

# Un hackathon dataviz pour aider une ONG

## De la théorie aux (bonnes) pratiques de la data visualisation

*Mars, 2019*

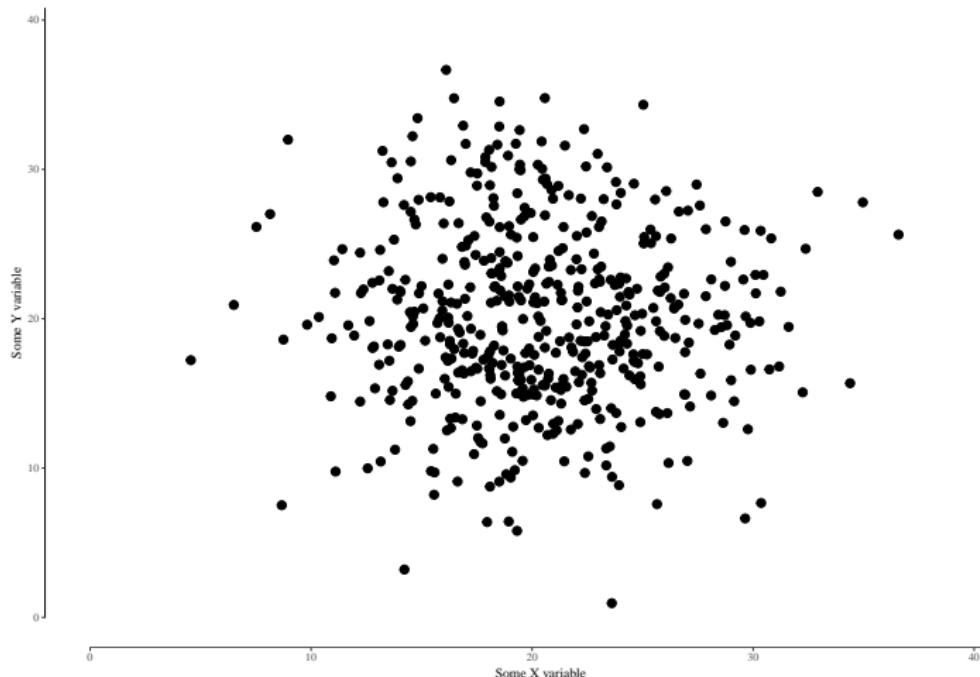
Christophe Bontemps  
*Toulouse School of Economics (INRA)*  
&  
Édith Maulandi  
*Viz For SocialGood*



# THE “VISUAL PERCEPTION” OF A GRAPHIC

What do you see ?

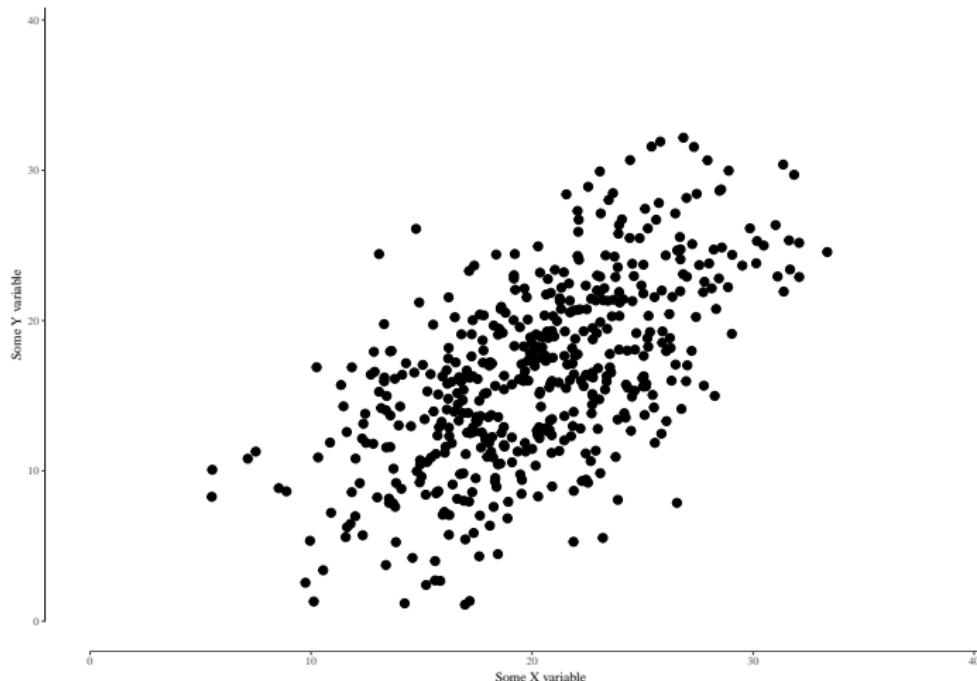
Some points (N = 500 )



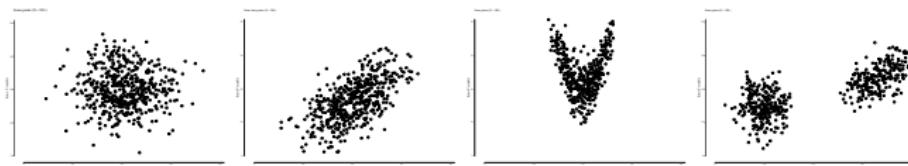
# THE “VISUAL PERCEPTION” OF A GRAPHIC

And here, what do you see ?

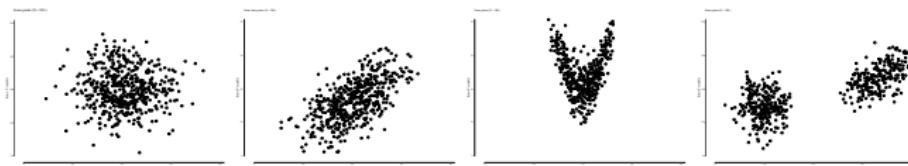
Some other points (N = 500 )



# “VISUAL PERCEPTION” AS A STATISTICAL TEST

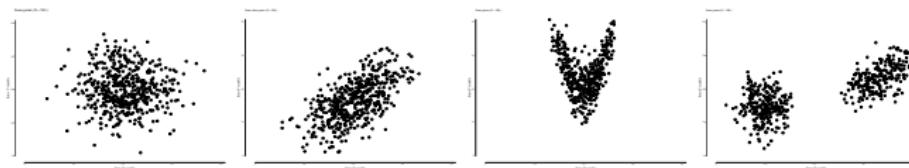


# “VISUAL PERCEPTION” AS A STATISTICAL TEST



*“The human eye acts is a broad feature detector and general statistical test”. Buja et al. (2009)*

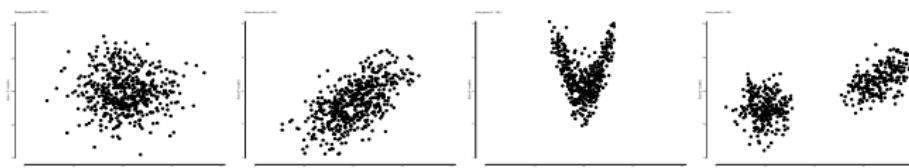
# "VISUAL PERCEPTION" AS A STATISTICAL TEST



*"The human eye acts is a broad feature detector and general statistical test". Buja et al. (2009)*

**Test :**  $H_0 : \{\text{There is "nothing"}\} = \{\text{No relation}\}$

# "VISUAL PERCEPTION" AS A STATISTICAL TEST



*"The human eye acts is a broad feature detector and general statistical test". Buja et al. (2009)*

**Test :**  $H_0 : \{\text{There is "nothing"}\} = \{\text{No relation}\}$

$H_1 : \{\text{There is "something"}\} = \{\text{There is some relation}$   
(Correlation, linearity, heterogeneity, groups..) }

# “VISUAL PERCEPTION” AS A COMPARISON



# “VISUAL PERCEPTION” AS A COMPARISON



- ▶ What do you see here ?

# “VISUAL PERCEPTION” AS A COMPARISON



- ▶ What do you see here ?



# “VISUAL PERCEPTION” AS A COMPARISON



- ▶ What do you see here ?



Difficult to see the maximum/minimum of each curve...

# “VISUAL PERCEPTION” AS A COMPARISON



- ▶ What do you see here ?



Difficult to see the maximum/minimum of each curve...

Idea shared by Gelman (2004) and Munzner (2014)

# ROAD MAP : A REVOIR

- ▶ What is data visualisation ?
  - ▶ What for ? What questions ?
  - ▶ Type of graphics, classics, Case studies
  - ▶ Does dynamic help ?
- ▶ Rules
  - ▶ Bertin, Tufte and Cleveland's rules
  - ▶ Tables *vs* graphics
  - ▶ Visual perception
- ▶ Visualizing complexity
  - ▶ Mixing variables types
  - ▶ Visualising many
  - ▶ *Maps and networks*
  - ▶ *Case studies*
- ▶ *Interactions*

# GOALS OF DATA VISUALISATION

Data visualisation serves at least two main purposes

- ▶ Data exploration

Graphs as visual tests, **comparisons** → **short** time to built and to read

# GOALS OF DATA VISUALISATION

Data visualisation serves at least two main purposes

- ▶ Data exploration

Graphs as visual tests, **comparisons** → **short** time to built and to read

- ▶ Data representation

Summaries, **comparisons**, storytelling → **long** time to build, short time to read

# GOALS OF DATA VISUALISATION

Data visualisation serves at least two main purposes

- ▶ Data exploration

Graphs as visual tests, **comparisons** → **short** time to built and to read

- ▶ Data representation

Summaries, **comparisons**, storytelling → **long** time to build, short time to read

The problem is that :

*“ Communicating implies **simplification**  
data exploration implies **exhaustivity**”*

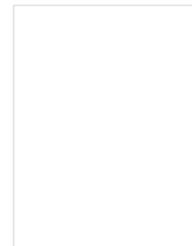
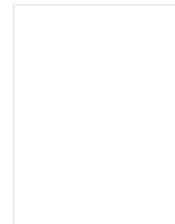
## THE BIG PICTURE : FROM "DATA TO VIZ"



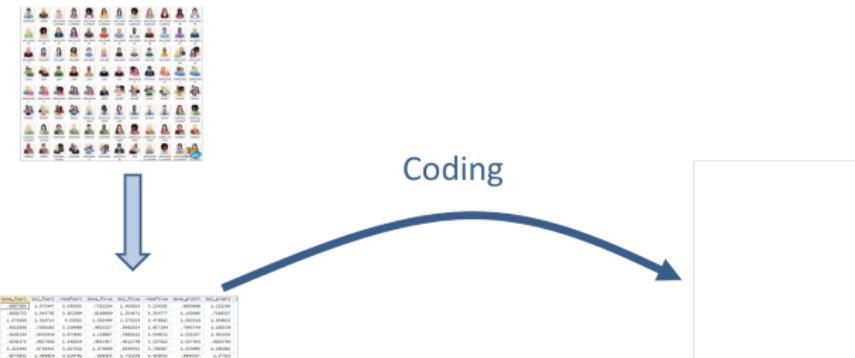
# THE BIG PICTURE : FROM "DATA TO VIZ"



# THE BIG PICTURE : FROM "DATA TO VIZ"

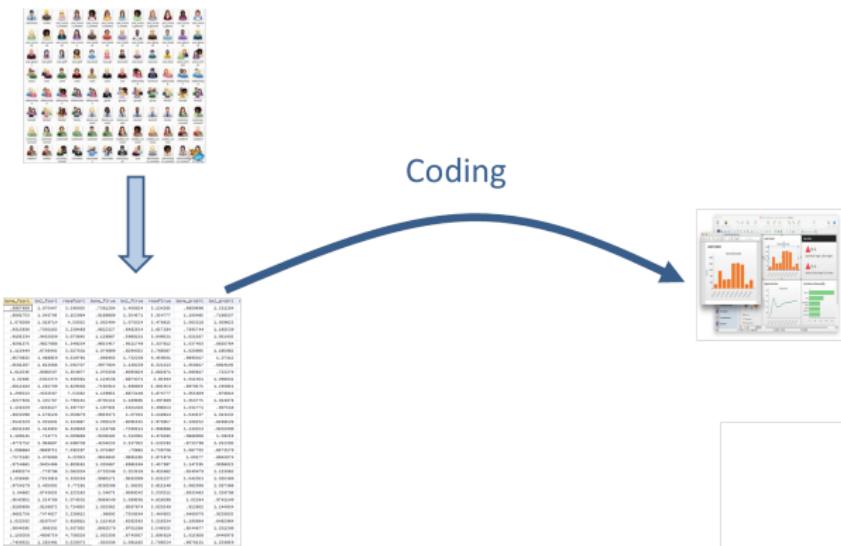


# THE BIG PICTURE : FROM "DATA TO VIZ"

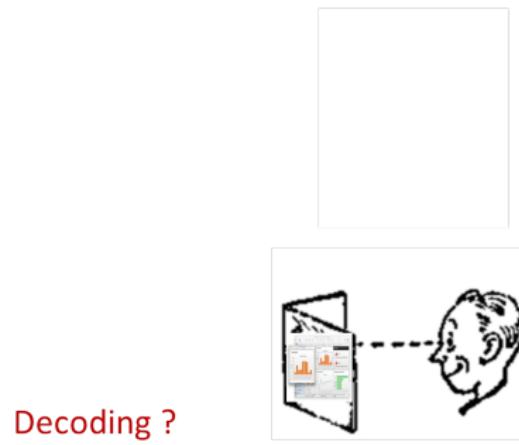


Row	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	Col 8	Col 9	Col 10
1	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
2	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
3	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
4	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
5	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
6	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
7	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
8	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
9	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000
10	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

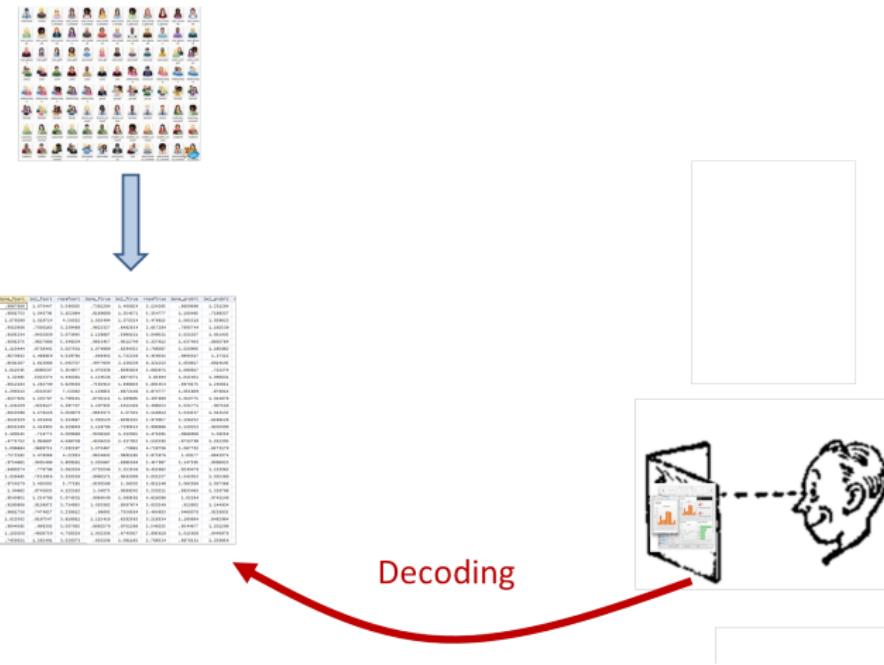
# THE BIG PICTURE : FROM "DATA TO VIZ"



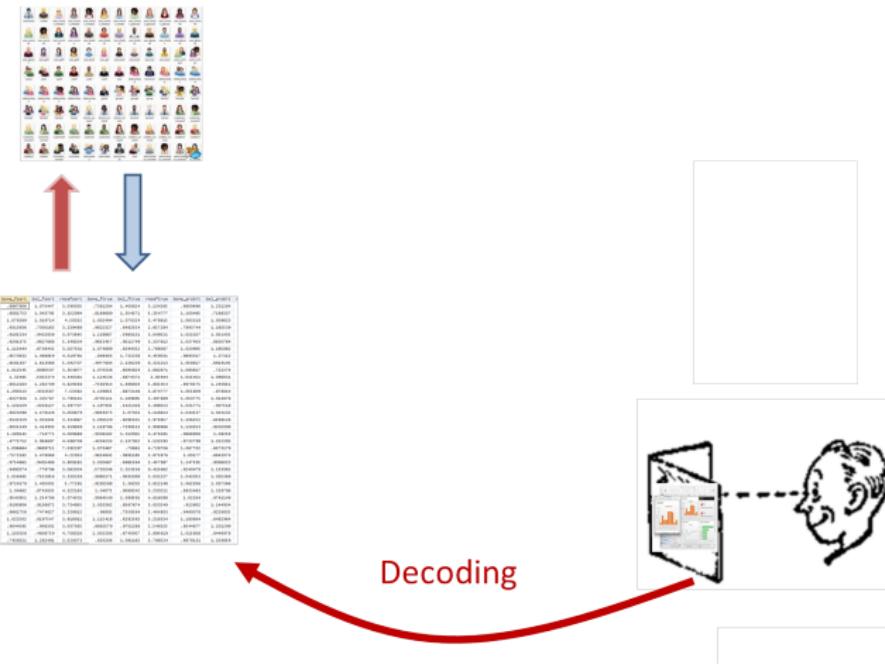
# THE BIG PICTURE : FROM "VIZ TO DATA"



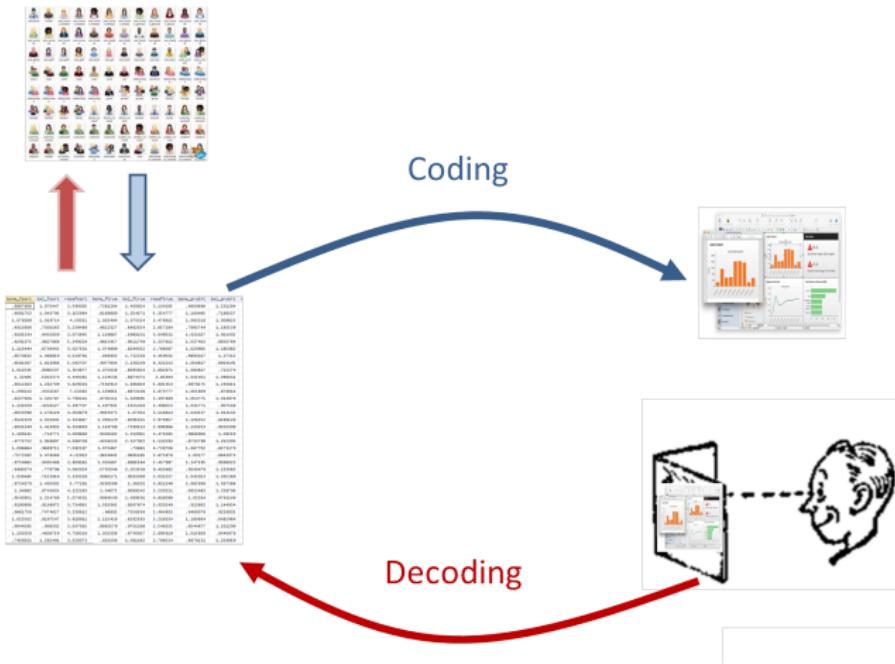
# THE BIG PICTURE : FROM "VIZ TO DATA"



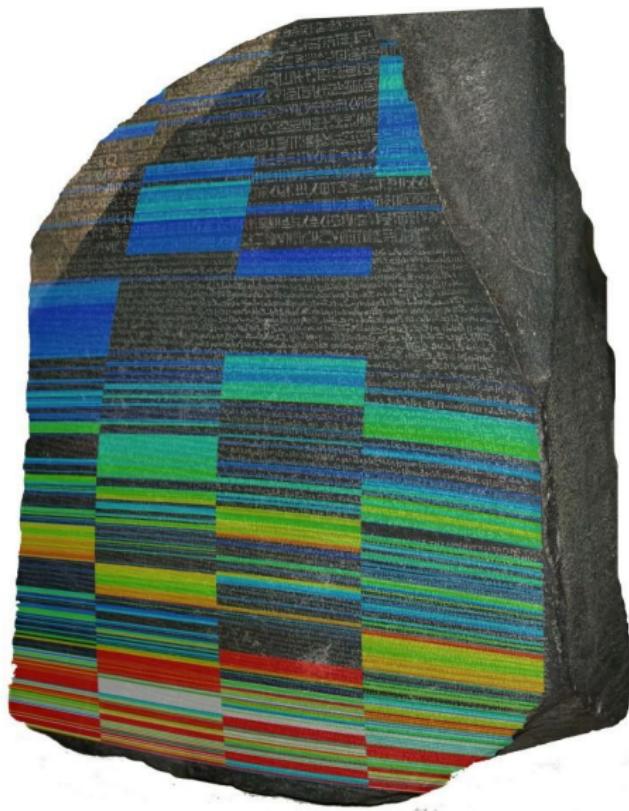
# THE BIG PICTURE : FROM "VIZ TO DATA"



# THE BIG PICTURE : FROM "VIZ TO DATA"



# **[ - DECODING VISUAL INFORMATION - ]**



# GRAPHICS IMMEDIATE TO UNDERSTAND



FIGURE – Where do people run in Paris

source : (N. Yau)

<http://flowingdata.com/2014/02/05/where-people-run/>

# GRAPHICS IMMEDIATE TO UNDERSTAND

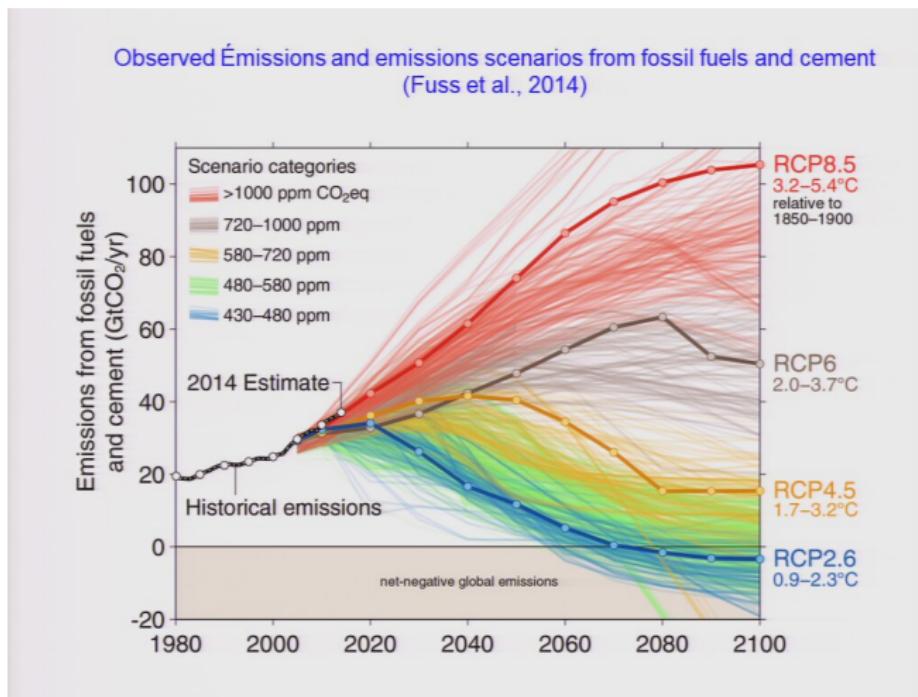


FIGURE – Climate forecast uncertainty (S. Planton)

# GRAPHICS THAT NEED EXPLANATIONS :

## Everyone

Sleeping, eating, working and watching television take up about two-thirds of the average day.

