

Tipos de Graficos

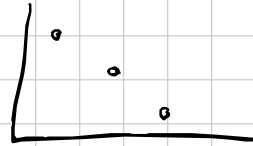
Qualitative data

PIE CHART



→ % total

Dot plot



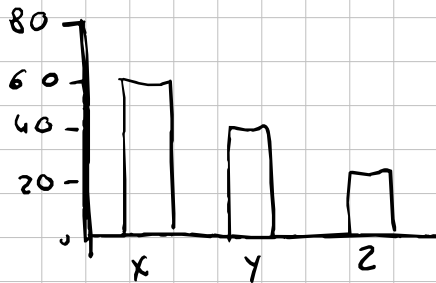
→ Categorias

Quantitative - Bar graph

→ Numbers

→ Distance

...



Histogram

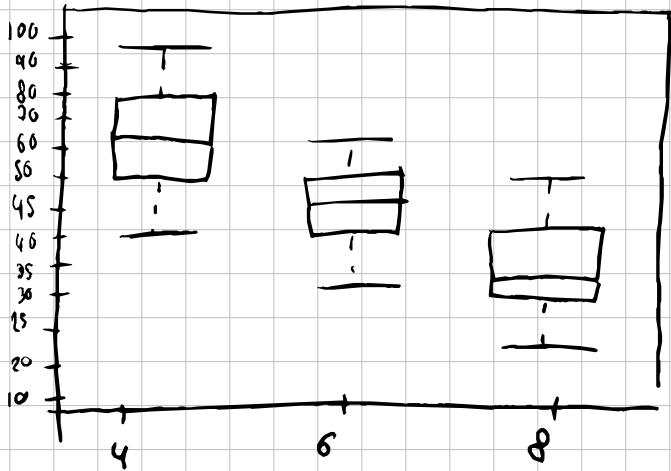
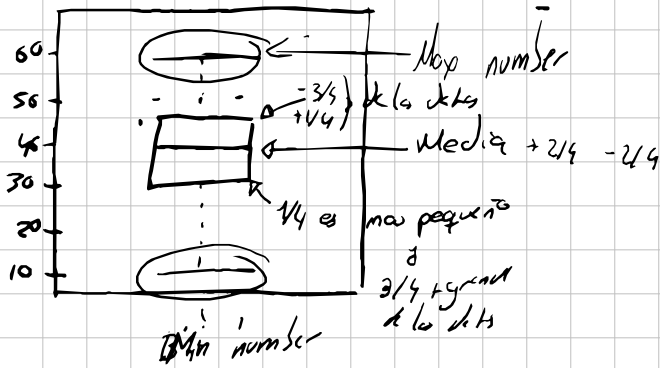
→ Permite adaptar el tamaño de los bloques a la frecuencia

→ Densidad: (altura de la barra)

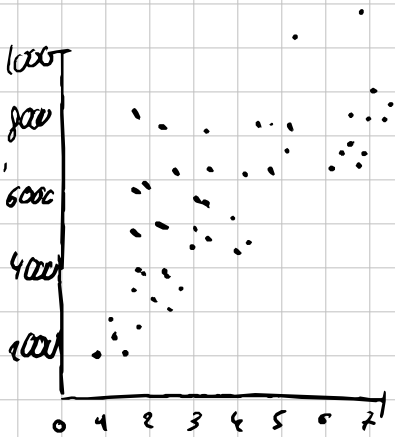
→ Porcentajes:



Box-Plot - Caja y bigote



Scatterplot - Dispersión



Uso principalmente
para visualizar
la relación entre dos
valores

Como por ejemplo
años de estudio y
salarios

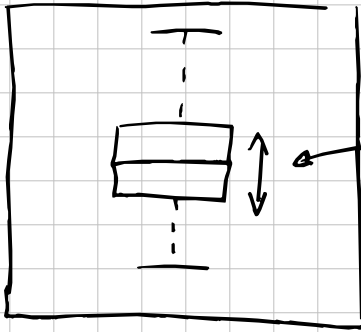
* Es importante dar todo el contexto necesario
, pero tampoco excesivo

