

MS-A0503 First course in probability and statistics

3B Statistical datasets

Jukka Kohonen

Department of mathematics and systems analysis
Aalto SCI

Academic year 2021–2022
Period III

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

What is statistics (as a science)?

- Applying and developing methods for studying **random or uncertain** phenomena in the *real world*.
 - The methods are based on the mathematical laws of probability.
 - **Sources of uncertainty** are many: physical randomness; unknown properties of the real world phenomenon; random sampling on purpose; measurement errors; missing data . . .
 - . . . generally, the same math applies.
- Roughly:
 - Probability theory tells: How a certain process **produces** data.
 - Statistics tells: What **was** the process that produced the data.
- Statistics is applicable whenever you have data; especially if there is any kind of uncertainty or randomness.
 - Most fields of engineering and business have data, so they can (and do) use statistics.

Two basic approaches of statistics

Descriptive statistics

Present and **describe** the data “as it is”, either fully, or in a summary way.

- Tables (“the raw data”)
- Graphs (visualization)
- Numerical summaries or “statistics” (e.g. average, minimum, maximum)

Statistical inference

Infer facts about the real phenomenon that lies “behind” the data.
Generalize, e.g. sample \rightarrow population; or measurements \rightarrow universal physical law.

- Stochastic models
- Parameter estimation
- Hypothesis testing

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

Statistical data

Typically (not always), statistical data is in a table, the **data frame**, where

- rows correspond to **units**, e.g. people
- columns are the **variables** observed for each unit, e.g. height

Caveat: Different fields of science/engineering use different words. E.g. “units” may be “objects”, “items”, “data points” (geometrically thinking), “records” (in databases)

Depending on number of variables, we may call our data *univariate*, *bivariate*, *multivariate*.

Levels of measurement = What kind of values

- **nominal scale** \approx **categorical**: just distinct classes
gender: {male, female}
piece of DNA: {**A**denine, **C**ytosine, **G**uanine, **T**hymine}
- **ordinal scale**: classes have meaningful order
cloth size: { XS < S < M < L < XL }
Likert scale: { str. disagree < disagree < neutral < agree < str. agree }
- **numerical**: values have arithmetic meaning
 - **interval scale**: differences $x - y$ are meaningful
calendar dates, Celsius temperature
 - **ratio scale**: also quotients x/y are meaningful
length, weight, distance, Kelvin temperature

Notes:

- all can be *represented* as numbers, e.g. adenine=1, cytosine=2, guanine=3, but arithmetic might not make sense.
- nominal sometimes called “qualitative”, but other meanings
- this is *not* the discrete/continuous distinction. Numerical data can be well discrete; e.g. **counts** (frequencies)

Data set (Data frame)

- **data set** = sequence of elements (units) of the same type, e.g. numbers, identifiers, or lists of values (one for each variable)
- Often arranged in a table; (R terminology) **"data frame"**
- Order of units often not meaningful, so we could treat it as a set (or *multiset*, if many identical observations possible)

E.g. course feedback: ((12345A, 5, 1, 5), (98759K, 1, 5, 2), (33312K, 4, 4, 3), (23453B, 4, 4, 3), (21453U, 3, 3, 3)), ((223344, 5, 5, 5)), ((98313A, 5, 5, 5))

7 units; 4 variables (one string and three numbers).

Student ID	General	Workload	Usefulness
12345A	5	1	5
98759K	1	5	2
33312K	4	4	3
23453B	4	4	3
21453U	3	3	5
223344	5	5	5
98313A	5	5	5

