

dati

dati

**informazioni raccolte in modo
sistematico**

in un contesto e per uno scopo

dati non è il plurale di aneddoto

osservazioni aneddotiche

25

Nicola

41

Gerard

blue

185

30

variabili

etichette per insiemi di dati

**in studi sperimentali o
osservativi**

dipendenti o indipendenti

variabili

**possono essere numeri o
categorie, between o within**

**all'interno di queste
distinzioni fondamentali ci
sono tantissime sfumature**

R

creare un dataframe

**esempio di uso di una
funzione: pie()**

**decifrare help(pie):
descrizione, uso, argomenti,
vedi anche, esempi e note**

la nota su help(pie)

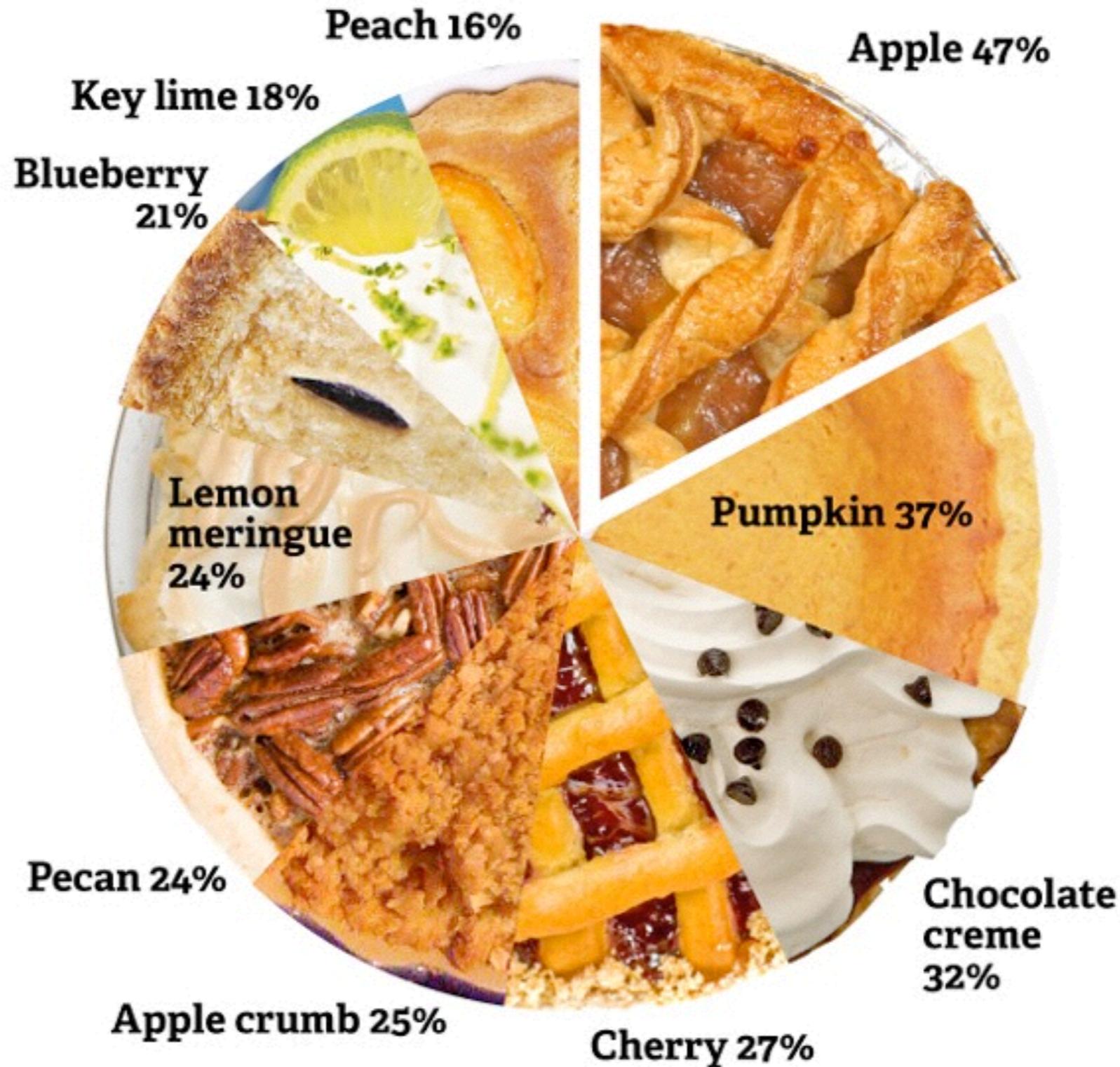
Note

Pie charts are a very bad way of displaying information. The eye is good at judging linear measures and bad at judging relative areas. A bar chart or dot chart is a preferable way of displaying this type of data.

Cleveland (1985), page 264: “Data that can be shown by pie charts always can be shown by a dot chart. This means that judgements of position along a common scale can be made instead of the less accurate angle judgements.” This statement is based on the empirical investigations of Cleveland and McGill as well as investigations by perceptual psychologists.



What are your three most favorite types of pie?



