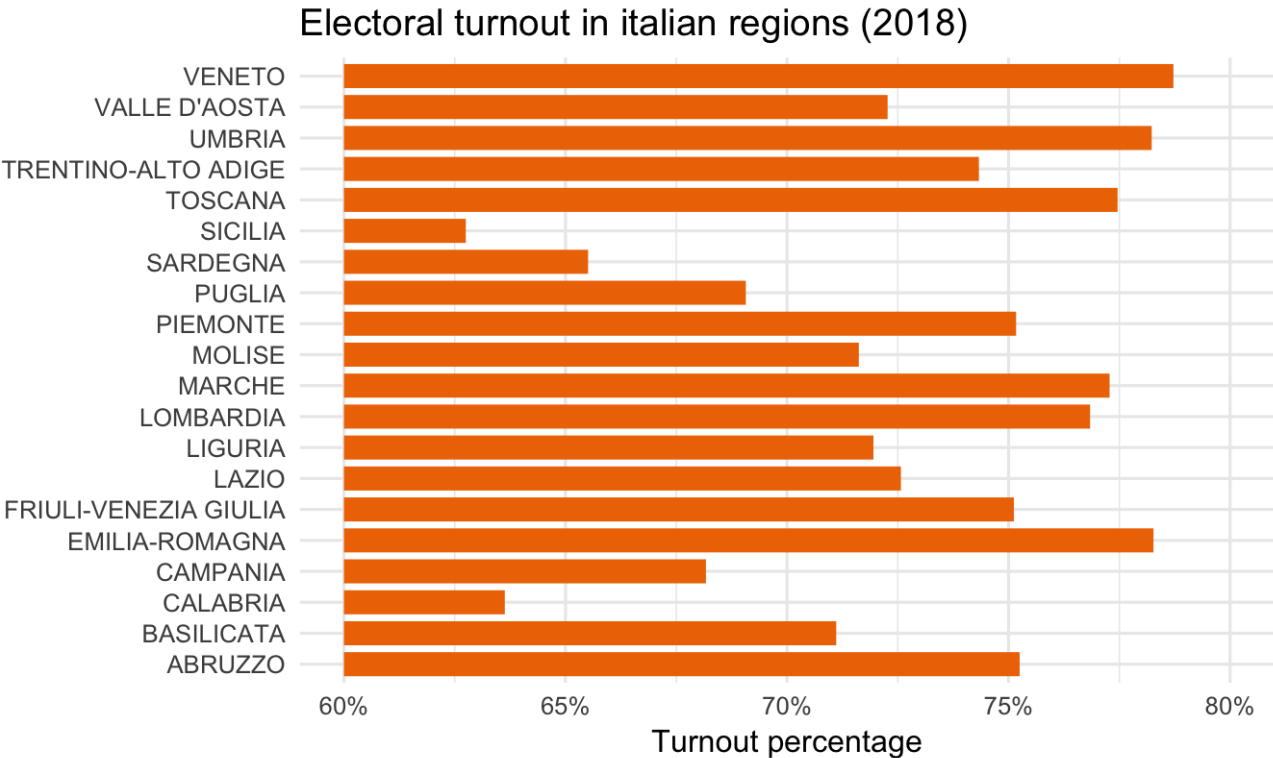- Quantitative values are encoded only as length of the bars
  - The axis is a supporting element
- Width of bars plays no role
  - Bars are just very thick lines
- Bars require a zero-based scale
  - See: Lie factor!

# Comparison – Barplot



Electoral turnout in italian regions (2018)

# Barplot (non zero based scale)



Electoral turnout in italian regions (2018)

# Barplot (non zero based scale)



Electoral turnout in italian regions (2018)

Proportionality:

$$LF = \frac{18.72/2.75}{78.72/62.75} = 5.43$$

# Barplot vertical labels



Electoral turnout in italian regions (2018)

# Bars Guidelines

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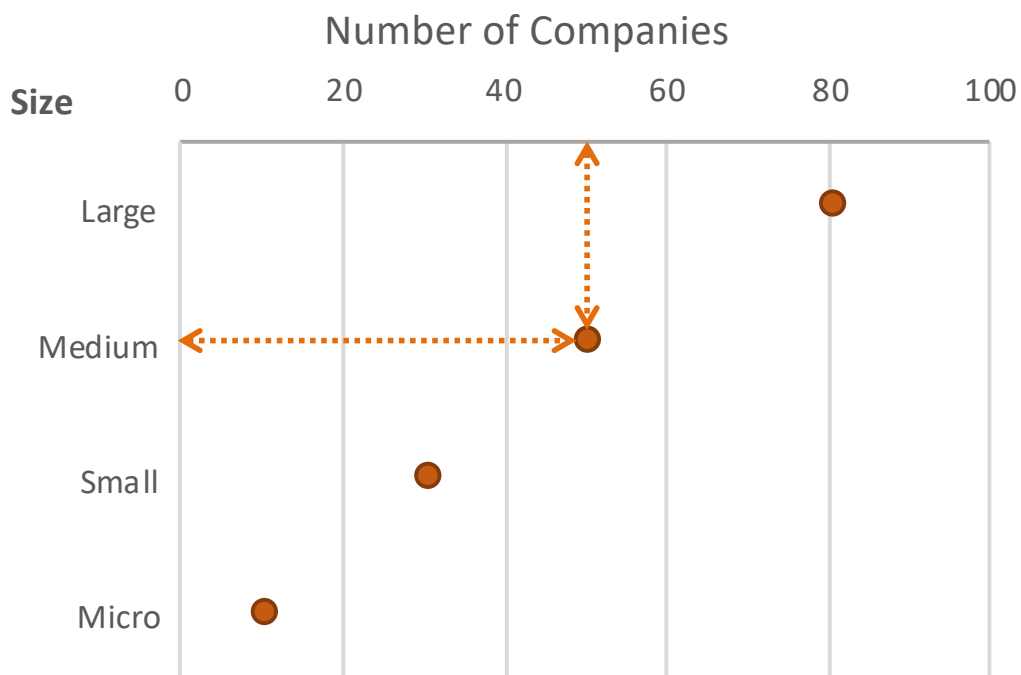  - No spacing between bars that are not labeled on the axis (legend categories)
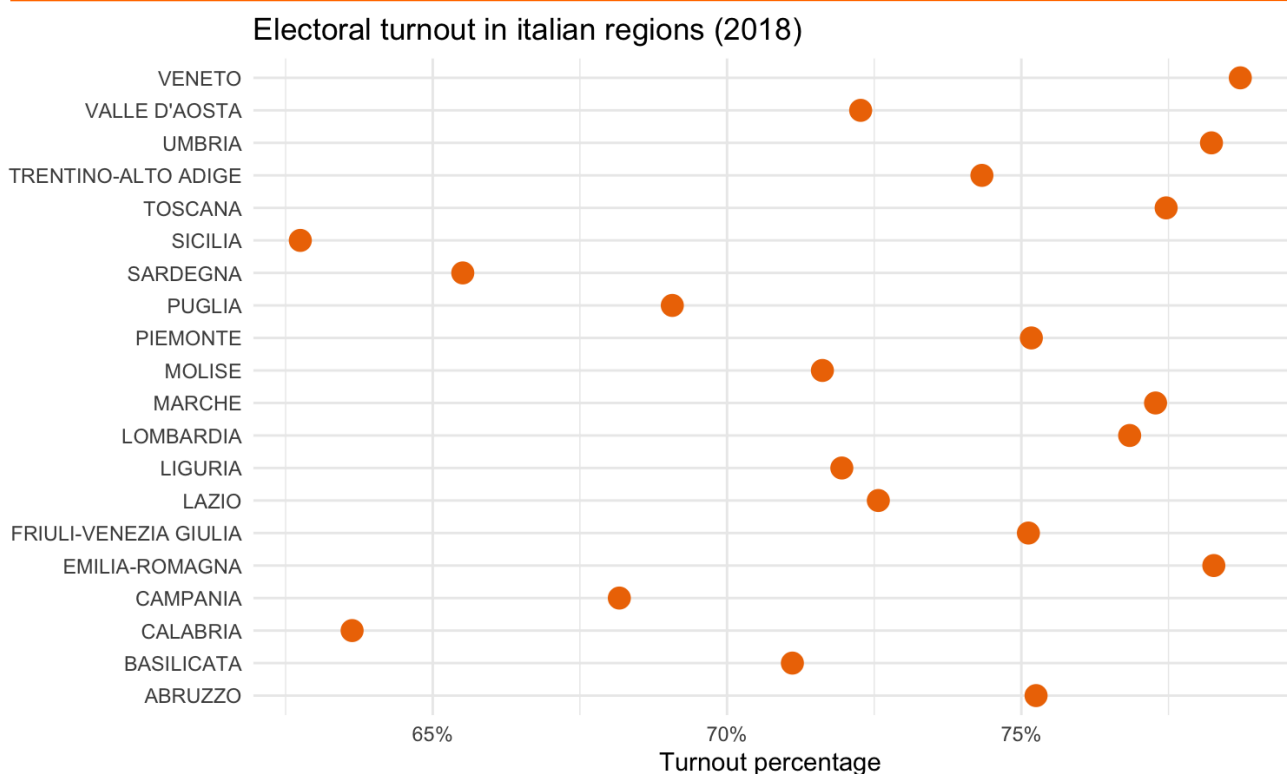  - No overlapping bars

# Position – Dots plot

# 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
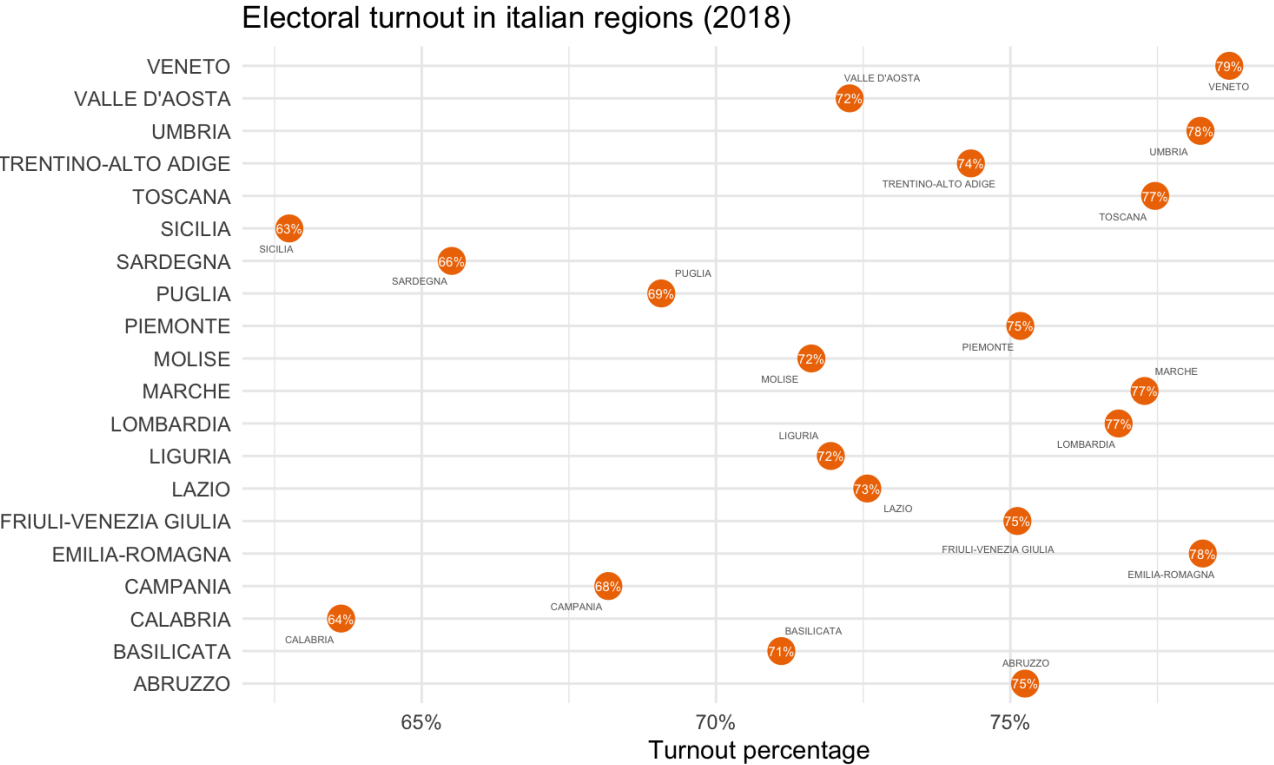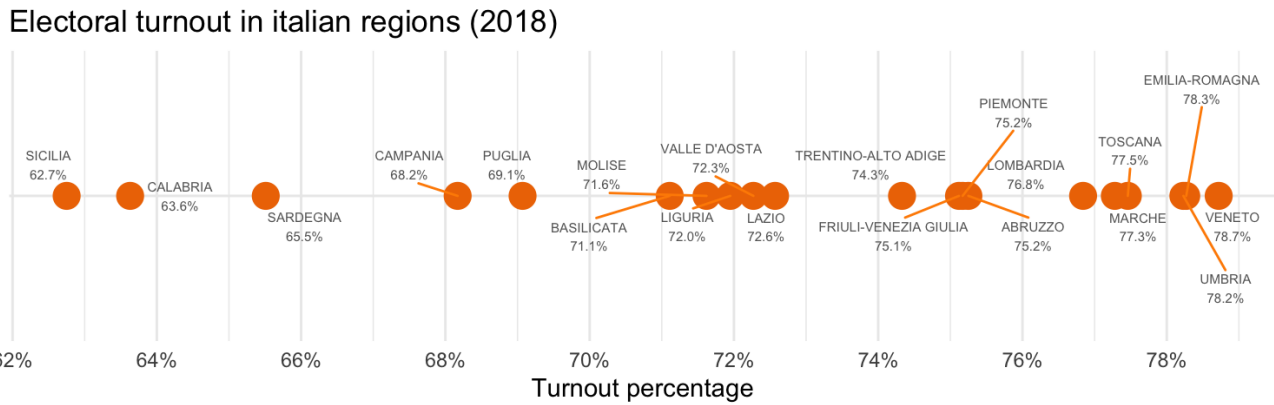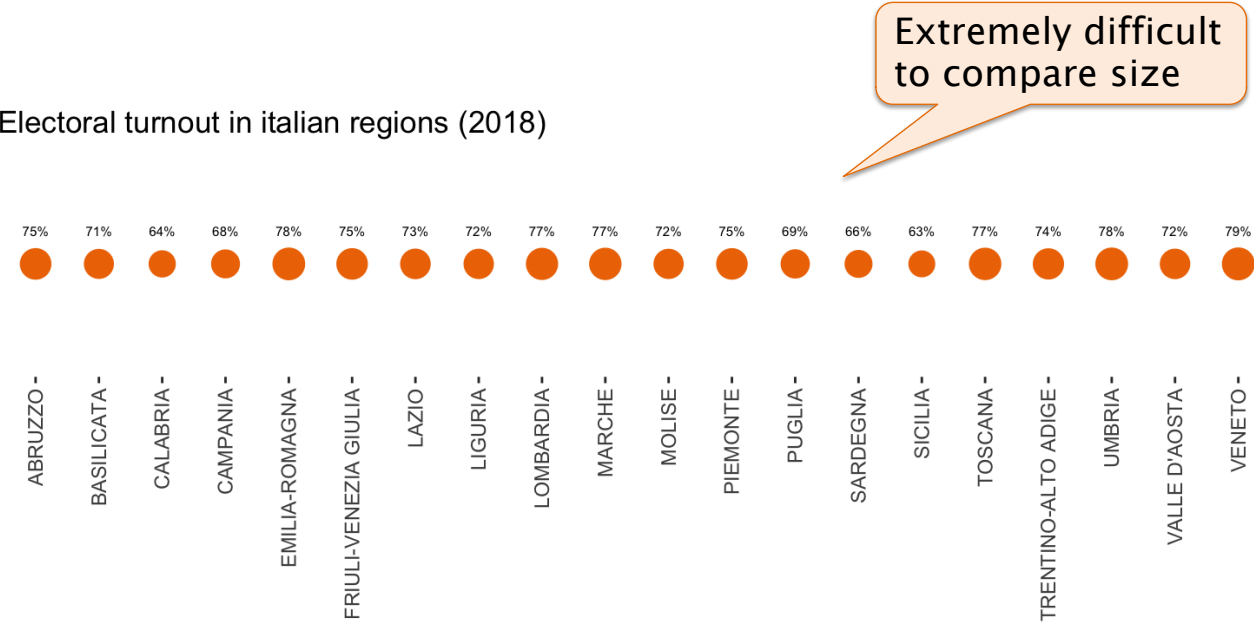
# Comparison – Dot plot