Everyone	Employed	White	Age 15-24	H. S. grads	No children
Men	Unemployed	Black	Age 25-64	Bachelor's	One child
Women	Not in lab...	Hispanic	Age 65+	Advanced	Two+ children

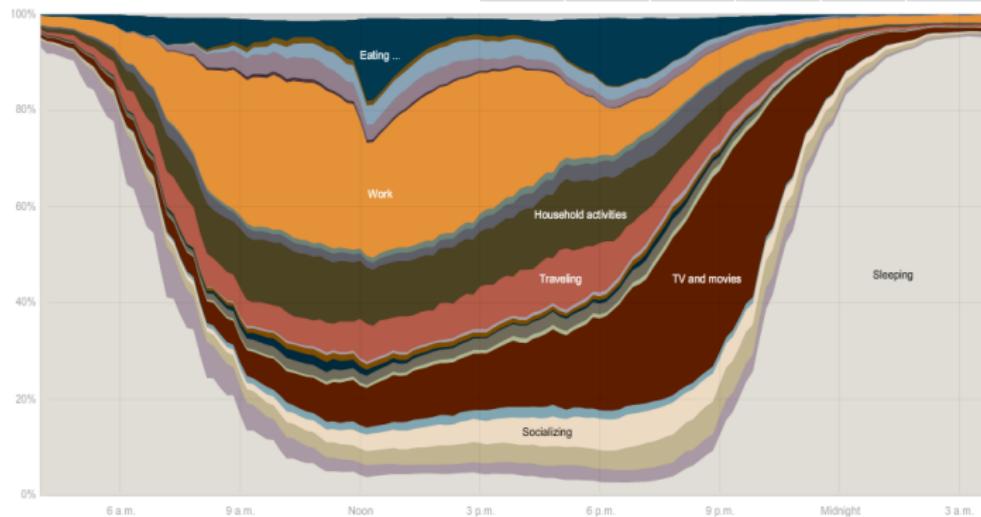


FIGURE – How people spend their days (NYT).

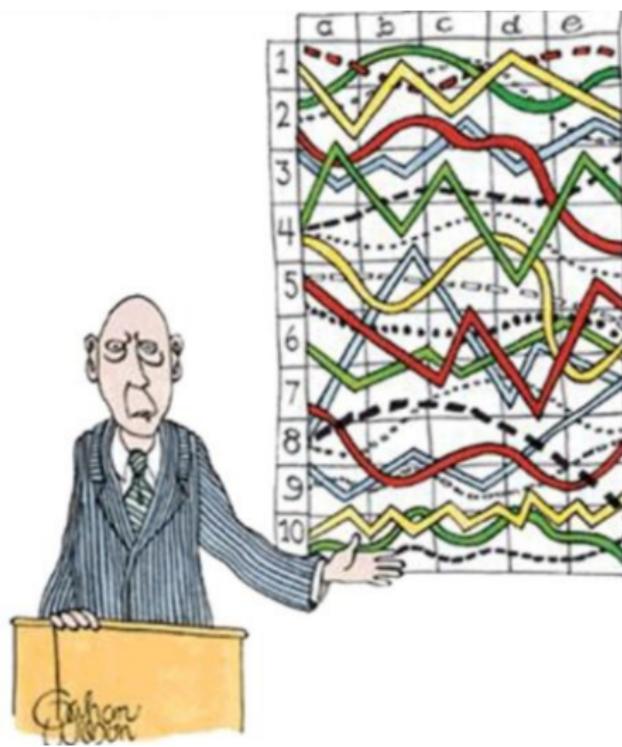
# SO WHAT ARE THE RULES ?

- ▶ Can you name some rules for a good graphic ?

# SO WHAT ARE THE RULES ?

- ▶ Can you name some rules for a good graphic ?
- ▶ Or maybe, rules for a "bad" graphic first..

# [- MISTAKES ! -]



**"I'll pause  
for a moment  
so you can let  
this information  
sink in."**

# USUAL MISTAKES

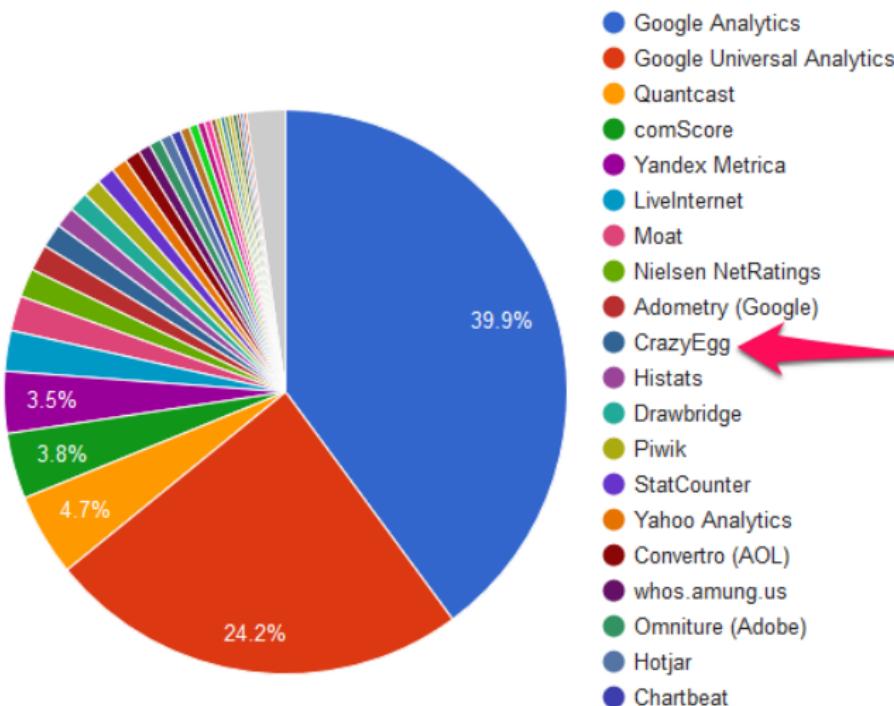


FIGURE – (source WTF Visualization)

# EVEN WORSE : 3D-PIECHART

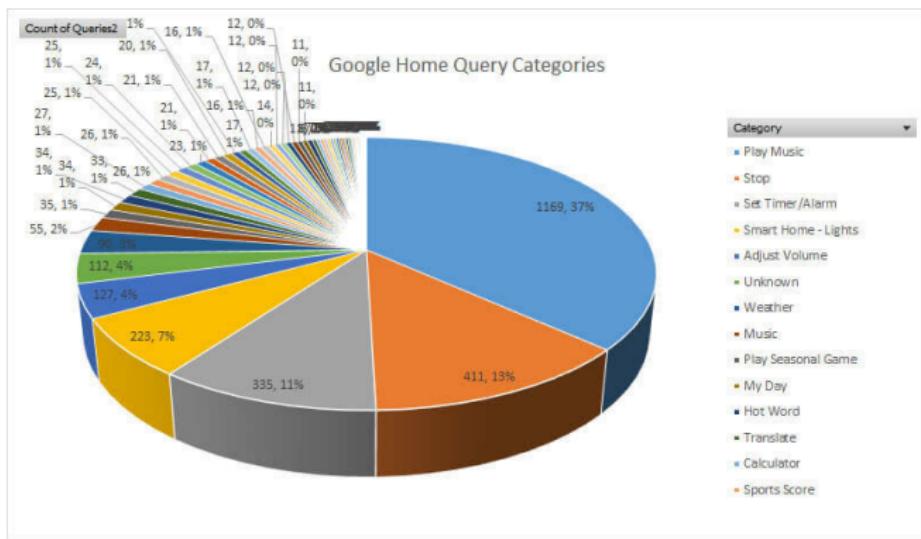


FIGURE – (source WTF Visualization)

# SOME PEOPLE USE IT !

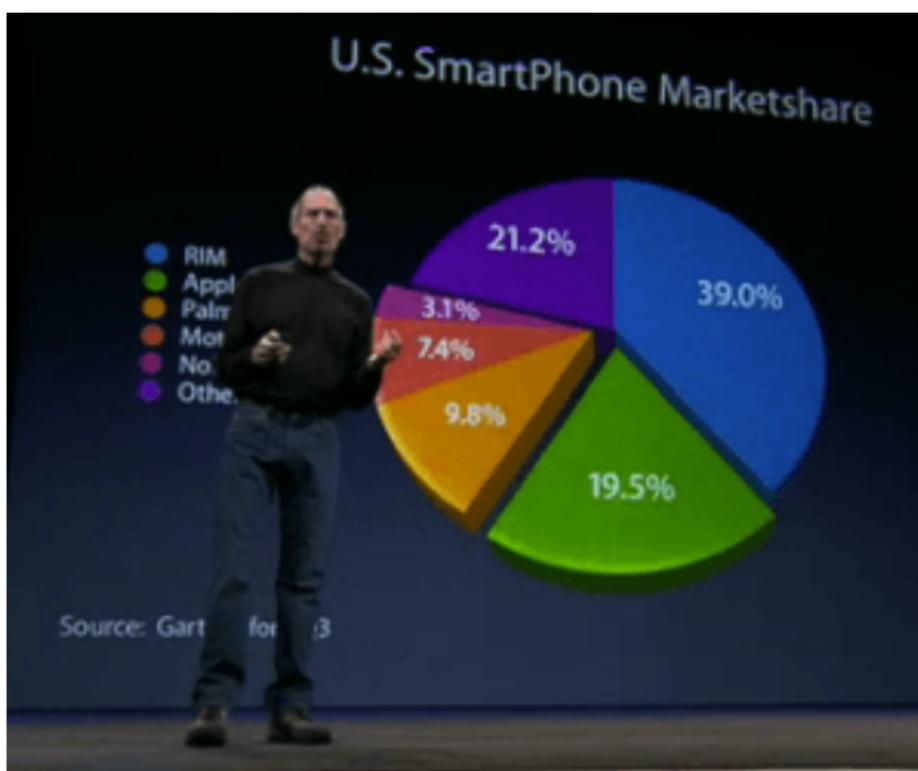


FIGURE – Source : 5 mistakes in Dataviz

# WHY IS IT A BAD PRACTICE ?

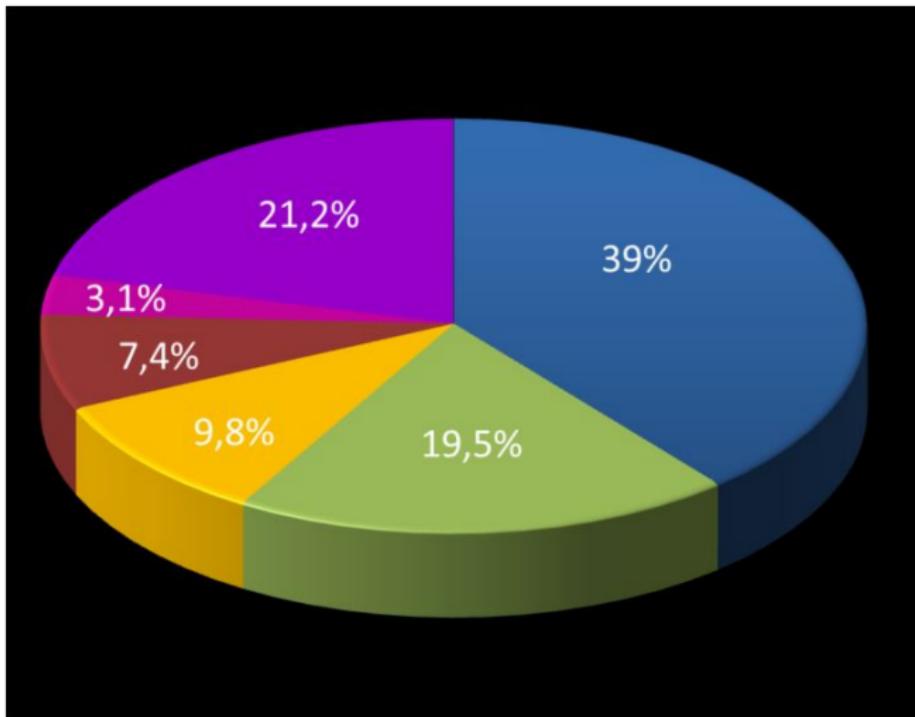


FIGURE – Source : reproducing Steve Job's dataviz (C. Bontemps)

# THE AREAS DO NOT REPRESENT THE DATA

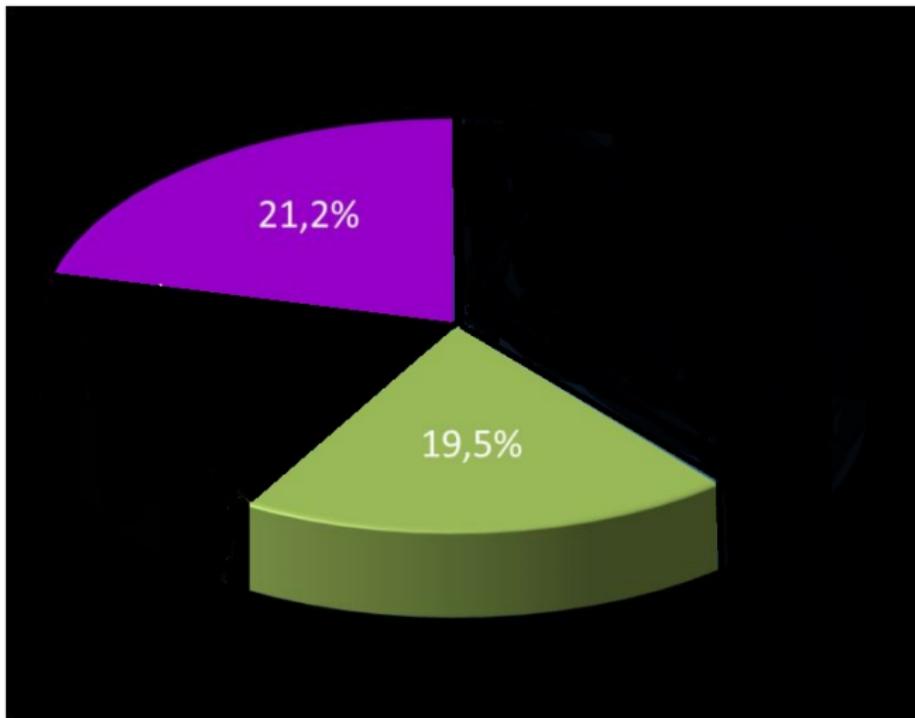


FIGURE – Source : reproducing Steve Job's dataviz (C. Bontemps)

TOP :  $20 \text{ cm}^2 = 21.2\%$ , BOTTOM :  $30 \text{ cm}^2 = 19.5\%$

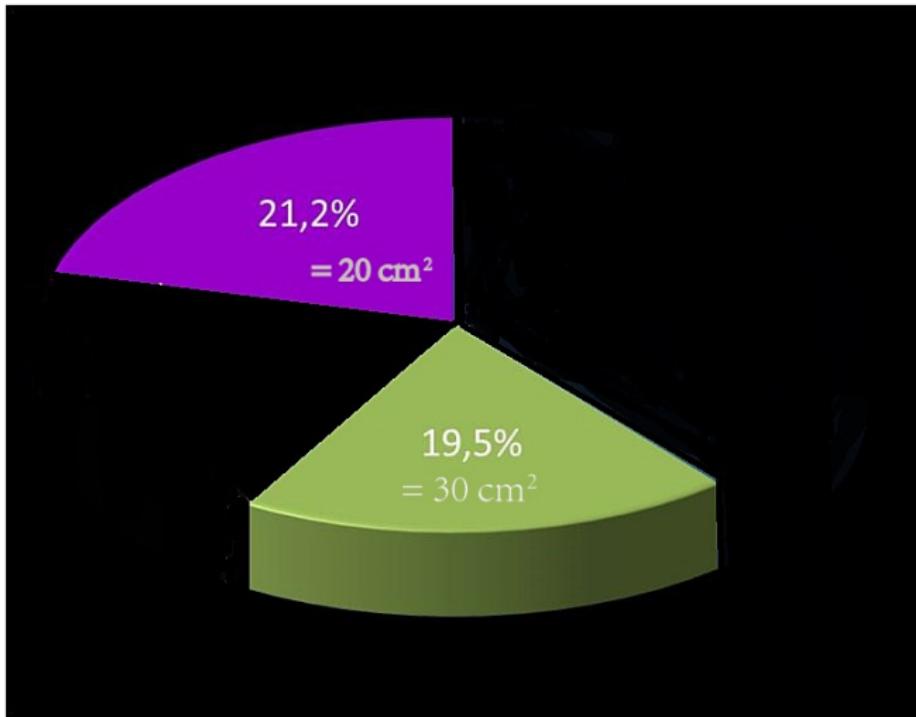


FIGURE – Source : reproducing Steve Job's dataviz (C. Bontemps)

# USUAL MISTAKES

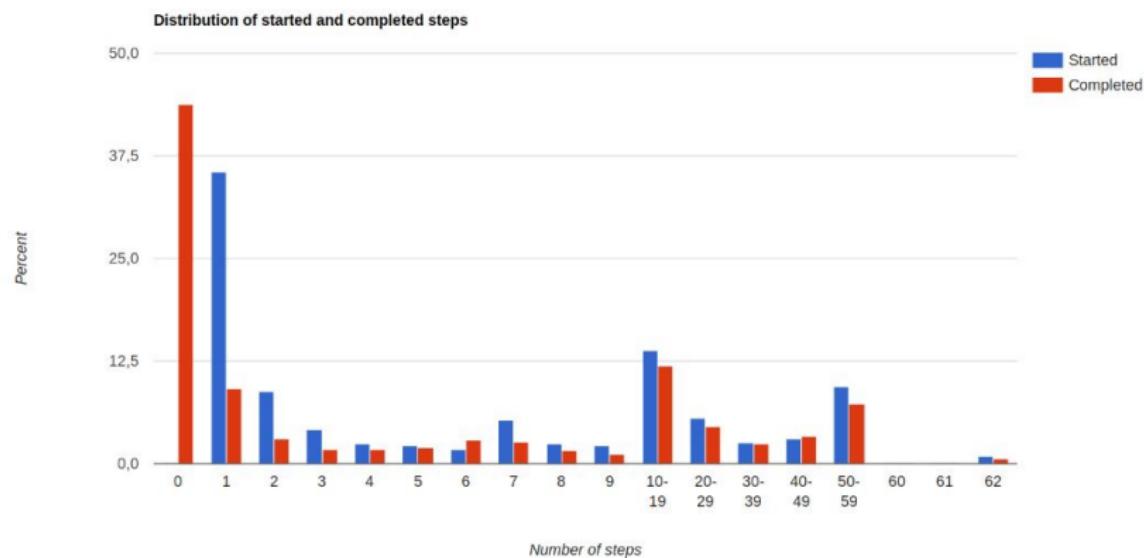


FIGURE – Percentage of student starting and completing a “step” in MOOC

Source : My 2017' students (Jan & Mohamed), but also recent researchers' presentations

# WHAT'S WRONG WITH THIS GRAPHIC ?

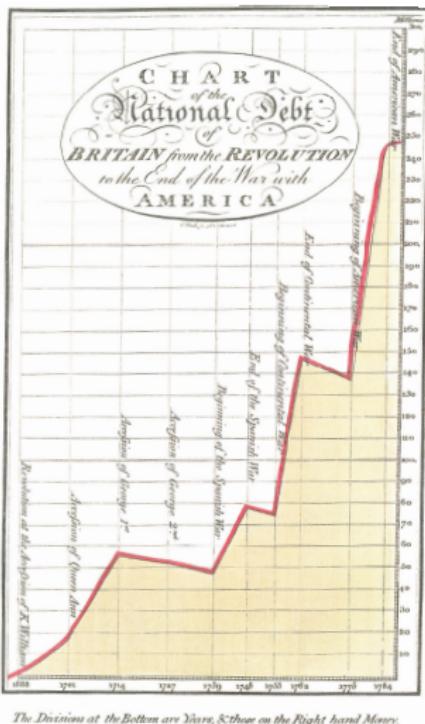


FIGURE – Government spending "Skyrocketing". Source : Tufte (2001) from Playfair(1786)

# SCALES ARE MISLEADING!

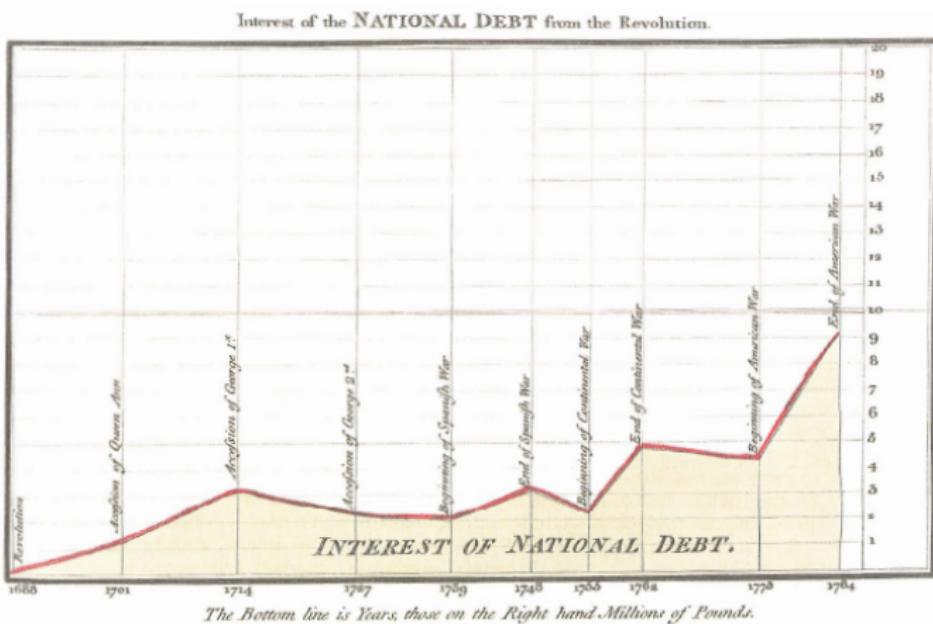


FIGURE – Government spending "Skyrocketing" (revisited)- Source : Tufte (2001) from Playfair(1786).

