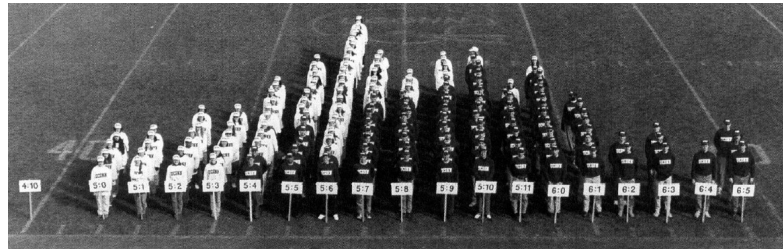


Introduction to Biostatistics



Arranged by Linda Strausbaugh (Genetics 147:5, 1997)

Introduction to Biostatistics

Lecture 1B and 2

Henrik Renlund

October 5-6, 2015



Contents of Lecture 1-2

Data and descriptive statistics

Probability theory and models

Sampling, SE(M) and CLT.

Visual tests (of normality)

What shall we learn today (and tomorrow)?

- Data description
 - Graphs
 - Tables and summary measures
- Probability Models
 - Glimpse at theory
 - Normal distribution and some properties
 - Some properties of samples and the Central Limit Theorem.

A data set contains one or more *variables*, typically arranged as:

Individual	Gender	Age	Albumin	Diabetes	Happiness
1	M	29	3.92	0	☺
2	F	37	4.12	0	☹
3	F	25	4.75	1	—
⋮	⋮	⋮	⋮	⋮	

Theoretical data categories:

- Numerical
 - discrete; typically integer valued 1, 2, . . . , like Age, or
 - continuous; i.e. any value in an interval, like Albumin.
- Categorical
 - nominal, e.g. Gender, or
 - ordinal, e.g. Happiness: ☹, —, ☺.

In practice data might be stored in the wrong format.

Check this prior to analysis! This is especially important if data has been transferred (between formats, different OS, etc).

Common problems:

- date- and categorical data suddenly stored as integers
- numerical values stored as text (due to ',./"')

Usually (bias alert!) one only needs to distinguish

- measurements (Age, Albumin are numerical), and
- categories (Gender, Diabetes, Happiness are categorical).

It is useful to provide a summary table of the variables you are working with. Choice of descriptive measures may be context dependent.

<i>variable</i>	Diabetes No		Diabetes Yes	
<i>value</i>	<i>mean</i>	<i>sd</i>	<i>mean</i>	<i>sd</i>
Age	32.0	15.9	32.5	14.1
Albumin	4.20	0.37	3.80	0.50
	<i>percent</i>	<i>n</i>	<i>percent</i>	<i>n</i>
Gender				
M	64%	27	52%	22
F	36%	15	48%	20
Happiness				
(61%	19	36%	15
—	23%	7	36%	15
)	16%	5	36%	12

- Table and caption should be self contained.
- Every table should be referred to in text.*
- Put captions *above* the table.
- Avoid excessive precision. (Mean age 65.63?)
- Use adequate measures of location and spread.

*The table should illustrate some part of the data relating to the hypothesis of the paper. Referring to a table should fit the narrative.

Generally bad: "Table 3 displays values of biomarkers in group A and B".

Better: "The mean level of measured biomarkers are higher in group A compared to B (Table 3)."

DESCRIPTIONS	MODELS	SAMPLING	VISUAL TESTS
○○○○●○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○	○○○○○○

What about missing data?

Missing data is usually a problem for observational studies.

Typically data is missing in explanatory variables.

Suppose we investigate death within one year of myocardial infarction with a model that includes gender, age, BMI (some missing) and smoking status (some missing).

The statistical software default is to include only those individuals with complete case data on all variables that you include.

Complete case analysis will only give an unbiased result if the reason that a variable is missing has nothing to do with the actual value.

If the missing mechanism is known this might be utilized in the analysis.

DESCRIPTIONS	MODELS	SAMPLING	VISUAL TESTS
○○○○○●○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○	○○○○○○

Solutions...ish

There is no trick that guarantees a non-biased analysis.

Single imputation: replace missing value with a "typical" value for that variable (e.g. mean or median). Better: try to model the variable and predict it from the other variables. Both methods underestimates the variance in the variable and will give overly optimistic results.

Multiple imputation: create multiple imputed data sets where the missing values are replaced differently in each, e.g. drawn at random from the non-missing values. Better: use (random) estimates from modeling the variable.

"It is not that multiple imputation is so good; it is really that other methods for adresssing missing data are so bad."

(Donald Rubin)

DESCRIPTIONS	MODELS	SAMPLING	VISUAL TESTS
○○○○○●○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○	○○○○○○

Visualization of (continuous) data

A sufficiently small data set might not need visualization.

The level of the protein albumin was recorded in a sample (of size 8) of mice (56 days old):

1.88 2.03 2.11 1.77 2.04 2.05 1.94 1.95

(measured in g/dL - this data set will return in lecture 3).