Electoral turnout in italian regions (2018)

# Comparison – Dot plot

### Electoral turnout in italian regions (2018)

# Comparison – Strip plot

### Electoral turnout in italian regions (2018)

# Comparison – Area – Bubbles

Electoral turnout in italian regions (2018)

Extremely difficult to compare size

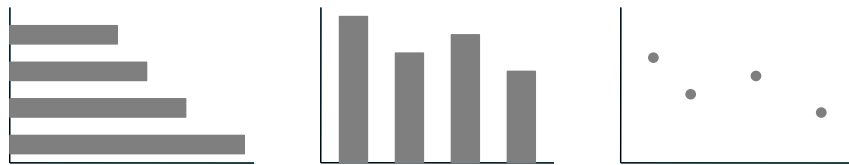| 75% | 71% | 64% | 68% | 78% | 75% | 73% | 72% | 77% | 77% | 72% | 75% | 69% | 66% | 63% | 77% | 74% | 78% | 72% | 79% |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ABRUZZO | BASILICATA | CALABRIA | CAMPANIA | EMILIA-ROMAGNA | FRIULI-VENEZIA GIULIA | LAZIO | LIGURIA | LOMBARDIA | MARCHE | MOLISE | PIEMONTE | PUGLIA | SARDEGNA | SICILIA | TOSCANA | TRENTINO-ALTO ADIGE | UMBRIA | VALLE D'AOSTA | VENETO |

# Count – Isotype

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

Literacy in England and Wales

Among 10 men

Illiterates — Literates

1841
1871
1901
1931

Among 10 women

1841
1871
1901
1931

ISOTYPE INSTITUTE

# Ranking

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

# Ranking

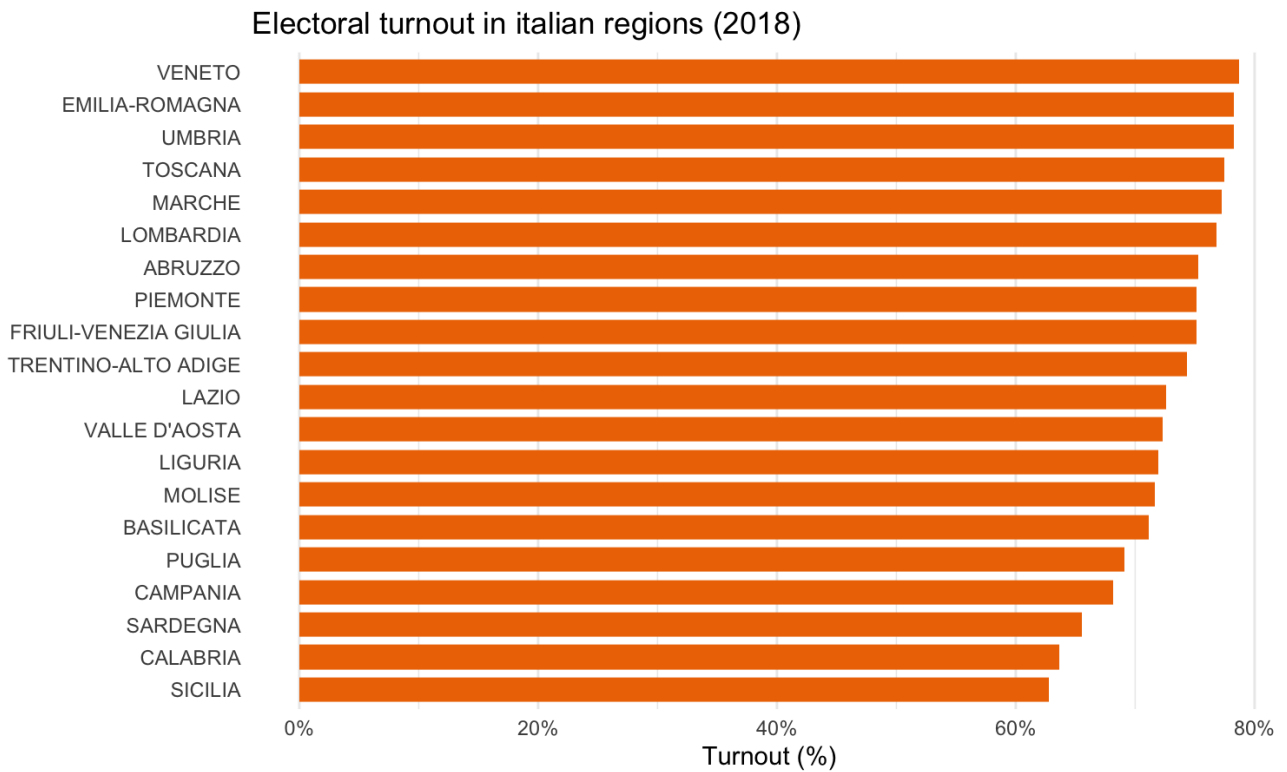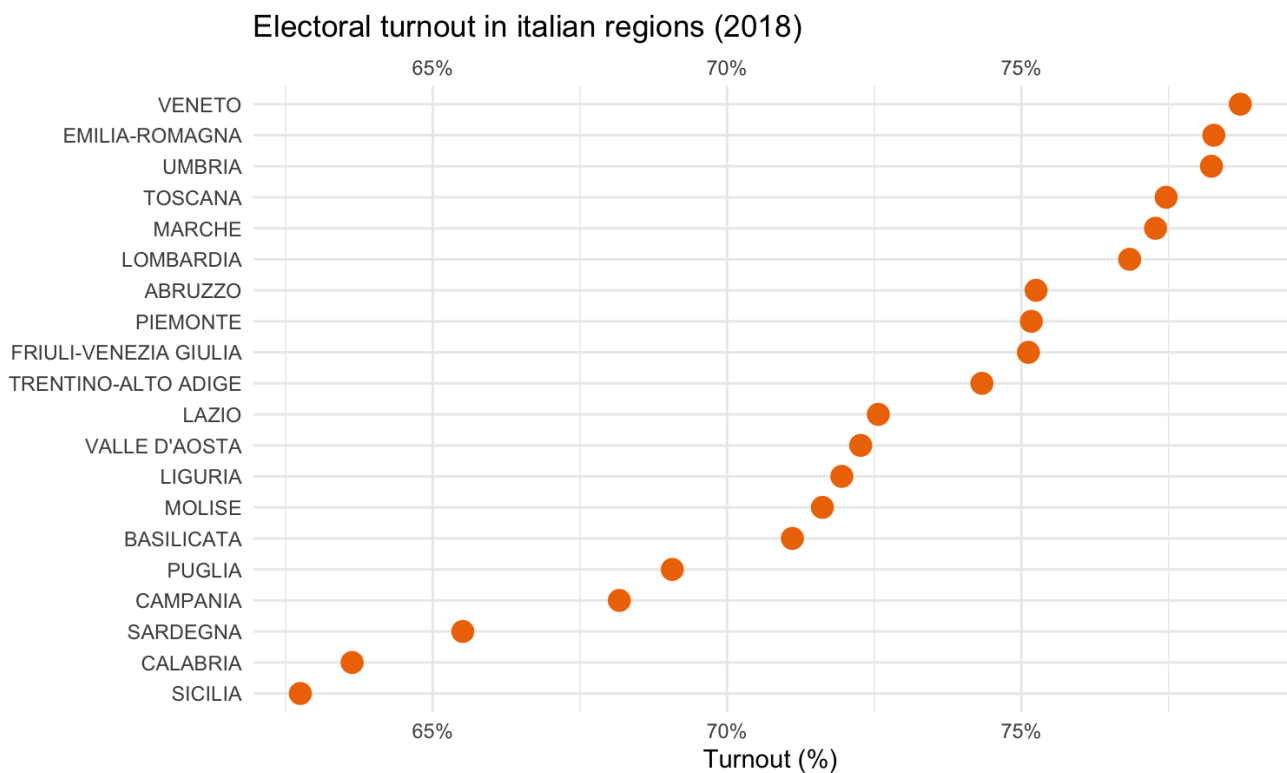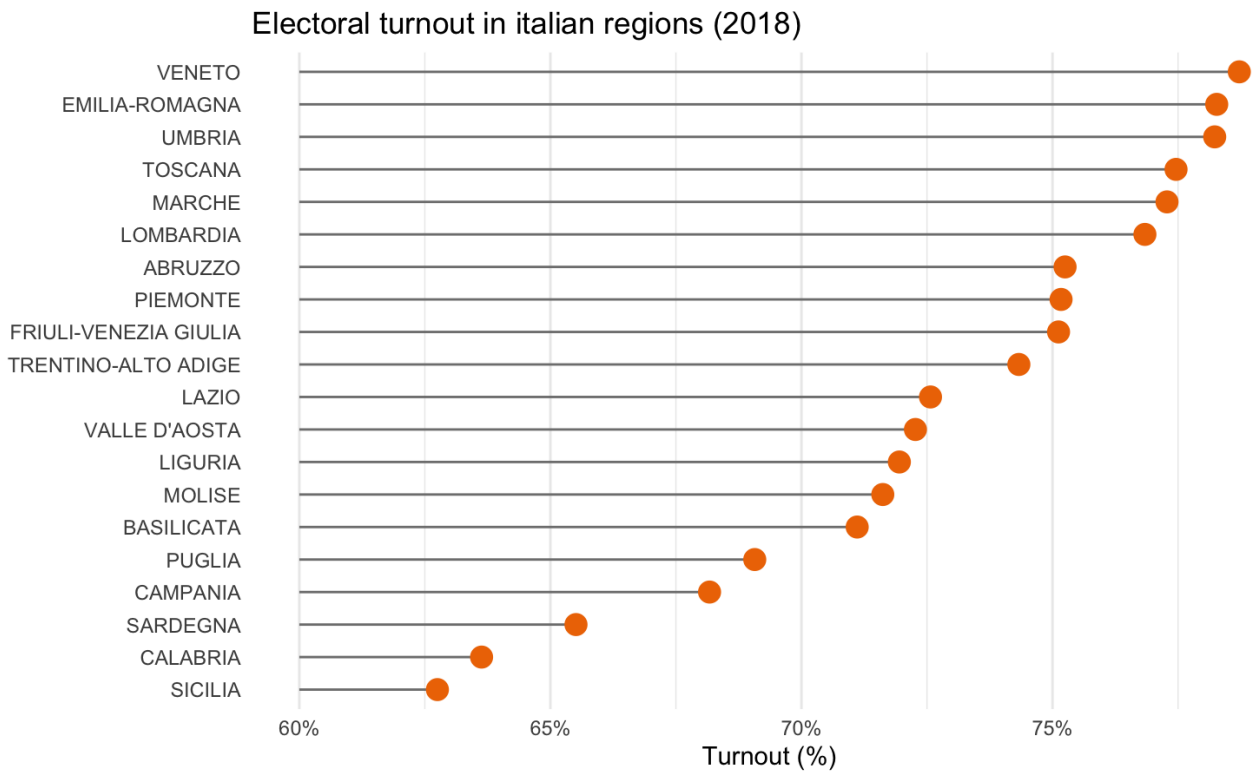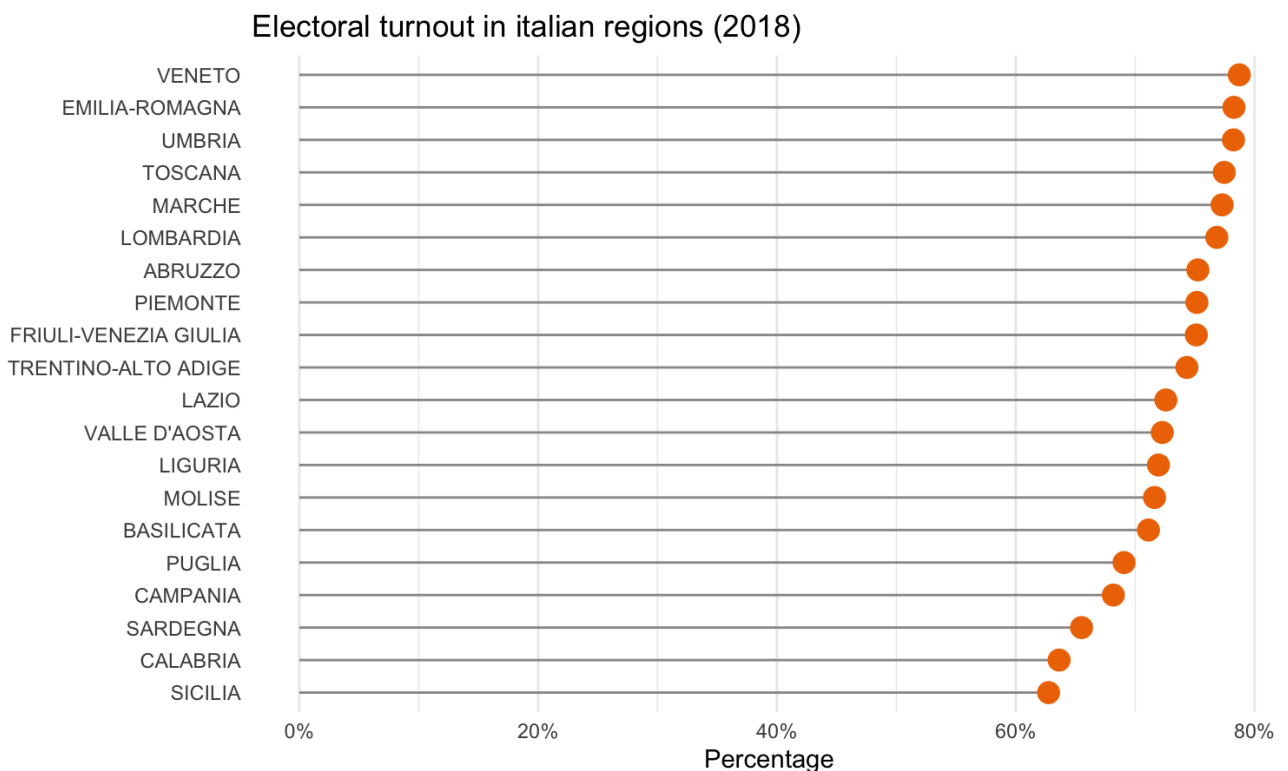| Purpose | Sort order | Chart orientation |
|---|---|---|
| Highlight the highest value | Descending | H: highest on top<br>V: highest on left |
| Highlight the lowest value | Ascending | H: lowest on top<br>V: lowest on left |

