

To analyse the data, we take the difference between the observations within each pair, and treat the differences as a single sample

To estimate the population mean difference.

We test:  $H_0: \mu = 0$

A  $(1 - \alpha) 100\%$  confidence interval for  $\mu$ :

$$\bar{X} \pm t_{\alpha/2} \cdot SE(\bar{X})$$

↓

Sample mean difference

$SE(\bar{X}) = \frac{s}{\sqrt{n}}$

↖ sample standard deviation of the differences

↖ number of differences

$$t = \frac{\bar{X} - 0}{SE(\bar{X})} \rightarrow \text{find p-value and draw a conclusion}$$

## Chi-Square Tests (for categorical data)

### Chi-square tests for one-way tables

In a one-way table, observations are classified according to a single categorical variable.

Example: famous genetics experiment in 1905 (investigation of inheritance in sweet peas)

One pure line of peas that had purple flowers and long pollen grains was crossed with another pure line that had red flowers and round pollen grains.

This first generation was then self-crossed

Results for 381 plants: phenotypes

purple / long	purple / round	red / long	red / round
284	21	21	55

If these two genes are inherited independently, a 9:3:3:1 ratio would be expected

$H_0$ : True ratio of phenotypes is 9:3:3:1

$H_a$ : True ratio of phenotype is not 9:3:3:1

→ in other words:

Expected counts

$P(\text{purple and long}) = \frac{9}{16}$	$\frac{9}{16} \cdot 381$	214,3
$P(\text{purple and round}) = \frac{3}{16}$	$\frac{3}{16} \cdot 381$	71,4
$P(\text{red and long}) = \frac{3}{16}$	$\frac{3}{16} \cdot 381$	71,4
$P(\text{red and round}) = \frac{1}{16}$	$\frac{1}{16} \cdot 381$	23,8

Are the observed counts significantly different from those expected under the null hypothesis?

	purple/long	purple/round	red/long	red/round
observed	284	21	21	55
expected	214.3	71.4	71.4	23.8

↓ differences seem quite large

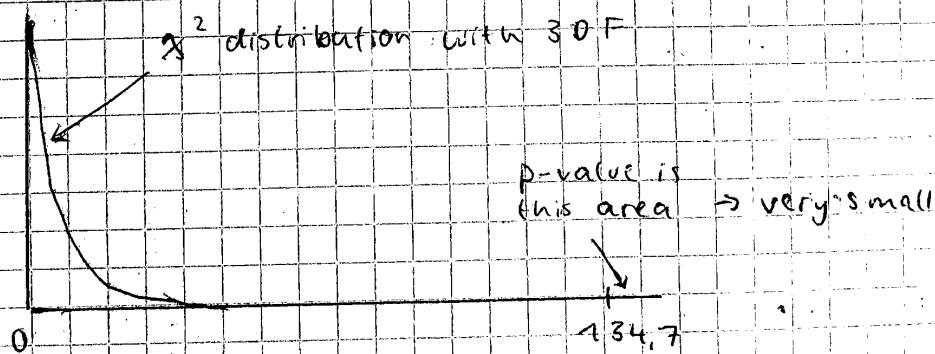
$$\text{test statistic: } \chi^2 = \sum_{\text{all cells}} \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} \Rightarrow \chi^2 = 134.7 \\ \text{DF} = 4-1=3$$

If the null hypothesis is true the chi-square test statistic is going to have approximately a chi-square distribution.

degrees of freedom = # of cells - 1  
(ie number of categories)

If the observed counts are close to the expected counts then  $\chi^2$  is going to be small.  
→ large values of  $\chi^2$  give us evidence against "H<sub>0</sub>"

When we get our p-value, we take the observed value of our test statistic  $\chi^2$ . The p-value is going to be the probability, under the null hypothesis of getting that value or sth even larger. (The area to the right of  $\chi^2$  under a chi-square distribution)



There is extremely strong evidence that these phenotypes do not occur in a 9:3:3:1 ratio (aka not genes are independent)

→ This study was significant in the discovery of genetic linkage, where genes close to each other on the same chromosome are more likely to be inherited together.

## (Chi-square Tests of Independence (two-way tables))

Example: A study of 11160 alcohol drinkers on university campuses revealed:

	never	occasional	frequent	total
trouble with police	71	154	398	623
no trouble with police	4992	2808	2737	10537
total	5063	2962	3135	11160

$H_a$ : They are not independent

$$\chi^2 = \sum_{\text{all cells}} \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} = 469.6$$

\* Expected count =  $\frac{\text{row total} \times \text{column total}}{\text{Overall total}}$

$$DF = (\# \text{ rows} - 1) \times (\# \text{ columns} - 1) = (2-1) \cdot (3-1) = 2$$

$$\rightarrow p\text{-value} < 2.2 \cdot 10^{-16}$$

$\rightarrow$  There is extremely strong evidence that binge drinking and trouble with police are not independent.