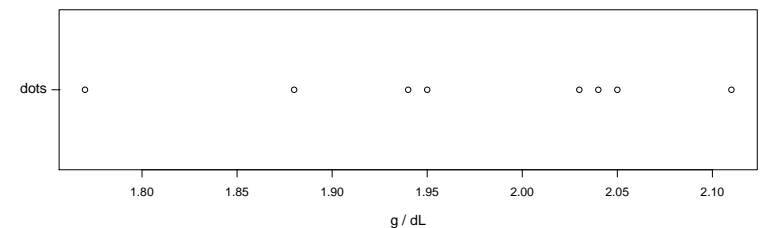
One simple way to get some handle on data is to order it:

1.77 1.88 1.94 1.95 2.03 2.04 2.05 2.11

DESCRIPTIONS	MODELS	SAMPLING	VISUAL TESTS
○○○○○●○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○○○○○○○	○○○○○○○○○○○○○○	○○○○○○

Dotplot of albumin data

A dotplot is a one dimensional plot of the data.



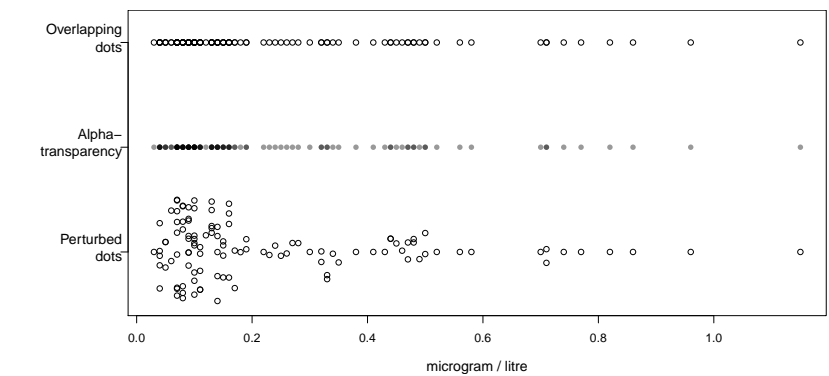
If there are non-unique (or close) points, the data set may appear smaller than it really is.

This can be alleviated by

- perturbation, or,
- alpha transparency.



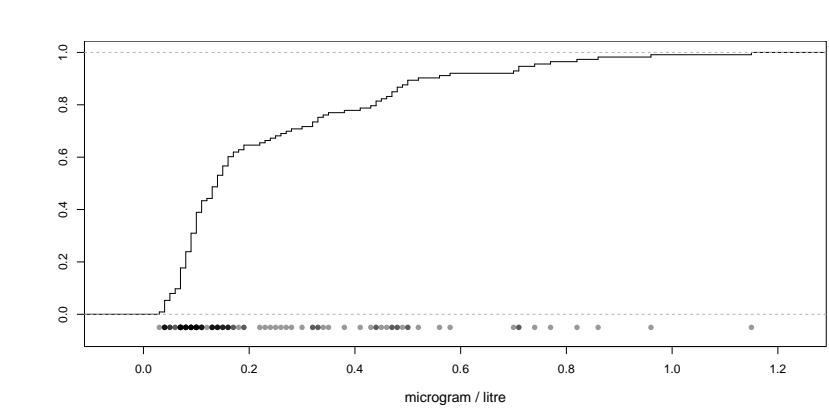
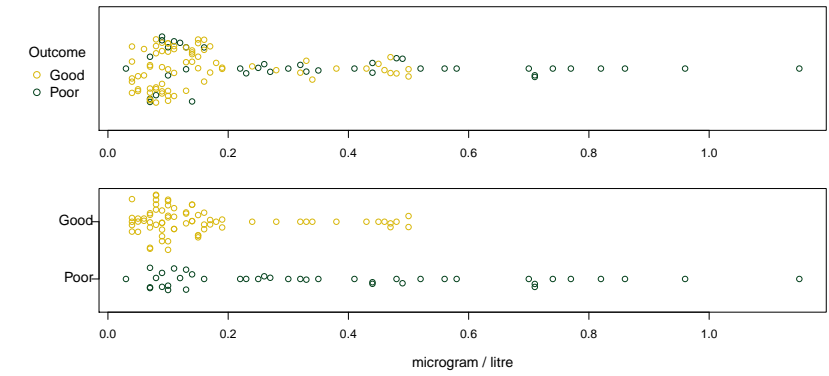
A biomarker - the protein S100 β - was measured for 113 individuals with aneurysmal subarachnoid hemorrhage. (To return in lecture 10.)



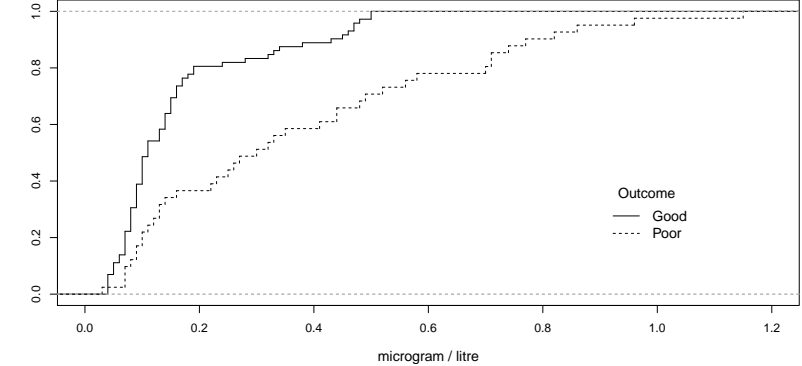
- The k th percentile is a value v such that k percent of your data lies below (or at) v . (Usually not uniquely defined.)
- The 50th percentile (the *median*) is the point which divides your ordered sample equally. (Only 'unique' if sample is odd, else use mean of the two midpoints.)
- *The Quartiles*: Q1 is the 25th percentile, Q2 is the 50th percentile and Q3 is the 75th percentile.
- We can describe all percentiles with the *empirical cumulative distribution function* (ecdf).
When the sample size is 113 the jumps in the ecdf will be multiples of $1/113 \approx 0.01$.