# NOWADAYS TOO :

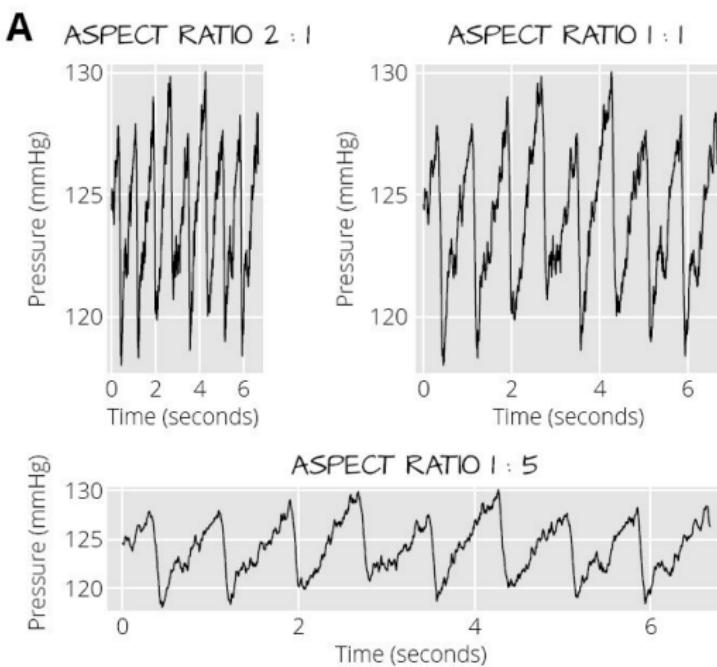


FIGURE – From Allen and Erhardt (2016b)

# NOWADAYS TOO :

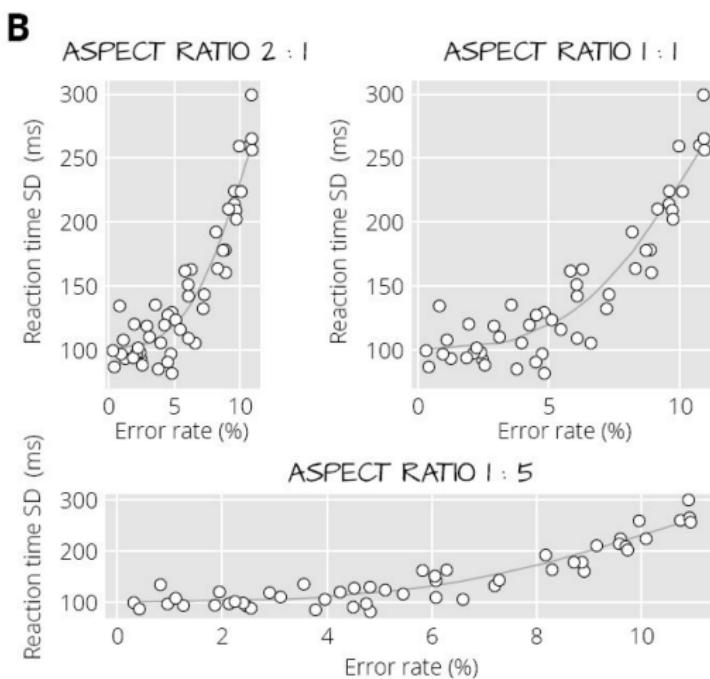
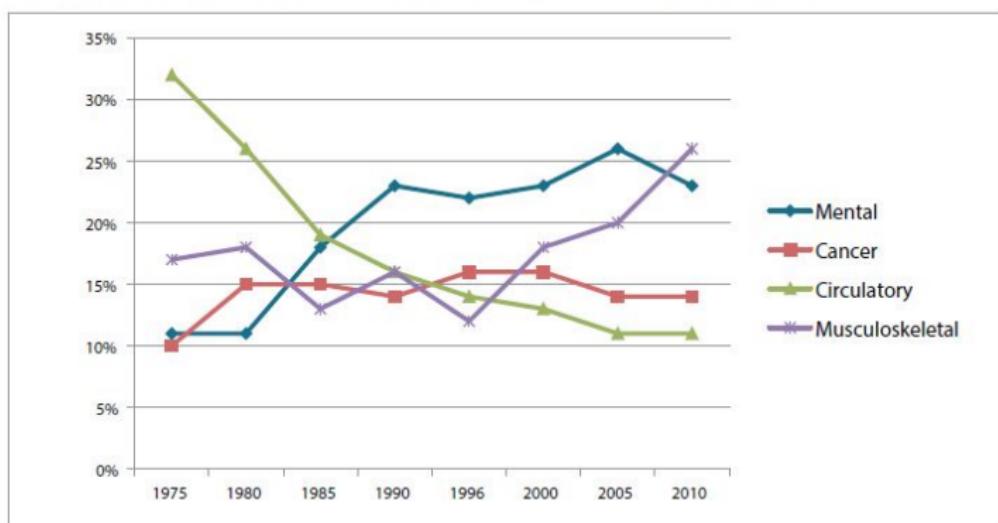


FIGURE – Cleveland (1994) recommends that curves follow a  $45^\circ$  angle. From Allen and Erhardt (2016b)

## WHAT'S WRONG WITH THIS GRAPHIC? (HARDER)

## 27. Initial DI Worker Awards by Major Cause of Disability—Calendar Years 1975–2010



*Source:* Social Security Advisory Board (2012).

FIGURE – Major Cause of Disability - 1975-2010 (J. Schwabish, 2014).

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)

Questions :

- in 2010, what is the major cause of disability ?

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)

Questions :

- ▶ in 2010, what is the major cause of disability ?
- ▶ in 1975, what was the major cause of disability ?

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)

Questions :

- ▶ in 2010, what is the major cause of disability ?
- ▶ in 1975, what was the major cause of disability ?
- ▶ **In the recent years, which causes have increased/decreased the most ?**

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)

Questions :

- ▶ in 2010, what is the major cause of disability ?
- ▶ in 1975, what was the major cause of disability ?
- ▶ In the recent years, which causes have increased/decreased the most ?
- ▶ ...

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)

Questions :

- ▶ in 2010, what is the major cause of disability ?
- ▶ in 1975, what was the major cause of disability ?
- ▶ In the recent years, which causes have increased/decreased the most ?
- ▶ ....
- ▶ You do not remember a damn thing of this graph !

# MORE (SMALLER) GRAPHS, ARE OFTEN BETTER

Initial DI Worker Awards by Major Cause of Disability—  
Calendar Years 1975–2010  
(Percent)

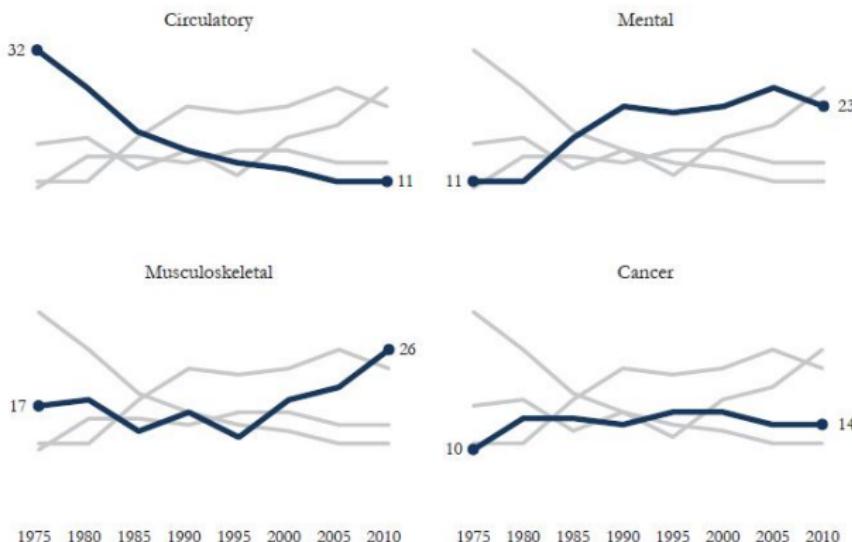
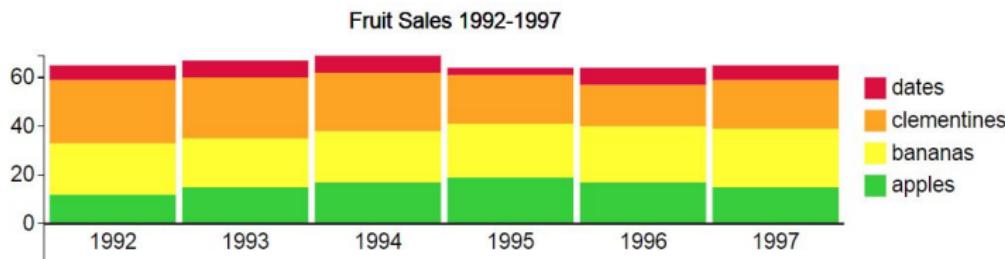


FIGURE – Major Cause of Disability- 1975-2010 (J. Schwabish).

Cf. "brushing" (ex : for parallel Coordinates plots)

# YOUR TURN : WHAT'S WRONG WITH THIS GRAPHIC ?



# "SMALL MULTIPLES" PROVIDE A BETTER VIEW

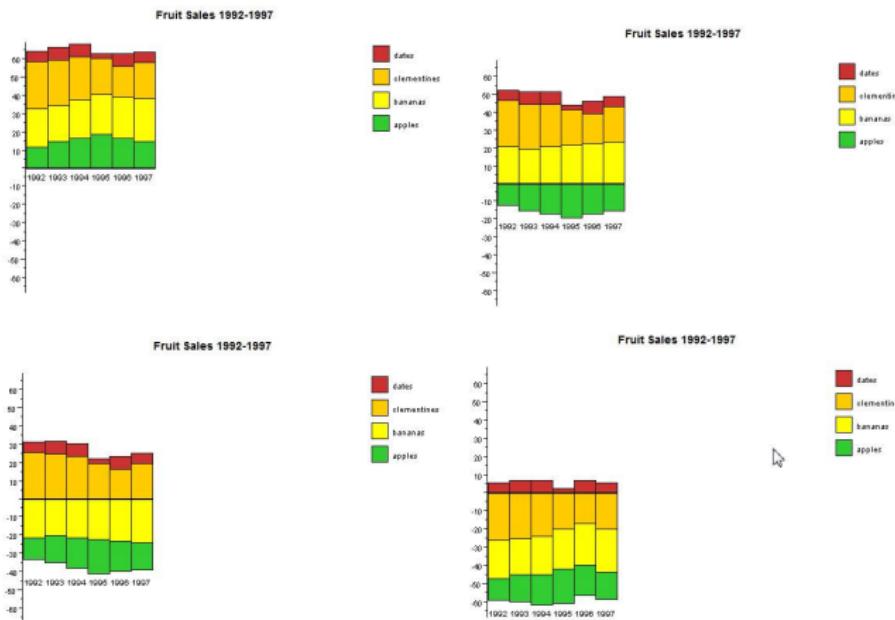


FIGURE – from Dix and Ellis (1998) example

# AND HERE COMES DYNAMIC DATA VISUALISATION !

Happily Jane is using an **interactive stacked histogram**.

Try it yourself.

Click on the coloured banana entry in the key, or one of the banana parts of any of the histogram bars.

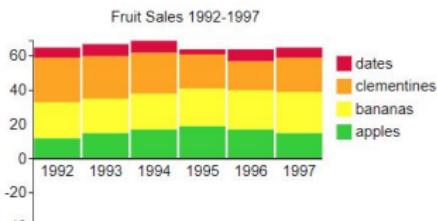
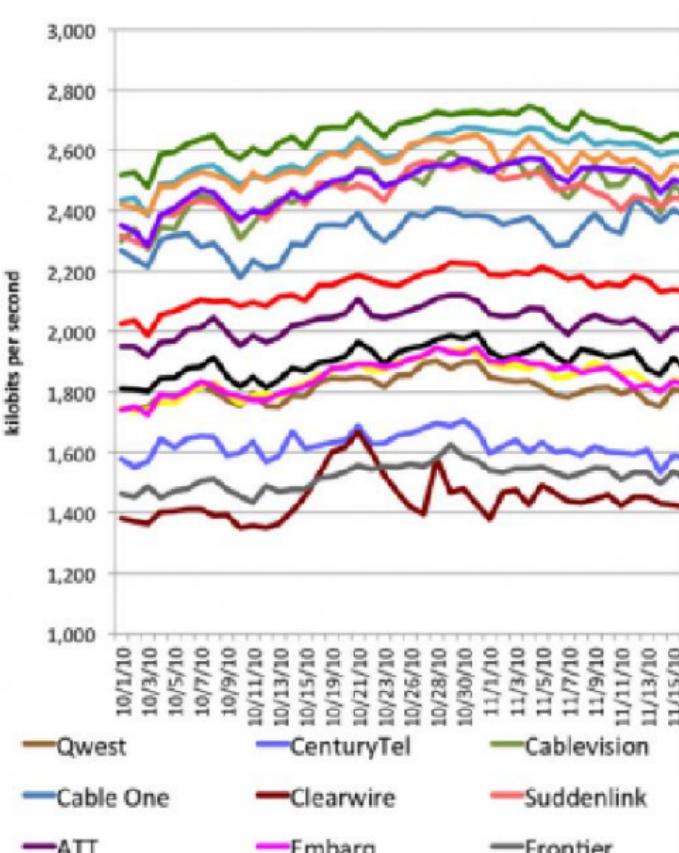


FIGURE – Dynamic fruit sales

from meandeviation.com

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)



# KEEP ALL YOUR AUDIENCE

Normal



Color-blind

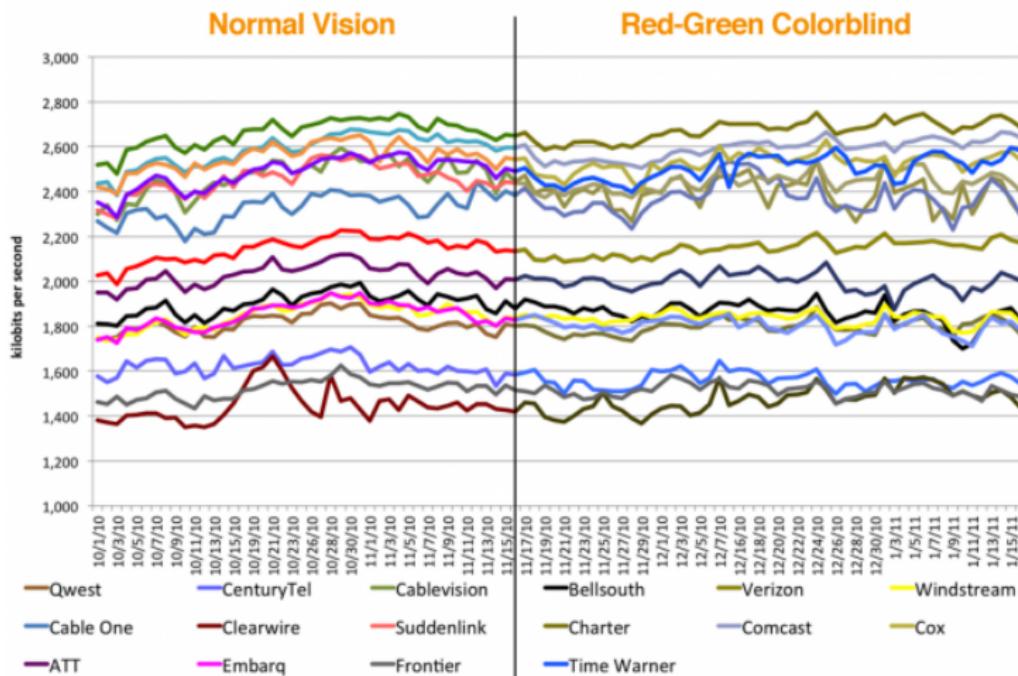


5%



0.35%

# WHICH MEANS THAT FOR 5 % OF MEN :



See also the ggplot option + `scale_colour_colorblind()`



# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...

# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...
- ▶ Axis and scale (starting at zero !)

# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...
- ▶ Axis and scale (starting at zero !)
- ▶ Context, labels

# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...
- ▶ Axis and scale (starting at zero !)
- ▶ Context, labels
- ▶ No multiple scales !

# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...
- ▶ Axis and scale (starting at zero !)
- ▶ Context, labels
- ▶ No multiple scales !
- ▶ Colors that are distinguishable

# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...
- ▶ Axis and scale (starting at zero !)
- ▶ Context, labels
- ▶ No multiple scales !
- ▶ Colors that are distinguishable
- ▶ A good message ! Unambiguous message

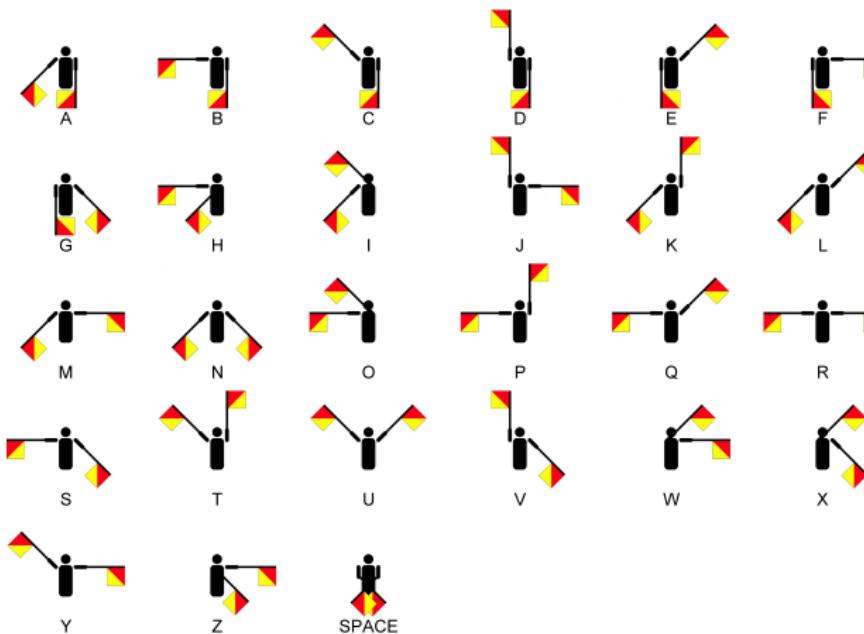
# SO WHAT ARE THE RULES ?

Can you name some rules for a good (*resp.* bad) graphic ?

- ▶ Cite some rules...
- ▶ Axis and scale (starting at zero !)
- ▶ Context, labels
- ▶ No multiple scales !
- ▶ Colors that are distinguishable
- ▶ A good message ! Unambiguous message
- ▶ Should not be overloaded

# [- WHY CODING ? -]

If decoding is hard, why coding?



## GRAPHICS reveal DATA : ANSCOMBE (1973) QUARTET

We use here 4 couples of random variables :  $(X_1, Y_1)$ ,  $(X_2, Y_2)$ ,  $(X_3, Y_3)$  and  $(X_4, Y_4)$ . All four data sets have the same descriptive statistics.

