Graph Construction

Data Management and Visualization







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 plot/chart elements
 - Composition of such common 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*) and their aesthetic attributes (*e.g.*, *color*, *position*)
- GUIDE: guides (e.g., axes, legends)

Specification for a scatter plot

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

SoftEng

Graph visual components

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

Visual Encoding

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

SoftEng

VISUAL RELATIONSHIPS

Data Visualization

Understanding



Relationships

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

Quantitative encoding



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 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!

SoftEng

Comparison – Barplot





Barplot (non zero based scale)



Barplot (non zero based scale)



SoftEng

Barplot vertical labels



Bars Guidelines

- Use horizontal bars when
 - A descending order ranking
 - Categorical label don't fit
- Proximity
 - Use a 1:1 bar:spacing ratio ±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



Area – Bubble plot

Electo	Electoral turnout in italian regions (2018)																		
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 -



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 V: highest on left
Highlight the lowest value	Ascending	H: lowest on top V: lowest on left

SoftEng

Ranking – Barplot





Ranking - Dot plot



SoftEng.

30

Lollypop (non zero based scale)



Lollypop (zero based scale)



Deviation

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



SoftEng

Angle + Position - Gauge



Length+Position- Bullet Graph



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



Pre-post variation

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

SoftEng

Slope chart



Dumbbell plot



Clustered bars





Proportion (Part-to-whole)

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



SoftEng

42

Length - Stacked Bar



Beware MS-Excel Default



Stacked bar graph



Stacked bars w/percentage



Area - Treemap





Area – Treemap

Fossil 78.3%	Nuclear 2.5% Renewable 19.2%	

SoftEng

48

Area + Count - Waffle / Grid





Area + Angle - Pie Chart





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



Pareto chart



SoftEng

Marimekko Chart



Distribution

- Two main types
 - 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

SoftEng

Histogram



SoftEng

Frequency polygon



Boxplot



Violin plot



Violin + Boxplot



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



https://unstats.un.org/unsd/genderstatmanual/Print.aspx?Page=Presentation-of-gender-statistics-in-graphspreak and the statistics-in-graph statis



Multiple Frequency polygons



Multiple Box plot

SOftEng



Violin plot



Multiple box plots





Multiple violin plots



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



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

lacroarea	N° Domanda	-50%	0%	50%	100%
Organizzazione del	Il carico di studio complessivo degli insegnamenti previsti nel periodo didattico è accettabile?	D1			
periodo didattico	2 L'orario degli insegnamenti del periodo didattico è ben organizzato?	D2			
	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			
Organizzazione di	5 Le conoscenze preliminari da me possedute sono risultate sufficienti per la comprensione della materia ?	D5			
questo insegnamento	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			
	9 Il docente rispetta gli orari di svolgimento dell'attività didattica?	D9			
	10 Il docente è disponibile a fornire chiarimenti e spiegazioni?	D10			
Efficacia del docente	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			
		-50%	0%	50%	100%

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



SoftEng

Line plot



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



Bars



SoftEng

79

Streamgraph



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 (+ O)
 - Balance size of points and graph
 - Use outlined shapes
- Lines must not obscure points

Scatter plot



SoftEng

Overplotting

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

SoftEng

Overplotting example



SoftEng

86

Overplotting – Small



SOftEng

88

Overplotting – Outlined



Overplotting - Transparent



Overplotting – Jittering



SoftEng

Points and Lines



SoftEng

92

Slope of lines



Slope of lines



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 2013 2018 VENETO EMILIA-ROMAGNA UMBRIA TOSCANA MARCHE LOMBARDIA ABRUZZO PIEMONTE FRIULI-VENEZIA GIULIA TRENTINO-ALTO ADIGE I AZIO VALLE D'AOSTA LIGURIA MOLISE BASILICATA PUGLIA CAMPANIA SARDEGNA CALABRIA SICILIA 75% 50% 25% 0% 25% 50% 75% Turnout 96

SOftEng

Categorical encoding attributes

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

Ordinal

Position



SoftEng

98

Color (hue)



Size



Point shape



Line style



SoftEng

Fill Texture



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

SoftEng

104

Heatmaps





Heatmaps

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

SoftEng

Heatmaps



106

SUPPORT ELEMENTS

SoftEng

108

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

SoftEng

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

110

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



SoftEng

112

Double scale (alternative)



Multiple graphs



SOftEng

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



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



121

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





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

128

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



Direct labeling



SoftEng

134

Direct labeling and color



Legend



SoftEng

136

Direct labeling

2003 Sales



SoftEng

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

SoftEng

138

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

SoftEng	140

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



PUC

Use the right measures

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



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

SoftEng

146

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



SoftEng

Avoid decorations

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



148

Α

В

Avoid decorations

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

SoftEng

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



SoftEng

152

Encoding $3D \rightarrow 2D$



Decorative 3D

Immatricol.



SoftEng

154

Decorative $3D \rightarrow 2D$



SoftEng
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 <u>http://ft.com/vocabulary</u>

References

- R.Olson. Revisiting the vaccine visualization
 - http://www.randalolson.com/2016/03/04/revisiting <u>-the-vaccine-visualizations/</u>
- Nathan Yau. 9 Ways to Visualize Proportions A Guide
 - http://flowingdata.com/2009/11/25/9-ways-tovisualize-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*
 - <u>http://graphics.cs.wisc.edu/Papers/2014/CG14/Preprint.pdf</u>

SoftEng

158