Dotplots can display groups.



Ecdf's can also display group differences.



An average value is a single value meant to be representative of the entire data set.

- **The mode:** is the single most common data point.
(Mostly meaningful for categorical data.)
- **The median:** is the midpoint of the ordered numerical sample.
- **The mean:** is the "center of gravity" of a data set.
Note: unlike the median, the mean value is sensitive to "extreme" values.

Ex: A small company has 5 employees, who earns 19, 21, 22, 24, 27 (K SEK) and a boss who earns 55.

Salaries	Employees	Entire company
Median	22	23
Mean	22.6	28

Some points:

- Small data sets might not need summary measures.
- Symmetric data has mean \approx median.
- (Easy enough to calculate both.)

- **Range** The difference between the maximum and the minimum value.
- **Interquartile range (IQR):** Q3-Q1.
- **Standard deviation (sd)** is given by the formula,

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

Where x_1, x_2, \dots, x_n is the sample ($n = 35$ in our example) and

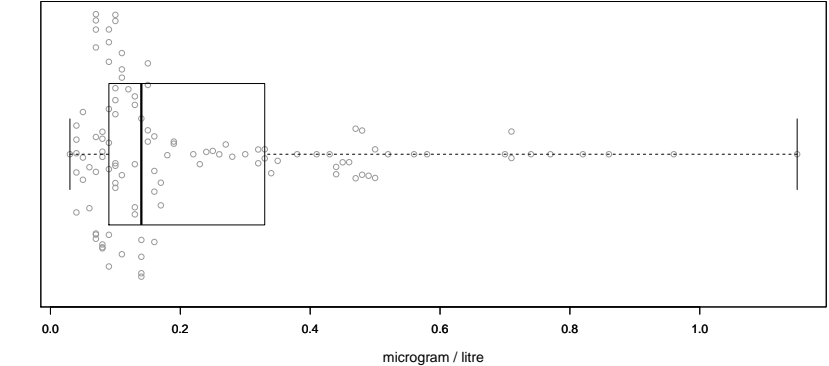
$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

is the (sample) mean.

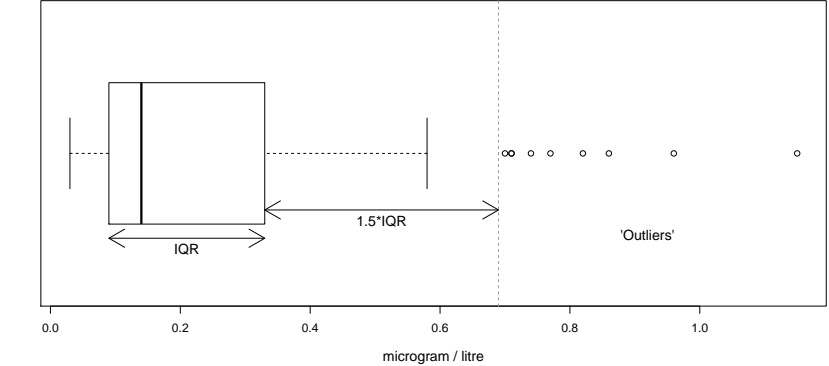
It is (approximately) the mean distance to the mean value.



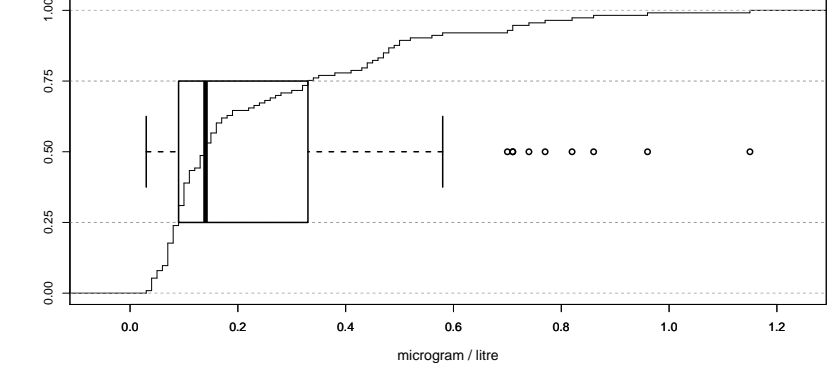
The boxplot usually show min, Q1, med, Q3 and max (the "5-point summary")...