Xs	Mean	Std. Dev.	Ys	Mean	Std. Dev.	$corr(X_i, Y_i)$	N
$X_1$	9	3.32	$Y_1$	7.5	2.03	0.8164	11
$X_2$	9	3.32	$Y_2$	7.5	2.03	0.8162	11
$X_3$	9	3.32	$Y_3$	7.5	2.03	0.8163	11
$X_4$	9	3.32	$Y_4$	7.5	2.03	0.8165	11

# ANSCOMBE (1973) QUARTET

All four data sets are described by the same linear model ( $Y_i = \alpha + \beta X_i + \epsilon_i$ ), revealing apparently the same relationships :

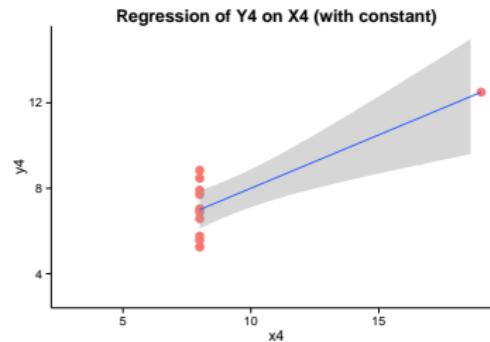
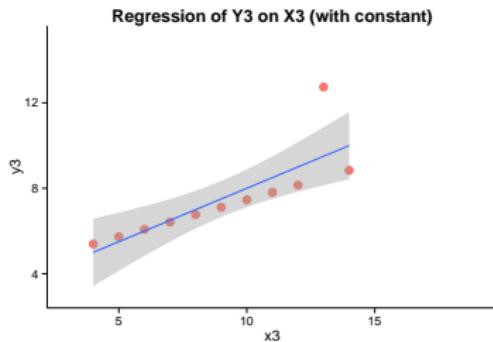
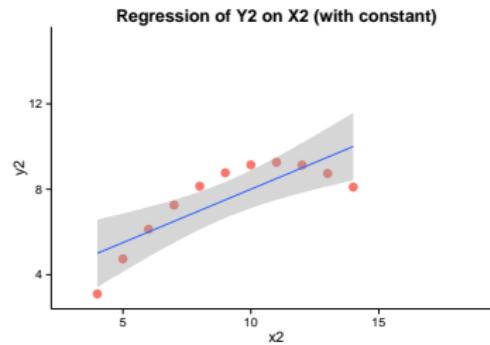
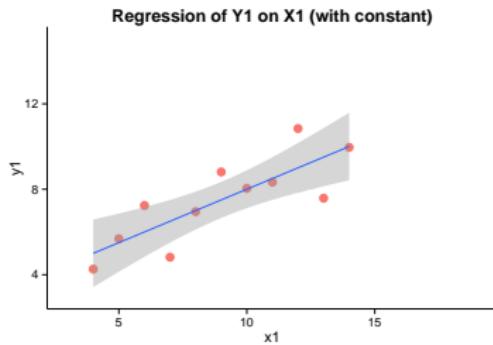
*Dependent variable :*

	$Y_1$	$Y_2$	$Y_3$	$Y_4$
Regressed on :				
$X_i, i=1,\dots,4$	0.500 ***	0.500 ***	0.500 ***	0.500 ***
Constant	3.000 **	3.001 **	3.002 **	3.002 **
$R^2$	0.667	0.666	0.666	0.667
Resid Std. Error	1.237	1.237	1.236	1.236
F Statistic	17.990***	17.966***	17.972***	18.003***

*Note : Data from Anscombe (1973). \* p < 0.1 ; \*\* p < 0.05 ; \*\*\* p < 0.01*

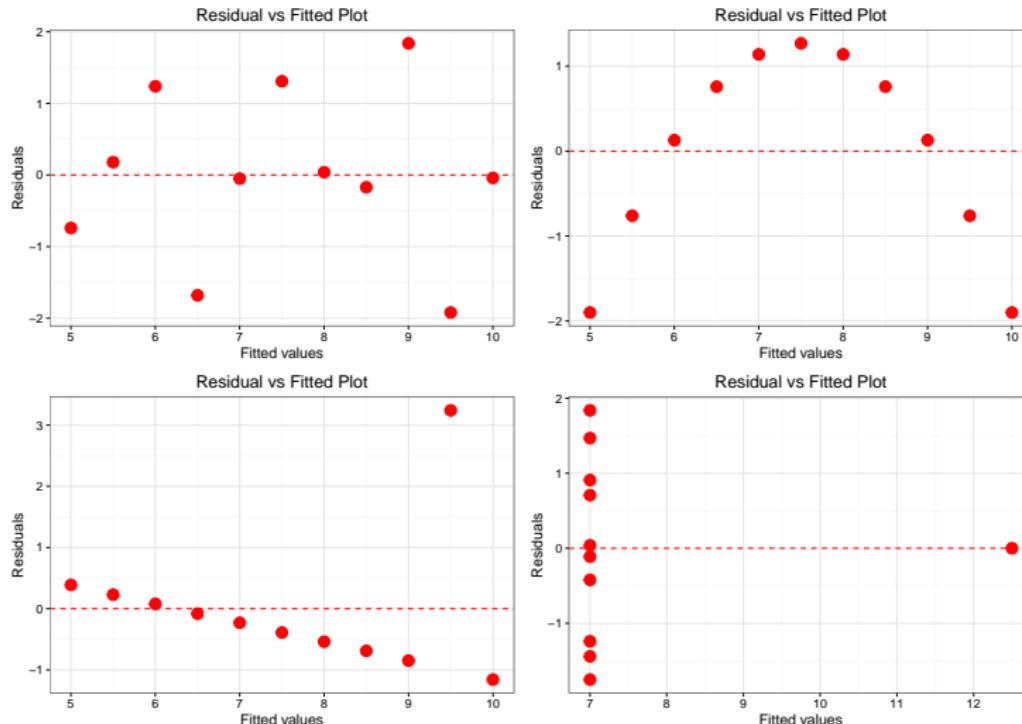
# ANSCOMBE (1973) QUARTET

A simple scatter plot (regression overlaid) shows something very different.

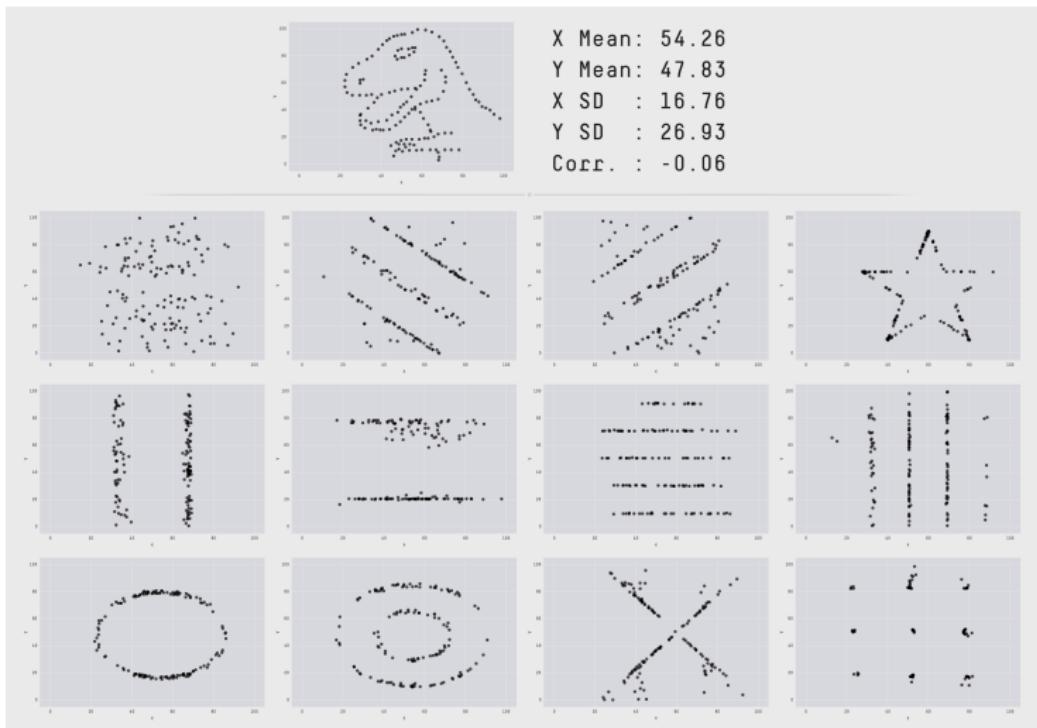


# ANSCOMBE (1973) QUARTET

NP : Plots of the residuals shows also same differences



## THE DATASAURUS...



Matejka and Fitzmaurice (2017), see also The Datasaurus website

# TABLES VS GRAPHICS

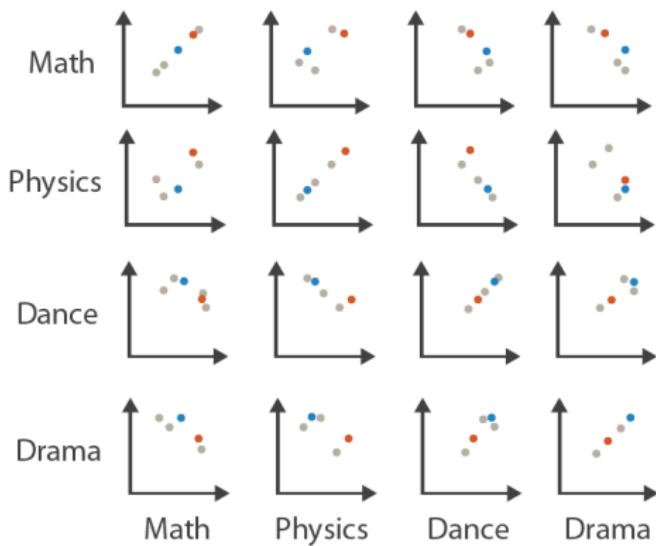
## Table

Math	Physics	Dance	Drama
85	95	70	65
90	80	60	50
65	50	90	90
50	40	95	80
40	60	80	90

From Munzner (2014)

## TABLES VS GRAPHICS

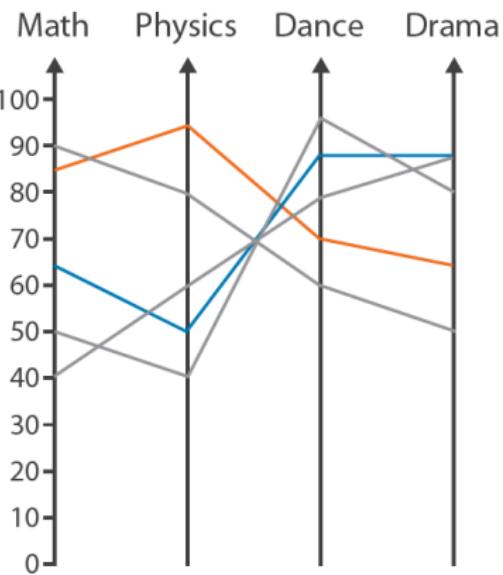
## Scatterplot Matrix



From Munzner (2014)

# TABLES VS GRAPHICS

Parallel Coordinates



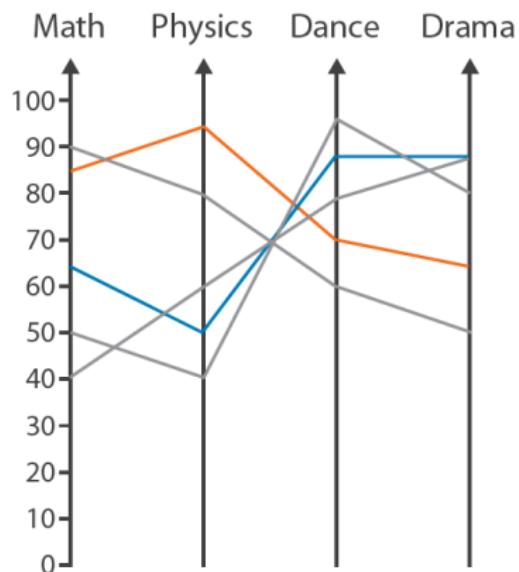
From Munzner (2014)

# TABLES VS GRAPHICS

Table

	Math	Physics	Dance	Drama
	85	95	70	65
	90	80	60	50
	65	50	90	90
	50	40	95	80
	40	60	80	90

Parallel Coordinates



From Munzner (2014)

## TABLES VS GRAPHICS : CASE STUDY

Data with many 0/1 variables :  
Facilities indicators for 16 towns (A-P)

Cities																	
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	#	Facilities
0	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0	1	School
0	1	1	1	0	0	1	0	0	0	0	1	0	0	1	0	2	Agricultural Coop
0	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0	3	Station
1	1	0	0	1	1	0	0	1	1	0	0	1	1	0	1	4	School (one class)
0	1	1	1	0	0	1	0	0	0	0	1	0	0	1	0	5	Vet
1	1	0	0	1	1	0	0	1	1	0	0	1	1	0	1	6	No doctor
0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	7	No running water
0	0	1	0	0	0	0	1	0	0	1	0	0	0	0	0	8	Police station
0	1	1	1	0	0	1	0	0	0	0	1	0	0	1	0	9	Land consolidation

from Bertin (1981)

# TABLES VS GRAPHICS :

"Ordering is visualising"



Using Bertin (1981) principles, The AVIZ lab has created the *Bertifier*. See also the R package Hahsler et al. (2008).

## TABLES VS GRAPHICS :

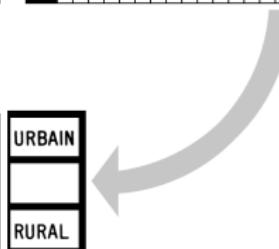
Data with many 0/1 variables :  
Facilities indicators for 16 towns (A-P)

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P		
																1	COLLÈGE
																2	COOPÉRATIVE AGRIC.
																3	GARE
																4	ÉCOLE CLASSE UNIQUE
																5	VÉTÉRINAIRE
																6	PAS DE MÉDECIN
																7	PAS D'ADDUCTION D'EAU
																8	GENDARMERIE
																9	REMENBREMENT

Bertin (1981)

## TABLES VS GRAPHICS

## LA MATRICE ORDONNABLE



# TABLES VS GRAPHICS



Bertin (1981)

# TABLES VS GRAPHICS

ordering is the key

#		→	z	a	e	u	—	M	p	b	d	g	l	o	c	x	h
1	HIGH SCHOOL																
3	RAILWAY STATION																
8	POLICE STATION																
2	AGRICULTURAL COOPERATIVE																
5	VETERINARY																
9	LAND REALLOCATION																
4	ONE ROOM SCHOOL																
6	NO DOCTOR																
7	NO WATER SUPPLY																

Using Bertin (1981) principles, The AVIZ lab has created the *Bertifier*. See also the R package Hahsler et al. (2008).

# GOOD GRAPHICS ?

It the excellent Handbook of data visualisation Chen et al. (2007), we find some good questions :

- ▶ What to Whom, How and Why ?

*A graphic may be linked to three pieces of text : its **caption**, a **headline** and an **article** it accompanies. Ideally, all three should be consistent and complement each other.*

# GOOD GRAPHICS ?

It the excellent Handbook of data visualisation Chen et al. (2007), we find some good questions :

- ▶ What to Whom, How and Why ?

*A graphic may be linked to three pieces of text : its **caption**, a **headline** and **an article** it accompanies. Ideally, all three should be consistent and complement each other.*

- ▶ Present or explore data ?

*Different purpose, different requirements !*

# GOOD GRAPHICS ?

It the excellent Handbook of data visualisation Chen et al. (2007), we find some good questions :

- ▶ What to Whom, How and Why ?

*A graphic may be linked to three pieces of text : its **caption**, a **headline** and **an article** it accompanies. Ideally, all three should be consistent and complement each other.*

- ▶ Present or explore data ?

*Different purpose, different requirements !*

- ▶ Choice of Graphical form ?

*Choice depends on the type of data to be displayed (e.g. univariate continuous data, bivariate categorical data, etc..) and on what is to be shown.*

# GOOD GRAPHICS ?

It the excellent Handbook of data visualisation Chen et al. (2007), we find some good questions :

- ▶ What to Whom, How and Why ?

*A graphic may be linked to three pieces of text : its **caption**, a **headline** and **an article** it accompanies. Ideally, all three should be consistent and complement each other.*

- ▶ Present or explore data ?

*Different purpose, different requirements !*

- ▶ Choice of Graphical form ?

*Choice depends on the type of data to be displayed (e.g. univariate continuous data, bivariate categorical data, etc..) and on what is to be shown.*

- ▶ Unique solution ?

*There is **not** always a unique optimal choice and alternatives can be equally good or good in different ways, emphasizing different aspects of the same data.*

# “GOOD” OR “BAD” GRAPHICS?

*“There are no “good” nor “bad” graphics (...), there are graphics answering legitimate questions and graphics that do not answer question at all ”*

Bertin (1981)

It is easy to criticize ... but are there some rules ?

# WHAT'S WRONG WITH THIS GRAPHIC ?

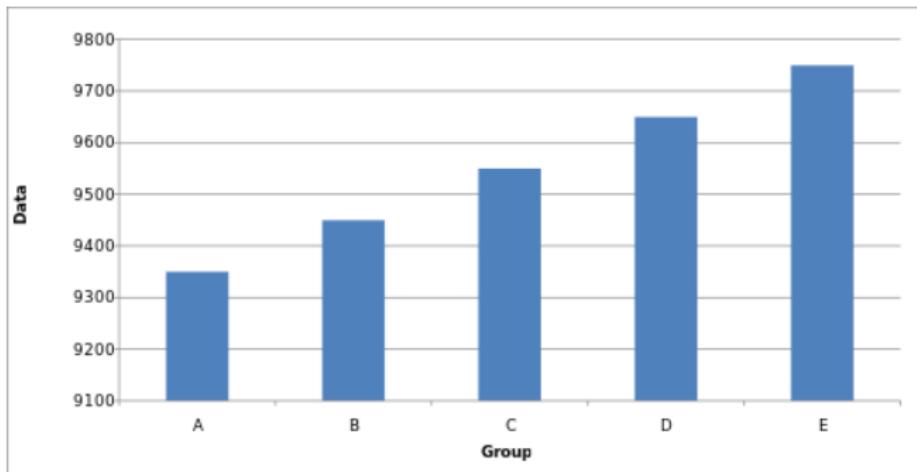


FIGURE – (source Wikipedia)

# WHAT'S WRONG WITH THIS GRAPHIC ?

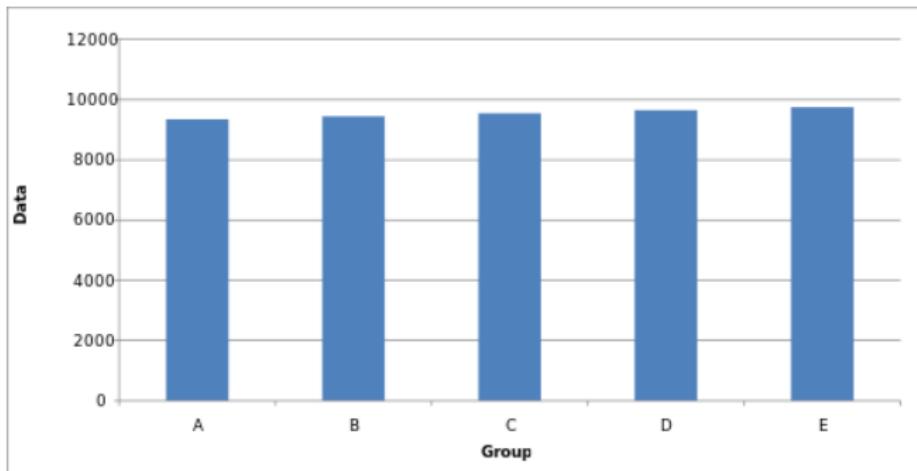


FIGURE – (source Wikipedia)

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)



FIGURE – Are you looking at the right thing ?

# WHAT'S WRONG WITH THIS GRAPHIC ? (HARDER)

## LIMITED SCOPE

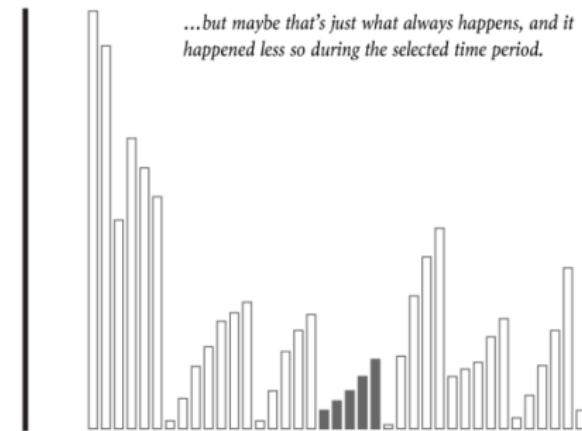
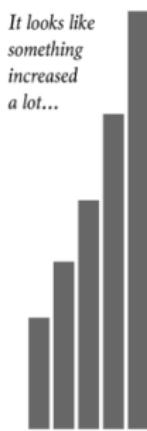


FIGURE – Are you looking at the right thing ?

from Flowing data

# EDWARD R. TUFTE'S RULES

In his seminal book, Tufte (2001) propose some principles for displaying quantitative information.

**Data :** *Above all, show the data*

# EDWARD R. TUFTE'S RULES

In his seminal book, Tufte (2001) propose some principles for displaying quantitative information.

**Data :** *Above all, show the data*

**Question :** *Induce the viewer to think about the substance rather than about methodology, graphic design. Encourage the eye to compare different piece of data.*

# EDWARD R. TUFTE'S RULES

In his seminal book, Tufte (2001) propose some principles for displaying quantitative information.

**Data :** *Above all, show the data*

**Question :** *Induce the viewer to think about the substance rather than about methodology, graphic design. Encourage the eye to compare different piece of data.*

**Data-ink ratio :** *Maximize the ink-data ratio. Erase all non data ink, Erase redundant information*

# EDWARD R. TUFTE'S RULES

In his seminal book, Tufte (2001) propose some principles for displaying quantitative information.

**Data :** *Above all, show the data*

**Question :** *Induce the viewer to think about the substance rather than about methodology, graphic design. Encourage the eye to compare different piece of data.*

**Data-ink ratio :** *Maximize the ink-data ratio. Erase all non data ink, Erase redundant information*

**Integrity :** *Avoid distorting what the data have to say*

# EDWARD R. TUFTE'S RULES

In his seminal book, Tufte (2001) propose some principles for displaying quantitative information.

**Data :** *Above all, show the data*

**Question :** *Induce the viewer to think about the substance rather than about methodology, graphic design. Encourage the eye to compare different piece of data.*

**Data-ink ratio :** *Maximize the ink-data ratio. Erase all non data ink, Erase redundant information*

**Integrity :** *Avoid distorting what the data have to say*

**General to specific :** *Reveal the data at different levels of detail (from broad picture to fine structure)*

# EDWARD R. TUFTE'S RULES

In his seminal book, Tufte (2001) propose some principles for displaying quantitative information.