Average and standard deviation

If we have univariate numerical data: $\vec{x} = (x_1, \dots, x_n)$

Average (sample mean) $m(\vec{x}) = \frac{1}{n} \sum_{i=1}^n x_i$

Variance $\text{var}(\vec{x}) = \frac{1}{n} \sum_{i=1}^n (x_i - m(\vec{x}))^2$

Standard deviation $\text{sd}(\vec{x}) = \sqrt{\text{var}(\vec{x})}$

Eg. $\vec{y} = (0, 0, 1, 1, 2, 2)$

$$m(\vec{y}) = \frac{1}{6} (0 + 0 + 1 + 1 + 2 + 2) = 1$$

$$\text{var}(\vec{y}) = \frac{1}{6} \left((0-1)^2 + (0-1)^2 + (1-1)^2 + (1-1)^2 + (2-1)^2 + (2-1)^2 \right) = \frac{2}{3}$$

$$\text{sd}(\vec{y}) = \sqrt{\frac{2}{3}} \approx 0.8165$$

Caveat: Sometimes $n - 1$ used as divisor for variance and sd, for technical reasons; more about this later (in parameter estimation).

Example

Calculate sample mean and standard deviation for the following data sets

$$\vec{x} = (1, 1, 1, 1, 1),$$

$$\vec{y} = (0, 0, 1, 1, 2, 2),$$

$$\vec{z} = (0, 2, 0, 2, 0, 2, 0, 2, 0, 2),$$

$$\vec{w} = (\underbrace{0, 0, 0, 0, \dots, 0, 0, 0, 0}_{666666 \text{ times}}, 1000000, \underbrace{0, 0, \dots, 0, 0}_{333333 \text{ times}}).$$

Dataset	Mean	SD
\vec{x}	1	0.0000
\vec{y}	1	0.8165
\vec{z}	1	1.0000
\vec{w}	1	999.9995

Average and standard deviation are *summaries*, they do not tell everything about the data. (Just like in probability distributions.)

Computing the summary statistics *from data*

Notation	Name	R	Matlab	Excel
$m(\bar{x})$	Average	mean()	mean()	AVERAGE()
$sd(\bar{x})$ $sd_s(\bar{x})$	Standard deviation (Corrected) std.dev.	sqrt(1-1/n)*sd() sd()	std(,1) std()	STDEV.P() STDEV.S()
$var(\bar{x})$ $var_s(\bar{x})$	Variance (Corrected) variance	(1-1/n)*var() var()	var(,1) var()	VAR.P() VAR.S()
$cov(\bar{x}, \bar{y})$ $cov_s(\bar{x}, \bar{y})$	Covariance (Corrected) covariance	(1-1/n)*cov() cov()	cov(,1) cov(1)	COVARIANCE.P() COVARIANCE.S()
$cor(\bar{x}, \bar{y})$	Correlation	cor()	corrcoef()	CORREL()
$q_{0.5}(\bar{x})$ $q_{0.25}(\bar{x})$ $q_{0.75}(\bar{x})$	Median Lower quartile Upper quartile	median() quantile(, .25) quantile(, .75)	median() quantile(, .25) quantile(, .75)	MEDIAN() PERCENTILE.INC(, .25) PERCENTILE.INC(, .75)

Caveat. Terminology and notation varies across sources. For technical reasons, many computer programs offer the so-called “unbiased” or “Bessel-corrected” *sample variance* and *sample standard deviation*, where the divisor is $n - 1$ instead of n . (Don’t worry now – we will talk about this in parameter estimation.)

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

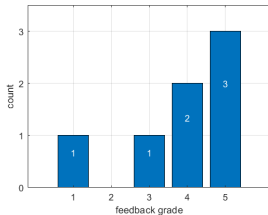
Counts = (absolute) frequencies of values

The **count**, or (absolute) **frequency** of a value x , in the univariate dataset \vec{x} , is

$$n_{\vec{x}}(x) = \#\{i : x_i = x\}$$

Course feedback, pick one variable “General” → univariate data (5, 1, 4, 4, 3, 5, 5). Frequency as a table and a bar chart:

x	1	2	3	4	5
$n_{\vec{x}}(x)$	1	0	1	2	3



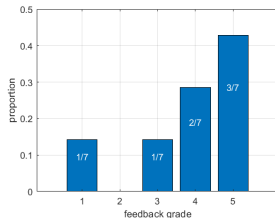
Proportions = relative frequencies

The **proportion**, or **relative frequency** of value x in dataset \vec{x} is

$$f_{\vec{x}}(x) = \frac{n_{\vec{x}}(x)}{n} = \frac{\#\{j : x_j = x\}}{n}$$

Course feedback, pick “General”, dataset (5, 1, 4, 4, 3, 5, 5), relative frequencies as a table and bar chart

x	1	2	3	4	5
$f_{\vec{x}}(x)$	$\frac{1}{7}$	0	$\frac{1}{7}$	$\frac{2}{7}$	$\frac{3}{7}$