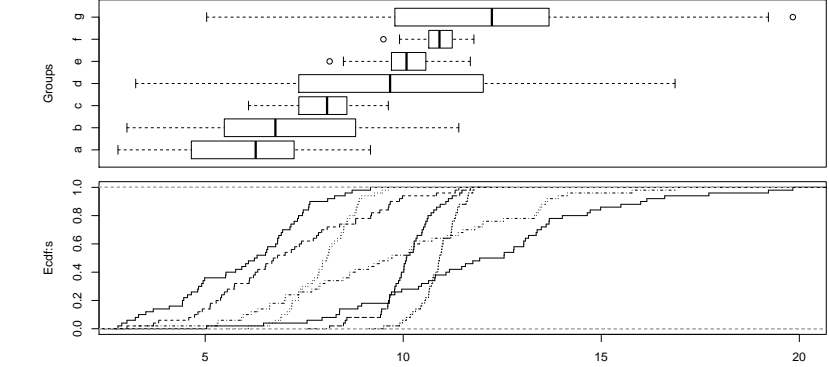
... but most software mark points that are more than 1.5 times the IQR away from 'the box'.



The boxplot contains less information.

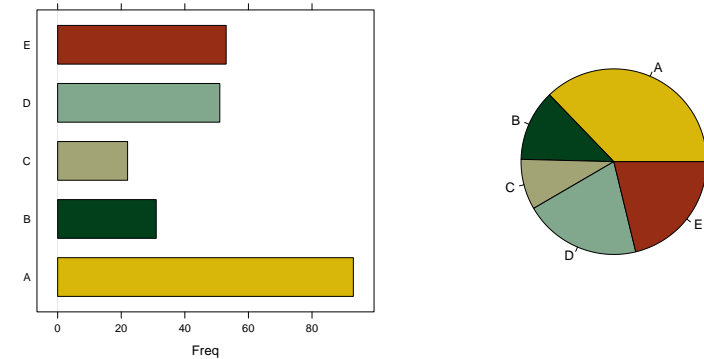


Here is fake data with 7 subgroups (a-g).



- Graph and caption should be self-contained.
- Every graph should be referred to in the text.
- 'Economy' Do not make a graph which is more easily expressed in text or a small table.
E.g. graph with a single boxplot.
- Pattern or detail?
E.g.
 - a ecdf gives a lot of detail of a data set.
 - graph with multiple boxplots can reveal pattern among subgroups.
- Avoid 2D graphs shown in 3D.
- Avoid pie charts? It depends...

Is it of interest to be able to compare summation of levels?
(Is $A + B + C$ smaller than $D + E$?)



- Probabilities are numbers between (and including) 0 and 1. (0% and 100%.)
- **Prob(A)** denotes "the probability of the event A".
- The basic rule of probability is that if A and B are mutually exclusive events the

$$\mathbf{Prob(A \text{ or } B) = Prob(A) + Prob(B).}$$

(Probability is a "measure".
But what does it measure?)

A yet undetermined random value is called a *random variable* (RV).

Suppose that in the population there are 49 % non-smokers, 20 % former smokers and 31% current smokers. Then the smoking status X of a person selected at random is a RV with a probability function

Value v	non	former	current
Prob($X = v$)	0.49	0.20	0.31

Integer-valued RV's are also described by probability functions.

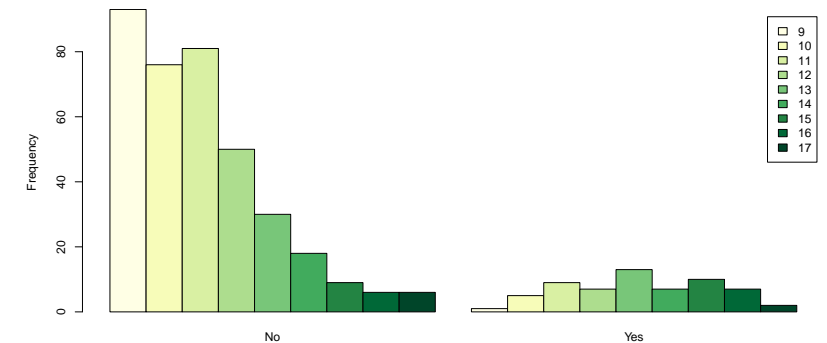
Some are easy to model realistically, let e.g. Z = 'the outcome of the throw of a die'. Then $\text{Prob}(Z = k) = 1/6$ for all $k = 1, 2, \dots, 6$, or, equivalently

Value k	1	2	3	4	5	6
$\text{Prob}(Z = k)$	1/6	1/6	1/6	1/6	1/6	1/6

Others are difficult. Let Y = 'the number of cigarettes a person selected at random smoked yesterday'.

It is difficult to know what the probability function of Y is.

Visualizing 'age' versus 'smoking'



FEV data set

430 children (9-17 years of age) had their age, forced expiratory volume in 1 second (FEV) and smoking status (!) recorded. (To be seen again in lecture 8.)

A barchart is a way to visualize a variable with a small number of unique values (often categorical). They are visual analogous of tables.

Ex: how do the ages distribute over smoking status?