## ANOVA

Analysis of variance, or in short ANOVA, is a technique from statistical inference that allows us to deal with several populations.

To see why we need ANOVA, we will consider an example. Suppose we are trying to determine if the mean weights of green, red, blue and orange M&M candies are different from each other. We will state the mean weights for each of these populations  $\mu_1, \mu_2, \mu_3$  and  $\mu_4$ , respectively. We may use the appropriate hypothesis test several times:

- $H_0: \mu_1 = \mu_2$  There are many problems with this kind of analysis. We will have six p-values. Our confidence in the overall process will be  $0.95 \cdot 0.95 \cdot 0.95 \cdot 0.95 \approx 0.74$ .
- $H_0: \mu_2 = \mu_3$  Thus the probability of a type I error has increased.
- $H_0: \mu_3 = \mu_4$  And we cannot compare these four parameters as a whole by comparing them two at a time.
- $H_0: \mu_1 = \mu_3$
- $H_0: \mu_2 = \mu_4$

## One-way analysis

One-way ANOVA is a statistical method that tests:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k \quad (\text{It is a test on means})$$

by comparing the variability between groups to the variability within groups

The assumptions are the same as those of the pooled-variance two-sample t-test  
(population variances are equal)

Let:

- $k$  represent the number of groups
- $X_{ij}$  represent the  $j$ th observation in the  $i$ th group
- $\bar{X}_i$  represent the mean of the  $i$ th group ( $k$  of them)
- $\bar{X}$  represents the overall mean
- $s_i$  represent the standard deviation of the  $i$ th group
- $n_i$  represent the number of observations in the  $i$ th group
- $n = n_1 + n_2 + \dots + n_k \rightarrow$  total number of observations

. The total sum of squares:

$$SS(\text{Total}) = \sum_{\text{all obs}} (X_{ij} - \bar{X})^2 \quad (\rightarrow s^2 = \frac{SS(\text{Total})}{n-1})$$

The total sum of squares is partitioned into two components:

$$\begin{array}{l} \text{Sum of square} \\ \text{treatment} \end{array} \quad SS(T) = \sum_{\text{groups}} n_i (\bar{X}_i - \bar{X})^2$$

$$\begin{array}{l} \text{sum of squares} \\ \text{error} \end{array} \quad SSE = \sum_{\text{within groups}} (n_i - 1) s_i^2 \quad - \text{if there is a lot of variability within the groups SSE tends to be large}$$

$$\Rightarrow SS(\text{Total}) = SST + SSE$$

$$DF(\text{Total}) = DFT + DFE$$

source	df	sum of squares	mean square	f
treatments	k - 1	SST	SST / (k - 1)	MST / MSE
error	n - k	SSE	SSE / (n - k)	-
total	n - 1	SS (Total)	-	-

$$MSE = s_p^2 \\ (\text{mean square error})$$

The test statistic is  $F = \frac{MST}{MSE}$

If  $H_0$  is true, this test statistic will have an F distribution and MST and MSE have to be the same quantity. (F close to 1 where its "mean" is)

numerator df: k - 1  
denominator df: n - k

If  $H_0$  is false, MST will tend to be bigger than MSE and the test statistic will tend to be large.

example: can self-control be restored during intoxication?

An experiment randomly assigned 11 males to each of four groups

- Group A received 0,62 mg/kg of alcohol
- Group AC received alcohol plus caffeine
- Group AR received alcohol and a monetary reward for performance
- Group P received a placebo

Scores on a word stem completion task involving "controlled memory processes" were recorded.

Higher scores were indicative of greater self-control

$$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$$

$H_a$ : The means are not all equal

Group A	Group B	Group AR	Group P
$\bar{x}_1 = 0,08$	$\bar{x}_2 = 0,266$	$\bar{x}_3 = 0,441$	$\bar{x}_4 = 0,4$
$s_1 = 0,194$	$s_2 = 0,170$	$s_3 = 0,182$	$s_4 = 0,167$

resulting ANOVA table is:

source	df	ss	ms	f	p-value
treatments	3	0,972	0,324	10,18	0,00004
error	40	1,272	0,0319	-	-
total	43	2,243	-	-	-

There is very strong evidence ( $p\text{-value} = 0,00004$ ) that the groups do not all have the same population mean.  
(the different treatments do not all have the same effect)

but how you could investigate further and calculate confidence intervals for all the pairwise differences ( $\mu_1 - \mu_2, \mu_3 - \mu_4, \text{etc.}$ )

## Two-way ANOVA

We use a two-way ANOVA when we want to know how two independent variables, in combination, affect a dependent variable.