**Data :** *Above all, show the data*

**Question :** *Induce the viewer to think about the substance rather than about methodology, graphic design. Encourage the eye to compare different piece of data.*

**Data-ink ratio :** *Maximize the ink-data ratio. Erase all non data ink, Erase redundant information*

**Integrity :** *Avoid distorting what the data have to say*

**General to specific :** *Reveal the data at different levels of detail (from broad picture to fine structure)*

**Context :** *Graphical display should be closely integrated with the statistical and verbal descriptions of the data set.*

# PRACTICAL EXAMPLE : DATA-INK RATIO

Let's start with a classical graph (R default - Boxplot )

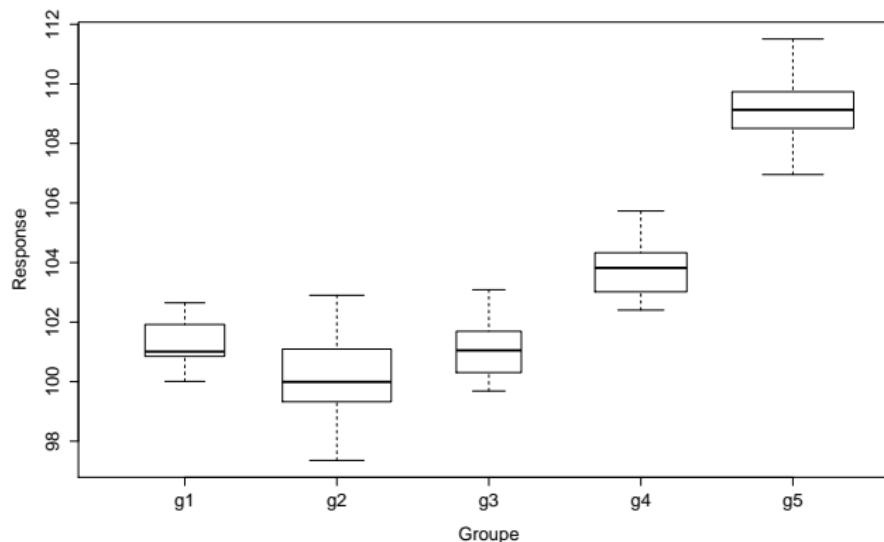


FIGURE – Distribution of a continuous variable on 4 groups

# ERASE ALL NON DATA INK

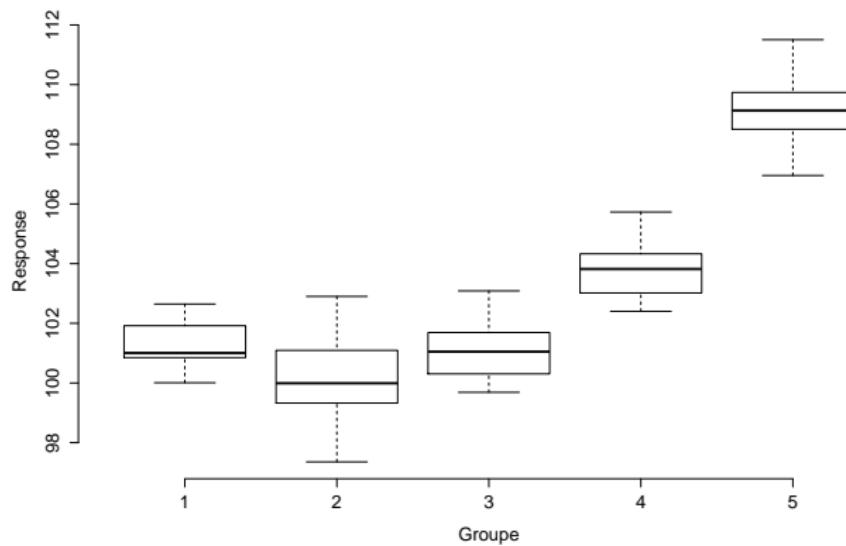


FIGURE – Distribution of a continuous variable on 4 groups

# ERASE ALL REDUNDANT !

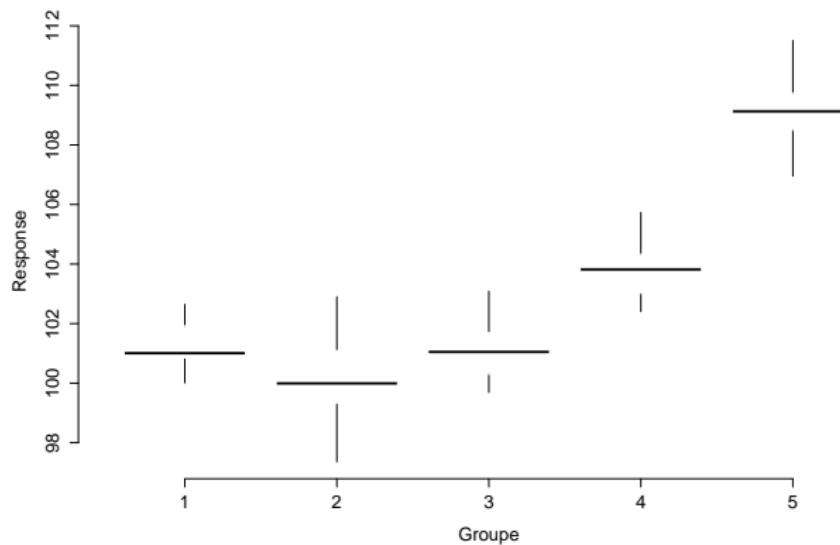


FIGURE – Distribution of a continuous variable on 4 groups

# GOING FURTHER...

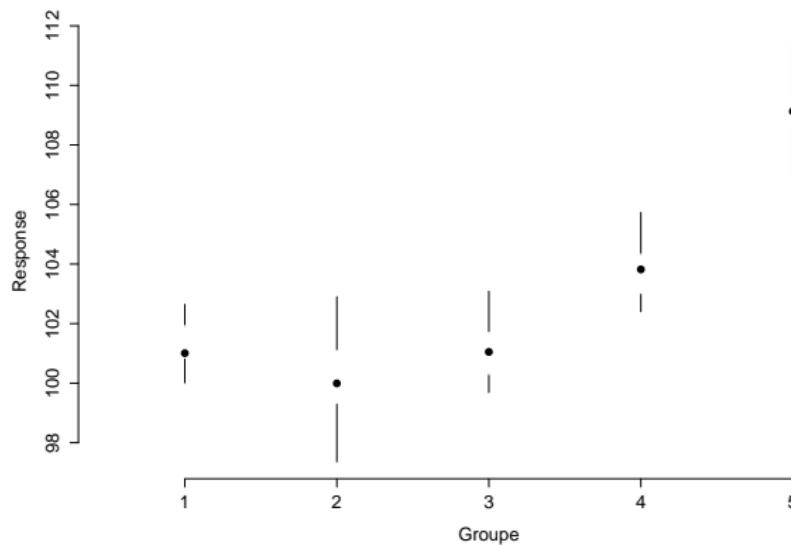


FIGURE – Distribution of a continuous variable on 4 groups

AND SHOW THE DATA...

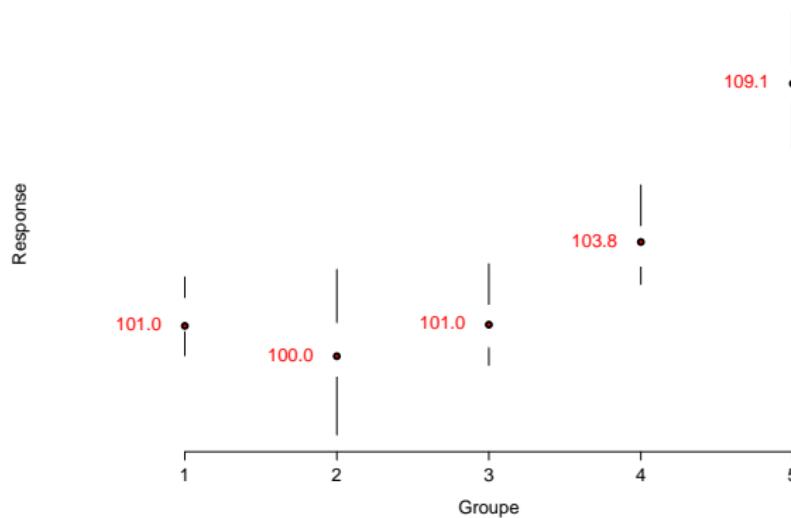


FIGURE – Distribution of a continuous variable on 4 groups

# HAVE WE LOST SOMETHING ?

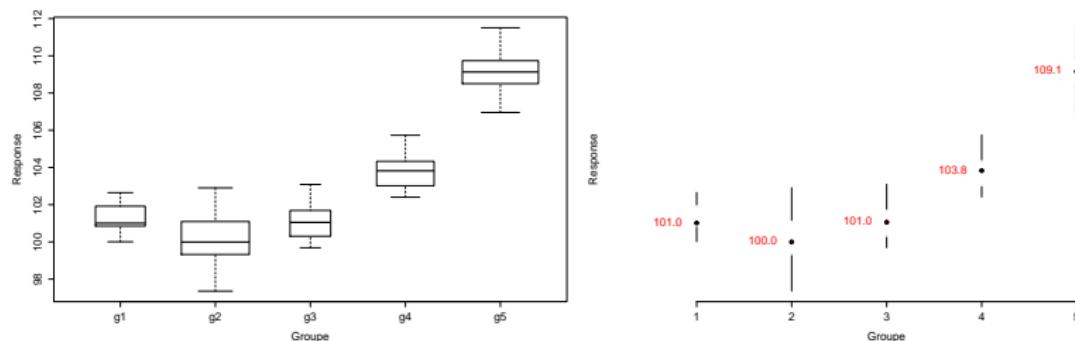


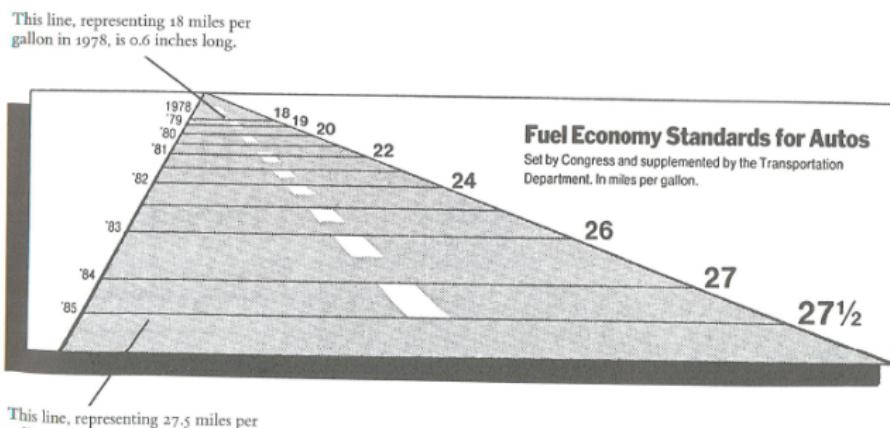
FIGURE – Distribution of a continuous variable on 4 groups

Did you noticed that group 1 and group 3 had the same median (101.0) ? see the `ggplot` theme + `theme_tufte()`

# INTEGRITY : THE LIE FACTOR

$$\text{LieFactor} = \frac{\text{Size of effect shown in graphic}}{\text{Size of effect in data}} \quad (1)$$

A Lie Factor  $\neq 1$  indicates a substantial distortion



New York Times, August 9, 1978, D-2.

FIGURE – Fuel economy standards. (E. Tufte - from NY Times 1978) ↗

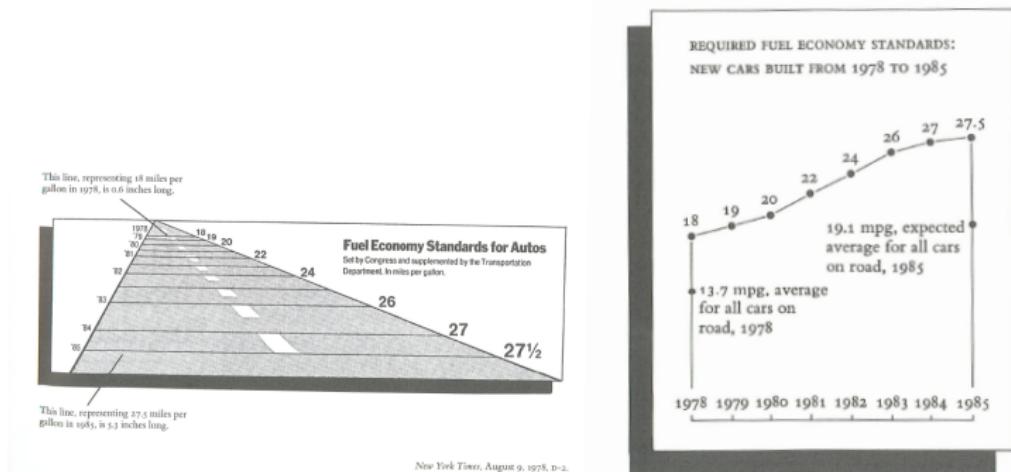


FIGURE – Fuel economy standards (revisited)

The "18 mpg" line measures 1.5 cm (in 1978); the "27.5 mpg" measures 13 cm (in 1985)  
 → Lie factor = 14.5%!!!

## PRACTICAL EXAMPLE : CONTEXT

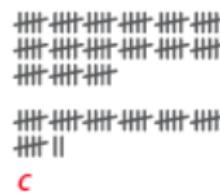
In the next 5 minutes, you will design “**the best way**” to **compare** two numbers

**75** and **37**.

From S. Ortiz

# CONTEXT MATTERS !

75, 37        
**a**                **b**



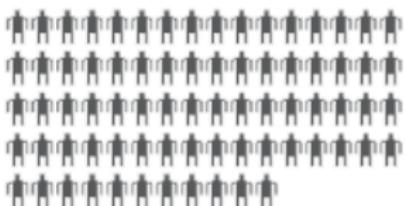
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



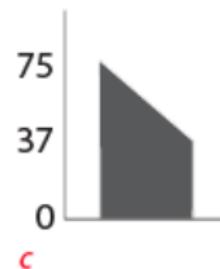
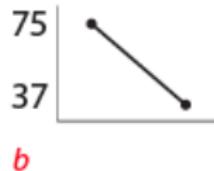
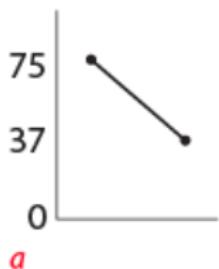
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !

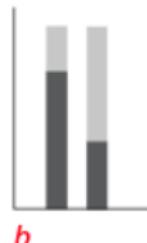


*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*a*



*b*

*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



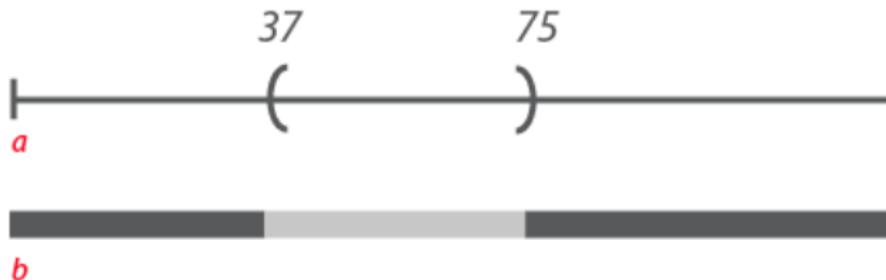
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



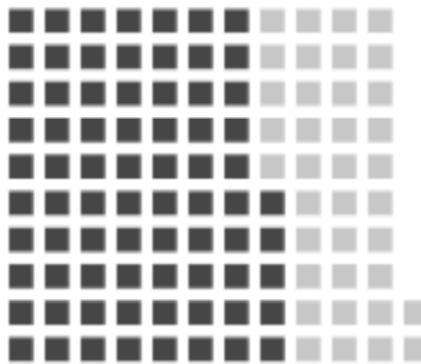
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



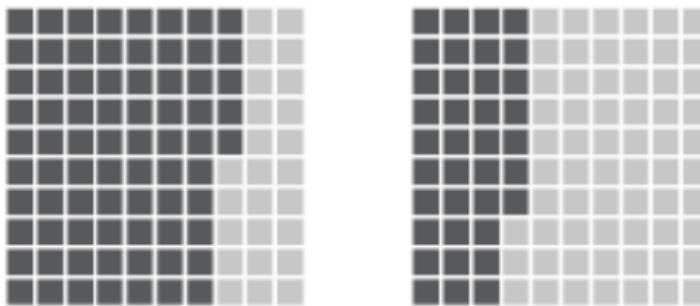
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

## CONTEXT MATTERS!



Solutions proposed by S. Ortiz

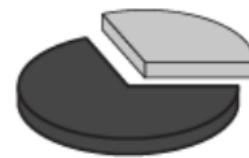
# CONTEXT MATTERS !



*a*



*b*



*c*

*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



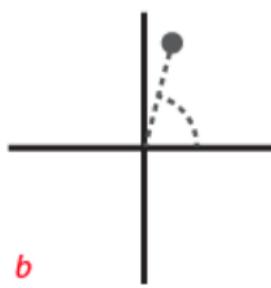
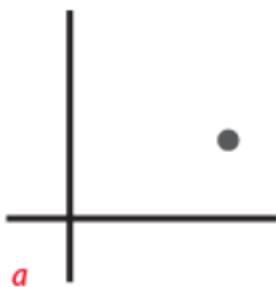
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*a*

*b*

*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



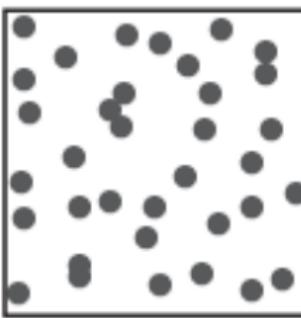
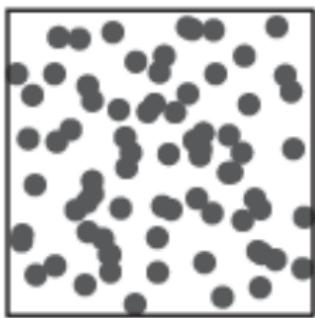
*a*



*b*

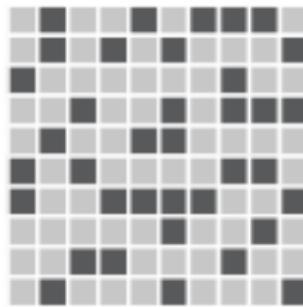
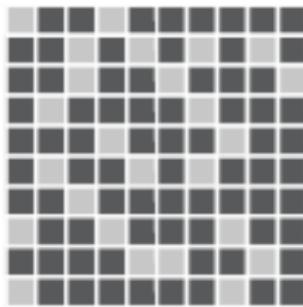
*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



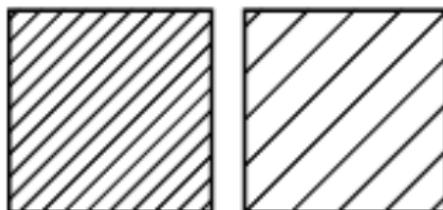
*Solutions proposed by S. Ortiz*

## CONTEXT MATTERS!



Solutions proposed by S. Ortiz

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS !



*Solutions proposed by S. Ortiz*

# CONTEXT MATTERS (FOLLOW-UP)

What are the problems you've encountered ?

# CONTEXT MATTERS (FOLLOW-UP)

What are the problems you've encountered ?

- ▶ Context?

# CONTEXT MATTERS (FOLLOW-UP)

What are the problems you've encountered ?

- ▶ Context ?
- ▶ Audience ?

# CONTEXT MATTERS (FOLLOW-UP)

What are the problems you've encountered ?

- ▶ Context ?
- ▶ Audience ?
- ▶ What to compare ?

# CONTEXT MATTERS (FOLLOW-UP)

What are the problems you've encountered ?

- ▶ Context ?
- ▶ Audience ?
- ▶ What to compare ?
- ▶ Units ?

# CONTEXT MATTERS (FOLLOW-UP)

What are the problems you've encountered ?

- ▶ Context ?
- ▶ Audience ?
- ▶ What to compare ?
- ▶ Units ?
- ▶ Multiple solutions ?

# [RULES :] BERTIN'S APPROACH (A VISUAL LANGUAGE)

If graphs are used to communicate, it is a form of language.

Any language has a grammar, "words" and logic. Let us study the science that deals with signs or sign language : "*The Semiology*".

TABLE – Bertin's definition of 8 visual variables

Position (x, y)

Size

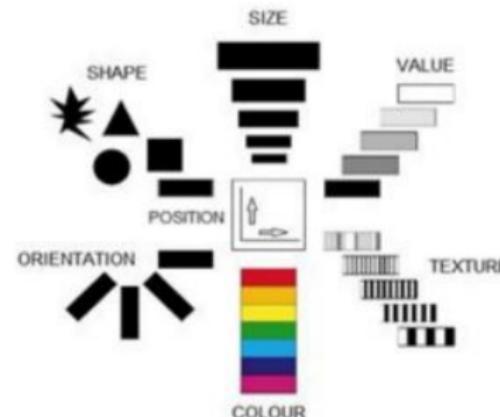
Value

Texture

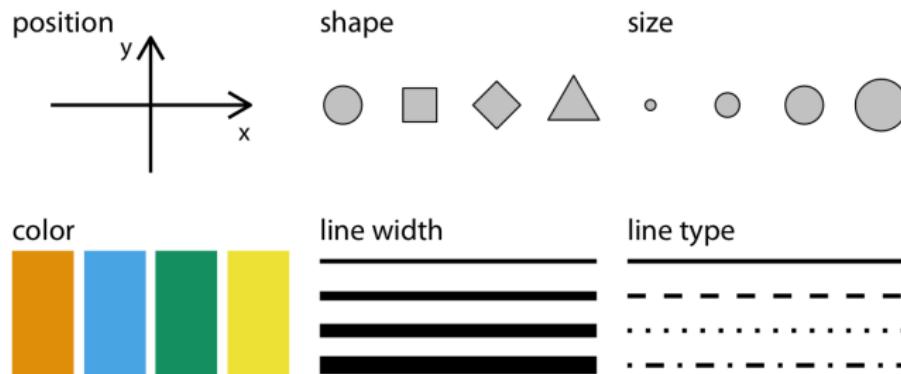
Colour

Orientation

Shape



# A VISUAL LANGUAGE



From *Claus O. Wilke*

# THESE VARIABLES SERVE DIFFERENT GOALS

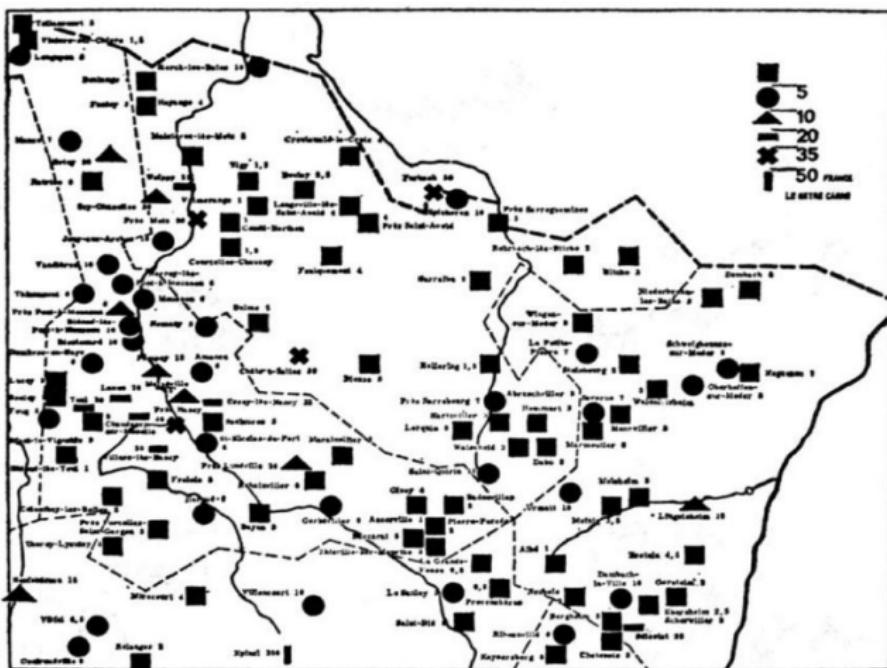
Visual variable *syntactics*, designating each visual variable as suited or not for levels of measurement :  
Equivalence, differences, order, proportions.