Observation: $\sum_x f_{\vec{x}}(x) = 1$, thus $f_{\vec{x}}(x)$ is a probability distribution!
It is the **empirical distribution** of the dataset \vec{x} .

Empirical distribution

Proposition

If an element X is chosen uniformly at random, from the dataset $\vec{x} = (x_1, \dots, x_n)$, then X is a discrete random variable, whose density corresponds to the empirical distribution: $f_X(x) = f_{\vec{x}}(x)$. Furthermore,

$$\mathbb{E}(X) = m(\vec{x}), \quad (1)$$

$$\text{SD}(X) = \text{sd}(\vec{x}), \quad (2)$$

$$\text{Var}(X) = \text{var}(\vec{x}). \quad (3)$$

Also, for any function g , we have

$$\mathbb{E}[g(X)] = \frac{1}{n} \sum_{i=1}^n g(x_i). \quad (4)$$

Example

For the dataset $\vec{y} = (0, 0, 1, 1, 2, 2)$, determine the empirical distribution, and its mean and standard deviation.

The relative frequencies are

y	0	1	2
$f_{\vec{y}}(y)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

If random variable Y has density $f_{\vec{y}}(y)$, then

$$\mathbb{E}(Y) = \sum_{y=0}^2 y f_{\vec{y}}(y) = 0 \times \frac{1}{3} + 1 \times \frac{1}{3} + 2 \times \frac{1}{3} = 1,$$

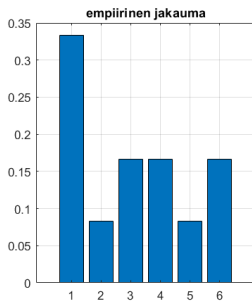
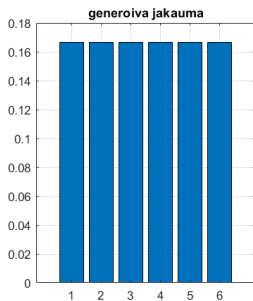
$$\text{Var}(Y) = \sum_{y=0}^2 (y-1)^2 f_{\vec{y}}(y) = (0-1)^2 \times \frac{1}{3} + (1-1)^2 \times \frac{1}{3} + (2-1)^2 \times \frac{1}{3} = \frac{2}{3}$$

$$\implies m(\vec{y}) = \mathbb{E}(Y) = 1$$

$$\implies \text{sd}(\vec{y}) = \sqrt{\text{var}(\vec{y})} = \sqrt{\text{Var}(Y)} = \sqrt{\frac{2}{3}} \approx 0.8165$$

Generating vs. empirical distribution

Rolling a fair die. Generating distribution (stochastic model) versus empirical distribution from 12 rolls that were (5,1,6,4,3,1,1,6,2,4,1,3).



- The empirical distribution *is* a probability distribution and you can use all tools of probability distributions.
- But *different distribution* from the generating one.
- For large n they are reasonably close (by LLN).

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

Binning (Grouping)

Eg. ages of all Finns 31.12.2015.