quantitative

Both of your independent variables should be categorical

Again, we use the F-test which compares the variance in each group to the overall variance in the dependent variable

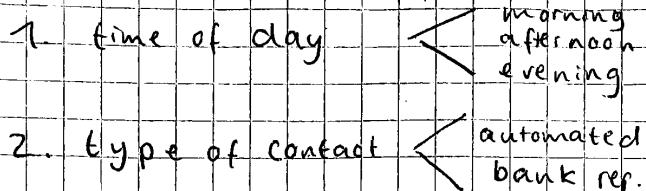
If the variance within groups is smaller than the variance between groups, the F-test will find a higher F-value, and therefore a higher likelihood that the difference observed is real and not due to chance.

assumption: the dependent variable should be normally distributed

example: customer satisfaction at friendly bank

		Factor 2: Type Contact		Row means
		Automated	bank repr.	
Factor 1: time of day	Morning	6,5,8,4 $\bar{x} = 5,75$	8,7,9,9 $\bar{x} = 8,25$	$\bar{x} = 7$
	afternoon	3,5,6,5 $\bar{x} = 4,75$	9,10,6,8 $\bar{x} = 8,25$	$\bar{x} = 6,5$
	evening	5,5,7,5 $\bar{x} = 5,5$	9,10,10,9 $\bar{x} = 9,5$	$\bar{x} = 7,5$
column means		$\bar{x} = 5,13$	$\bar{x} = 8,175$	$\bar{x} = 7$
* same with these				→ but instead with two-way ANOVA we can perform the test at the same time
→ two-way ANOVA studies average customer satisfaction regarding:				1. the interaction of contact type and time of day 2. type of contact 3. time of day

## Variables (factors)



n = number of items in cell  
r = number of rows  
i = number of columns

## Analysis of variance response

Source	DF	SS	MS
Contact type	1	66,67	66,67
Contact time	2	4,00	2,00
Interaction	2	2,33	1,17
Error	18	29,00	1,61
Total	23	102,00	

degrees of freedom

row factor = r - 1

column factor = i - 1

interaction = (r-1)(i-1)

error = n · (n-1)

Total = n · pi - 1

$$MS = \frac{SS}{d.f.}$$

$$\text{test statistic: } F = \frac{MS_{\text{interaction}}}{MS_{\text{error}}}$$

Critical value:  $\alpha = 0,05$ , d.f.<sub>r</sub>, d.f.<sub>i</sub>

↓                          ↓  
numerator denominator

1.  $H_0$ : There is no interaction

$H_1$ : There is interaction

$$\text{test statistic: } F = \frac{MS_{\text{interaction}}}{MS_{\text{error}}} = \frac{1,17}{1,61} = 0,73$$

Critical value:  $\alpha = 0,05$  d.f.<sub>r</sub> = 2 d.f.<sub>i</sub> = 18

$$\rightarrow F = 3,55$$

→ The test statistic is so small that we don't get out far enough to get into the rejection region, so there is no interaction between contact type and contact time

2.  $H_0$ : No difference in customer satisfaction between the two types of contact

$H_1$ : There is a difference

$$\text{test statistic: } F = \frac{MS_{\text{contact type}}}{MS_{\text{error}}} = \frac{66,67}{1,61} = 41,41$$

Critical value: d.f.<sub>r</sub> = 1 d.f.<sub>i</sub> = 18  $\rightarrow F = 4,41$

→ Certainly the 41,41 puts us far enough out to where we "know" we're in the rejection region so we reject  $H_0$  and there is a difference between contact type

3.  $H_0$ : There is no difference in customer satisfaction among the times of contact.

$H_1$ : There is a difference.

Test Statistic:  $F = \frac{MS_{\text{time}}}{MS_{\text{error}}} = \frac{2}{1,61} = 1,24$

Critical value:  $\alpha = 0,05$  d.f.<sub>N</sub> = 2 d.f.<sub>0</sub> = 18

$$\rightarrow F_0 = 3,55$$

$\rightarrow$  1,24 is not in the critical region, therefore we accept the  $H_0$ , meaning times of contact doesn't affect customer satisfaction

## Correlation & Regression

Objectives of regression:

- find out whether there is a relationship between two variables
- forecast observations

Regression analysis explores the relationship between a quantitative response variable and one or more explanatory variables

If there is only one explanatory variable, we call it simple regression and if there are more than one explanatory variable, we call it multiple regression.

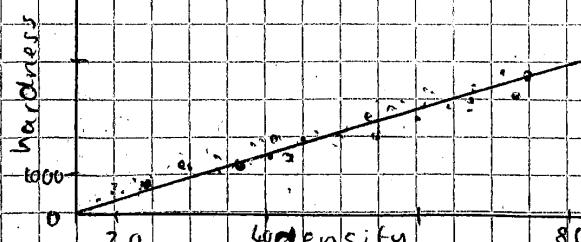
Simple linear regression:

Example: Can the density of Australian timber be used to predict its Janka hardness?