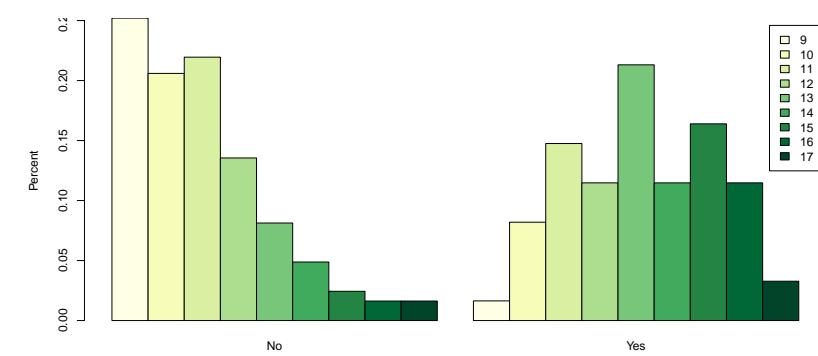
Smoking	Age								
	9	10	11	12	13	14	15	16	17
No	93	76	81	50	30	18	9	6	6
Yes	1	5	9	7	13	7	10	7	2

If the groups (smokers/non-smokers) aren't balanced it is difficult to compare the distributions.

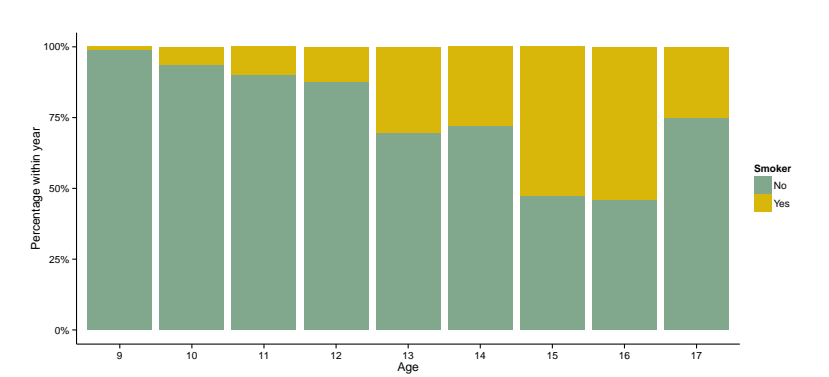
Tabulate/plot the percentages within groups:

Age	Smoking (%)	
	No	Yes
9	0.252	0.016
10	0.206	0.082
11	0.220	0.148
12	0.136	0.115
13	0.081	0.213
14	0.049	0.115
15	0.024	0.164
16	0.016	0.115
17	0.016	0.033
Sum	1.000	1.000

Visualizing 'age' versus 'smoking' as percentages within smoking groups



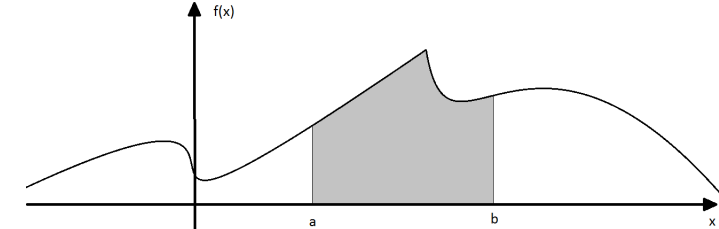
Visualizing 'smoking' versus 'age' as percentages within age groups



Probability model for continuous data

Recall that a RV is a yet undetermined value among several possible numbers.

A continuous RV is described by its *density function*.

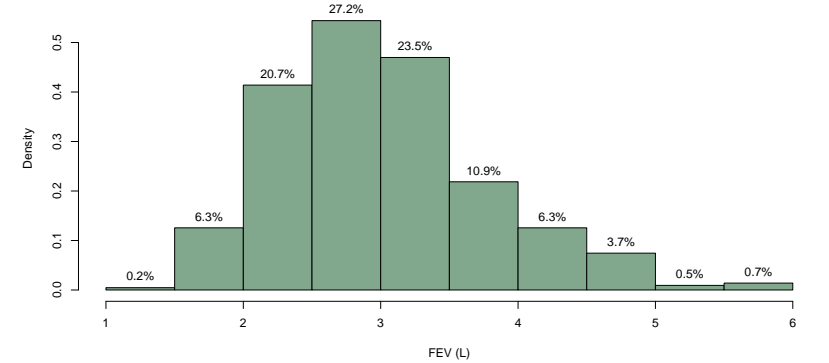


If X has density function f as above, then we compute probabilities as

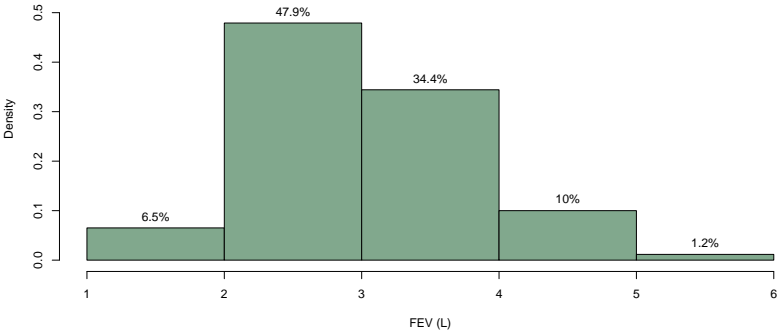
$$\text{Prob}(a \leq X \leq b) = \text{Area}(a,b).$$