30%

GLI INTERVISTATI
CHE NON RITENGONO
L'UNITÀ NAZIONALE
UN BENE IRRINUNCIABILE

55%

GLI INTERVISTATI
PER I QUALI A DIVIDERE
IL PAESE SONO
I PARTITI POLITICI



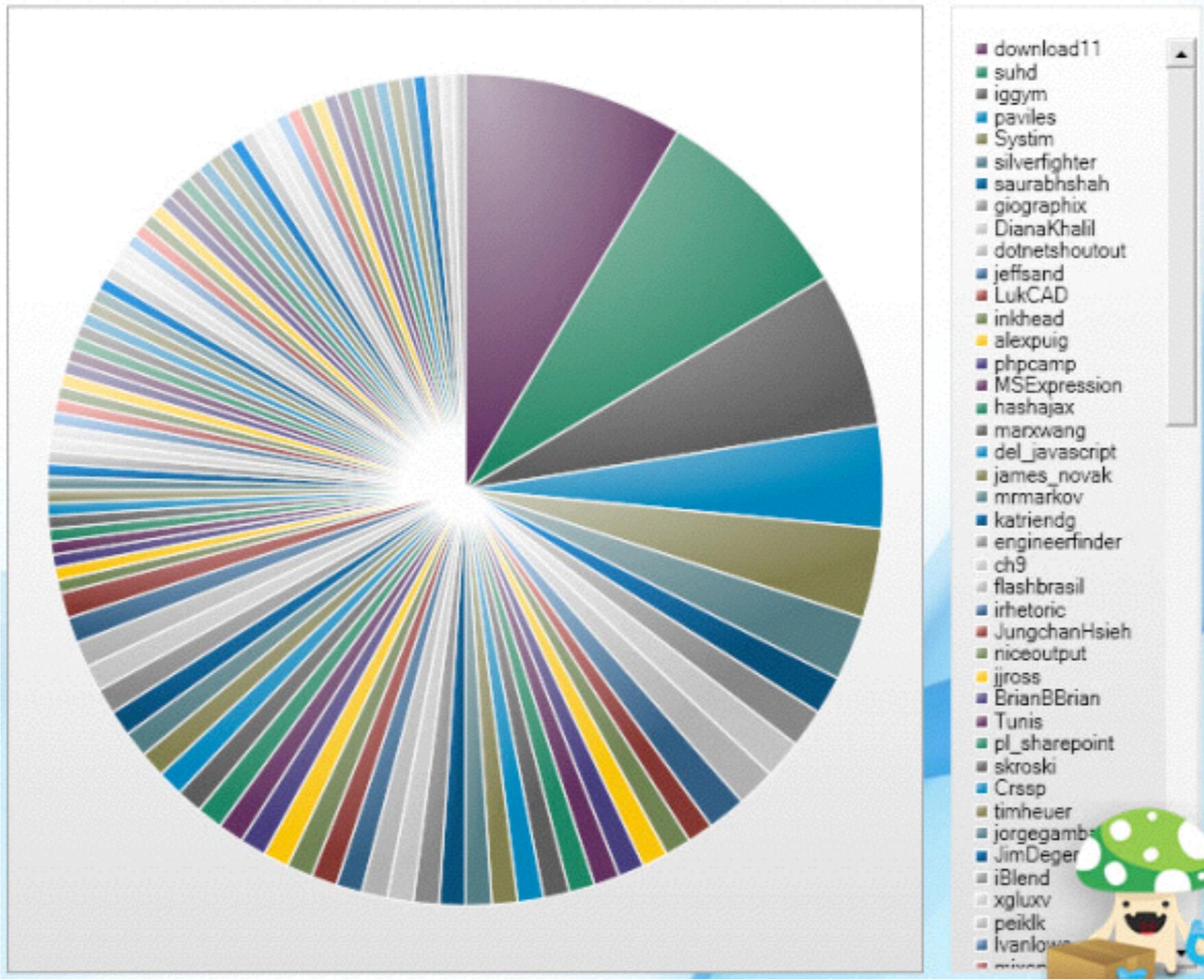
40%

GLI INTERVISTATI CHE
CONSIDERANO LA MAFIA
IL NOSTRO MAGGIOR
MOTIVO DI VERGOGNA

30%

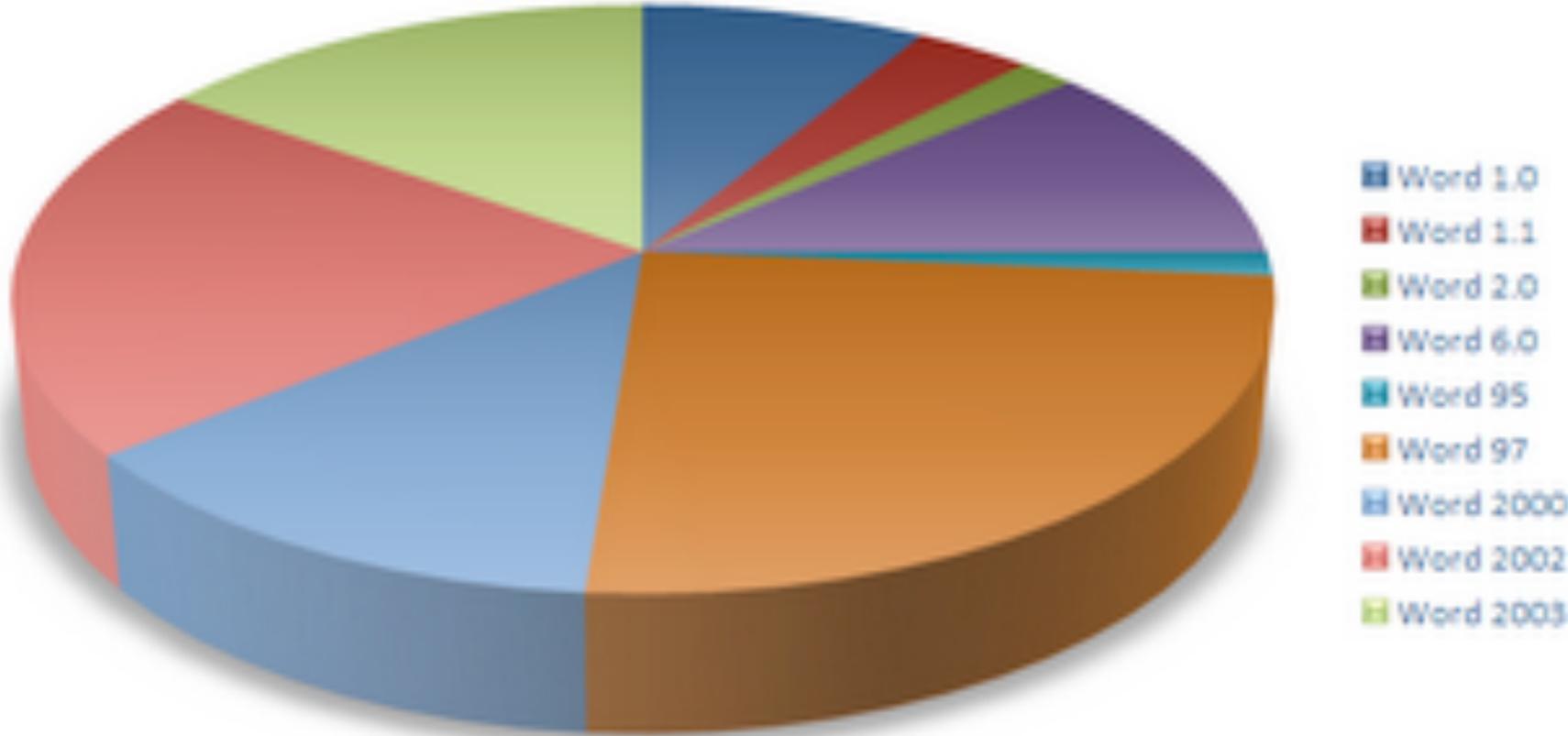
GLI INTERVISTATI
SECONDO I QUALI L'ARTE
DI ARRANGIARSI È LA NOSTRA
MIGLIORE QUALITÀ

100 Most Active Tweeters