The first 4 observations (of 36):

explanatory variable (independent)	x density	24,7	39,4	53,4	24,8
response variable (dependent)	y hardness	484	1210	1880	427

We call the variable we are predicting y and the variable using to predict it x variable



We attempt to fit a line through those points and then we use that line for prediction

We will assume a linear relationship between  $Y$  and  $X$ :

$$E(Y|X) = \beta_0 + \beta_1 X + \epsilon$$

$\downarrow$   
 $\rightarrow \text{unknown}$        $\nwarrow \text{slope}$        $\nearrow \text{y intercept}$        $\nwarrow \text{irreducible error}$

The observed values of  $Y$  vary along the line. We account for that variability with a random error component  $\epsilon$ .

We usually use the method of least squares to estimate  $\beta_0$  and  $\beta_1$ .

### The Least Squares Regression Line

We use sample data to find the estimated regression line:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\text{Residual} = \text{Observed} - \text{Predicted}$$

$$e_i = Y_i - \hat{Y}_i$$

The residuals are the vertical distances between the points and the line points. Points below the line have negative residuals and points above the line have positive residuals.

We minimize the sum of the absolute value of the residuals

$\hat{\beta}_0$  and  $\hat{\beta}_1$  are chosen to minimize:

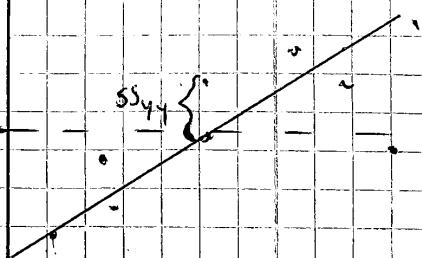
$$\begin{aligned}\sum e_i^2 &= \sum (Y_i - \hat{Y}_i)^2 \\ &= \sum (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2\end{aligned}$$

$$SS_{XX} = \sum (X_i - \bar{X})^2 \quad s_x^2 = \frac{SS_{XX}}{n-1}$$

$$SS_{YY} = \sum (Y_i - \bar{Y})^2 \quad s_y^2 = \frac{SS_{YY}}{n-1}$$

$$SP_{XY} = \sum (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\text{cov}(x, y) = \frac{SP_{XY}}{n-1}$$



$$\rightarrow \hat{\beta}_0 = Y - \hat{\beta}_1 X$$

$$\hat{\beta}_1 = \frac{SP_{XY}}{SS_{XX}} = \frac{\text{cov}(X, Y)}{\text{var}(X)}$$

Now if we took all distances and added them up that would add to zero.

### Exercise

X	-1	1	2	4	6	7	n = 6
Y	-1	2	3	3	5	8	

Does there exist a linear relationship?

x	y	$x^2$	$y^2$	xy
-1	-1	1	1	1
1	2	1	4	2
2	3	4	9	6
4	3	16	9	12
6	5	36	25	30
7	8	49	64	56
$\Sigma$	19	107	112	107

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$\hat{\beta}_1 = \frac{(\sum x)(\sum y) - n \sum xy}{(\sum x)^2 - n \sum x^2} = \frac{19 \cdot 20 - 6 \cdot 107}{(19)^2 - 6 \cdot 107} = \frac{-262}{-281} \approx 0,93$$

$$\hat{\beta}_0 = \frac{(\sum x)(\sum y) - (\sum y)(\sum x^2)}{(\sum x)^2 - n \sum x^2} = \frac{19 \cdot 107 - 20 \cdot 107}{19^2 - 6 \cdot 107} = \frac{-107}{-281} \approx 0,3808$$

$$\rightarrow \hat{y} \approx 0,9324x + 0,3808$$

### Correlation calculation:

$$s_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} = \sqrt{\frac{\sum x^2 - \frac{1}{n} (\sum x)^2}{n-1}} = \sqrt{\frac{107 - \frac{1}{6} \cdot 19^2}{5}} \approx 3,0605$$

$$s_y = \sqrt{\frac{112 - \frac{1}{6} (20)^2}{5}} \approx 3,0111$$

$$\text{Correlation } r = \frac{\sum xy - \frac{1}{n} (\sum x)(\sum y)}{(n-1) s_x s_y} \approx 0,9477$$

↓  
strong correlation

## Multiple Regression

General parametric equation:

$$y = f(x) + \epsilon$$

\ depends on statistical method

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p$$

This regression type is more relevant in real life since it can take multiple factors (e.g. height, age to determine marital status) into account.