Numerical Summary Measures

- ▷ Mean = average = media
- ▷ Median: punto medio, la mitad de los datos más Mediana alto y mitad más bajos.
- ▷ Mode: más repetido

- ▷ 1st Quartile o cuartil: 25% más bajo
- ▷ 2nd Quartile: media
- ▷ 3rd Quartile: 75% más bajo

Percentiles, the five Numbers Summary, and Standard Deviation



Interquartile range
3rd quartile - 1st quartile
it measures how spread the data is

Standard Deviation: Same with the values

\bar{x} average of the numbers $x_0 \dots x_n$

formula
$$s = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

Takes the square root of the average of each number minus the average

EXAMING OF HOW TO PRODUCE DATA

STATISTICAL INTERFERENCE

- Pick a random sample of the population to use it to calculate everything
- Population: whole dataset
- Parameter: quantity / value / data we are interested from the population
- Sample - muestra
- Statistic (estimate): parameter solo medido en el sample

Simple Random vs Stratified Random

SAMPLING

► Sample of convenience:

No suele ser una buena manera de sacar datos. Un ejemplo de este tipo sería coger solo clientes de un barrio en lugar de toda España.

Bias. Sample que favorece un resultado concreto.

Selection Bias: sample of convenience hace más probable uno u otro sobre otro.

Non-responsive bias: la gente que responde puede ser diferente a la que no responde.

voluntary response: generalmente los reviews solo vienen de clientes con o muy buenas o muy malas experiencias.

► Simple random sample

• Select subject without replacement

► Stratified Random Sample

Dividir el population en grupos de topicos similares "strata".

Bias and Chance Error

$$\text{Estimate} = \text{parameter} + \underbrace{\text{Bias}}_{\substack{\text{bias / error} \\ \text{partly caused} \\ \text{by the} \\ \text{sample}}} + \text{chance error}$$

Observation Studies

Son los resultados de comparación de datos
atribuir de una variable aunque no tiene
porque sea directamente la que altera la otra
pero se predice que una variable este ligada
a otra determinada.

Ej. + red meat =? + cancer

Pero realmente no es la carne roja porx es que la
gente que come carne roja suele hacer menos
ejercicio y beber mas alcohol

A este alcohol + ejercicio se le llama confounding variable

Para poder asegurar un vinculo es necesario un
experimento para asegurar su efecto.

Similar a los estudios medicos y los
grupos medicados y de control con el placebo

Un experimento debe ser double-blind ni los pacientes ni los examinadores deben saber
en que grupo esta el paciente

Randomization

▷ \mathcal{R}

Probability

When A & B dependent $P(A \text{ or } B) = P(A) + P(B)$

Independent A & B

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

$$P(B|A) = \frac{P(A \& B)}{P(A)}$$

$$P(A \& B) = P(A) \cdot P(B|A)$$

• Useful in independent P

Bayes rule

$$\begin{aligned} P(B|A) &= \frac{P(A \& B)}{P(A)} = \frac{P(B \& A)}{P(A)} = \\ &= \frac{P(A|B) \cdot P(B)}{P(A)} \end{aligned}$$

$$\left\{ \begin{array}{l} P(D|+) = \frac{P(+|D) P(D)}{P(+)} \Rightarrow \\ \frac{P(+|D) P(D)}{P(+|D) P(D) + P(+|no D) P(no D)} \end{array} \right.$$