	Visual variable syntactics :			
	Variable suited for :			
Position (x, y)	$\equiv$	$\neq$	O	$\propto$
Size	$\equiv$	$\neq$	O	$\propto$
Value	$\equiv$	$\neq$	O	$\propto$
Texture	$\equiv$	$\neq$	O	
Colour	$\equiv$	$\neq$		
Orientation	$\equiv$	$\neq$		
Shape	$\equiv$			

$\equiv$  : Equivalence  $\neq$  : Differences ; O : Order ;  $\propto$  : Proportions

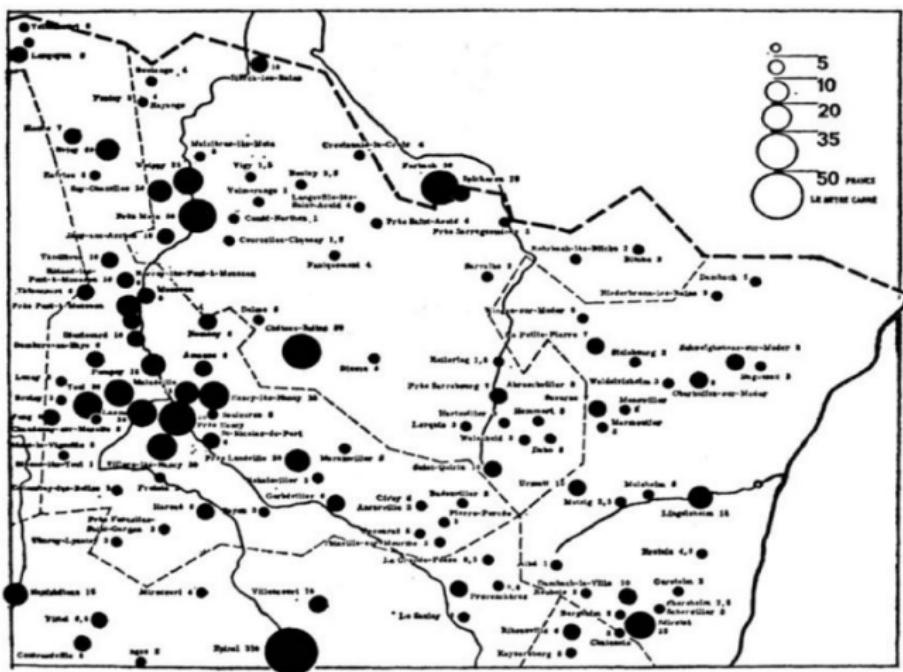
SHAPE IS NOT SUITABLE FOR PROPORTIONALITY

## Price of land in the East of France Bertin (1970)



## SIZE IS SUITABLE FOR PROPORTIONALITY

## Price of land in the East of France Bertin (1970)



# WHAT ELSE ?

Is **color** suitable for  
proportionality ?

# A NOTE ON COLORS

“Colors” are not suited for ordering nor for proportionality !  
Try putting the following hues in order from low to high.



## A NOTE ON COLORS

These colors are easy to order from low to high.

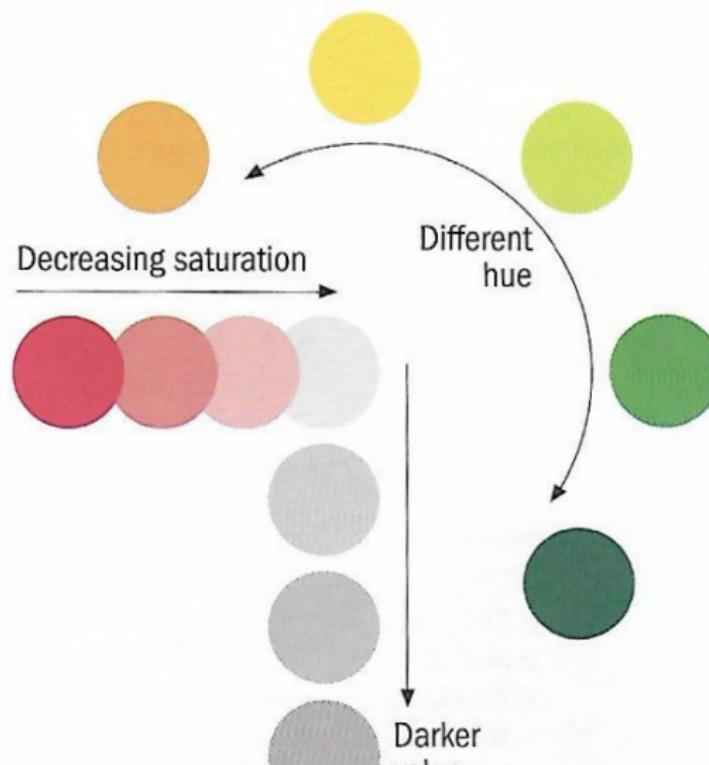
Few (2008) provides meaningful solutions for choosing palettes of colors, for example for heatmaps.



See also the *ggplot* theme `theme_few()`

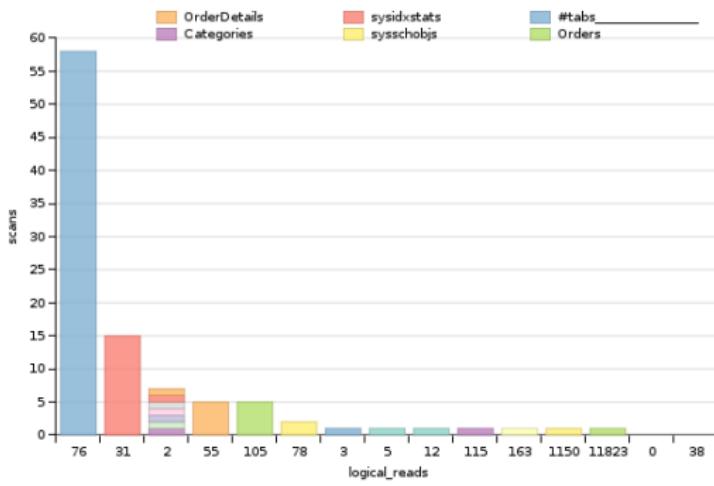
# A NOTE ON COLORS

What we call "colors" are three dimensional objects



# A NOTE ON COLORS (FINAL)

Colors are sometimes a graphic puzzle Tufte (2001).



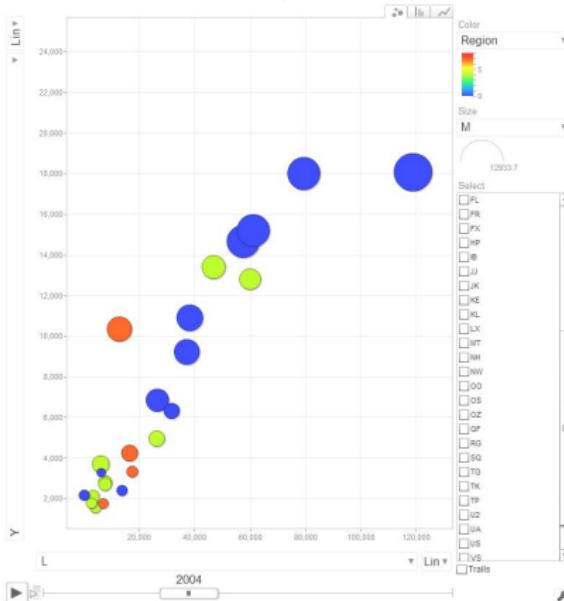
scans

Your eyes will go back and forth from the graph to the legend...

Source <http://viz.wtf/image/135265269618>)

# CONJUNCTION OF COLOURS AND PROPORTIONALITY

## Productivity of Airlines



(Demo with googleVis)

## A QUIZ ON PROPORTIONALITY :

If 100% of the US prisoners are represented by the big green square...what is the percentage for each group ?

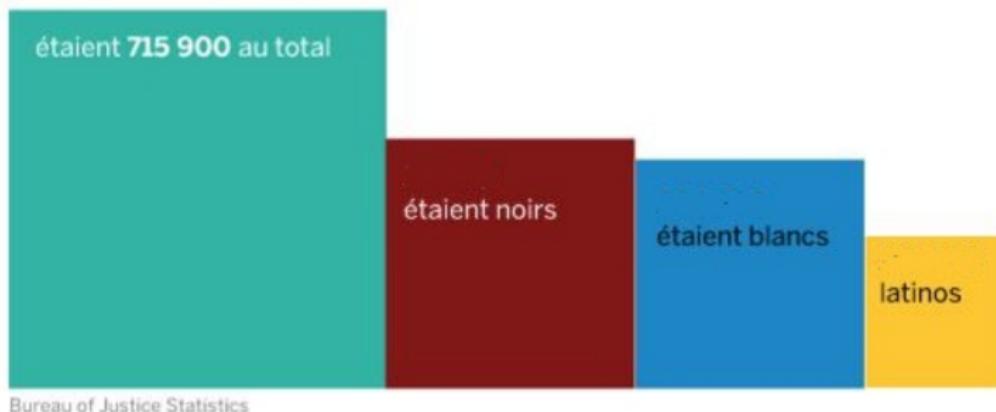


FIGURE – Ethic composition of prisoners in Jail in 2008 in the USA.  
(*Le Monde* 5/12/2014)

# NOT SO SIMPLE...

If 100% of the US prisoners are represented by the big square...what is the percentage for each group ?

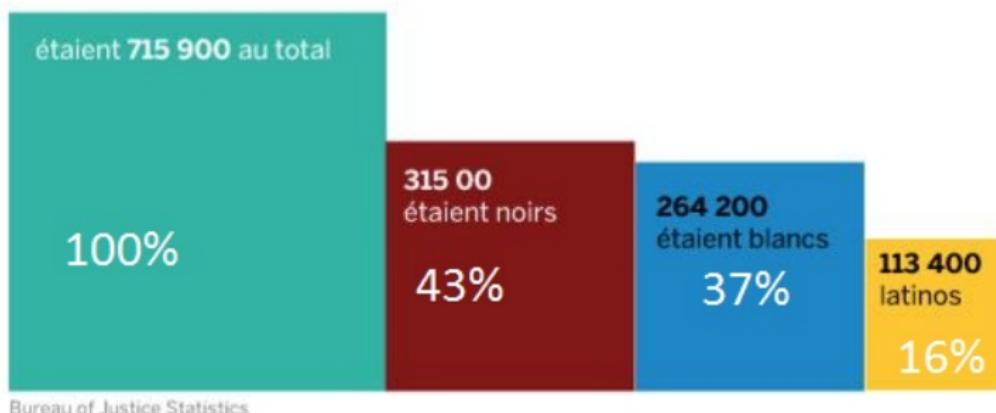
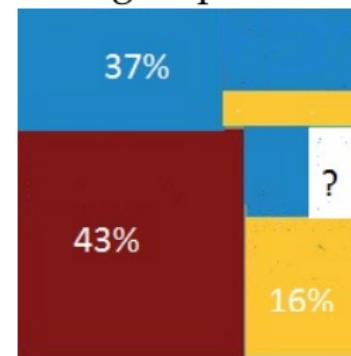


FIGURE – Ethic composition of prisoners in Jail in 2008 in the USA.  
(Le Monde 5/12/2014)

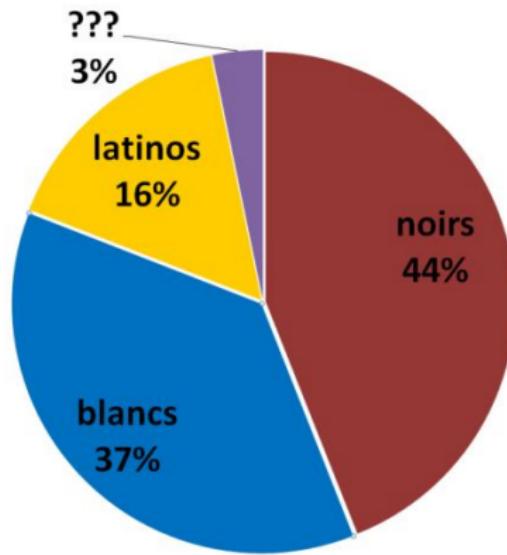
# VERIFICATION

If 100% of the US prisoners are represented by the big square...what is the percentage for each group ?



OR...

If 100% of the US prisoners are represented by the big square...what is the percentage for each group ?



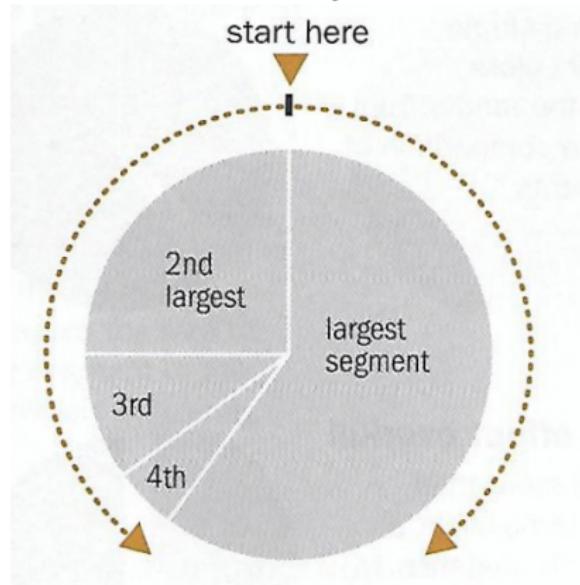
# A NOTE ON PIE CHARTS

Pie charts are graphics that can be used if :

- ▶ Few categories
- ▶ "good" distribution of shares
- ▶ "good" choice of colors

# A NOTE ON PIE CHARTS (CONTINUED)

One should cleverly use of the reference line



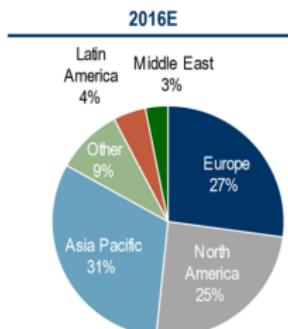
from Wong (2010)

# A NOTE ON PIE CHARTS (FINAL)

One should **not** use pie chart for comparisons

## APAC Has Become the Largest Travel Market<sup>(1)</sup>

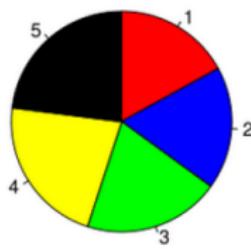
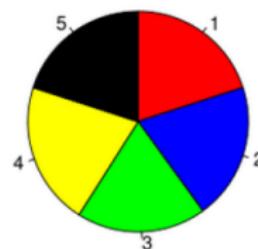
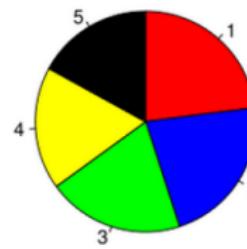
Regional Share of Total Tourism Contribution to Global GDP



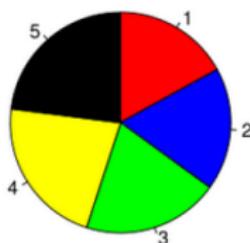
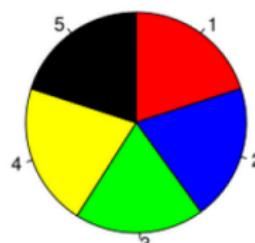
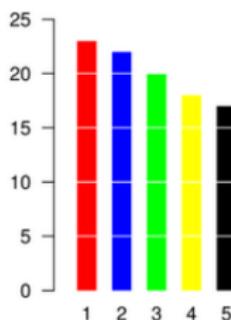
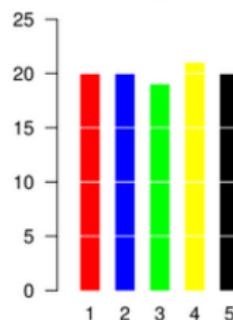
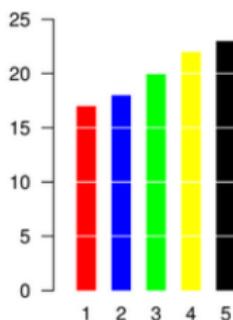
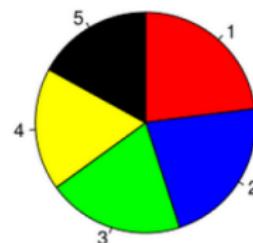
see SWD challenge

<sup>(1)</sup> World Travel & Tourism Council

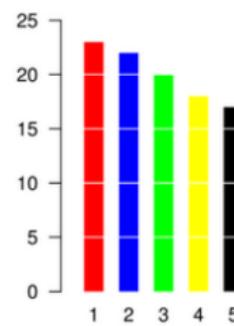
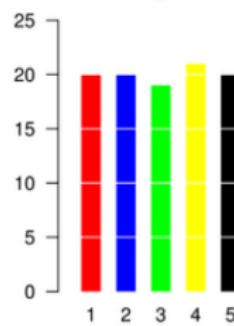
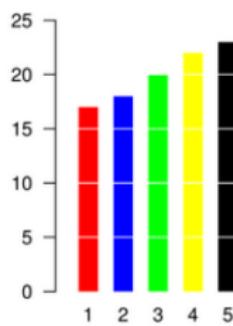
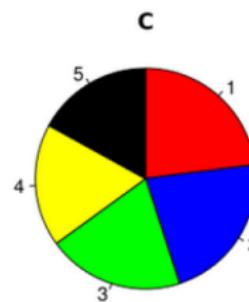
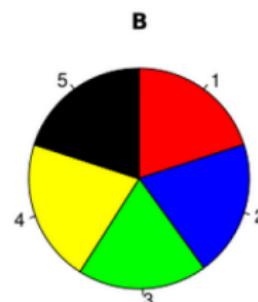
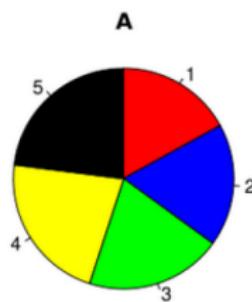
# BE CAREFUL WITH PIE CHARTS !

**A****B****C**

# BE CAREFUL WITH PIE CHARTS !

**A****B****C**

# BE CAREFUL WITH PIE CHARTS !



Source <https://twitter.com/freakonometrics/status/612742330160951296>

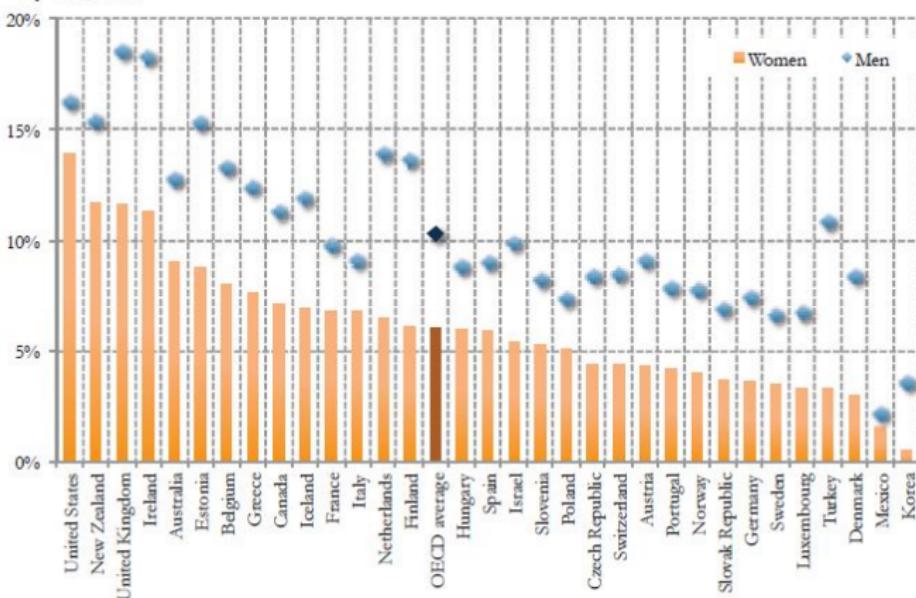
# “GOOD” OR “BAD” GRAPHICS?

*“There are no “good” nor “bad” graphics (...), there are graphics answering legitimate questions and graphics that do not answer question at all ”*

Bertin (1981)

# SCHWABISH (JEP, 2014) BEFORE-AFTER

Percentage of Employed Who Are Senior Managers,  
by Sex, 2008



Source: Author, based on OECD (no date) and Rampell (2013).

FIGURE – An Unbalanced Chart - Original

# SCHWABISH (JEP, 2014) BEFORE-AFTER

Percentage of Employed Who Are Senior Managers, by  
Gender, 2008  
(percent)

● Women ● Men

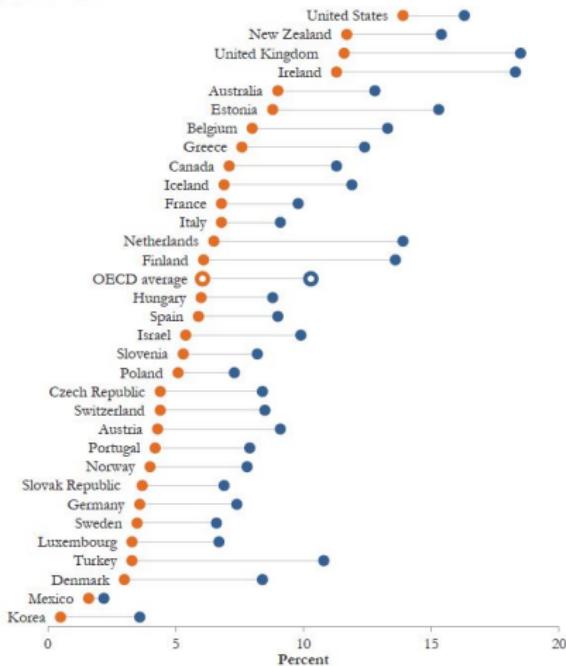
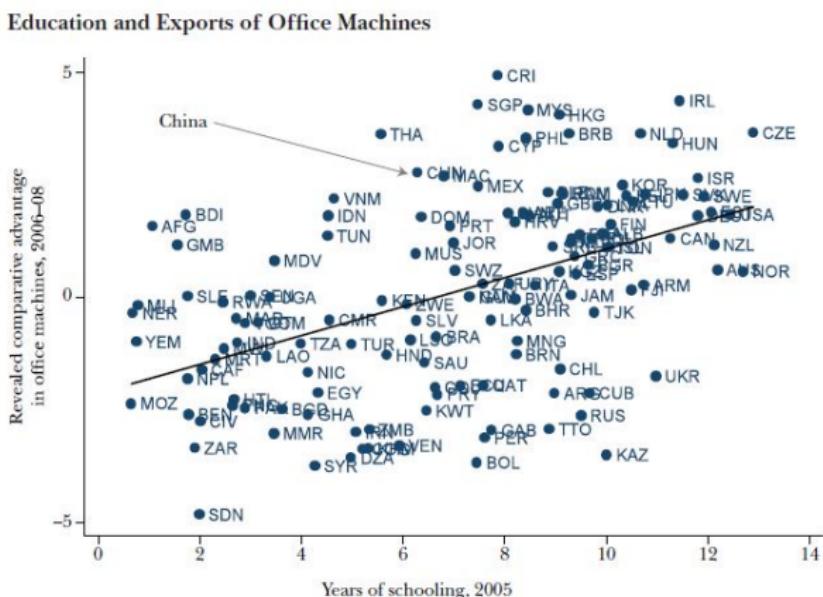


FIGURE – An Unbalanced Chart - Revised

SCHWABISH (JEP, 2014) BEFORE-AFTER



*Source:* Hanson (2012).