## Bayesian Statistics

We have covered the Bayes' theorem already but as a quick review:

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

where A and B are events,  $P(A|B)$  is the conditional probability that event A occurs given that event B has already occurred.

example: picking a card from a pack of traditional playing cards

There are 52 cards in the pack, 26 of them are red and 26 of them are black.

What is the probability of the card being a 4 given that we know the card is red?

A: card is a 4      B: card is red

$$P(B|A) = P(\text{red}|4) = \frac{1}{2} \quad (\text{2 red cards left})$$

$$P(A) = P(4) = \frac{4}{52} = \frac{1}{13}$$

$$P(B) = P(\text{red}) = \frac{1}{2} \quad (\text{half the cards})$$

$$P(A|B) = P(4|\text{red}) = \frac{\frac{1}{2} \cdot \frac{1}{13}}{\frac{1}{2}} = \underline{\underline{\frac{1}{13}}}$$

# Prior probability and Posterior probability

example: diagnostic testing on HIV

early HIV testing in the military

- first screen with ELISA
- if positive, two more rounds of ELISA
- if either positive, test with two Western blot assays
- only if both positive, determine HIV infection

data:

ELISA

- sensitivity (true positive): 93%
- specificity (true negative): 99%

Western blot

- sensitivity: 99,9%
- specificity: 99,1%

prevalence:  $\frac{1,48}{1000}$  (adult Americans were HIV positive by 1980)

What is the probability that a recruit who tested positive in the first ELISA has HIV?

$$P(\text{has HIV} | \text{ELISA}+) = ?$$

$$P(\text{HIV}) = 0,00148$$

$$P(+ | \text{HIV}) = 0,93 \quad P(- | \text{no HIV}) = 0,99$$

prior probability ( $\hat{=}$  probability that an observation will fall into a group before you collect your data - best rational assessment of the probability of an outcome based on the current knowledge before an experiment is performed)

→ prior to any testing what probability should be assigned to a recruit having HIV?

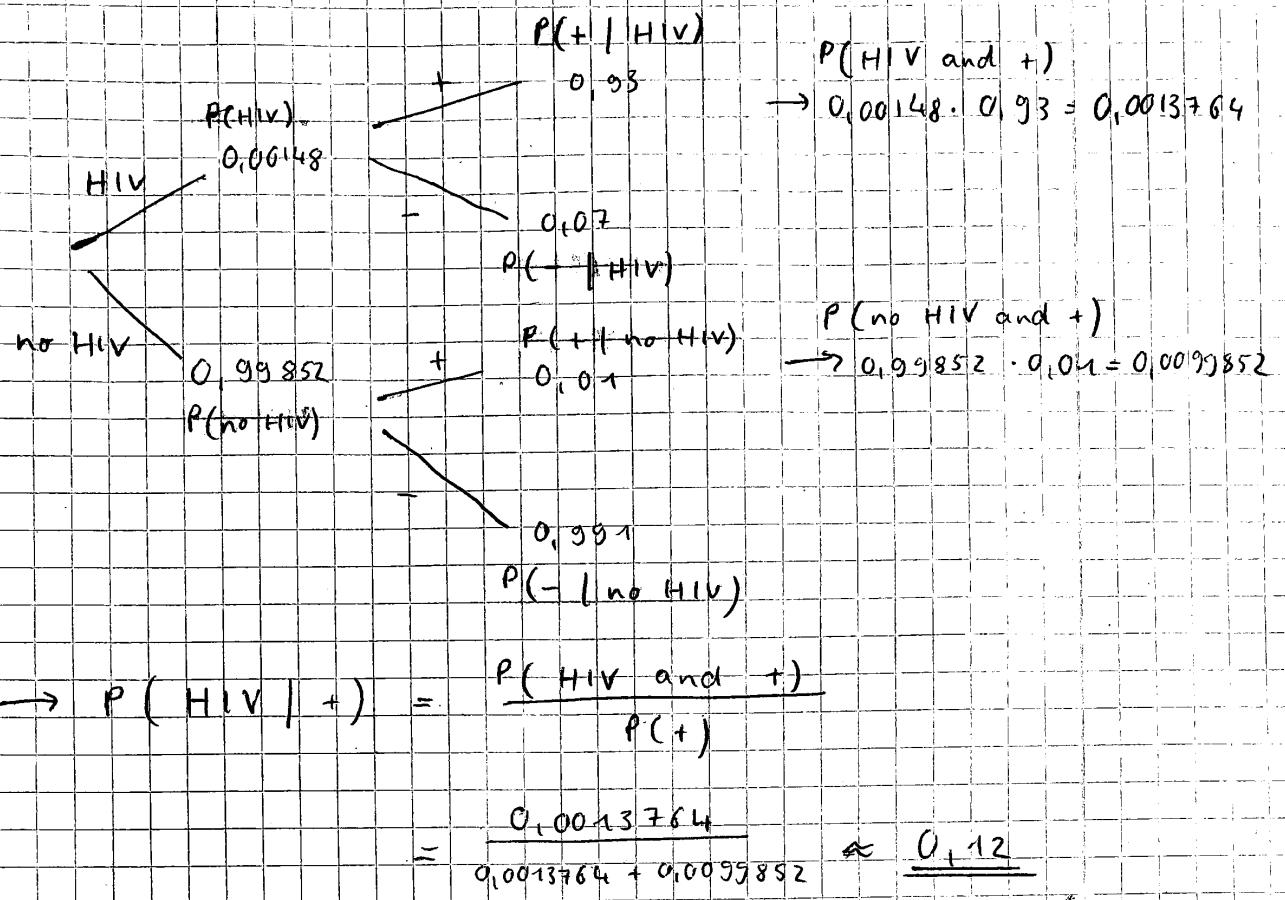
→ prevalence of the disease in the population