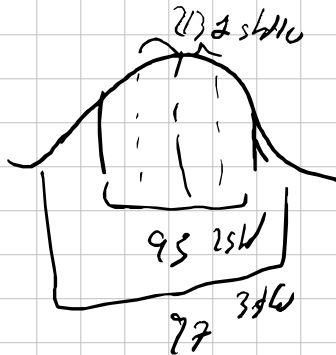
Normal curve - bell shape

Empirical rule

1 - 2/3 data in 1 Std dev

95% in 2 standard

99.7 in 3 standard



Standardizing data

$$Z = \frac{\text{data} - \bar{x}}{s} = Z\text{-score (No Un. 1)}$$

~~t, k, k~~

mean 0

std dev 1

Normal Approximation

Calculate % of height between

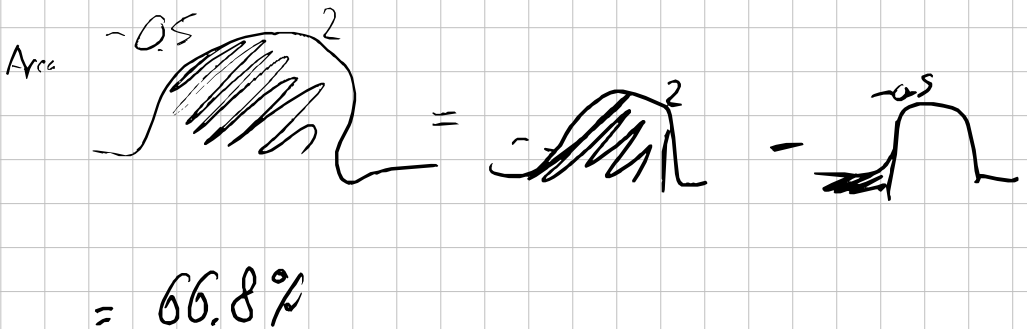
67.4 & 71.9

mean = 68.3

$s = 2.8$

$$\frac{67.4 - 68.3}{2.8} = -0.5$$

$$\frac{71.9 - 68.3}{2.8} = 2$$



Calcula % normal data

$$Z = \frac{\text{valor} - \bar{x}}{s}$$

$$\text{valor} = \bar{x} - Zs$$

Binomial Settings and Coefficient

2 in 3 in a 49% setting

$$P(2 \text{ in } 3) = P(AAB \text{ or } ABA \text{ or } BAA) =$$

$$= P(A) \cdot P(A) \cdot P(B) +$$

$$P(A) \cdot P(B) \cdot P(A) +$$

$$P(B) \cdot P(A) \cdot P(A) =$$

$$= \underline{3} \times (0.49)(0.49)(0.51) =$$

Independent
repetitions = m

for 2 in 5

$$5 \cdot (0.49)(0.49)(0.51)(0.51)(0.51)$$



Binomial Coefficient

$$m = \frac{n!}{k!(n-k)!}$$

Binomial Formula

$$P(k \text{ success in } n \text{ experiments}) = \frac{n!}{k! \cdot (n-k)!} p^k (1-p)^{n-k}$$

Expected Value and SE for the sum

$$S_n = \text{sum } n \text{ draws} \quad S_n = n \bar{x}_n$$

$$E(S_n) = n\mu \quad SE(S_n) = \sqrt{n}\sigma$$

The more n the more the SE goes up

$$\begin{array}{l} \text{SE} \\ \text{more } n \end{array} \left[\begin{array}{l} + SE(S_n) \\ - SE(n) \end{array} \right]$$

$$\begin{array}{l} \text{less } n \\ \text{SE} \end{array} \left[\begin{array}{l} + SE(n) \\ - SE(S_n) \end{array} \right]$$

Simulating Values: X has k outcomes

$$\mu = \sum_{i=1}^k x_i P(\bar{X} = x_i)$$

$$\sigma^2 = \sum_{i=1}^k (x_i - \mu)^2 P(\bar{X} = x_i)$$

X has ∞ outcomes, density f
(E_i which follow normal curve) ~~to~~ we will not use this

$$\mu = \int_{-\infty}^{\infty} x f(x) dx \quad \sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

Square root Law