Histogram (with density) of FEV

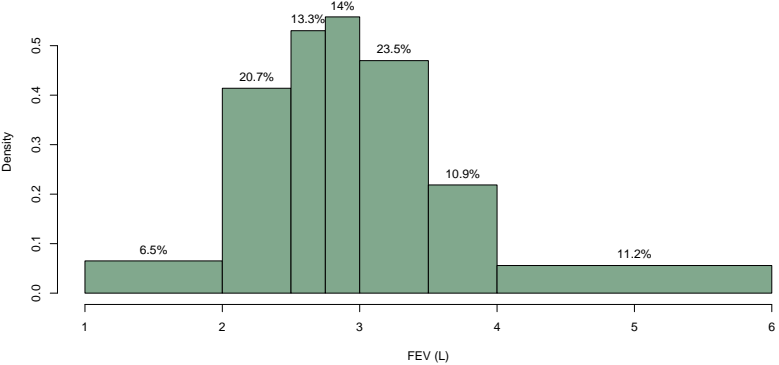
The histogram (with density on the y-axis) estimates the density function.



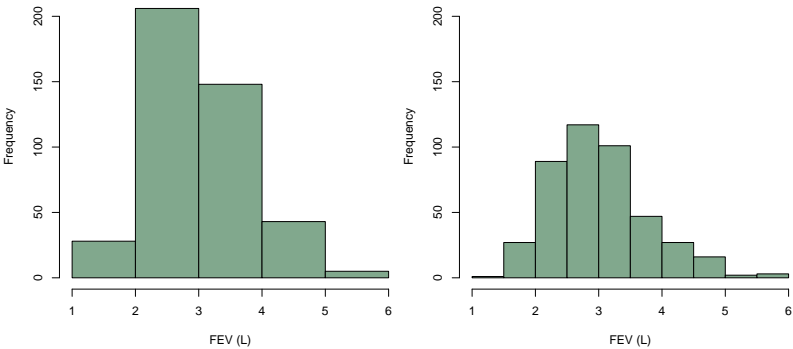
The interval width is arbitrary...



... and could be varying.

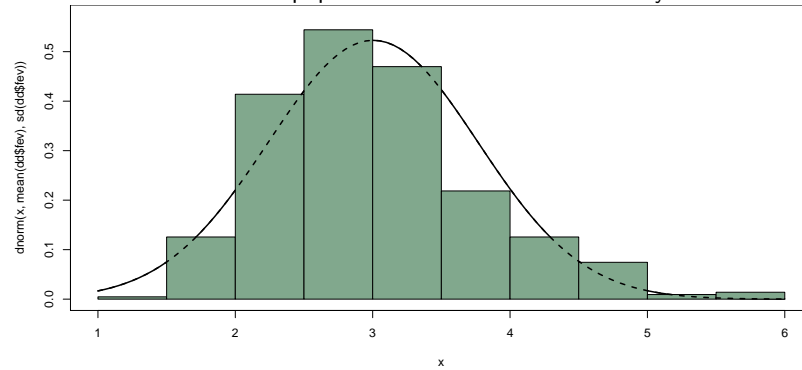


The scale of histograms with counts depends on chosen interval width and sample size. Below is the FEV data set with different widths.



Graph	Summary Measure	Theory
Ecdf	percentiles	(Cdf)
Boxplot	min, Quartiles, max	
Bar charts		Probability functions (discrete RV)
Histograms		Density functions (continuous RV)
	median, IQR	any distribution
	mean, s.d.	symmetrical distribution (\approx Normal distribution)

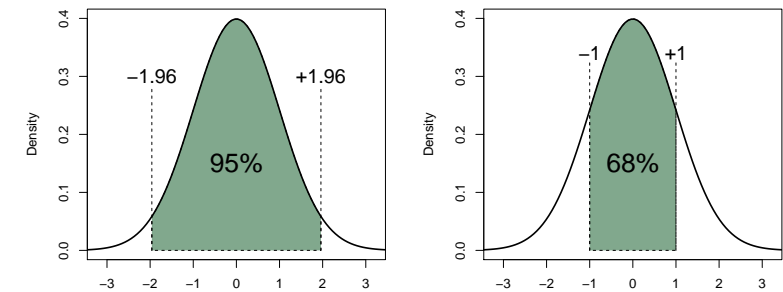
Often we assume that the population follows a Normal density curve.



The standard Normal distribution has a standard deviation of 1.

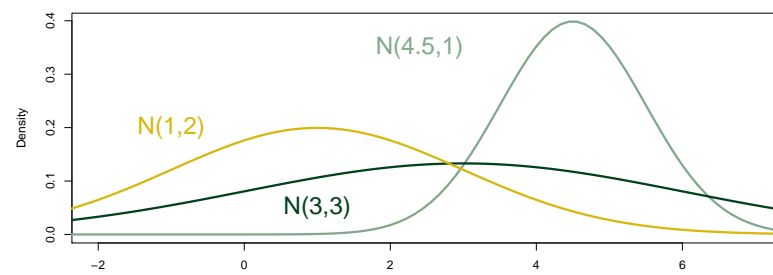
68% of the 'probability mass' lies within ± 1 .

95% of the 'probability mass' lies within ± 1.96 .