# Ranking – Barplot

Electoral turnout in italian regions (2018)

# Ranking – Dot plot

Electoral turnout in italian regions (2018)

# Lollypop (nonzero based scale)

Electoral turnout in italian regions (2018)



Turnout (%)

# Lollypop (zero based scale)

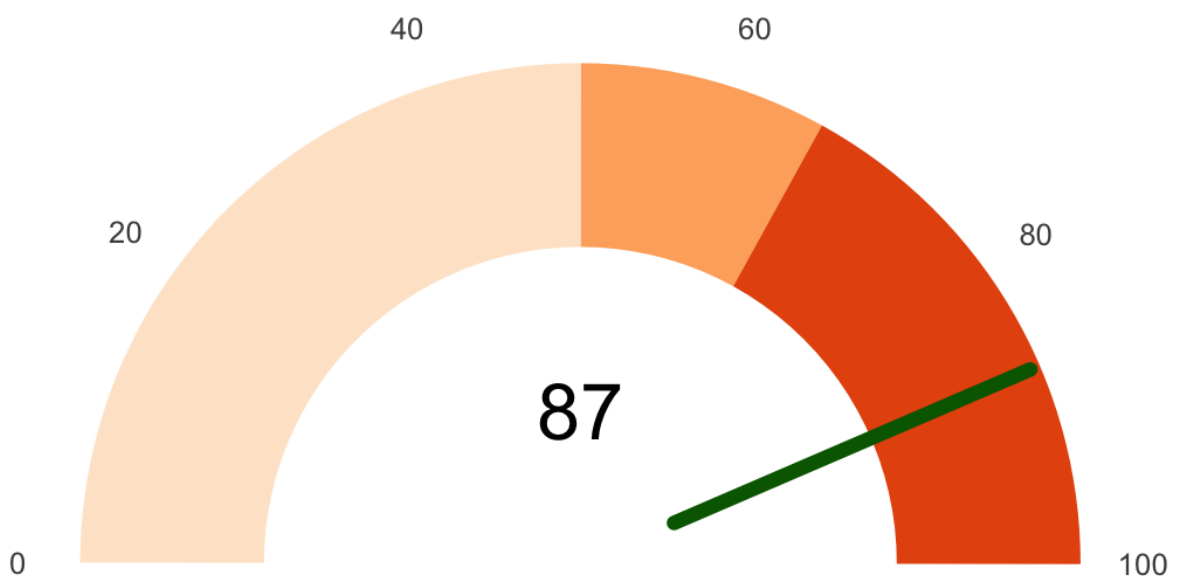Electoral turnout in italian regions (2018)



Percentage

# Deviation

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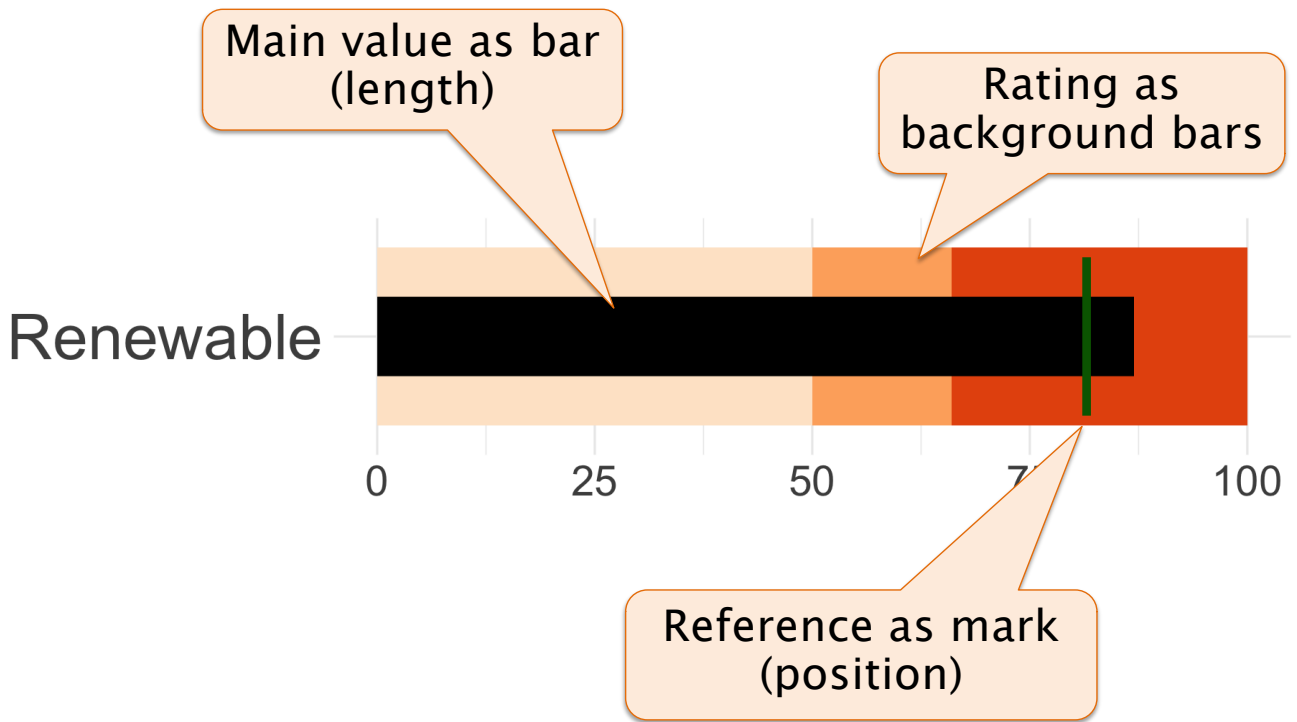  - Points (dots)
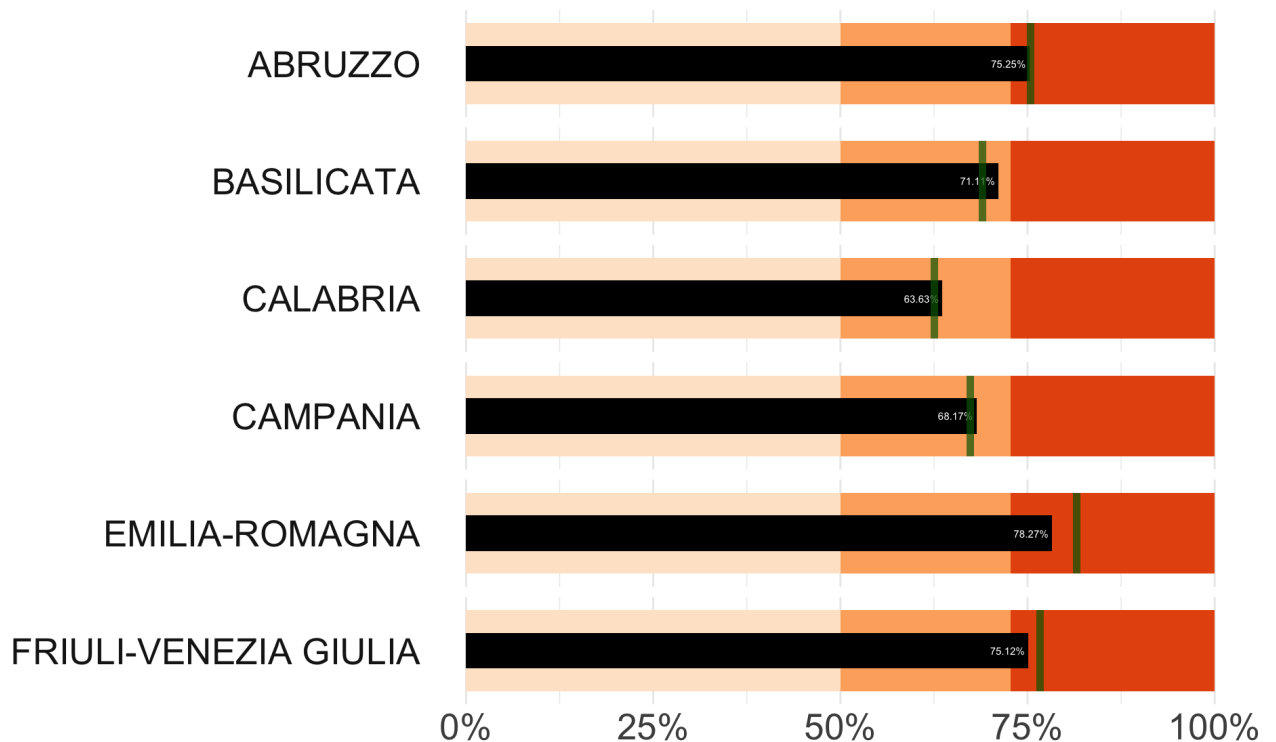  - Gauge
  - Bars
  - Bullet