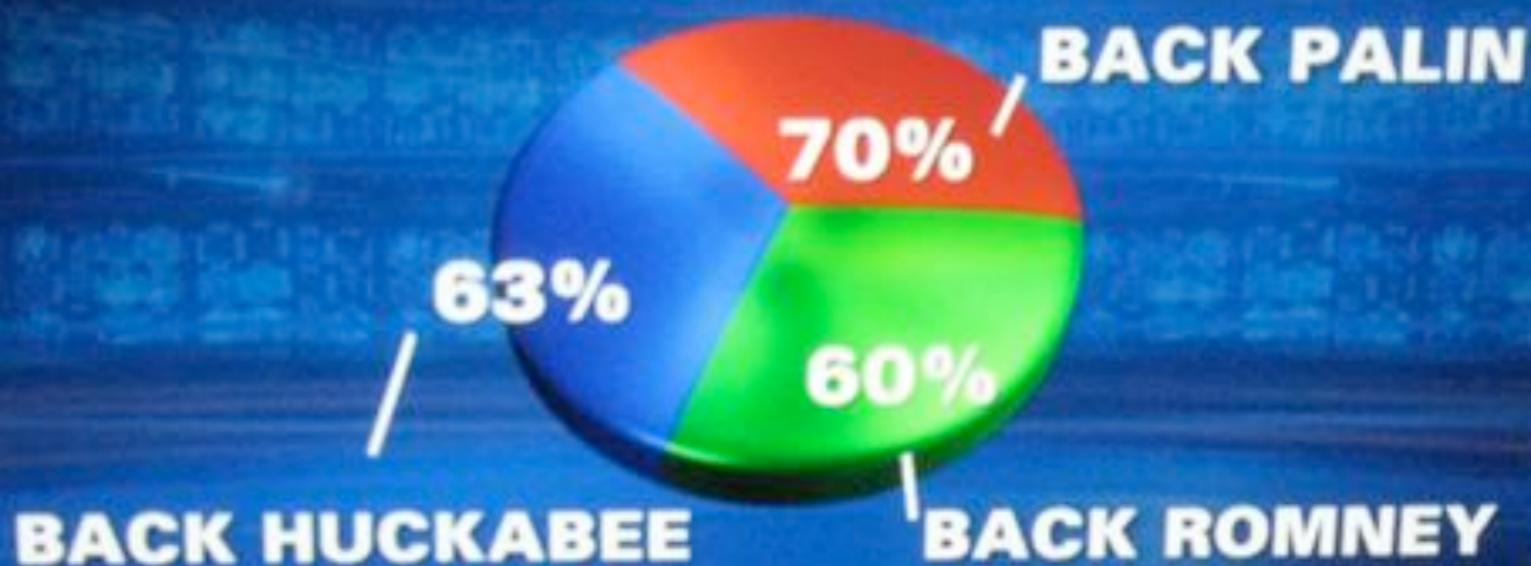
The Archivist, a Microsoft Desktop app to archive tweets

Microsoft Word Features By Version Added



2012 PRESIDENTIAL RUN

GOP CANDIDATES



FOX

47°

SOURCE: OPINIONS

DYNAMIC

misurare

measuring

a **process**

assigns numbers to things

**relations between things must
correspond to relations between
numbers**

things

numbers

A 

10

B 

20

B longer than A

20 > 10

L(A conc. B) = L (A) & L(B)