$$P(\text{HIV}) = 0,00148$$

posterior probability ( $\hat{=}$  the updated probability of an event occurring after taking into consideration new information)

When a recruit goes through HIV screening there are two competing claims: recruit has HIV and recruit doesn't have HIV.

→ If the ELISA yields a positive result, what is the probability this recruit has HIV?

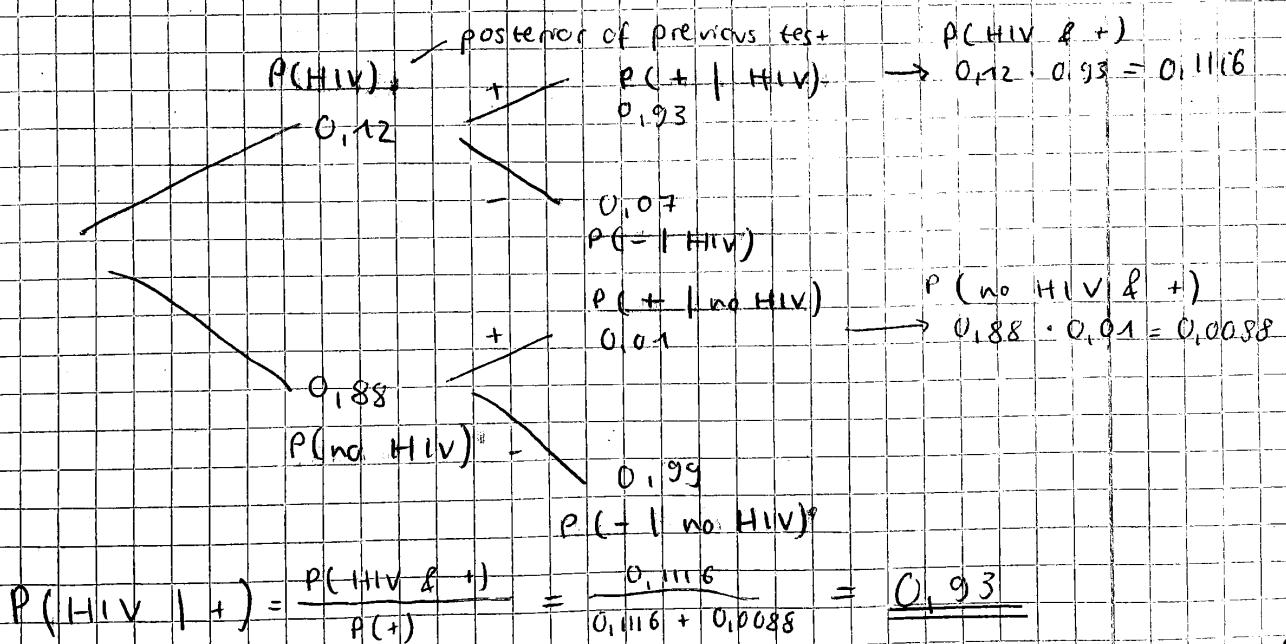


The probability that a recruit who tested positive in the first ELISA has HIV is  $\underline{\underline{0,12}}$ . ↑ posterior probability

## Bayes Updating

Retesting: Since a positive outcome in the ELISA doesn't necessarily mean that the recruit actually has HIV, they are retested.

What is the probability of having HIV if this second ELISA also yields a positive result?



→ Due to a second ELISA test we could increase the probability that a recruit who tested positive has HIV by 81%.

! The difference between prior and posterior probabilities characterizes the information we have gotten from the experiment or measurement.

→ We learned that a posterior probability can subsequently become a prior for a new updated posterior probability as new information arises and is incorporated into the analysis.