# Angle + Position – Gauge

87

# Length+Position– Bullet Graph



Main value as bar (length)

Rating as background bars

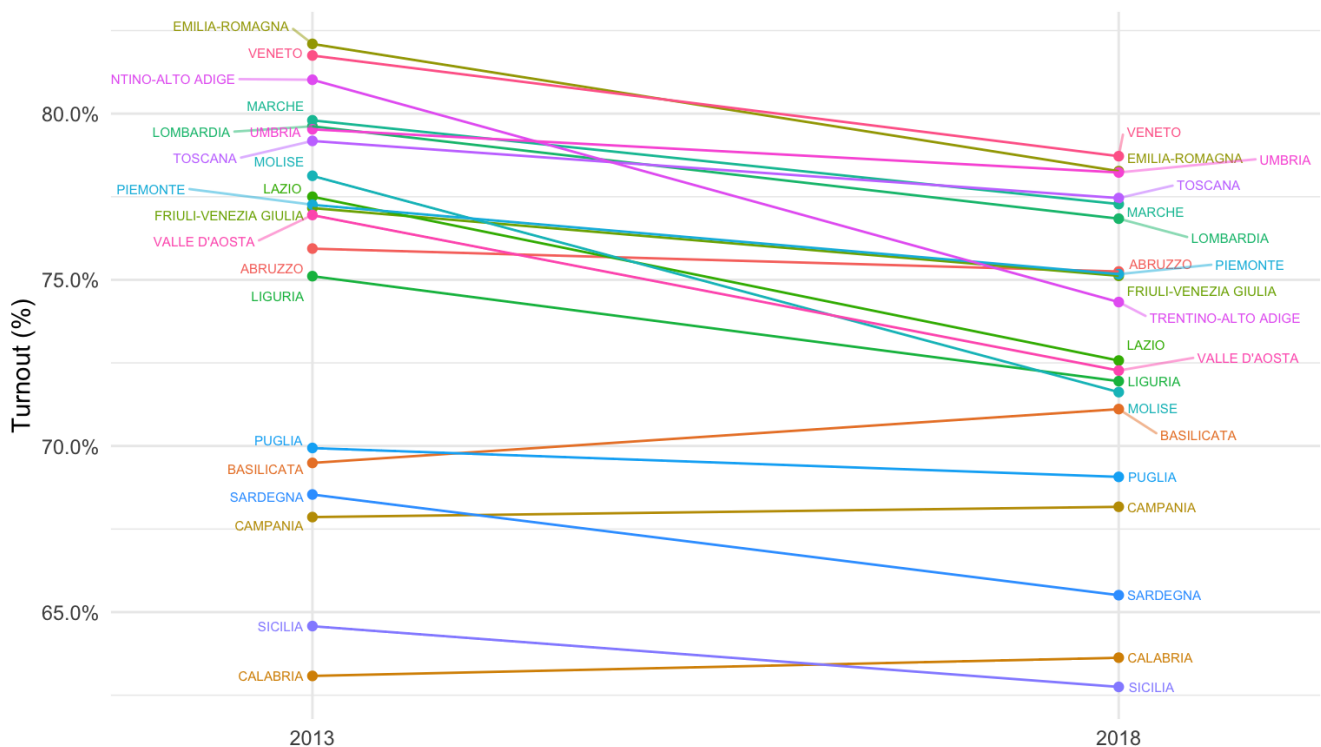Reference as mark (position)

# Length+Position– Bullet Graph

# Pre-post variation

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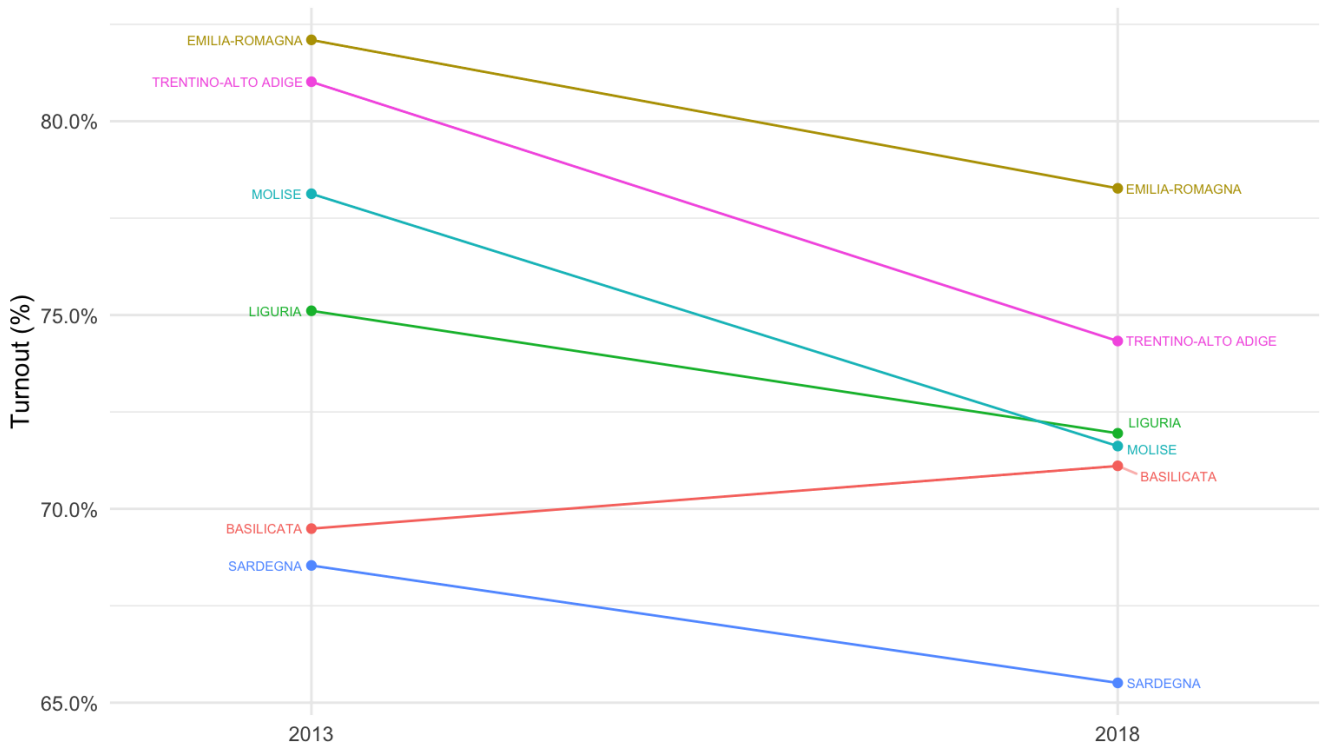  - ◆ With vs. without
  - ◆ …

# Slope chart



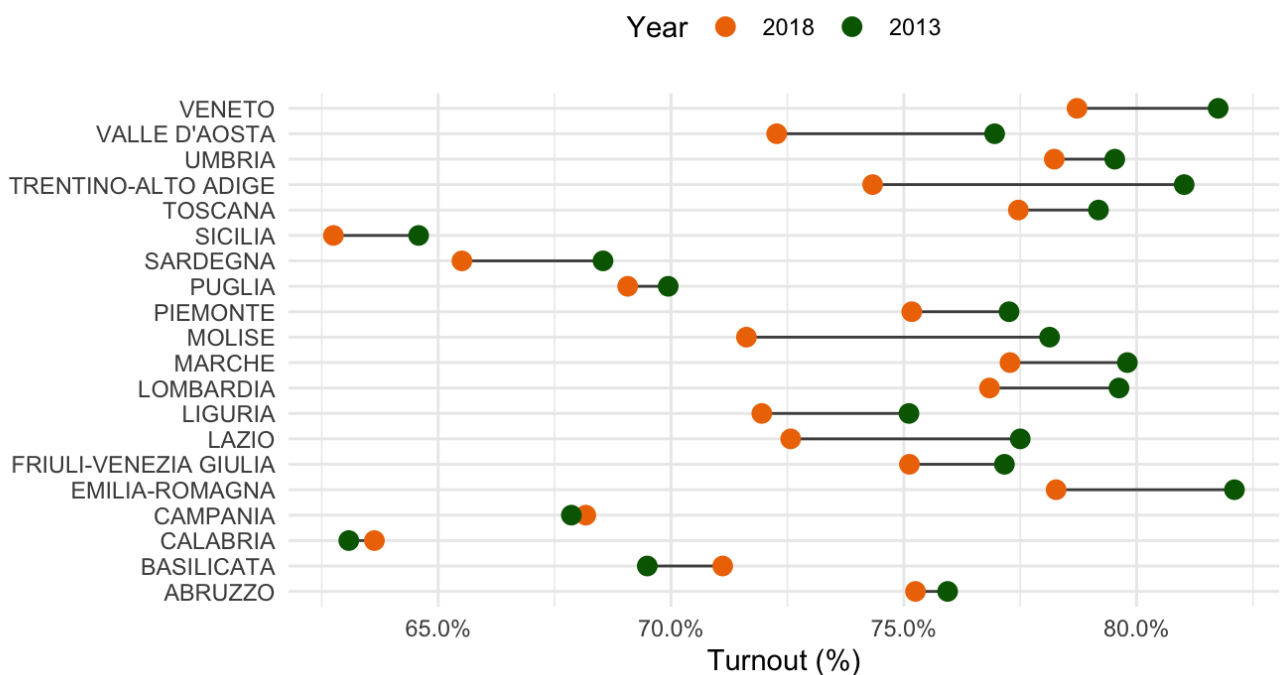Change in electoral turnout for italian regions (2018 vs. 2013)

# Slope chart

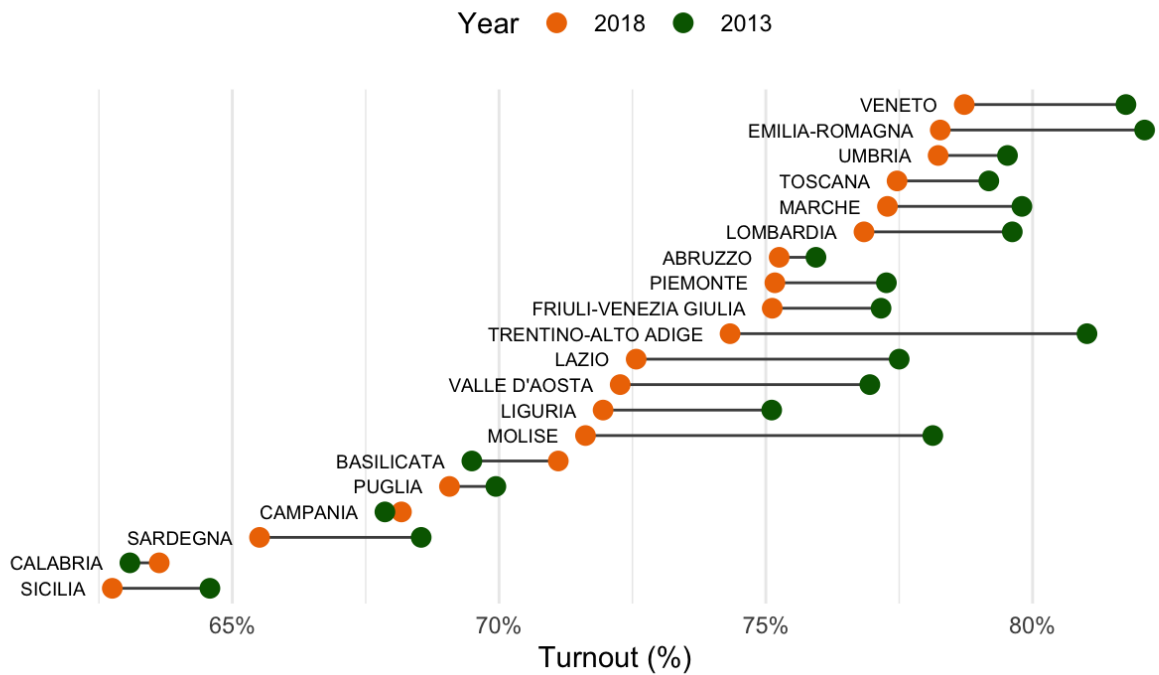Change in electoral turnout for italian regions (2018 vs. 2013)