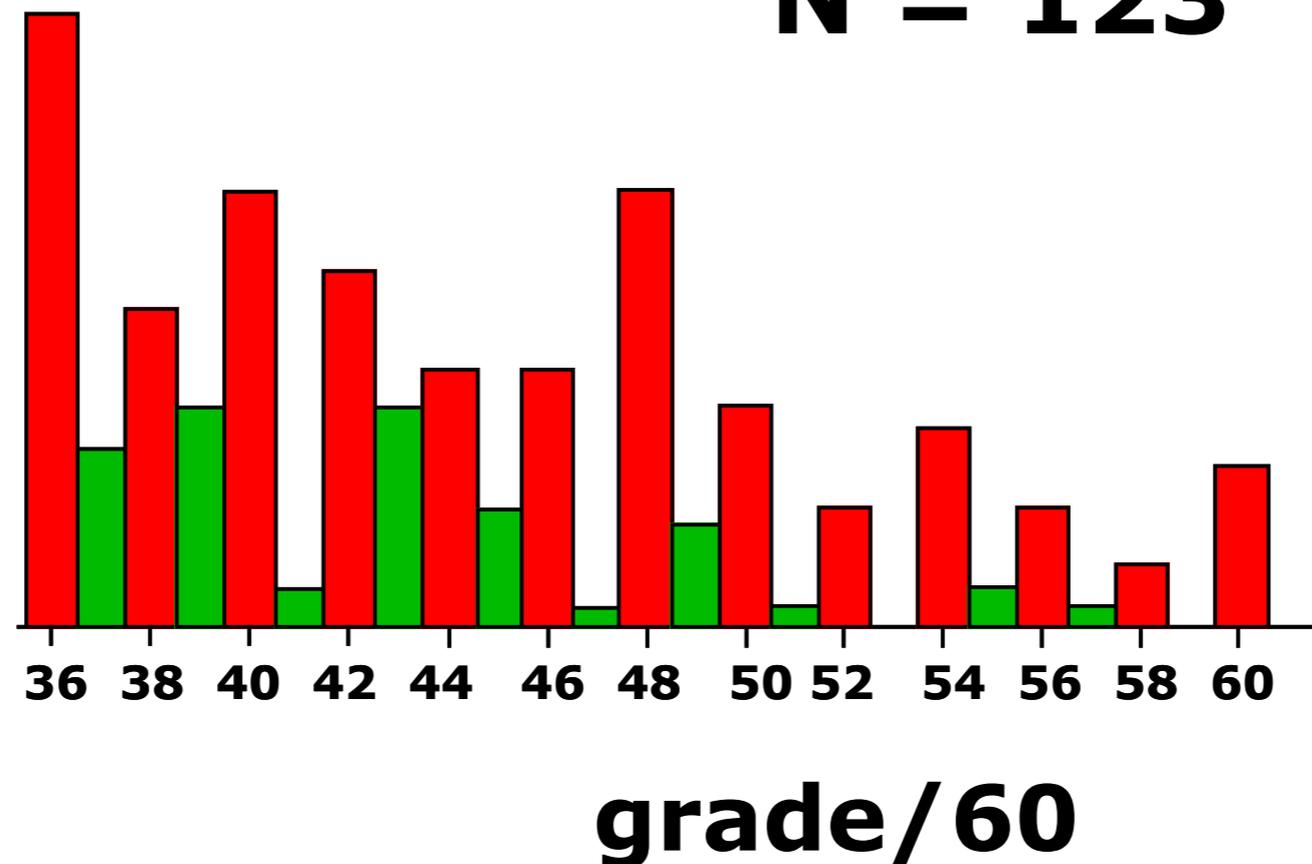
30 = 10 + 20

relations

**to understand the
measure, you must
understand the
process**

Italian school grades 1993/94

N = 123



distinzione utile

validità

cosa stiamo misurando?

il "significato"

affidabilità

stiamo misurando bene?

la "replicabilità"

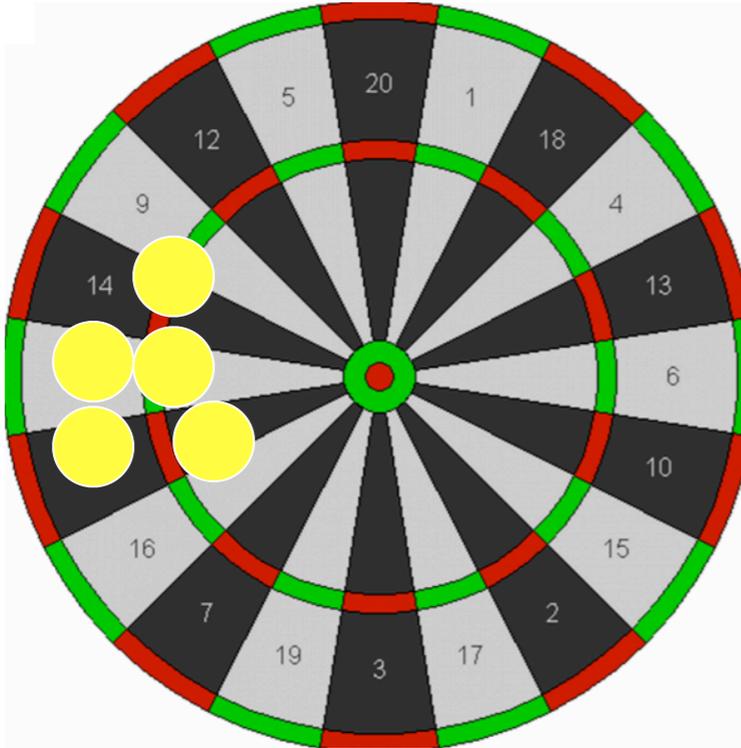
errors

variable error (random variation)