$n = 5\,487\,308$ units (data points)

Not probably good idea to draw as individual points (especially if ages are expressed in 1-day precision)

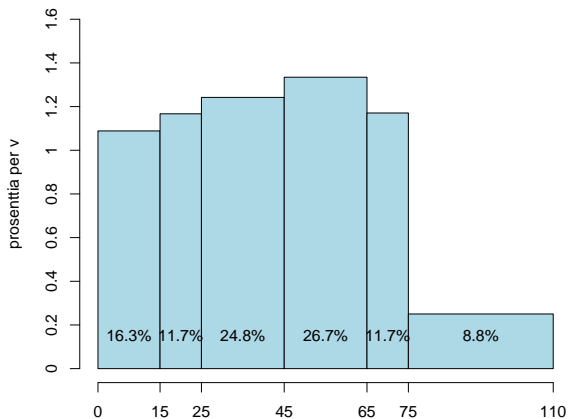
Let us **group** the data into **bins** at predefined boundaries.

Age (yr)	Frequency
0–14	896 023
15–24	640 387
25–44	1 363 155
45–64	1 464 640
65–74	642 428
75–	480 675

In fact we have transformed our real-valued variable (age) into a discrete one (group or bin index).

Example: Histogram with unequal widths

Finnish age distribution 31.12.2015 [Source: Tilastokeskus]



Age (yr)	Frequency
0-14	896 023
15-24	640 387
25-44	1 363 155
45-64	1 464 640
65-74	642 428
75-	480 675

The bars are an *approximation* of the true density function.

Could you find the proportion of Finns in the 1-year interval $[13, 14)$?

What about the the interval $[109, 110)$ years? Would it be accurate?

How to draw a histogram (allowing unequal bin widths)

- One bar for each bin (interval of possible values)
- Bar width = width of the interval
- Bar height = relative frequency *divided* by width

Example: Age distribution, first bar:

- Represents Finns with ages in interval $[0, 15)$
Note: values *strictly* smaller than 15; age in whole years 0–14
- Bar width = 15 years
- Frequency = 896023, relative frequency $896023/5487308 \approx 16.3\%$
- Bar height = $16.3/15 \approx 1.09$ (unit: % per year).
- Then bar *area* is the relative frequency.

(Typically, we use equal-width intervals, but not always.)

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

Bivariate data

Bivariate data = sequence (or multiset) of *pairs*

$$\vec{x}\vec{y} = ((x_1, y_1), \dots, (x_n, y_n)).$$

Alternatively, a pair (\vec{x}, \vec{y}) , where $\vec{x} = (x_1, \dots, x_n)$ and $\vec{y} = (y_1, \dots, y_n)$ are univariate data (note: as ordered sequences, so we know which x_i and y_i belong together).

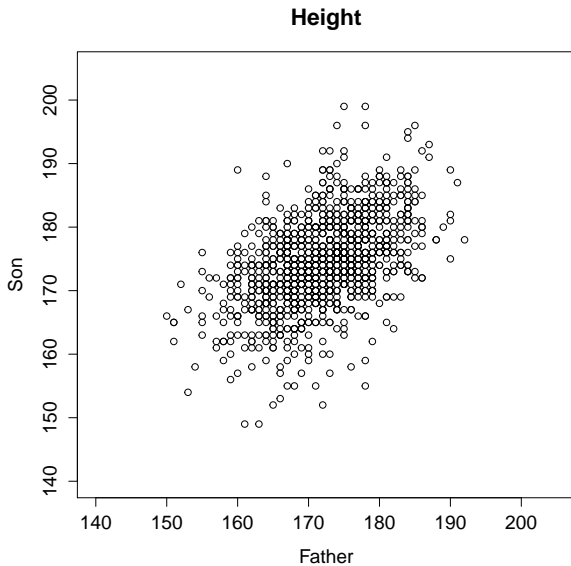
Course feedback: Two variables “General” and “Usefulness” composed as a bivariate dataset $((5,5), (1,2), (4,3), (4,3), (3,3), (5,5), (5,5))$

Univariate statistics $m(\vec{x}), m(\vec{y}), sd(\vec{x}), sd(\vec{y})$ are surely useful, but they tell nothing about the dependence between variables. Covariance and correlation tell (some aspects of) dependence.