# Dumbbell plot

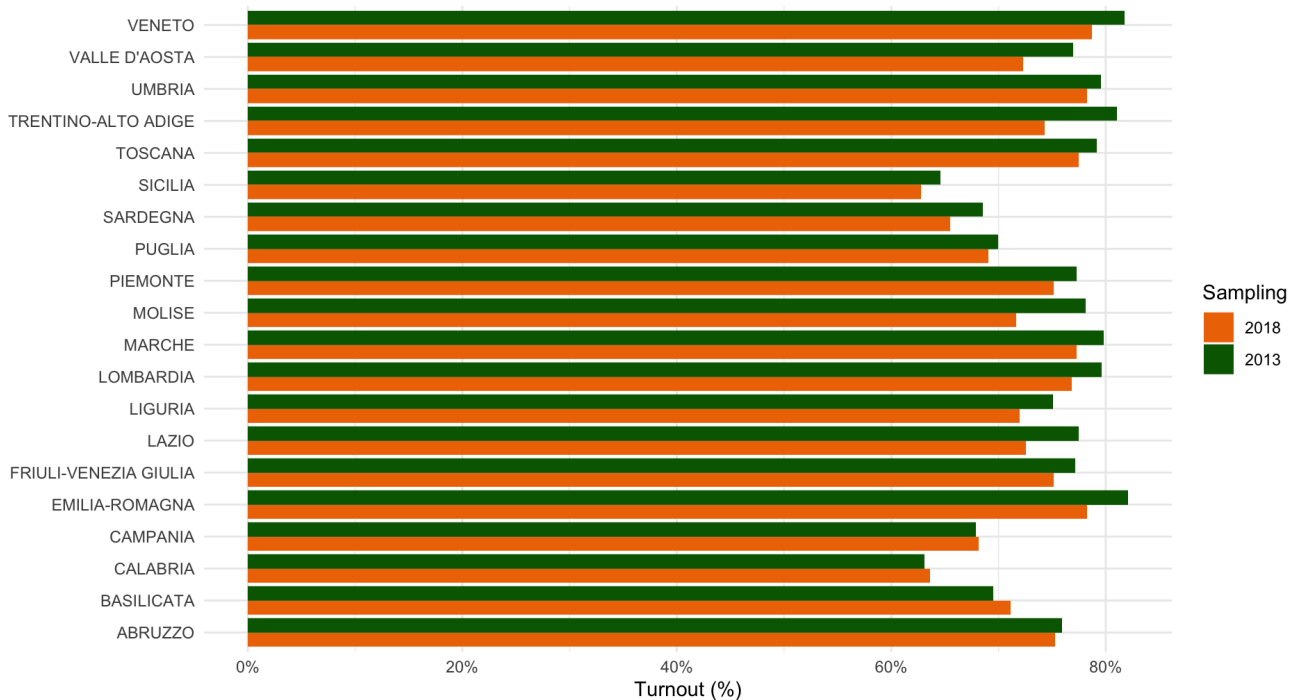Change in electoral turnout for italian regions (2018 vs. 2013)

# Dumbbell plot (sorted)

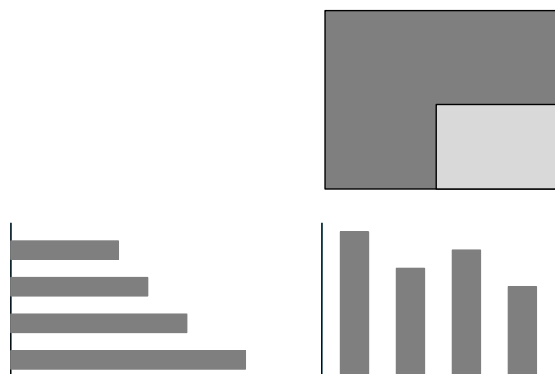Change in electoral turnout for italian regions (2018 vs. 2013)

# Clustered bars

# Proportion (Part-to-whole)

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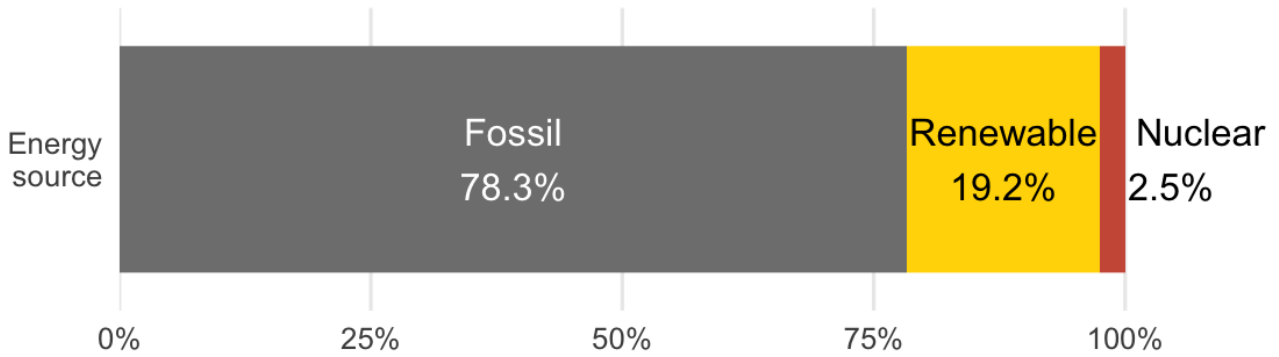- Can be used to compare frequency distribution across different categories sharing the same levels

# Proportion (Part-to-whole)

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

# Length – Stacked Bar



Energy source

| Fossil 78.3% | Renewable 19.2% | Nuclear 2.5% |

# Beware of MS–Excel Defaults

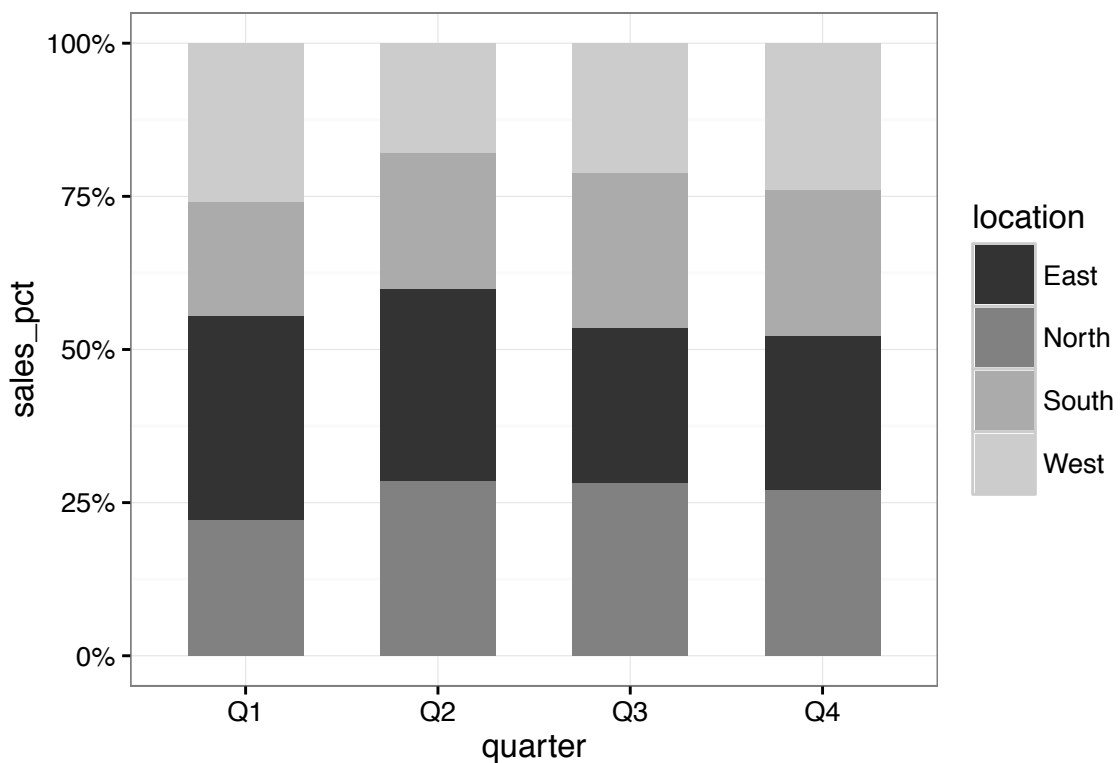|   | A | B |
|---|---|---|
| 1 | YES | 99% |
| 2 | NO | 1% |

# Stacked bar graph

# Stacked bars w/percentage

# Area – Treemap

# Area – Treemap

Fossil
78.3%

Nuclear
2.5%

Renewable
19.2%

# Area + Count – Waffle / Grid



Fossil
Renewable
Nuclear

# Area + Angle – Pie Chart



Nuclear
2.5%

Renewable
19.2%

Fossil
78.3%

# Pie charts

# Pie charts



■ sky

■ shady side of pyramid

■ sunny side of pyramid

# Pies

# Pies vs. Bars

# 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)

# Area+Angle+Length – Donut

# Marimekko Chart



Contribution to Overall Sales by Brand & Category (in US$)
(2011-12)

# Pareto chart

# Distribution

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

# 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

# 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

# Multiple distribution

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

# Paired diverging bargraph



Age group

| | Women % | Men % |
|---|---|---|
| 80 + | | |
| 75 - 79 | | |
| 70 - 74 | | |
| 65 - 69 | | |
| 60 - 64 | | |
| 55 - 59 | | |
| 50 - 54 | | |
| 45 - 49 | | |
| 40 - 44 | | |
| 35 - 39 | | |
| 30 - 34 | | |
| 25 - 29 | | |
| 20 - 24 | | |

Per cent women

Per cent men

# Multiple Frequency polygons



Time
- Noon
- 7pm
- 11pm

Turnout

Frequency

# Box plot

# Multiple Box plot

# Violin plot

# Multiple box plots

# Multiple violin plots

# Just dots for mean values

# Confidence intervals

# 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*

# Interval may be Asymmetric



Figure 5. Mean files per changeset.

> It is physically impossible to modify −6 files

# 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

# Diverging stacked bars

| Macroarea | N° | Domanda | | |
|---|---|---|---|---|
| Organizzazione del periodo didattico | 1 | Il carico di studio complessivo degli insegnamenti previsti nel periodo didattico è accettabile? | D1 | |
| | 2 | L'orario degli insegnamenti del periodo didattico è ben organizzato? | D2 | |
| Organizzazione di questo insegnamento | 3 | Le regole d'esame, gli obiettivi e il programma dell'insegnamento sono stati resi noti in modo chiaro? | D3 | |
| | 4 | L'insegnamento è stato svolto in maniera coerente con quanto dichiarato sul portale della didattica? | D4 | |
| | 5 | Le conoscenze preliminari da me possedute sono risultate sufficienti per la comprensione della materia ? | D5 | |
| | 6 | Il carico di studio richiesto da questo insegnamento è proporzionato ai crediti assegnati? | D6 | |
| | 7 | Il materiale didattico, indicato o fornito, è adeguato per lo studio della materia? | D7 | |
| | 8 | Le attività didattiche integrative (esercitazioni, lab, seminari, visite, ecc.) sono utili per l'apprendimento della materia? | D8 | |
| Efficacia del docente | 9 | Il docente rispetta gli orari di svolgimento dell'attività didattica? | D9 | |
| | 10 | Il docente è disponibile a fornire chiarimenti e spiegazioni? | D10 | |
| | 11 | Il docente interagisce efficacemente con gli studenti, stimolando l'interesse verso la materia? | D11 | |
| | 12 | Il docente espone gli argomenti in modo chiaro? | D12 | |
| Infrastrutture | 13 | Le aule in cui si svolgono le lezioni sono adeguate? | D13 | |
| | 14 | I locali e le attrezzature per le attività didattiche integrative sono adeguati? | D14 | |
| Interesse e soddisfazione | 15 | Sono interessato agli argomenti di questo insegnamento? (indipendentemente da come è stato svolto) | D15 | |
| | 16 | Sono soddisfatto di come è stato svolto questo insegnamento? | D16 | |
| | 17 | Al fine dell  apprendimento, la frequenza alle attività didattiche è utile? | D17 | |

# 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

# Time series

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

# Line plot

Italian Public Debt as percentage of GDP

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

# Bars

## Italian Public Debt as percentage of GDP



Bar chart values by year: 121%, 127%, 128%, 129%, 123%, 119%, 118%, 117%, 114%, 114%, 117%, 115%, 110%, 113%, 126%, 124%, 117%, 135%, 143%, 155%, 157%, 154%, 152%, 147%