## FIGURE – A Clutterplot Example - Original

# SCHWABISH (JEP, 2014) BEFORE-AFTER

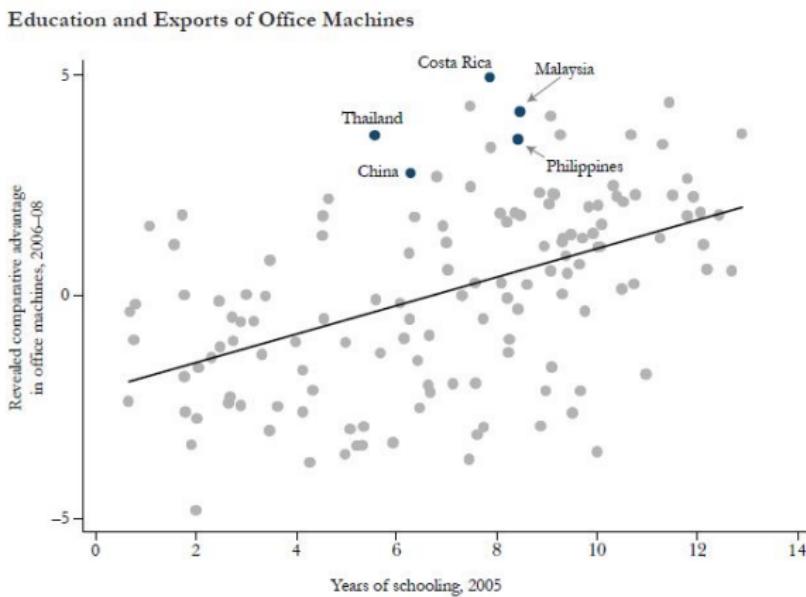
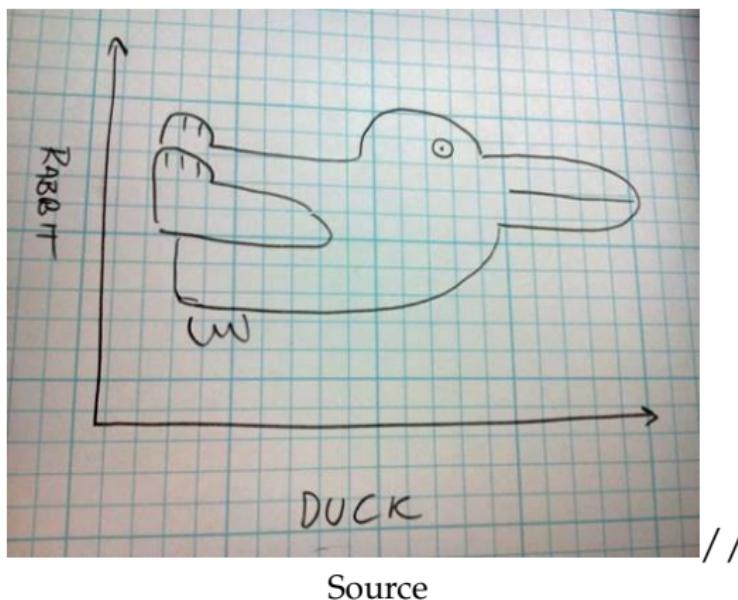


FIGURE – A Clutterplot Example - Revised

# A NOTE ON PERCEPTION

A bird (Duck, Toucan ?) on the X axis, a rabbit on the Y axis !



<http://flowingdata.com/2014/06/25/duck-vs-rabbit-plot/>

# "PREATTENTIVE" VARIABLES

How many "3" in that sequence ? (from Ware (2012))

# "PREATTENTIVE" VARIABLES

How many "3" in that sequence ? (from Ware (2012))

45929078059772098775972655665110049836645  
27107462144654207079014738109743897010971  
43907097349266847858715819048630901889074  
25747072354745666142018774072849875310665

# "PREATTENTIVE" VARIABLES

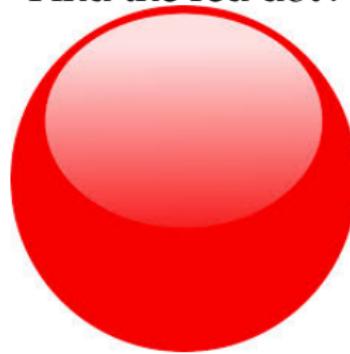
How many "3" in that sequence ? (from Ware (2012))

45929078059772098775972655665110049836645  
27107462144654207079014738109743897010971  
43907097349266847858715819048630901889074  
25747072354745666142018774072849875310665

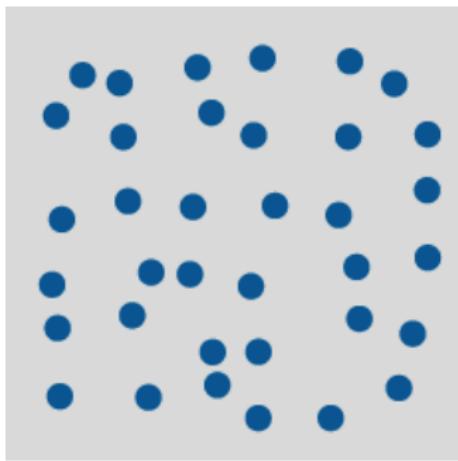
459290780597720987759726556651100498**3**6645  
271074621446542070790147**3**8109743897010971  
**4**3907097**3**49266847858715819048630901889074  
25747072**3**54745666142018774072849875**3**10665

# AND NOW...

Find the red dot!



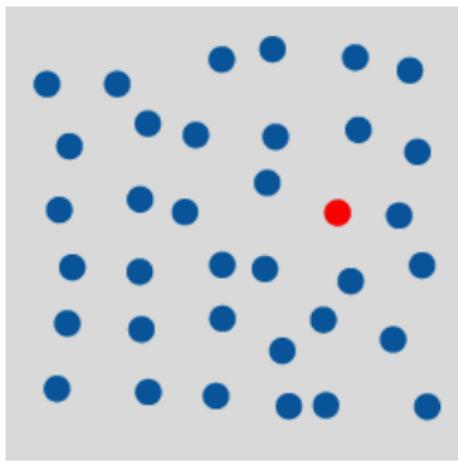
# TEST : FIND THE RED DOT !





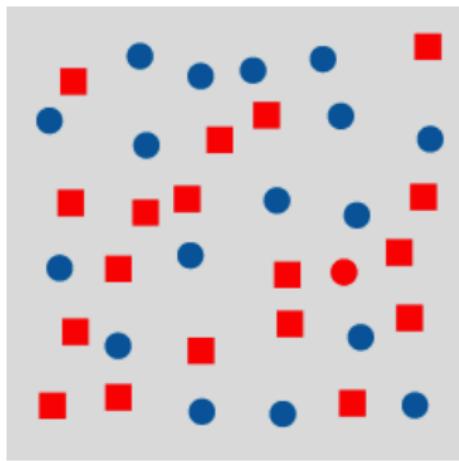
# TEST : FIND THE RED DOT !

# TEST : FIND THE RED DOT !



# TEST : FIND THE RED DOT !

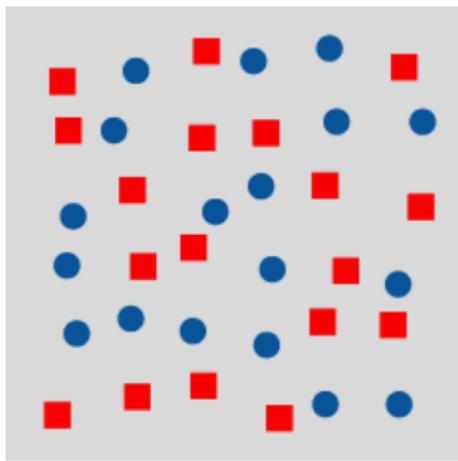
# TEST : FIND THE RED DOT !





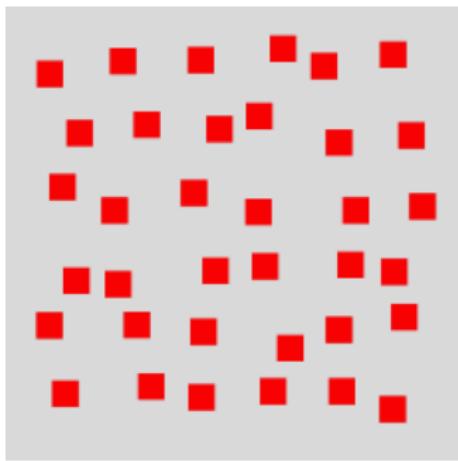
# TEST : FIND THE RED DOT !

# TEST : FIND THE RED DOT !



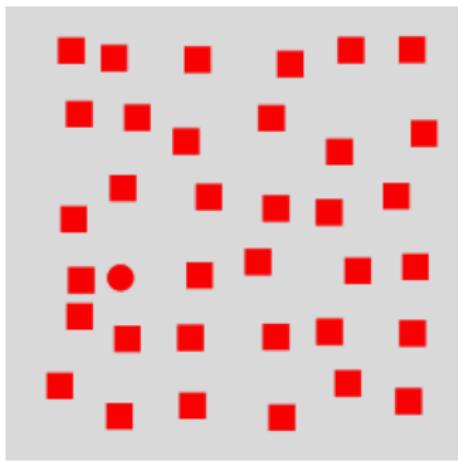
# TEST : FIND THE RED DOT !

# TEST : FIND THE RED DOT !



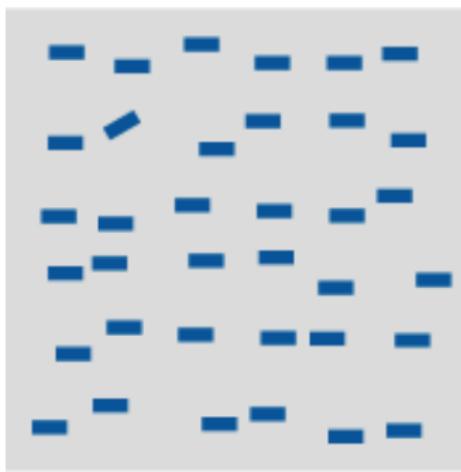
# TEST : FIND THE RED DOT !

# TEST : FIND THE RED DOT !



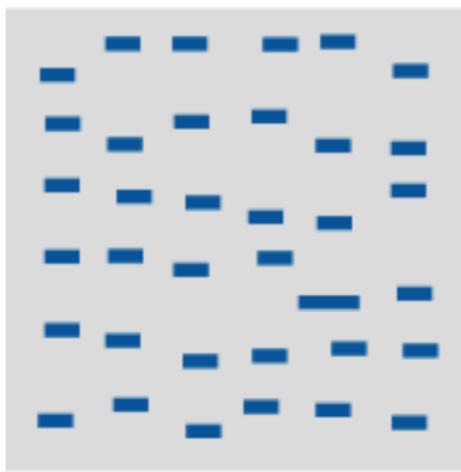
# HARDER : IS THERE A "STRANGER" ?

# HARDER : IS THERE A "STRANGER" ?



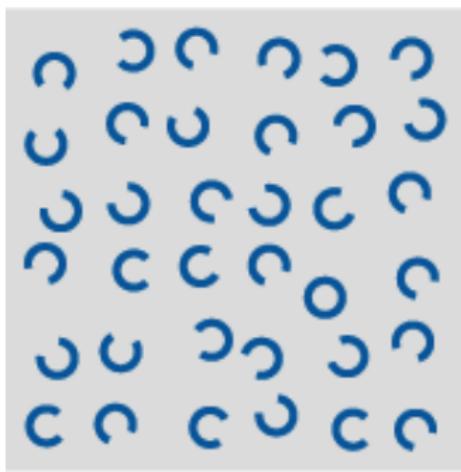
# HARDER : IS THERE A "STRANGER" ?

# HARDER : IS THERE A "STRANGER" ?



# HARDER : IS THERE A "STRANGER" ?

# HARDER : IS THERE A "STRANGER" ?



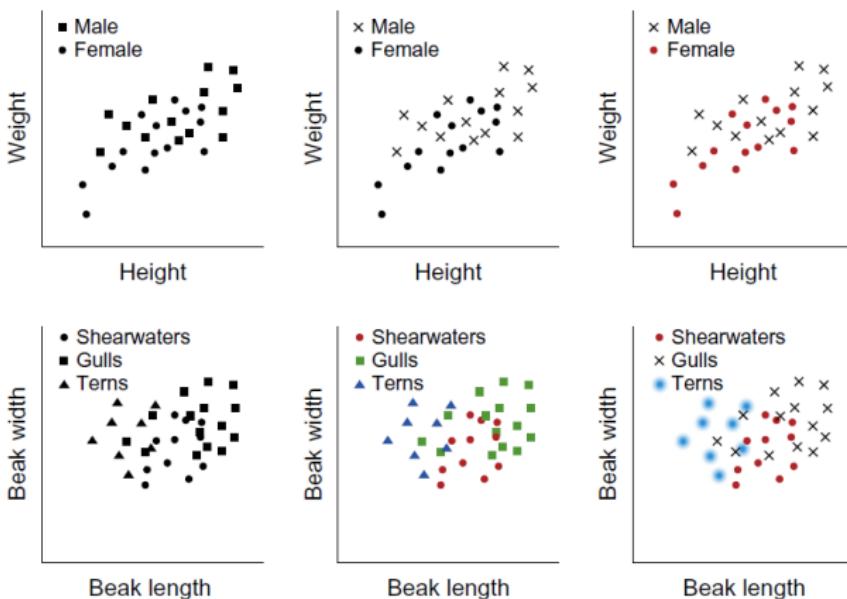
# HARDER : IS THERE A "STRANGER" ?

# THAT WASN'T EASY

- ▶ Preattentive concept, Treisman (1985)
- ▶ Some visual elements or patterns are detected immediately
- ▶ But there may be interferences (colour and form)
- ▶ Very useful (detection, explanatory and presentation)

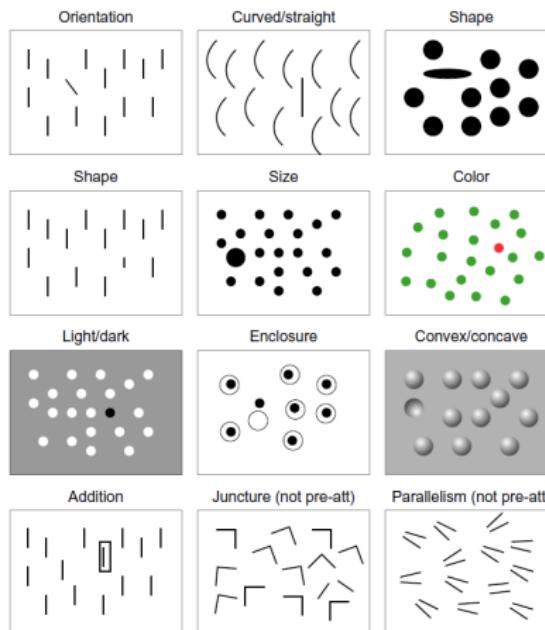
**Helpful to highlight a message !**

# TOO MUCH VARIATION DOESN'T HELP



From Ware (2012)

# MOST PREATTENTIVE VISUAL VARIABLES

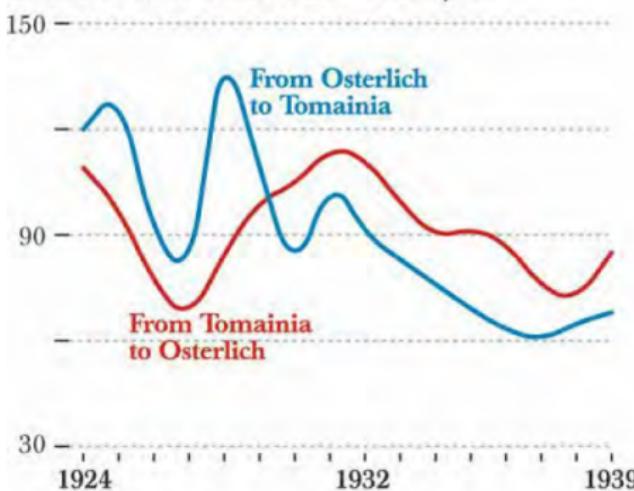


From Ware (2012)

# VISUAL PERCEPTION AND LINES

## Exports between Tomainia and Osterlich

In millions of Tomainian reichsmarks a year

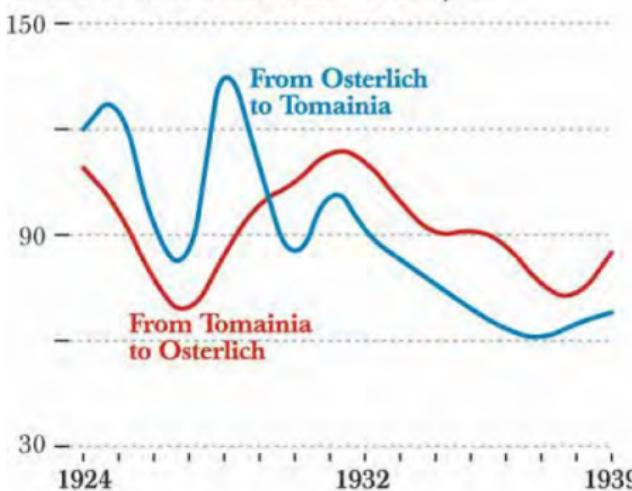


From Cairo (2012)

# VISUAL PERCEPTION AND LINES

## Exports between Tomainia and Osterlich

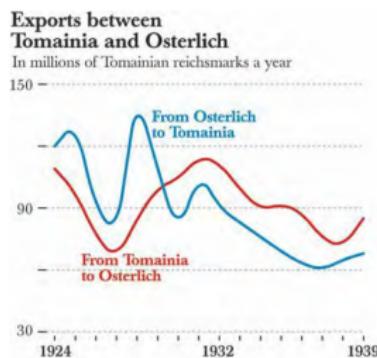
In millions of Tomainian reichsmarks a year



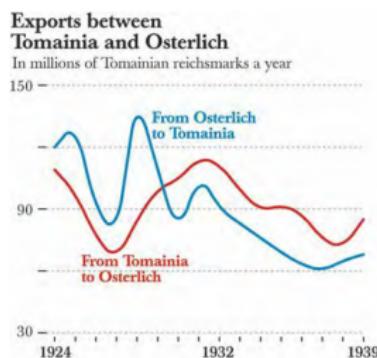
When was the biggest negative (positive) difference ?

From Cairo (2012)

# VISUAL PERCEPTION AND LINES

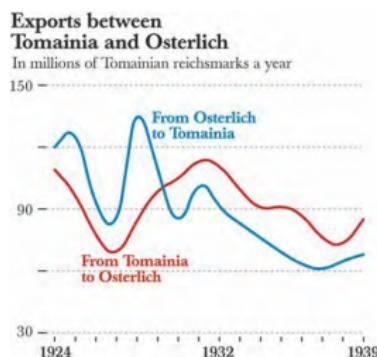


# VISUAL PERCEPTION AND LINES



When was the biggest negative (positive) difference ?

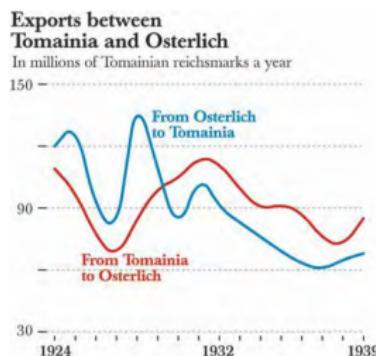
# VISUAL PERCEPTION AND LINES



When was the biggest negative (positive) difference ?



# VISUAL PERCEPTION AND LINES

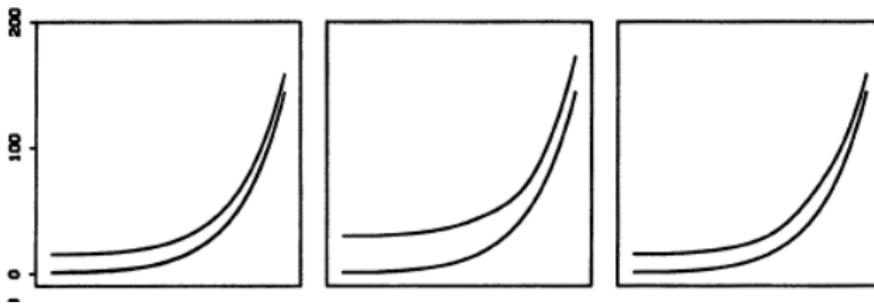


When was the biggest negative (positive) difference ?