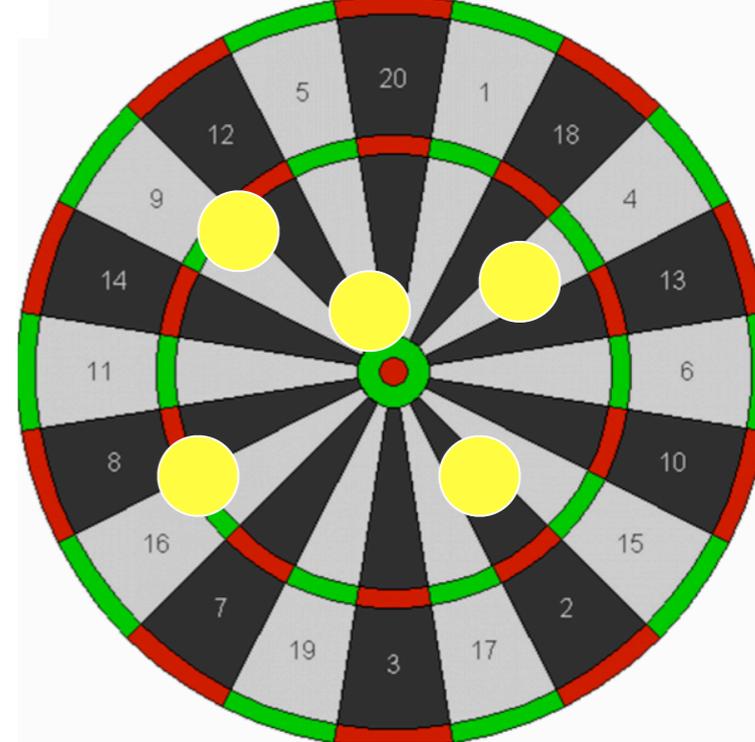
constant error (systematic bias)

precision: inverse of VE

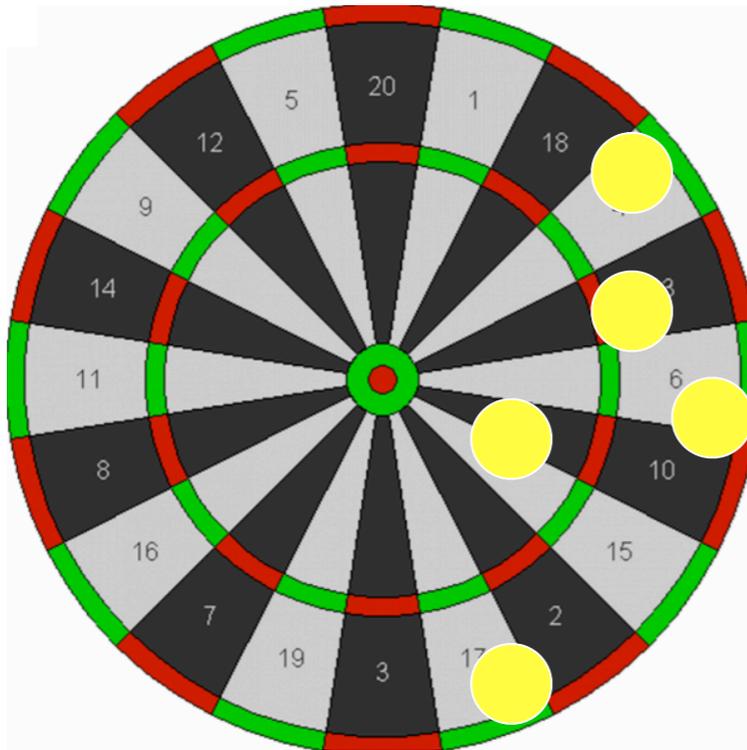
accuracy: inverse of CE



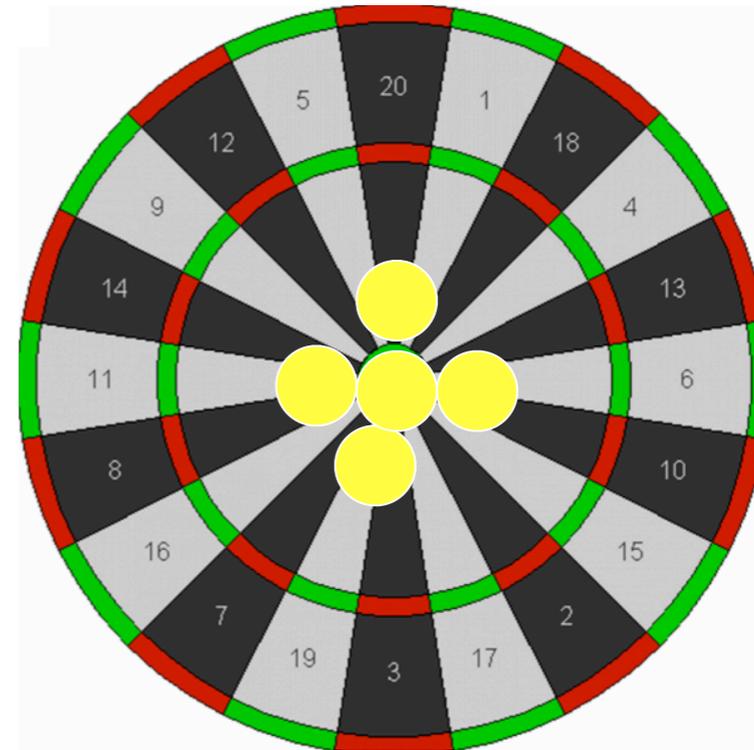
high precision
low accuracy



low precision
high accuracy



**low precision
low accuracy**



**high precision
high accuracy**

R

simulare una serie di misure

la funzione `rnorm()`

la funzione `hist()`

**scale di misura e
tipologie di dati**



Stanley S. Stevens
(1906 - 1973)

Stevens, S. S. (1946).
On the Theory of Scales
of Measurement.
Science, **103** (2684),
677–680.