X-axis: 1995, 2000, 2005, 2010, 2015

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

# Streamgraph



Twitter Topic Stream for Top Users

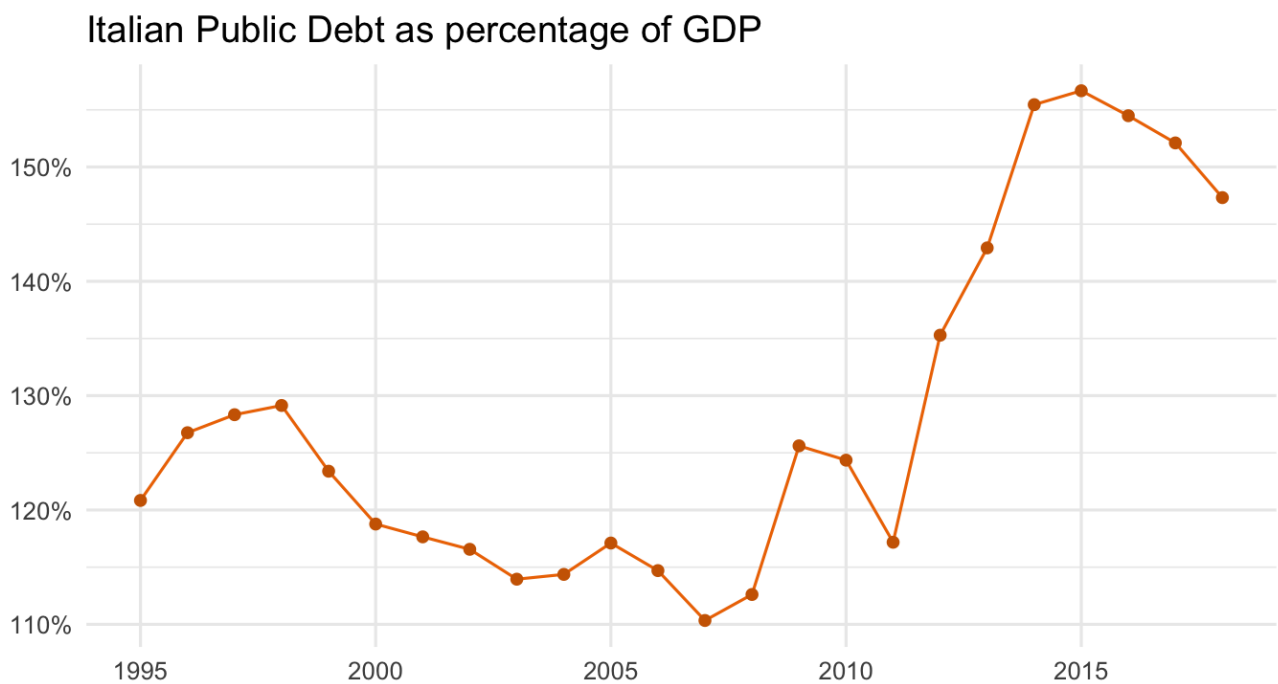Labels: TV, Google, Yahoo, QIK, Mahalo, Air, RSS, Flickr, Apple, Scoble, Pownce, Youtube, Problogger, Friendfeed, Microsoft, Mac, Facebook, Seesmic, Obama, Justin, Sxsw, Neoformix

X-axis: Dec 2006, Jan, Feb, Mar, Apr, May, Jun, Jul 2007, Aug, Sep, Oct, Nov, Dec, Jan, Feb, Mar 2008, Apr

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

# Correlation

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

# Points

# Points Guidelines

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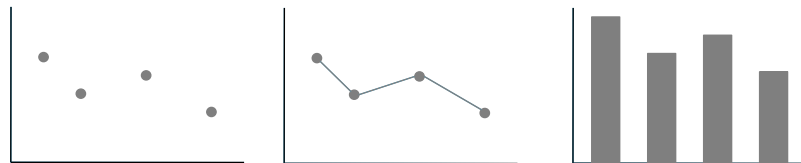  - ◆ Select radically distinct shapes (✚ ⭕)
  - ◆ Balance size of points and graph
  - ◆ Use outlined shapes
- **Lines must not obscure points**

# Scatter plot

Voter turnout in Italian general elections

# Overplotting

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

# Overplotting example

# Overplotting – Small

# Overplotting – Outlined

# Overplotting – Transparent

# Overplotting – Jittering

# Points and Lines



The slope encodes the amount of change.
**Warning**: non linear!

# Slope of lines

# Slope of lines



Trend line
Line of best fit
The slope encodes the regression coefficient

# Lines

- Easy perception of trends and overall shape of data

- Best suited for time series

- Variation encoded as slope
  - Clear direction
  - Approximate magnitude

# Paired diverging bars



Voter turnout in Italian general elections

# Categorical encoding attributes

- Encoding of categorical levels
  - ◆ Position (along an axis)
  - ◆ Size
  - ◆ Color
    - – Intensity

    Ordinal

    - – Saturation
    - – Hue
  - ◆ Shape
  - ◆ Fill pattern
  - ◆ Line style

# Position

Number of companies



| | Large | Medium | Small | Micro |
|---|---|---|---|---|

## Position ×
# Color (hue)

**2003 Sales**

■ Direct  ■ Indirect

# Size



Size corresponds to size

# Point shape

# Line style

# Fill Texture

# Fill Gray Levels

# 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

# Heatmaps

**Measles**



Vaccine introduced

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
Alaska, Ark., Calif., Conn., Del., Ga., Iowa, Ill., Kan., La., Md., Mich., Mo., Mont., N.D., N.H., N.M, N.Y., Okla., Pa., S.C., Tenn., Utah, Vt., Wis., Wyo.

1930  1940  1950  1960  1970  1980  1990  2000  2010

0k  1k  2k  3k  4k

---

# Heatmaps

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

# Heatmaps

Measles cases per US State per year

# Rainbow palette



Wind Chill / Heat Index (°F)
11:00 AM February 15, 2021 CST
Data provided by NOAA's National Weather Service. Created 11:10:34 AM February 15, 2021 CST. © Copyright 2021

# Gradient palette



WASH.

MONT.

N.D.

MINN.

MAINE

ORE.

IDAHO

S.D.

WIS.

MICH.

VT.

N.H.

N.Y.

MASS.

Between 0°
and 32°F

WYO.

IOWA

CONN.

R.I.

NEV.

NEB.

ILL.

IND.

OHIO

PA.

N.J.

UTAH

COLO.

**Low temperatures
below 0°F**

MO.

KAN.

W.VA.

MD.

DEL.

VA.

CALIF.

ARIZ.

KY.

N.C.

OKLA.

TENN.

S.C.

Above 32°F

N.M.

ARK.

MISS.

ALA.

GA.

TEXAS

Between 0°
and 32°F

LA.

Above 32°F

FLA.

Lowest temperatures forecast
Sunday through Tuesday

−10° 0°   10°      30°      50°      70°F

# SUPPORT ELEMENTS

# Support elements

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

# 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

# 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

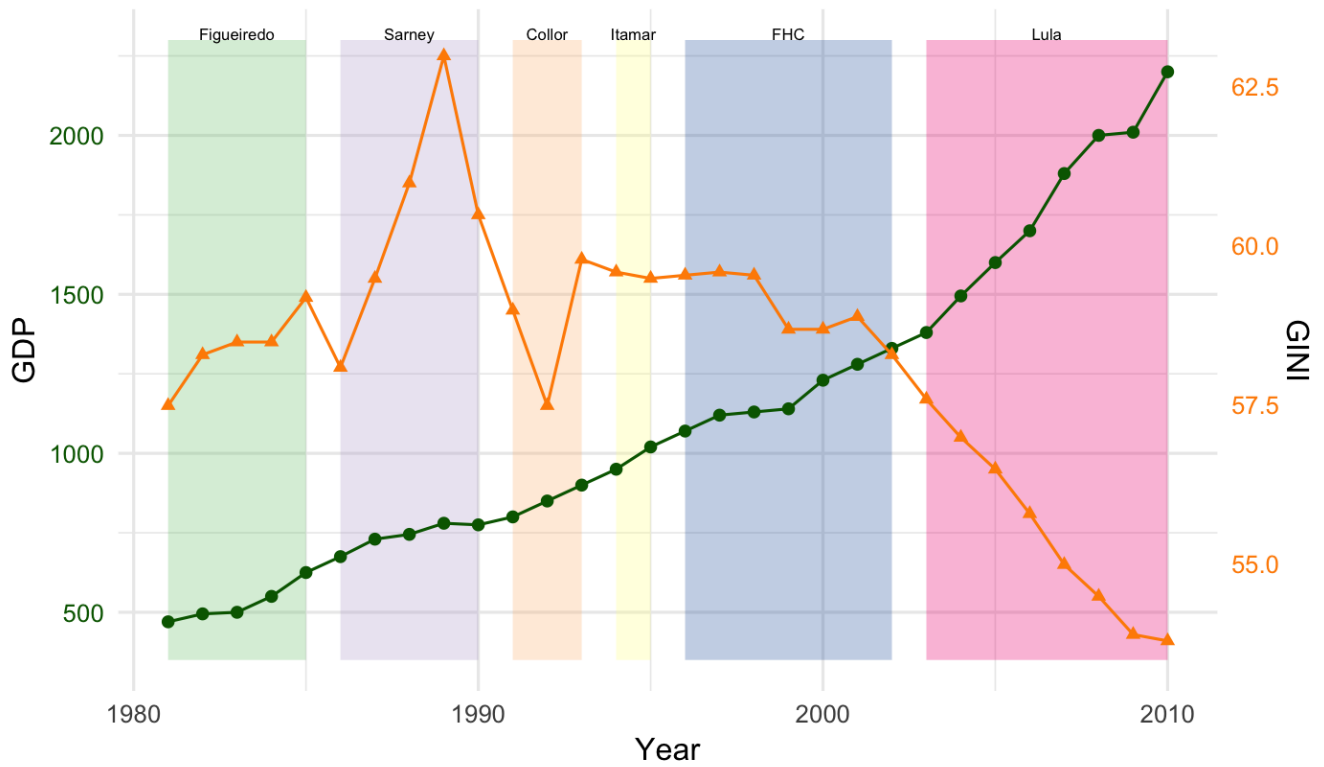# Multiple variables

- 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

# 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

# Small multiples



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

# Small multiples

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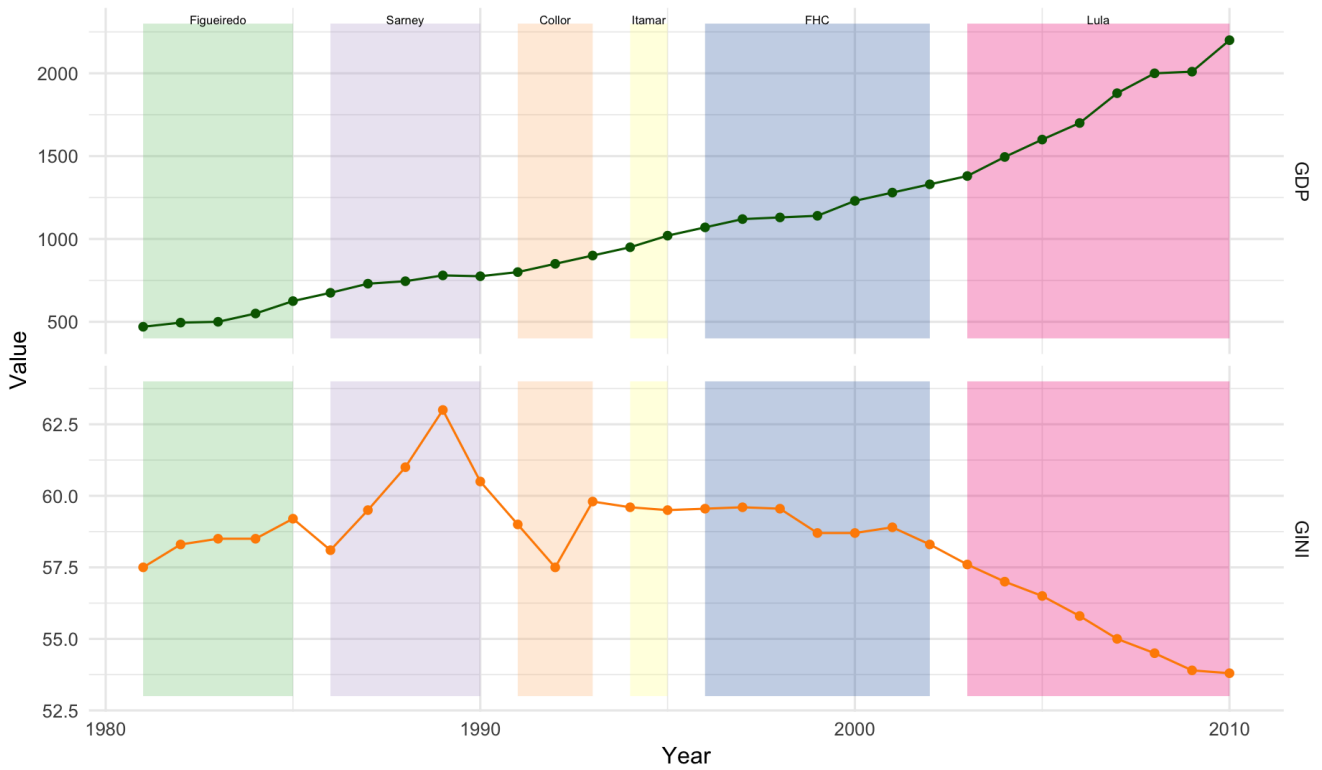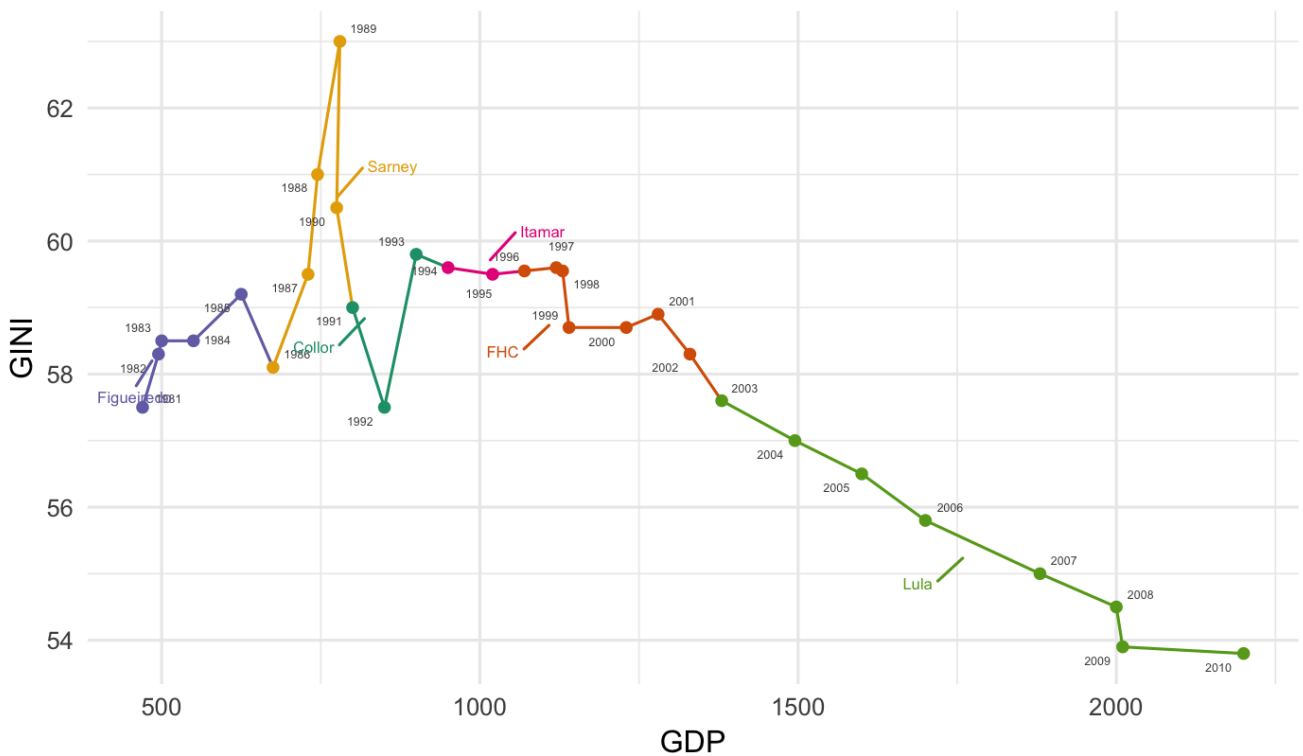  - Same ordering of categorical levels
- Arrangement
  - Align axis that involve comparison
    - Possibly along a matrix