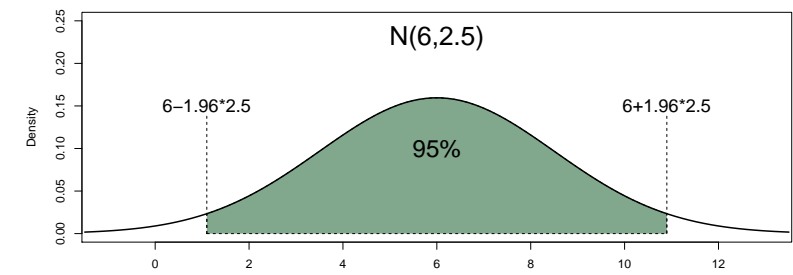


is determined by its mean (μ) and standard deviation (σ).

If X is $N(0,1)$ then $\mu + \sigma X$ is $N(\mu, \sigma)$.



If X is $N(\mu, \sigma)$, then 95% of observations will be between $\mu - 1.96\sigma$ and $\mu + 1.96\sigma$.



Properties of samples and sample means

Suppose we have a sample of size 10 from the population.

We want to know the *population mean* (μ).

We can estimate with the *sample mean* (μ^*).

But how good a guess is μ^* ?

Well... if the sample is drawn at random¹ then μ^* is *unbiased* (on average correct).

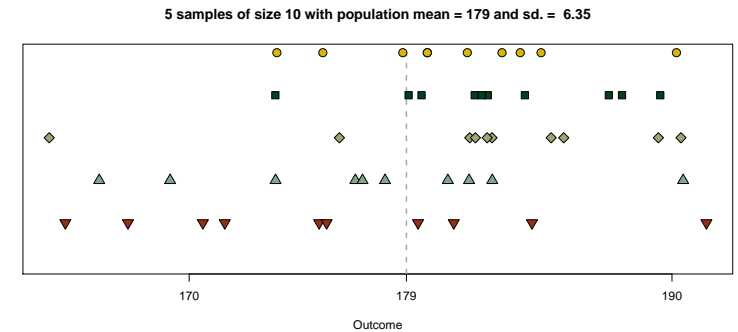
If we *knew* that the sd of $\mu^* = 180$ was, say 3.0, then there is an (approximate!) interpretation.

The estimated value 180 (cm) is on average off by 3.0 (cm).

But what *is* the standard deviation of a sample mean?

Consider several iterations of the procedure of drawing a sample of size 10.

¹All possible samples are equiprobable



Properties of sample means

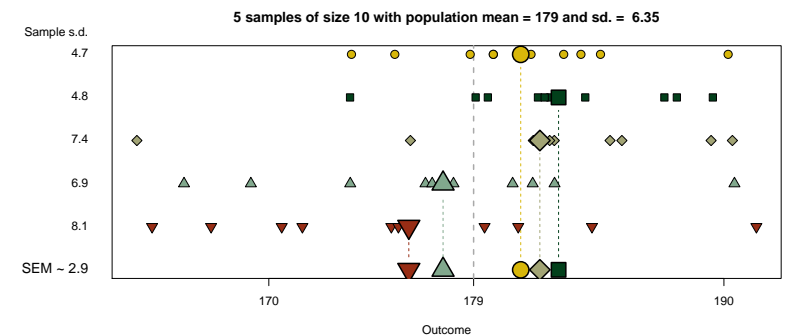
We can see that samples and their means and standard deviations vary (of course).

Standard Error of the Mean (SEM)

The standard deviation of the sample mean is referred to as SEM and a measure of (approx.) the average distance the population mean. Thus a measure of 'how good' the sample mean is as an estimator of the population mean.

In a Normal model the SEM is the population s.d. divided by \sqrt{n} .

The next slide tries to visualize this.



- a Normal population curve, and
- the Normal curve for the sample mean (for sample size 10).

Properties of the Normal distribution will aid further analysis using the mean of the sample.

What happens if the population follows some non-Normal curve?

Attempt to visualize CLT for the sample mean

Figure 1 is a line graph illustrating the density of the sample mean for different sample sizes. The x-axis is labeled 'Outcome' and ranges from 0.0 to 2.5. The y-axis is labeled 'Density' and ranges from 0.0 to 1.5. A vertical grey line is drawn at Outcome = 1.0. The legend indicates five sample sizes: 1 (Population), 2, 5, 10, and 20. The population density (solid yellow line) is centered at 0.0. The sample mean densities (dashed lines) are centered at 1.0. As the sample size increases, the sample mean densities become narrower and taller, illustrating the Central Limit Theorem.

	Outcome
--	---------

Quantiles divides your data into (roughly) equal piles.

- the median is the 2-quantile
- the tertiles are the 3-quantiles (the $33\frac{1}{3}$ percentile and the $66\frac{2}{3}$ percentile)
- the quartiles (Q1, Q2 and Q3) are the 4-quantiles.
- ... and so on.

13 patients had their peak expiratory flow (PEF, l/min) recorded once after inhaling each of two different asthma drugs (the order of which were random).