## Bayesian vs. frequentist

### definitions of probability

- frequentist: relative frequency in large number of trials

$$P(E) = \lim_{n \rightarrow \infty} \frac{n_E}{n} \quad (\text{proportion of times event occurs in } n \text{ trials, when } n \text{ goes to infinity})$$

- bayesian:

example: indifferent between winning

> \$1 if event E occurs  
or

> winning \$1 if you draw a blue chip from a box with  $1000 \times p$  blue chips and  $1000 \times (1-p)$  white chips

This means you're equating the probability of events E,  $P(E)$ , to the probability of drawing a blue chip from this box,  $p$ .

$$\rightarrow P(E) = p$$

→ It is an interpretation or estimate of probability as a personal judgement or degree of belief about how likely a particular event is to occur, based on the state of knowledge and available evidence.

### definitions of statistics:

- frequentist:

Frequentist Statistics tests whether an event occurs or not. It calculates the probability of an event in the long run of the experiment. With a frequentist approach the result of an experiment is dependent on the number of times the experiment is repeated.

## → flaws in frequentist statistics.

- p-values measured against a sample (fixed size) statistic changes with change in the stopping intention and sample size. So if two people work on the same data and have different stopping intention, they may get two different p-values for the same data, which is undesirable.

e.g. Person A chooses to stop tossing a coin after 100 tosses while B stops at 1000. For different sample sizes, we get different t-scores, thus different p-values

- Confidence intervals are not probability distributions therefore they do not provide the most probable value for a parameter.

e.g. "We are 95% confident that 60% to 64% of Americans think the federal government (poll on 1500 adults) does not do enough for middle class people"

means 95% of random samples of 1500 adults will produce confidence intervals that contain the true proportion of Americans who think the federal government does not do enough for middle class people.

but not: The true population proportion is in this interval 95% of the time

→ In other words for a single given confidence interval of a sample the true parameter is either in it or not.

## - Bayesian statistics

A mathematical procedure that applies probabilities to statistical problems, where probability expresses a degree of belief in an event.

The degree of belief may be based on prior knowledge about the event, for instance results of previous experiments or on personal beliefs about the event.

Bayesian statistical methods use Bayes' theorem to compute and revise probabilities after acquiring new data.

## Bayesian inference

As we already know, statistical inference is the process of deducing properties about a population or probability distribution from data.

Bayesian inference specifically, uses the Bayes' theorem.

example: probability of ice-cream sale with regards to the type of weather.

Bayes' theorem gives us the tools to use prior knowledge about the likelihood of selling ice cream on any other type of day (rainy, windy, snowy).

Let A represent the event that we sell ice cream and B the event of the weather.

Then we might ask what is the probability of selling ice-cream on any given day given the type of weather?

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

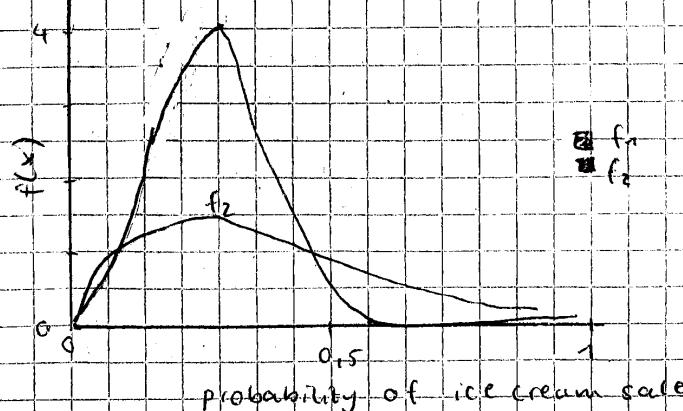
$P(A)$  is our prior, so in our case the probability of selling ice cream regardless of the type of weather outside. We might know this by looking at data that said 30 people out of a potential 100 actually bought ice-cream at some shop somewhere.

So  $P(A) = \frac{30}{100} = 0.3$ , prior to knowing anything about the weather.

(note: the prior belief can also be made up, though the resulting calculation will be affected by this choice)

If  $P(A)=0.3$  was just our best guess using a distribution of our prior belief is more appropriate, which is known as the prior distribution.

e.g.



The peak of both curves occur around 0.3, which as we said, is our best guess. The fact that  $f(x)$  is non-zero at other values of  $x$  shows that we're not completely certain that 0.3 is the true value.  $f_1$  shows that it's likely to be between 0 and 0.5 whereas  $f_2$  shows that it's likely to be between 0 and 1, meaning the prior prob expressed by  $f_2$  is less certain.

## Model form of Bayes' Theorem

Instead of event A, we usually use  $\theta$  called theta. Theta is of our interest, it represents the set of parameters.