# 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

# Log scale

- Reduce visual difference between quantitative data sets with significantly wide ranges

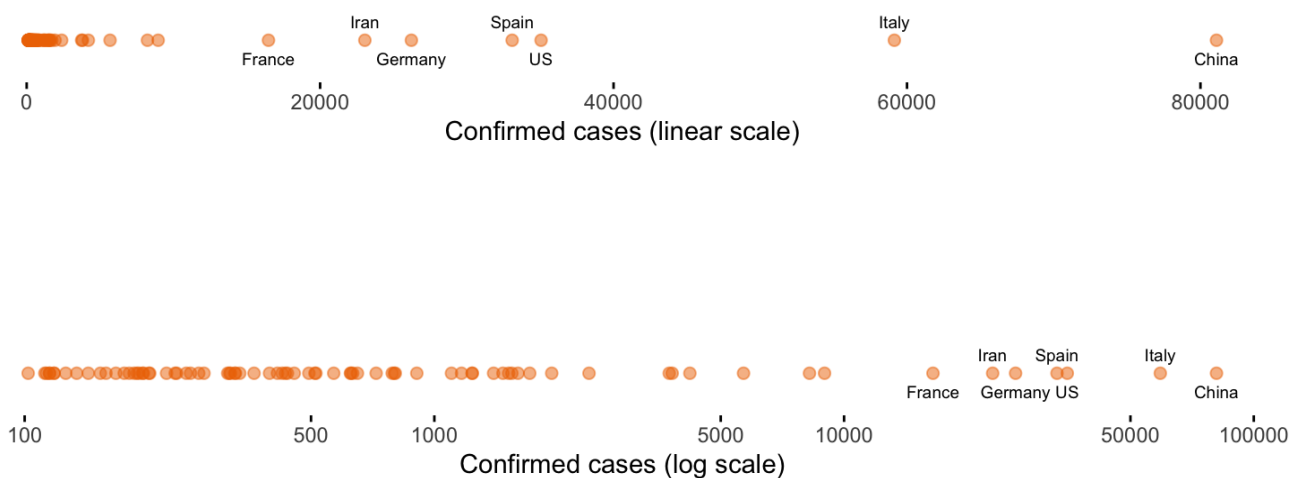- Differences are proportional to percentages

- Constant percentage increase correspond to a line

# Log scale – wide range

# Log scale – differences



+80

+144

400
300
200
100
0

324
180
100

A        B        C

Same **absolute gains** correspond to same distance

1000

100

10

1

+80%    +80%

100    180    324

A        B        C

Same **percentage gains** correspond to same distance

# Log scale – variation



140000
120000
100000
80000
60000
40000
20000
0

Q1    Q2    Q3    Q4

North
South

Parallel lines for same absolute gains

1000000
100000
10000
1000
100
10
1

Q1    Q2    Q3    Q4

North
South

Parallel lines for same percentage gains

# 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 / @jburnmurdoch
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

# Labels

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

# 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
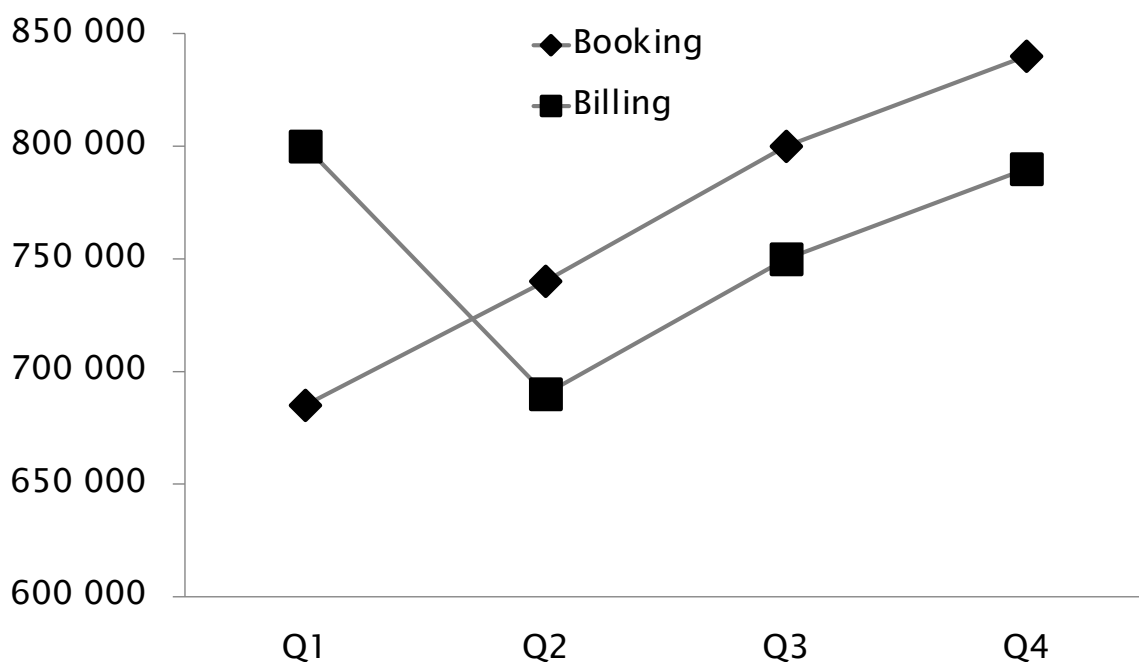
# Legends

- 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
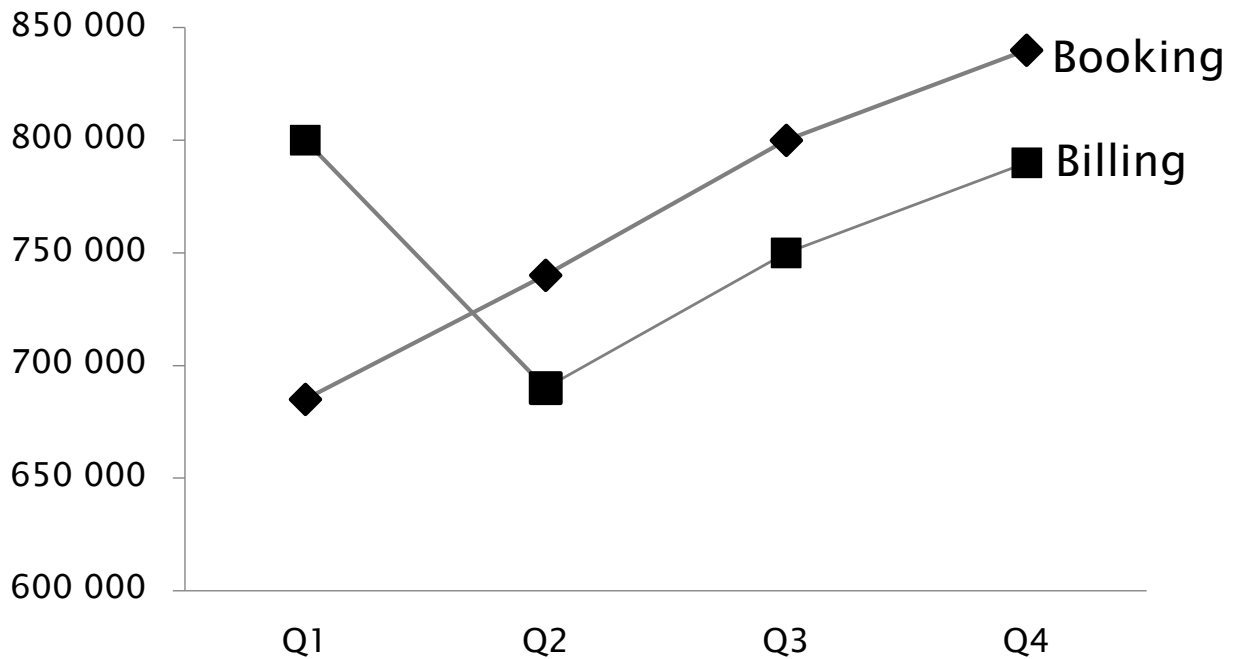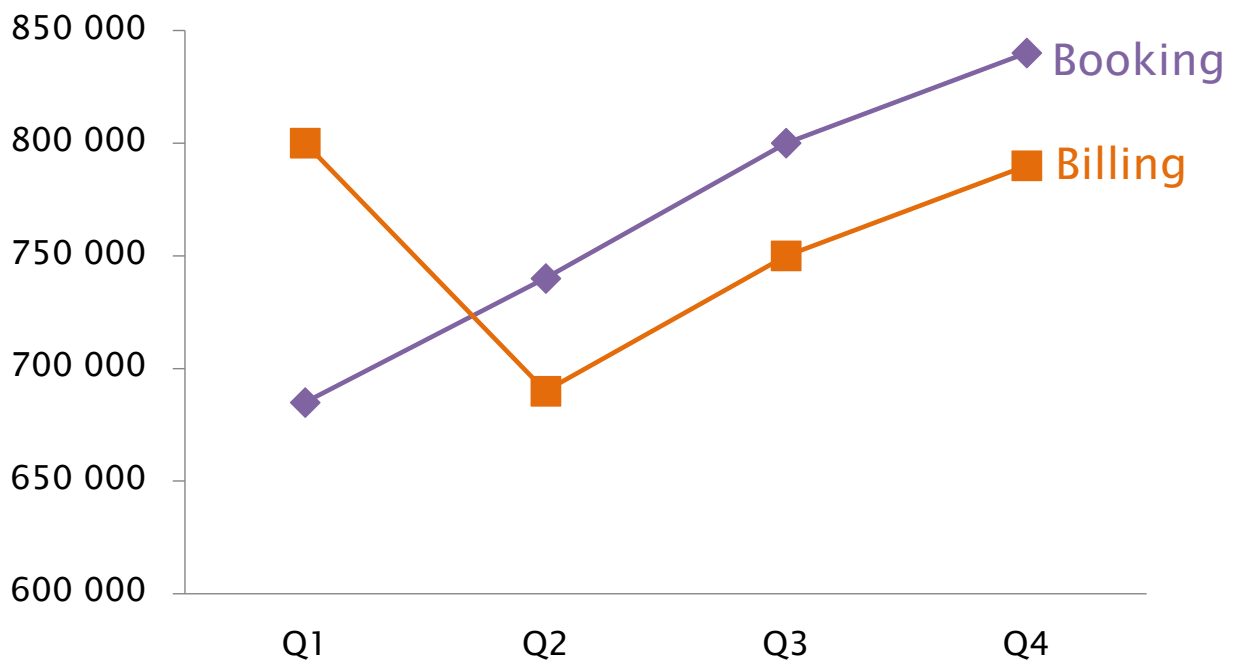  - Joint use of attributes ease discrimination
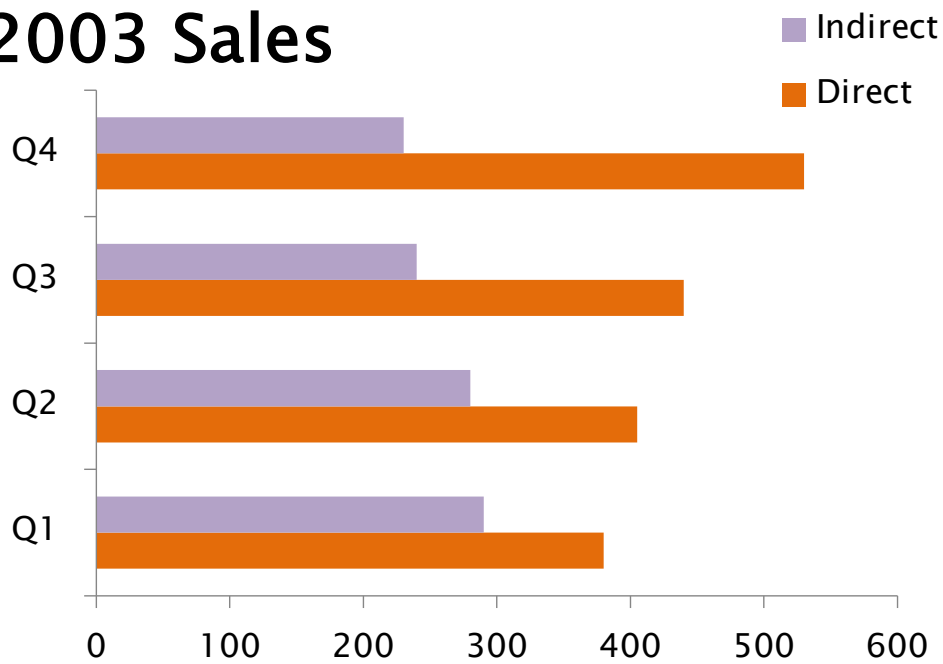
# Legend