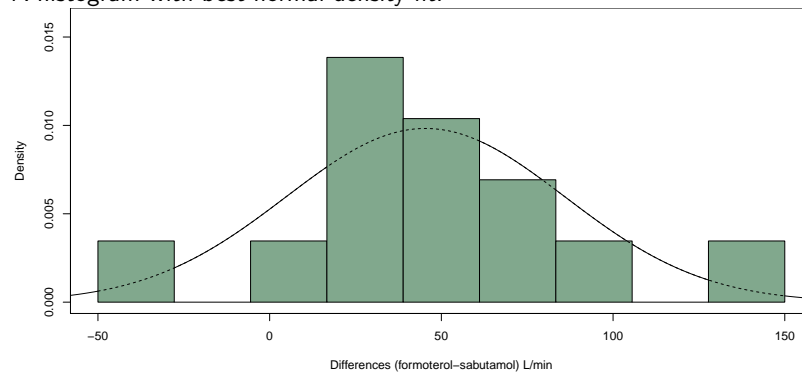
In *paired* data one usually look at the 13 differences as a measurement of differences in individual effect size.

Data:

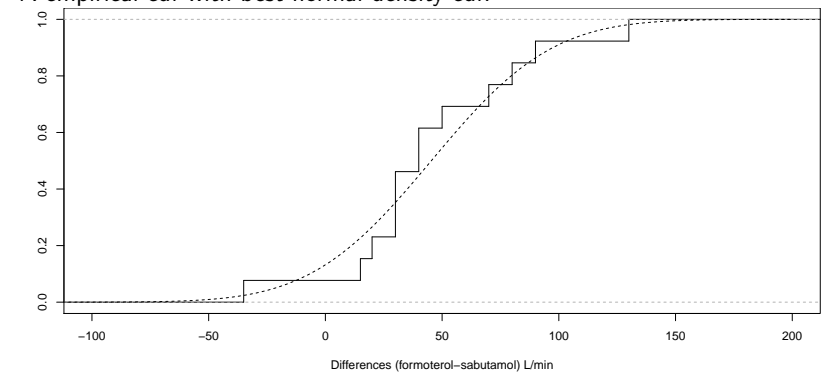
40, 50, 70, 20, 40, 30, -35, 15, 90, 30, 30, 80, 130

Is the normal distribution a good model for these 13 numbers?

A histogram with best normal density fit.

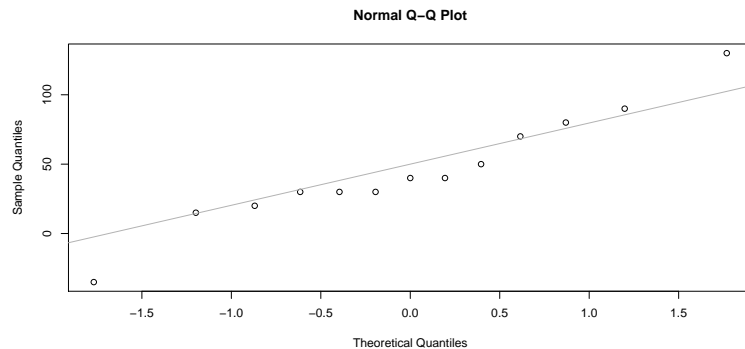


A empirical cdf with best normal density cdf.



The Quantile-Quantile plot

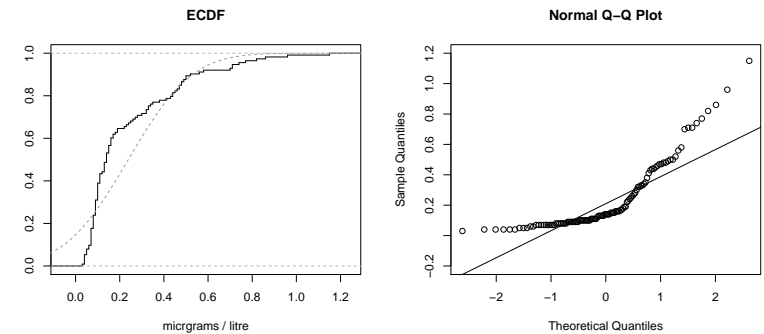
If the differences in effect size is sampled from a Normally distributed population its QQ-plot should be a straight line (approximately). A QQ-plot plots the sample of size n against the (slightly shifted) n -quantiles of the (standard) Normal distribution.



These visual test require some training.

There are also formal tests of normality (e.g. Shapiro-Wilks)

The S100 β measurements from the subarachnoidal bleeding example is certainly not normally distributed.



Second Summary

95% of observations from a Normal population lies within 1.96 multiples of the (population) s.d. from the (population) mean.

Means and in particular s.d. must be distinguished on three levels

- population,
- sample, and
- estimate.

The s.d. of the latter is called the Standard Error
(Standard Error of the Mean if the mean is used as estimate.)

(The sample s.d. is an estimate of the population s.d.)

The CLT explains why many estimates ('statistics') are (approx.) Normally distributed even though the population may not be.

References

- Chapters 1-8, 10: Petrie & Sabin. *Medical Statistics at a Glance*, Wiley-Blackwell (2009).
- Puhan et al. *More medical journals should inform their contributors about three key principles of graph construction*, Journal of Clinical Epidemiology, **59** (2006) 1017-1022.
- Franzblau & Chung. *Graphs, Tables, and Figures in Scientific Publications: The Good, the Bad, and How Not to Be the Latter*, American Society for Surgery of the Hand, **37A** (2012) 591-596.
- Kelleher & Wagener. *Ten guidelines for effective data visualization in scientific publications*, Environmental Modelling & Software **26** (2011) 822-827.