Scale types with their properties according to Stanley Smith Stevens				
	Nominal scale	Ordinal scale	Interval scale	Ratio scale
Logical/ math operations	×	X	X	✓
	÷	X	X	✓
	+	X	X	✓
	-	X	X	✓
	<	X	✓	✓
>	X	✓	✓	✓
=	✓	✓	✓	✓
≠	✓	✓	✓	✓
Examples: <i>Dichotomous and non-dichotomous</i> Variable name (data values)	<i>Dichotomous:</i> Gender (male vs. female) <i>Non-dichotomous:</i> Nationality (American/Chinese/etc)	<i>Dichotomous:</i> Health (healthy vs. sick), Truth (true vs. false), Beauty (beautiful vs. ugly) <i>Non-dichotomous:</i> Opinion ('completely agree/ 'mostly agree/ 'mostly disagree/ 'completely disagree')	Date (from 1457 BC to AD 2013) Latitude (from +90° to -90°)	Age (from 0 to 99 years)

Nominal, Ordinal, Interval, and Ratio Typologies are Misleading

Paul F. Velleman

Cornell University and Data Description, Inc.

Leland Wilkinson

SYSTAT, Inc. and Northwestern University

The American Statistician (1993),
47:1, 65-72.

It is relatively easy to construct situations in which the scale type of data depends on its interpretation or on what additional information is available.

At a reception sponsored by the ASA Section on Statistical Computing and the Section on Statistical Graphics, consecutively numbered tickets, starting with “1”, were allotted at the door as people entered so that a raffle could be held.

As a winning number, 126, was selected and announced, one participant compared it to her ticket to see if she had won, thus interpreting the “126” correctly as a nominal value.

She then immediately looked around the room and remarked that “It doesn't look like there are 126 people here”, now interpreting the same value, again correctly (but using the additional information that tickets had been allotted consecutively starting with 1), as a ratio-scale value.

One of the authors compared his ticket number (56) to the winning value and realized that he had arrived too soon to win the prize, thus interpreting the values ordinally.

If additional data about the rate and regularity of arrivals had been available, he might have tried to estimate by how much longer he should have delayed his arrival from the 70-ticket difference between his ticket and the winner, thus treating the ticket number as an interval-scaled value.

A common dataset reports facts about automobiles. One of these facts is the number of cylinders in the engine. In some analyses, this is a nominal category supporting such questions as “Are there significant differences among the gas mileages of cars with 8-cylinder, 6-cylinder, and 4-cylinder engines?” Of course, these categories are clearly ordered, so ordinal-based statistics would also be appropriate. But one might also ask about the average number of cylinders in, say, U.S. cars, and wonder whether this average had declined in recent years. This requires us to consider these data values (all of them integers) as interval-scale values — which they can certainly be. Finally, we might consider the size of each cylinder and compute the ratio of each car’s displacement to the number of its cylinders — a completely appropriate operation (for ratio- scale data).

The point of these examples, of course, is that the assertion, common to many traditional statistics texts, that “data values are nominal, ordinal, interval, or ratio” simplifies the matter so far as to be false. Scale type, as defined by Stevens, is not an attribute of the data, but rather depends upon the questions we intend to ask of the data and upon any additional information we may have. It may change due to transformation of the data, it may change with the addition of new information that helps us to interpret the data differently, or it may change simply because of the questions we choose to ask.

<i>Left</i>		<i>Right</i>	Row Sum	
		1 1 1 1	4	<i>More</i>
	1 1 1 0	0 1 1 1	3	
1 1 0 0	0 1 1 0	0 0 1 1	2	
1 0 0 0	0 1 0 0	0 0 1 0	1	
	0 0 0 0		0	<i>Less</i>

In this scale, the horizontal dimension comprises a qualitative (nominal) scale of attributes and the vertical dimension measures a quantitative (ordinal, interval, or ratio) scale. For example, each profile might be the presence or absence of each of four symptoms in a patient. In this case, the vertical scale might be related to severity of illness and the horizontal scale might be related to different syndromes.