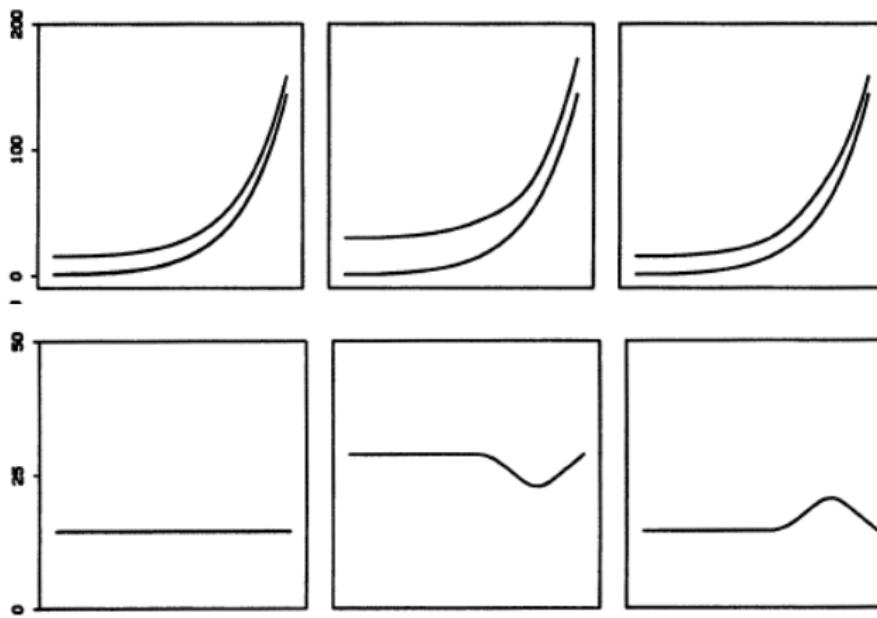


From Cairo (2012)

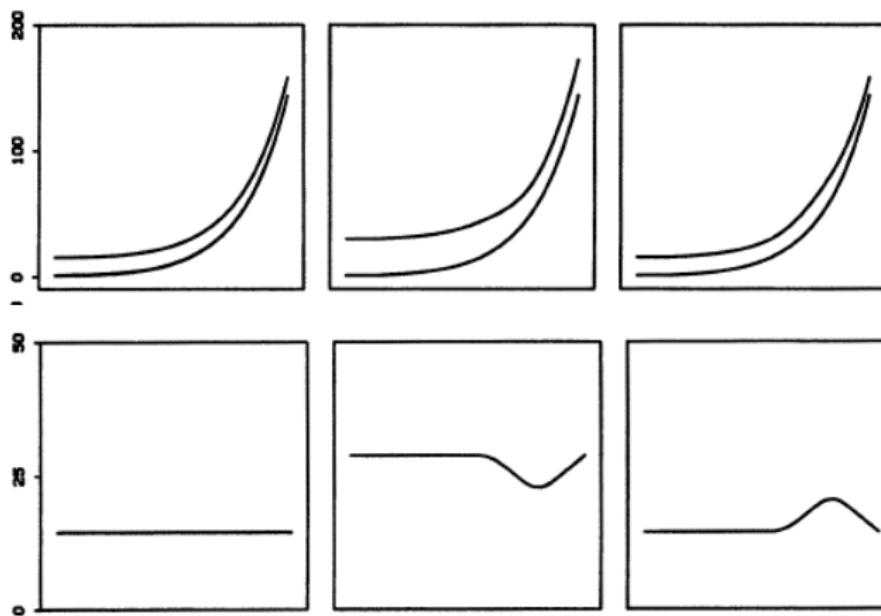
# THE CLEVELAND-MCGILL EFFECT



# THE CLEVELAND-MCGILL EFFECT

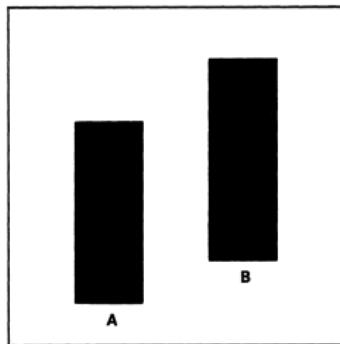


# THE CLEVELAND-MCGILL EFFECT

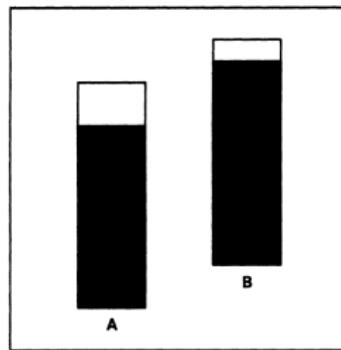
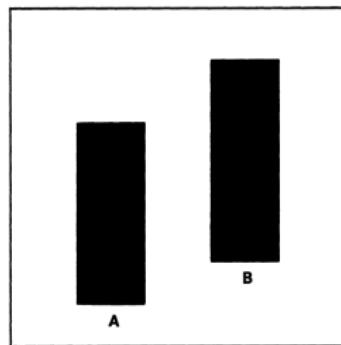


From Cleveland and McGill (1984)

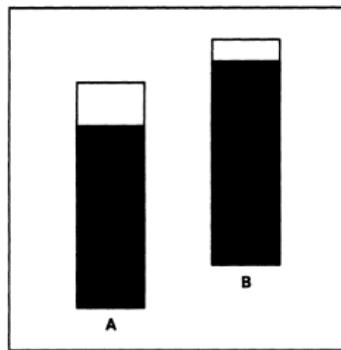
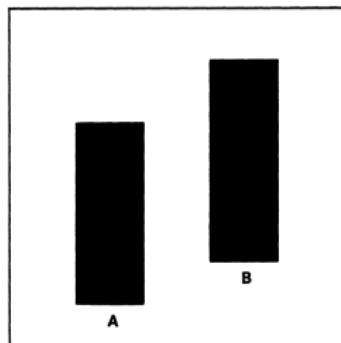
# WEBER'S LAW AND FRAMED BOXES



# WEBER'S LAW AND FRAMED BOXES

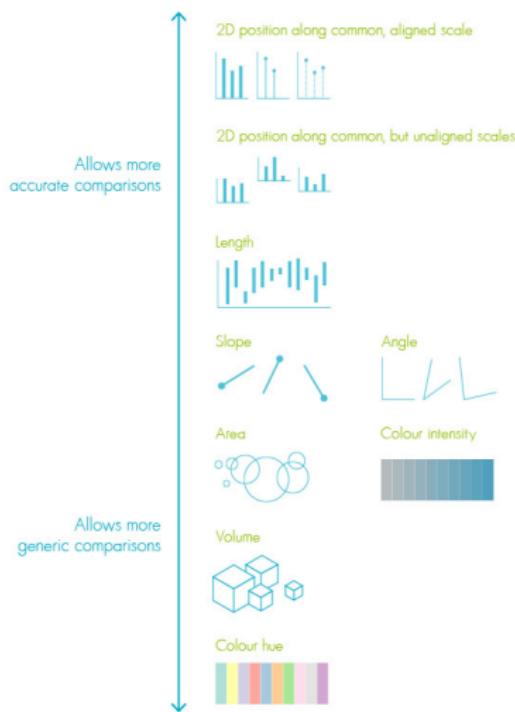


# WEBER'S LAW AND FRAMED BOXES



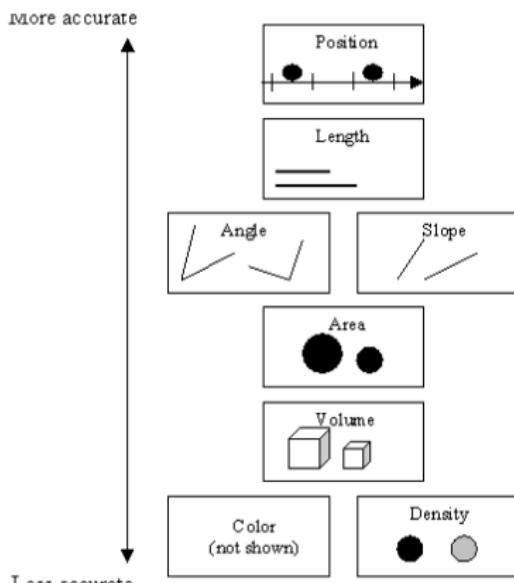
From Cleveland and McGill (1984)

# THE CLEVELAND-MCGILL SCALE



Cleveland and McGill (1984)

# THE CLEVELAND-MCGILL SCALE



Source <http://hcil2.cs.umd.edu/trs/99-20/99-20.html>

# THE DATAVIZ PROCESS

## STEP 1 : THINK

### THINK BEFORE YOU INK

#### DATA CHARACTERISTICS

Sample size: 25 subjects

Dimensions: 3

- Genotype (nominal, independent)
- Task difficulty (ordinal, independent)
- Reaction time (continuous, dependent)

#### HYPOTHESIS

Mean reaction time profiles over task difficulty depend on subject genotype.

# THE DATAVIZ PROCESS

## STEP 1 : THINK

### THINK BEFORE YOU INK

#### DATA CHARACTERISTICS

Sample size: 25 subjects

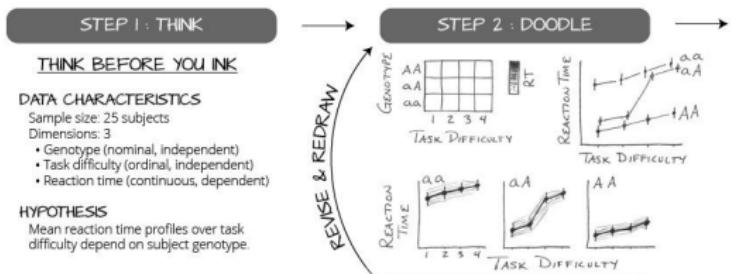
Dimensions: 3

- Genotype (nominal, independent)
- Task difficulty (ordinal, independent)
- Reaction time (continuous, dependent)

#### HYPOTHESIS

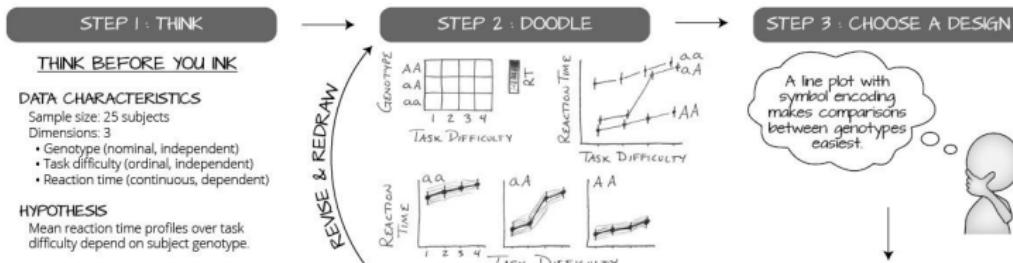
Mean reaction time profiles over task difficulty depend on subject genotype.

# THE DATAVIZ PROCESS



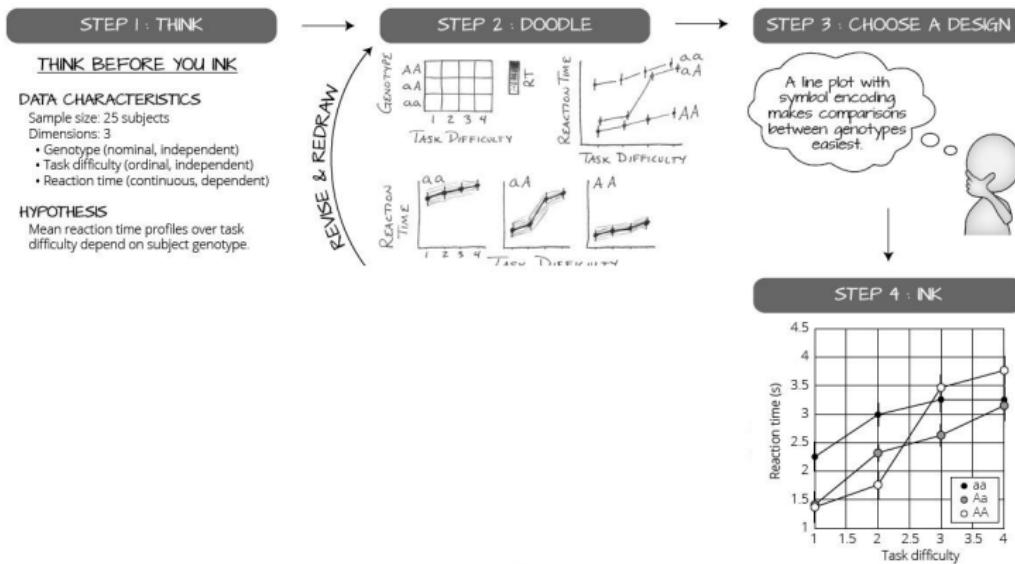
*From Allen and Erhardt (2016a)*

# THE DATAVIZ PROCESS



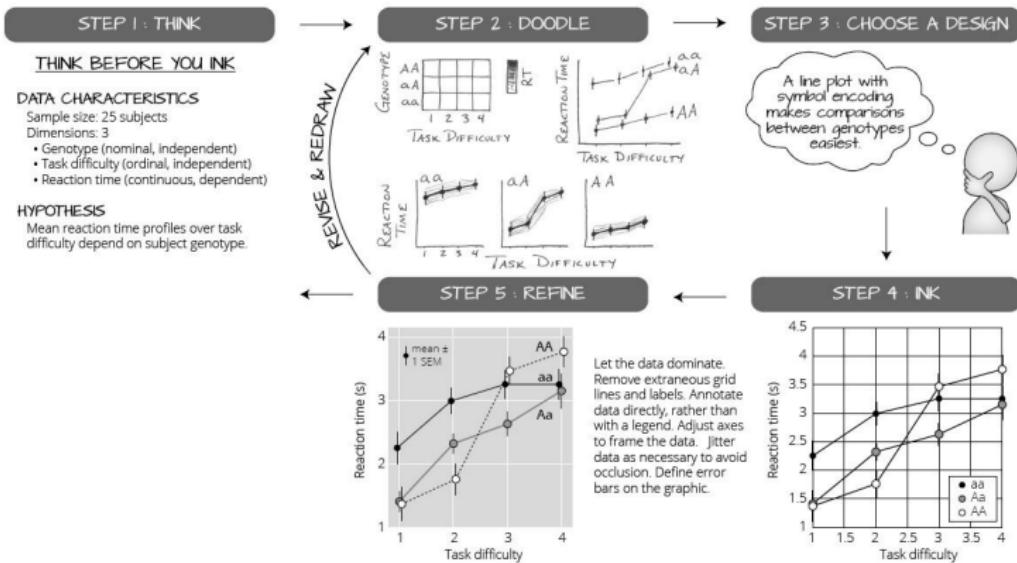
From Allen and Erhardt (2016a)

# THE DATAVIZ PROCESS



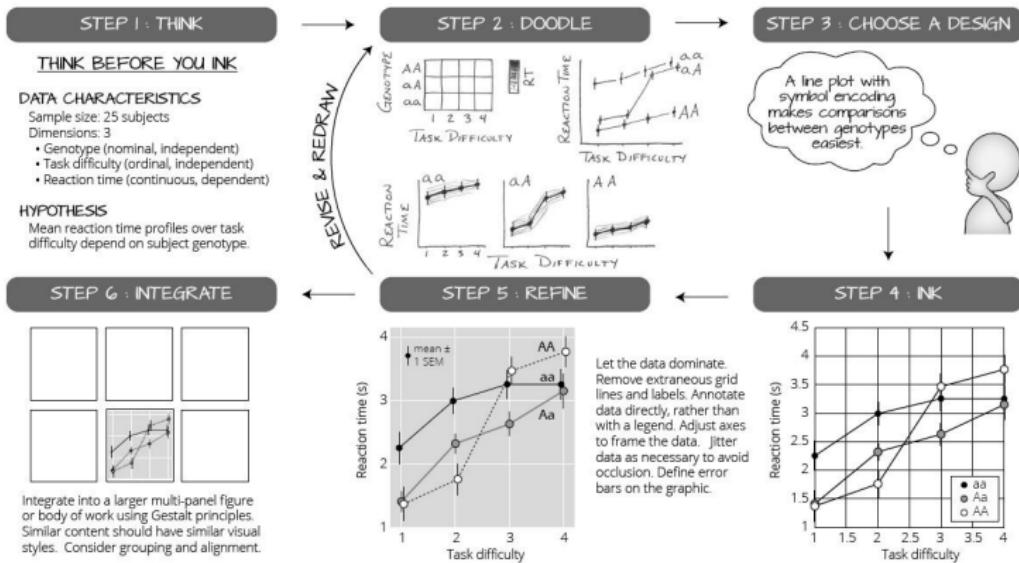
From Allen and Erhardt (2016a)

# THE DATAVIZ PROCESS



From Allen and Erhardt (2016a)

# THE DATAVIZ PROCESS



From Allen and Erhardt (2016a)

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

1. Show the data clearly

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

1. Show the data clearly
2. Use simplicity in design

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

1. Show the data clearly
2. Use simplicity in design
3. Use good alignment on a common scale for quantities to be compared

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

1. Show the data clearly
2. Use simplicity in design
3. Use good alignment on a common scale for quantities to be compared
4. Keep visual encoding transparent

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

1. Show the data clearly
2. Use simplicity in design
3. Use good alignment on a common scale for quantities to be compared
4. Keep visual encoding transparent
5. Use graphical forms consistent with those principles

# PARTIAL CONCLUSION

Gordon and Finch (2014) gives some nice principles

1. Show the data clearly
2. Use simplicity in design
3. Use good alignment on a common scale for quantities to be compared
4. Keep visual encoding transparent
5. Use graphical forms consistent with those principles

We may add some others (use preattentive elements, integrity, ...)

# CHECKLIST

Gordon and Finch (2014) provides a checklist

## A checklist for good graphical practice

How clear is your purpose in communication?

- What relationships or patterns can you identify in the graph?
- Are these the relationships or patterns you intended to represent?
- Can the viewer identify the patterns you wish to illustrate?
- Are the important comparisons you wish to show salient?

Make clarity a high priority.

- Does the graph have a clear title?
- Are the axes labelled?
- Are the units of the variables measured defined?
- Are the units of observation clear?
- Is the graph large enough?
- Would ordering groups or variables plotted improve the graph?
- Are all the graph labels horizontal?
- Use points to plot estimates (e.g. means, proportions) rather than bars.

Choose standard forms fit for your purpose.

- Use bars around points to indicate the precision of the estimates.
- Plot the estimates of interest (e.g. mean differences with confidence intervals) rather than standard summary statistics (group means).
- Plot inferences to support stories about models.
- Plot data to support stories about distributions and variation.

Consider detection issues.

- Can all the data points be seen?
- Are patterns in the data clear?

# CHECKLIST

Gordon and Finch (2014) provides a checklist

- Are the fonts large enough?
- Would it help to use jittering, or another form of representing multiple, identical values?

Would panels help?

- Are there grouping variables that can be used to panel the graph?
- Do the grouping variables correspond to the variation of interest?
- Would additional panels help?

Align quantities to be compared on a common scale.

- Has distortion of the data been avoided by using the same scales for the same measurement?
- Are measurements made on the same scale plotted on the same scale?
- Would transposition improve the graph?

Does the graph have grid lines?

- Light grey grid lines will help with accurate interpretation.

Are all the elements of the graph defined?

- What do points on the graph correspond to?
- Are estimates (e.g. means, proportions) plotted on the graph clearly defined?
- Are bars around points on the graph clearly defined?

How much decoding work does the viewer have to do?

- Is it easy for someone unfamiliar with your data to interpret your graph?
- Does the graph stand alone?
- Try it on a friend!

# REFERENCES I

- Allen, E. A. and Erhardt, E. B. (2016a). *Handbook of Psychophysiology*,. Cambridge University Press, 4th edition edition.
- Allen, E. A. and Erhardt, E. B. (2016b). Visualizing Scientific Data. In Cacioppo, J. T., Tassinary, L. G., and Berntson, G. G., editors, *Handbook of Psychophysiology*, pages 679–697. Cambridge University Press, 4 edition.
- Anscombe, F. J. (1973). Graphs in statistical analysis. *The American Statistician*, 27(1) :17–21.
- Bertin, J. (1970). La graphique. *Communications*, 15(1).
- Bertin, J. (1981). Théorie matricielle de la graphique. *Communication et langages*, 48(1) :62–74.
- Buja, A., Cook, D., Hofmann, H., Lawrence, M., Lee, E.-K., Swayne, D. F., and Wickham, H. (2009). Statistical inference for exploratory data analysis and model diagnostics. *Philosophical Transactions of the Royal Society of London A : Mathematical, Physical and Engineering Sciences*, 367(1906) :4361–4383.
- Cairo, A. (2012). *The Functional Art : An introduction to information graphics and visualization*. Voices That Matter. Pearson Education.

## REFERENCES II

- Chen, C.-h., Härdle, W. K., and Unwin, A. (2007). *Handbook of data visualization*. Springer Science & Business Media.
- Cleveland, W. S. (1994). *The Elements of Graphing Data*. Hobart Press, Summit : NJ, 2 edition.
- Cleveland, W. S. and McGill, R. (1984). Graphical perception : Theory, experimentation, and application to the development of graphical methods. *Journal of the American Statistical Association*, 79(387) :531–554.
- Dix, A. and Ellis, G. (1998). Starting simple - adding value to static visualisation through simple interaction. In Eds. T. Catarci, M. F. Costabile, G. S. and Tarantino, L., editors, *Proceedings of Advanced Visual Interfaces*, pages 124–134. L’Aquila, Italy, ACM Press.
- Few, S. (2008). Practical rules for using color in charts. *Visual Business Intelligence Newsletter*, (11).
- Gelman, A. (2004). Exploratory data analysis for complex models. *Journal of Computational and Graphical Statistics*, 13(4).
- Gordon, I. and Finch, S. (2014). Statistician heal thyself : Have we lost the plot ? *Journal of Computational and Graphical Statistics*, 24(4) :1210–1229.

## REFERENCES III

Hahsler, M., Hornik, K., and Buchta, C. (2008). Getting things in order : an introduction to the r package seriation. *Journal of Statistical Software*, 25(3) :1–34.

Matejka, J. and Fitzmaurice, G. (2017). Same stats, different graphs : generating datasets with varied appearance and identical statistics through simulated annealing. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 1290–1294. ACM.

Munzner, T. (2014). *Visualization Analysis and Design*. AK Peters Visualization Series. A K Peters/CRC Press, 1 edition.

Treisman, A. (1985). Preattentive processing in vision. *Computer Vision, Graphics, and Image Processing*, 31(2) :156–177.

Tufte, E. R. (2001). *The Visual Display of Quantitative Information*. Graphics Press, 2 edition.

Ware, C. (2012). *Information visualization : perception for design*. Elsevier.

Wong, D. M. (2010). *The Wall Street Journal guide to information graphics : The dos and don'ts of presenting data, facts, and figures*. WW Norton.