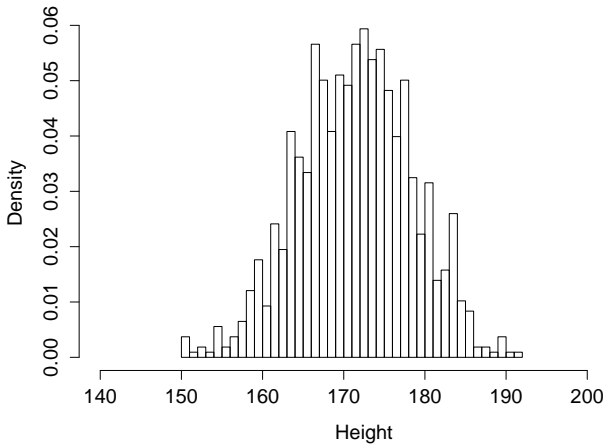
$$\text{cov}(\vec{x}, \vec{y}) = \frac{1}{n} \sum_{i=1}^n (x_i - m(\vec{x}))(y_i - m(\vec{y}))$$

$$\text{cor}(\vec{x}, \vec{y}) = \frac{\text{cov}(\vec{x}, \vec{y})}{sd(\vec{x}) sd(\vec{y})}$$

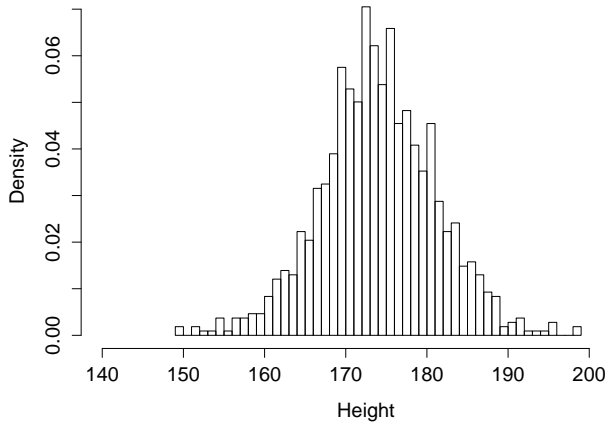
Scatterplot (scatter diagram)



Histogram of Fathers



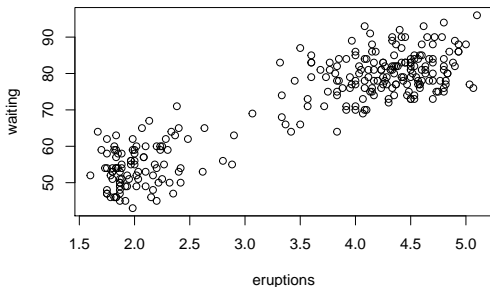
Histogram of Sons



Example: Eruptions of Old Faithful geysir

Scatterplot of 272 eruptions of *Old Faithful* (Yellowstone).

Two variables: eruption length and waiting time to next eruption.



Already a visual inspection (“eyeballing”) reveals interesting patterns (that are not captured by mean and standard deviation).

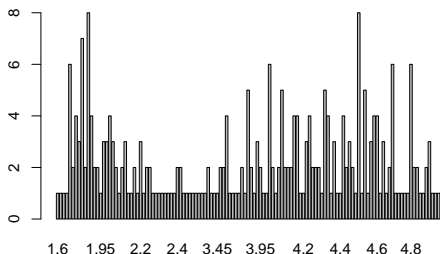
You can find the data in R, try `faithful` and `help("faithful")`

Old Faithful: bar chart of one variable...

We could try listing *all different values* of eruption length, and collect their frequencies (within $n = 272$)

x	1.6	1.667	1.7	1.733	1.75	...	5.1
$n_{\bar{x}}(x)$	1	1	1	1	6	...	1

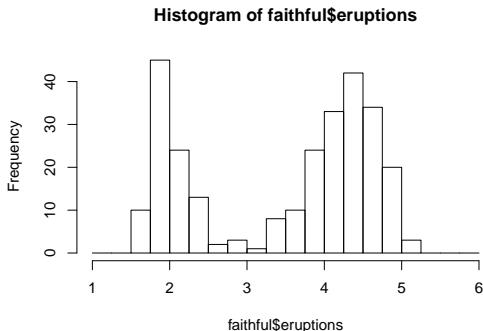
and draw a bar chart



Neither the table or the bar chart seems very informative.