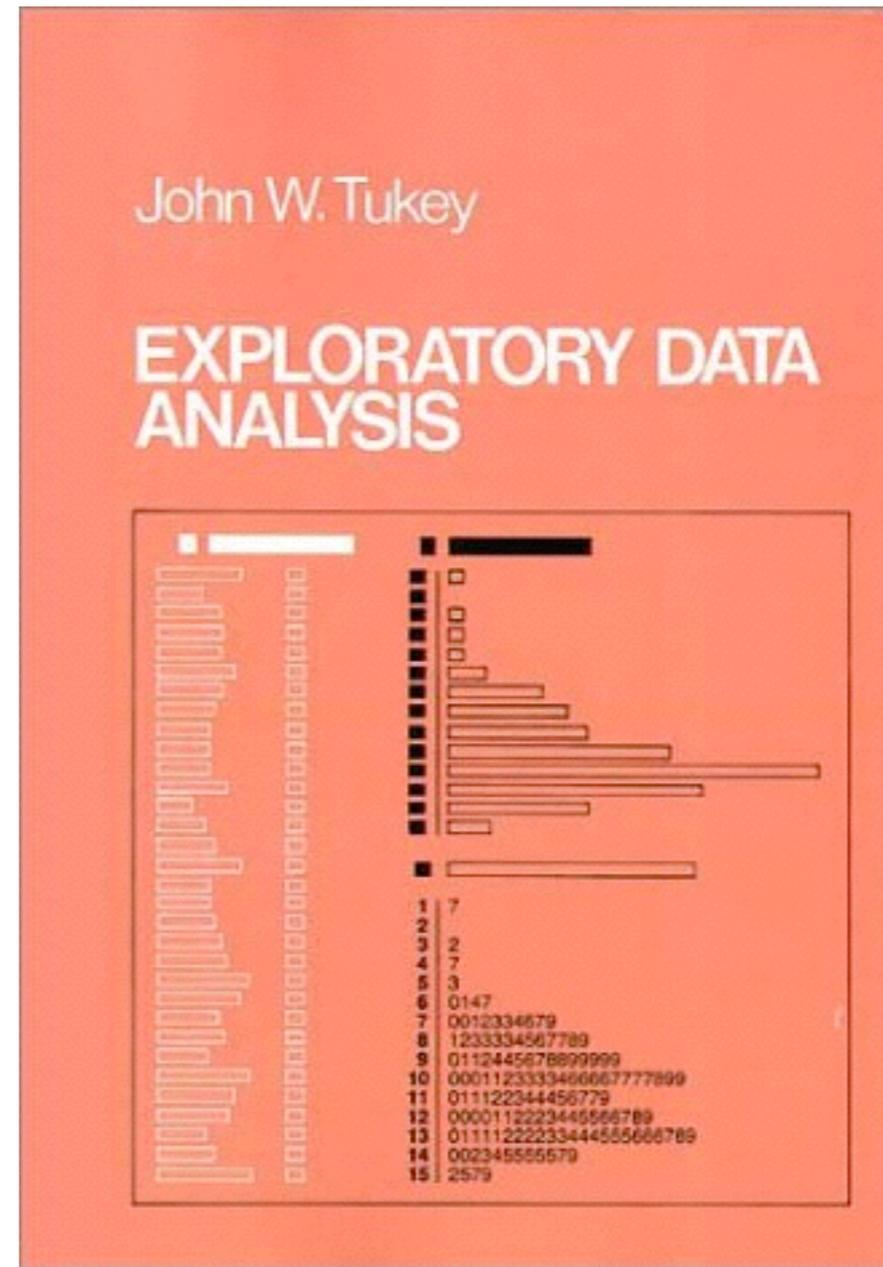
Conclusion

Measurement theory is important to the interpretation of statistical analyses. However, the application of Stevens's typology to statistics raises many subtle problems.

Statistics programs based on Stevens's typology suggest that doing statistics is simply a matter of declaring the scale type of data and picking a model. Worse, they assert that the scale type is evident from the data *independent of the questions asked of the data*. They thus restrict the questions that may be asked of the data. Such restrictions lead to bad data analysis and bad science.