# Direct labeling

# Direct labeling and color

# Legend

## 2003 Sales

# Direct labeling

## 2003 Sales

# 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

# DASHBOARD

# 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

# Dashboard

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

# Provide context for data

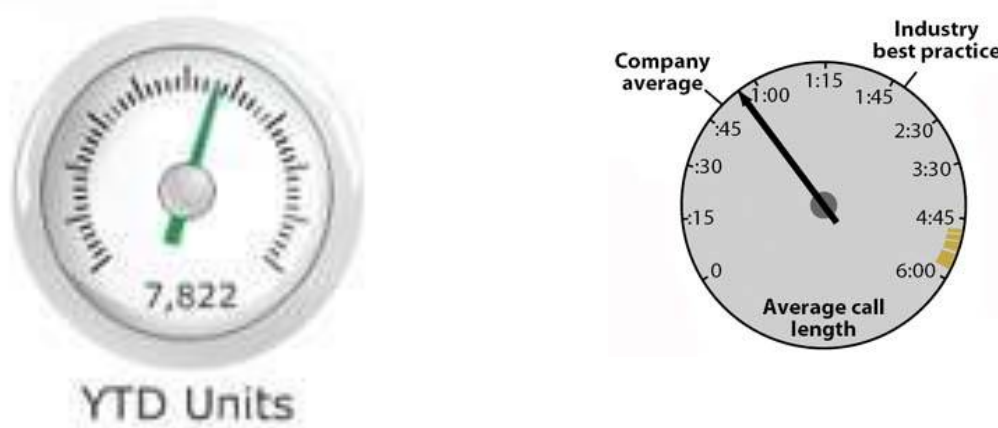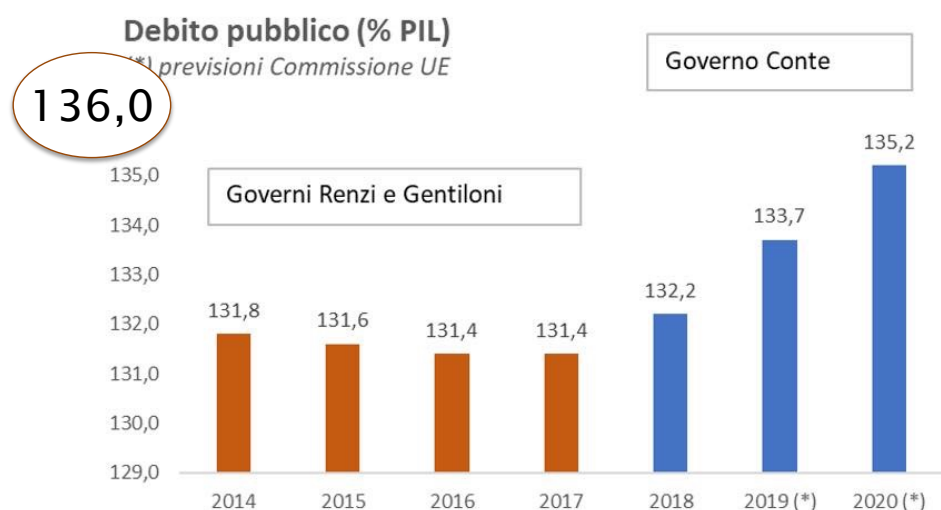- References allow judging the data

# Use appropriate detail

- Typical counterexamples
  - Dates with seconds detail
  - Decimals

# 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

| Year | Value |
|------|-------|
| 2015 | -0.20 |
| 2016 | -0.20 |
| 2017 | 0.00 |
| 2018 | 0.80 |
| 2019 | 1.50 |
| 2020 | 1.50 |

# Use appropriate visualization

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

# Visualization instruments

- Tables
  - ◆ Textual information

- Graphs
  - ◆ Visual information

# Avoid decorations

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

# Avoid decorations

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

# Avoid decorations

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

A

B

# Avoid decorations

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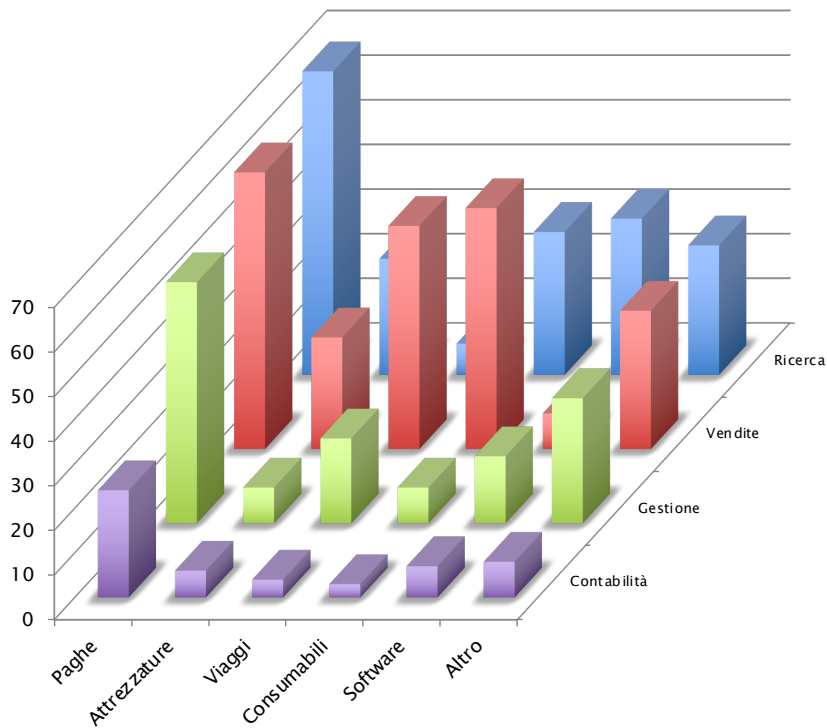  any measure
  - ◆ **Typically color**

# 3D diagrams

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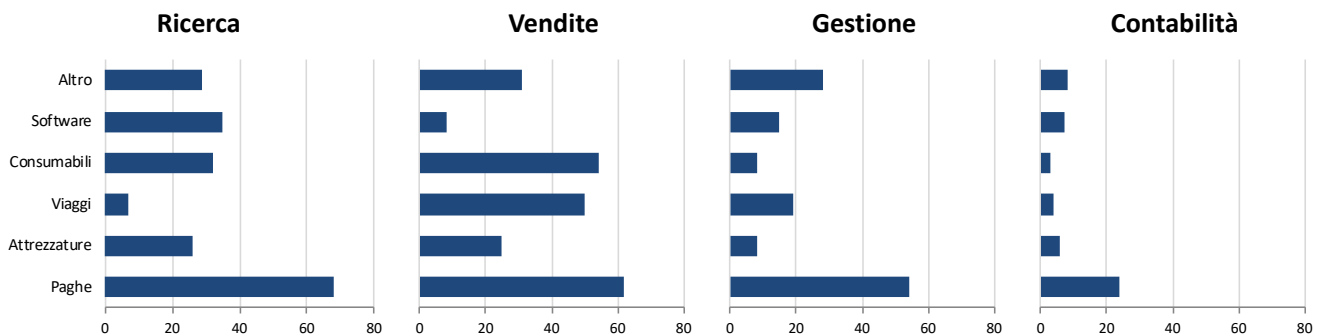  - ◆ Depth or height distract and make
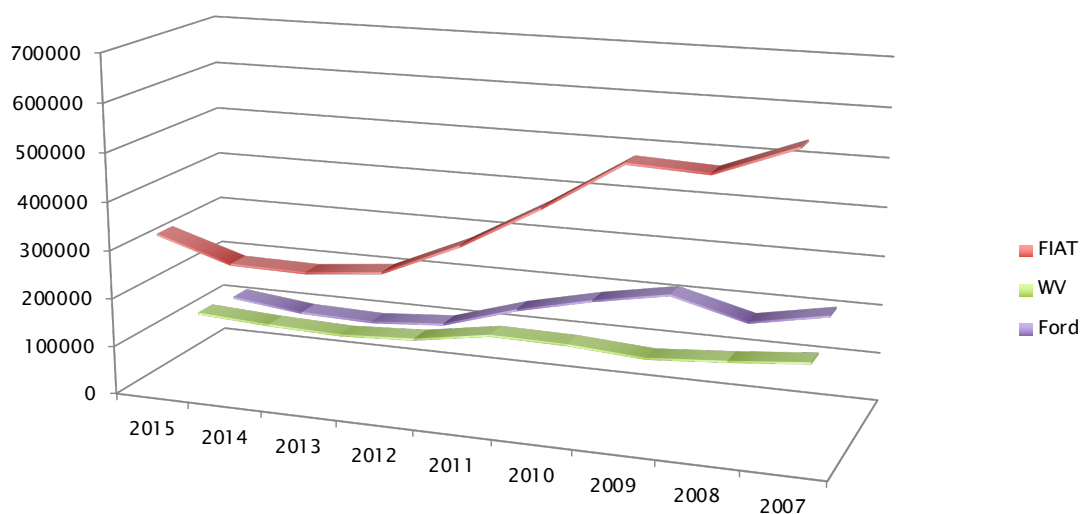    comparison more difficult

# Encoding 3D
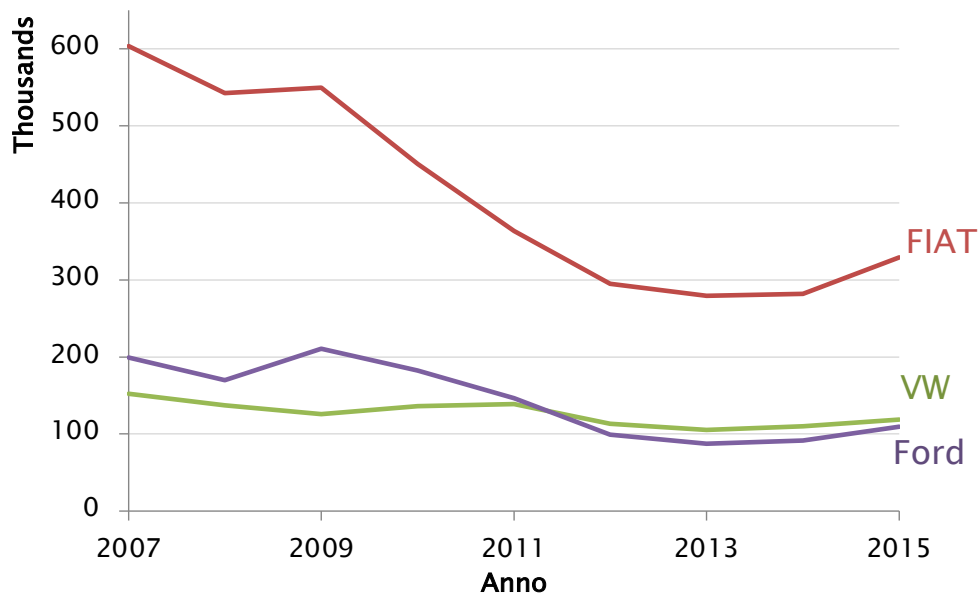
# Encoding 3D → 2D

# Decorative 3D

Immatricol.

# Decorative 3D → 2D

## Immatricolazioni auto per marchio sul mercato italiano

Immatricol.

# 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.

# 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

# 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.Correll, 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