Old Faithful: histogram

But if we group the data into 0.25-minute intervals such as $[2.00, 2.25)$, and plot the counts, we have a better picture of the distribution.



Try this on your own! What happens if you use more (finer) intervals? What if you use less (coarser) intervals?

Cross-tabulation (Contingency table)

The frequency of the pair (x, y) , in the data set, is

$$n_{\bar{x}\bar{y}}(x, y) = \#\{i : x_i = x \text{ and } y_i = y\}$$

Course feedback: “General” and “Usefulness” as bivariate dataset $((5,5), (1,2), (4,3), (4,3), (3,3), (5,5), (5,5))$ has this contingency table:

	<u>y</u>					
<u>x</u>	1	2	3	4	5	Sum
1	0	1	0	0	0	1
2	0	0	0	0	0	0
3	0	0	1	0	0	1
4	0	0	2	0	0	2
5	0	0	0	0	3	3
Sum	0	1	3	0	3	

Cross-tabulation of relative frequencies

The **relative frequency** of the pair (x, y) is

$$f_{\vec{xy}}(x, y) = \frac{\#\{i : x_i = x \text{ and } y_i = y\}}{n}$$

x	y					Sum
	1	2	3	4	5	
1	0	$\frac{1}{7}$	0	0	0	$\frac{1}{7}$
2	0	0	0	0	0	0
3	0	0	$\frac{1}{7}$	0	0	$\frac{1}{7}$
4	0	0	$\frac{2}{7}$	0	0	$\frac{2}{7}$
5	0	0	0	0	$\frac{3}{7}$	$\frac{3}{7}$
Sum	0	$\frac{1}{7}$	$\frac{3}{7}$	0	$\frac{3}{7}$	

$\sum_{x,y} f_{\vec{xy}}(x, y) = 1$, thus $f_{\vec{xy}}(x, y)$ is a probability distribution.
 $f_{\vec{xy}}(x, y)$ is the **empirical joint distribution** of the dataset \vec{xy} .

Empirical joint distribution

Proposition

If a pair (X, Y) is chosen uniformly at random, from the dataset $\vec{xy} = ((x_1, y_1), \dots, (x_n, y_n))$, it is a discrete random variable whose (joint) density is the empirical distribution $f_{X,Y}(x, y) = f_{\vec{xy}}(x, y)$ and also

$$\begin{aligned}\mathbb{E}(X) &= m(\vec{x}), & \mathbb{E}(Y) &= m(\vec{y}), \\ \text{SD}(X) &= \text{sd}(\vec{x}), & \text{SD}(Y) &= \text{sd}(\vec{y}), \\ \text{Var}(X) &= \text{var}(\vec{x}), & \text{Var}(Y) &= \text{var}(\vec{y}),\end{aligned}\tag{5}$$

and

$$\text{Cor}(X, Y) = \text{cor}(\vec{x}, \vec{y}),\tag{6}$$

$$\text{Cov}(X, Y) = \text{cov}(\vec{x}, \vec{y}).\tag{7}$$

Also, for any two-argument function g , we have

$$\mathbb{E}[g(X, Y)] = \frac{1}{n} \sum_{i=1}^n g(x_i, y_i).\tag{8}$$

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

Quantiles/Percentiles of data

Suppose data can be ordered from smallest to largest.

(OK for numerical or ordinal data; not for nominal data.)

If $0 < p < 1$, then the p -quantile (or $100p$ -percentile) $Q(p)$ is roughly the point x such that *proportion* p of the data is smaller than x , and $1 - p$ is greater.

- $Q(0.25)$ is **lower (first) quartile**; 25% of data is below
- $Q(0.5)$ is **median** or second quartile; 50% of data is below
- $Q(0.75)$ is **upper (third) quartile**; 75% of data is below

Note that half of data is between lower and upper quartiles.