John Tukey



tipologia

nomi

gradazioni

ranghi

frazioni *counted*

conteggi

quantità

bilanci

caratteristiche

etichette

categorie ordinate

classifica 1 - n

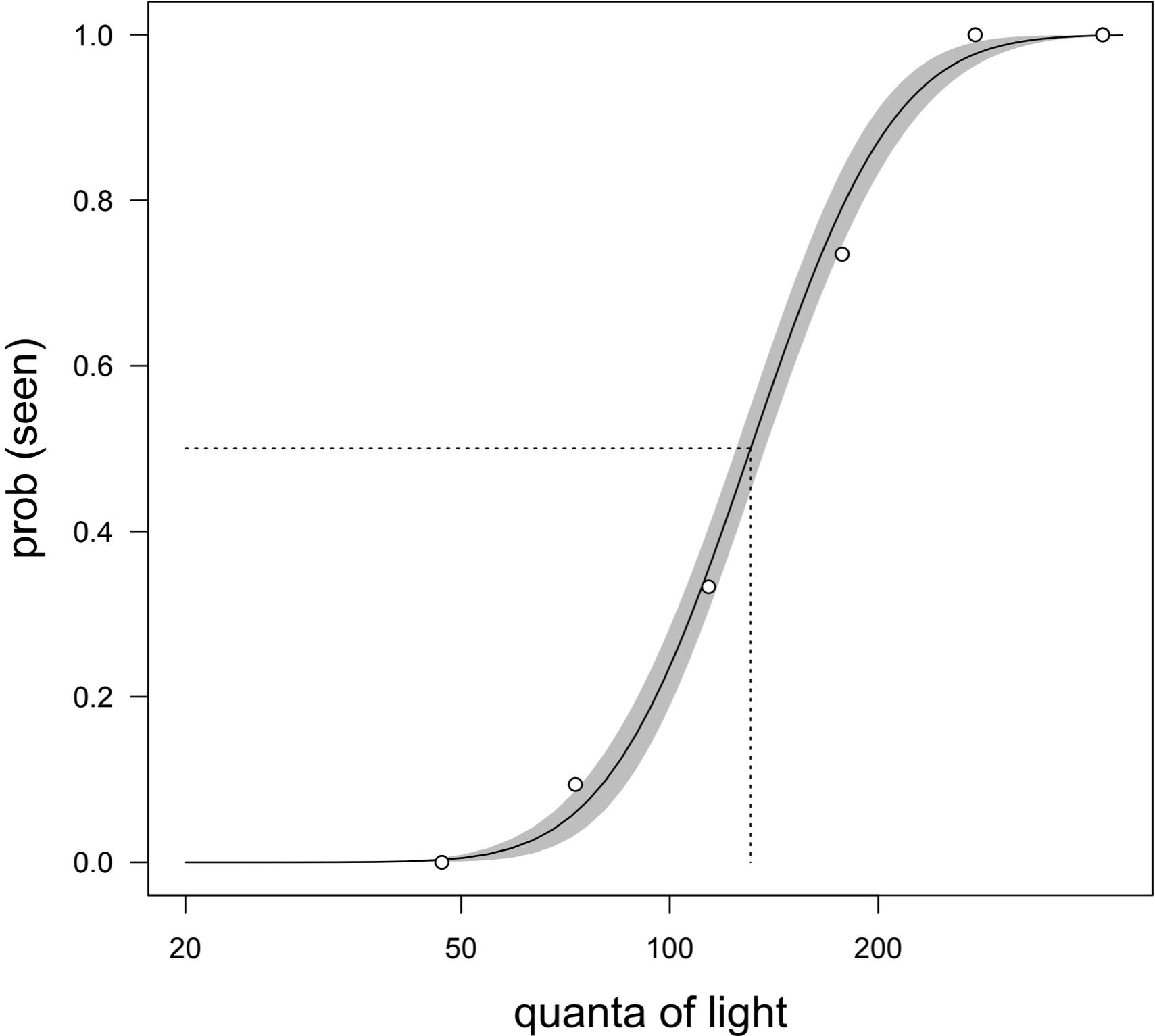
proporzioni 0 -1 o %

interi non negativi

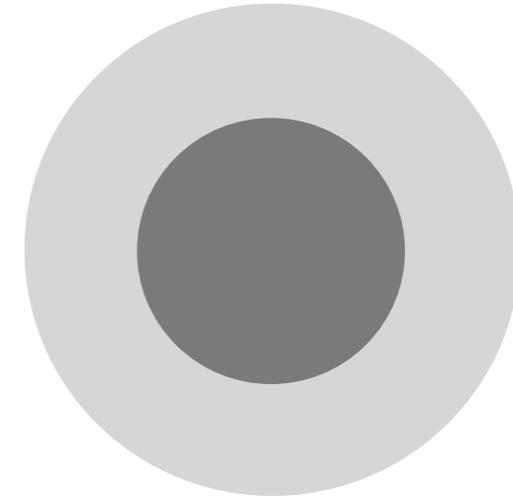
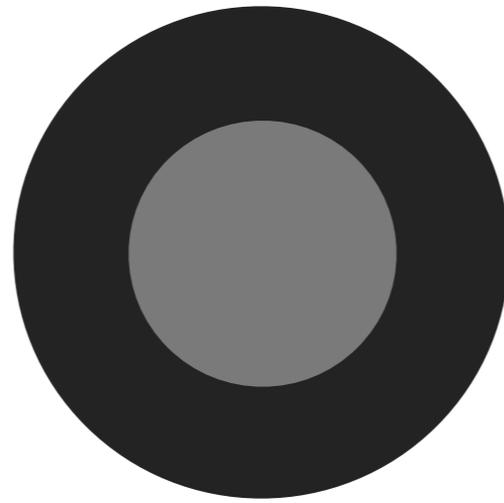
reali non negativi

reali negativi e positivi

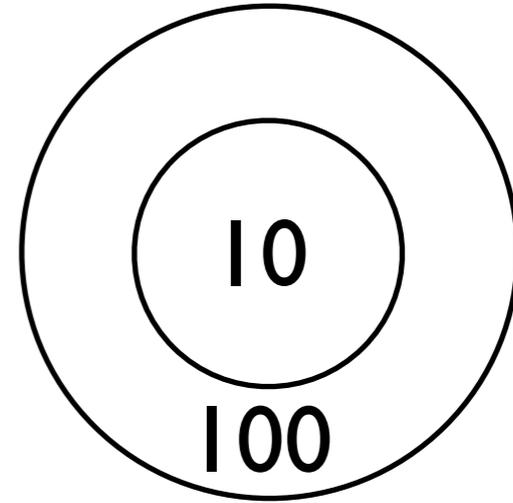
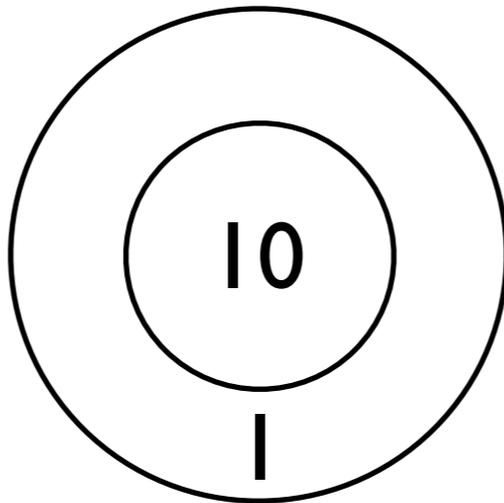
observer S.H.



configurazioni
disco-anello



luminanze



rapporti fra
luminanze

$$10/1 = 10$$

$$10/100 = 0.1$$

bilancio:
distanze

$$\log_{10}(10) = 1$$

$$\log_{10}(0.1) = -1$$