R: `quantile(x,p)`, `summary(x)`, `median(x)`

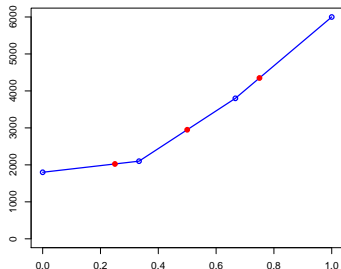
The “roughly” is because in finite data, you may not find exact quarters. There are some (varying) conventions for this.

Quantile function

One way to define the **quantile function** of dataset (x_1, \dots, x_n) :

- Order the data as $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$
- Divide the horizontal unit interval $[0, 1]$ into equal parts, at points $p_k = (k - 1)/(n - 1)$, $k = 1, \dots, n$
- Plot the points $(p_k, x_{(k)})$ and connect with lines

Example. Four **salaries** (eur/month): 3800, 1800, 2100, 6000



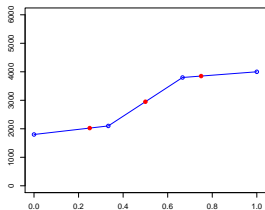
Quartiles = Evaluate the quantile function at 0.25, 0.50, 0.75

Example: Three small datasets

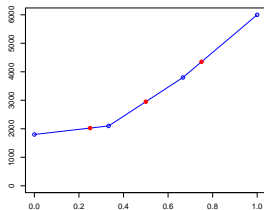
$$\vec{x} = (1800, 2100, 3800, 4000)$$

$$\vec{y} = (1800, 2100, 3800, 6000)$$

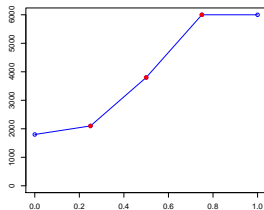
$$\vec{z} = (1800, 2100, 3800, 6000, 6000) \quad (n \text{ not divisible by four})$$



$$Q_x(0.50) = 2950,$$
$$m(x) = 2925$$



$$Q_y(0.50) = 2950,$$
$$m(x) = 3425$$



$$Q_z(0.50) = 3800,$$
$$m(x) = 3940$$

Contents

Introduction

Descriptive statistics

Empirical distribution

Histogram

Two-variable data (bivariate data, paired data)

Quantiles / Percentiles

From data to population

Sample, population and “population”

The finite dataset you have, the “sample”, is often thought to “represent” qualities of a larger “population”.

sample	population
Pearson's 1000 fathers and sons	All father-son pairs in ... ?
1000 poll responses	Opinions of 5 million Finns now
272 eruptions of Old Faithful	All its eruptions (past? future? potential?)
Drug effect on 30 patients	Drug effect on future patients
100 rolls of a loaded die	Potential infinite sequence of rolls

Population is statistical jargon for

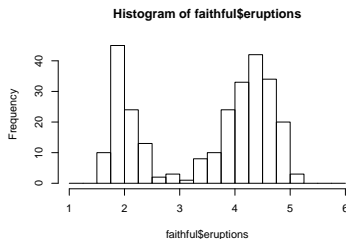
- where your data came from
(**data-generating mechanism**; **data source**)
- what you are trying to understand by looking at the data

Hence, terms such as *sample mean* and *population mean*.

The “population” may be quite concrete, or a figure of speech.

Old Faithful, once again

We have a *sample* of 272 eruption lengths. The real physical mechanism may be complicated, but perhaps we can think the lengths are **as if** they come from one particular **distribution** f . But what distribution?



One eruption length is a **random variable** X from some **generating distribution** or **underlying distribution** or **true distribution**. (“Population” if you want.)

Empirical distribution (shown above) approximates the generating distribution. Why? Think of the event $\{2.0 \leq X < 2.25\}$. We know by LLN that its *relative frequency* (long-run average) \approx its *probability* (in generating distribution).

Next lecture is about inference.

From the data, we will estimate some parameters concerning the reality “behind” the data.