So if we're trying to estimate the parameter values of a Gaussian distribution then  $\theta$  represents the mean  $\mu$  and the standard deviation  $\sigma$ .

$$\rightarrow \theta = \{\mu, \sigma\} \quad \begin{array}{l} \text{It usually represents a proposition} \\ (\text{e.g. "a coin lands on heads 50% of the time"}) \end{array}$$

Instead of event B, we'll see data or  $y = \{y_1, y_2, \dots, y_N\}$

These represent the data, i.e. the set of observations of measurements we have, we can also see it as our evidence.

$$\rightarrow P(\theta | \text{data}) = \frac{P(\text{data} | \theta) \cdot P(\theta)}{P(\text{data})}$$

$P(\theta | \text{data})$  is our posterior distribution, representing our belief about the parameter values after we have calculated everything on the right hand side, taking the observed data into account.

$P(\text{data} | \theta)$  is our likelihood distribution/function, which can be interpreted as the prob. of the evidence given  $\theta$  is true.

$\rightarrow$  It measures the goodness of fit of a statistical model to a sample of data for given values of parameters. The "likelihood" is equal to the probability that a particular outcome is observed when the true value of parameter is  $\theta$ .

$\rightarrow$  The Bayes' theorem implies that the posterior probability is proportional to the likelihood times the prior probability.

Models are the mathematical formulation of the observed events. Parameters are the factors in the models affecting the "observed" data.

example: fairness of a coin ( $\theta$ ).

Our question could be:

Given an outcome data what is the probability of the coin being fair ( $\theta = 0.5$ )

Thus  $P(\theta)$  is the prior, i.e. the strength of our belief in the fairness of the coin before the toss.

$P(\text{data} | \theta)$  is the likelihood of observing our result given our distribution for  $\theta$ .

$P(\text{data})$  is the probability of data determined by summing across all possible values of  $\theta$ , weighted by how strongly we believe in those particular values of  $\theta$ .

$P(\theta | \text{data})$  is the posterior belief of our parameters after observing the evidence, i.e. the number of heads.

### Bernoulli likelihood function

To recap: our likelihood function is the probability of observing a particular number of heads in a particular number of flips for a given fairness of coin.

1 represents heads. - 0 represents tails

$$P(y_1 = 1 | \theta) = \theta^y \quad (\text{if coin is fair, prob of heads is } 0.5)$$

$$P(y_1 = 0 | \theta) = (1-\theta)^{1-y}$$

$$\rightarrow P(y | \theta) = \theta^y \cdot (1-\theta)^{1-y} \quad y = \{0, 1\} \quad \Theta = (0, 1)$$

When we want to see a series of flips or heads:

$$P(y_1, y_2, \dots, y_n | \theta) = \prod P(y_i | \theta)$$

$$= \prod_i \theta^{y_i} \cdot (1-\theta)^{1-y_i}$$

Furthermore, if we are interested in the probability of number of heads  $z$  turning up in  $N$  number of flips:

$$P(z, N | \theta) = \theta^z \cdot (1-\theta)^{N-z}$$

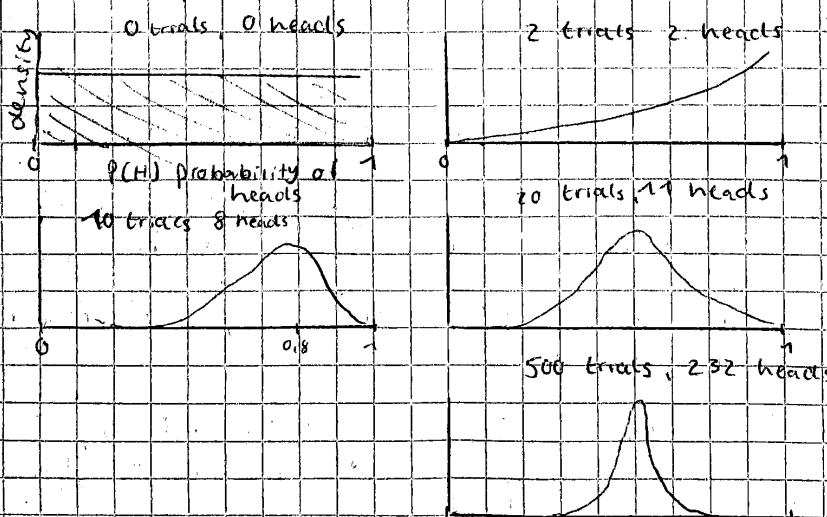
### Prior belief distribution

The mathematical function to represent prior beliefs is known as beta distributions and has the following form:

$$x^{\alpha-1} \cdot (1-x)^{\beta-1} / B(\alpha, \beta) \quad \theta = \frac{\alpha}{\alpha + \beta}$$

$\alpha$  and  $\beta$  are called the shape decoding parameters for the density function.

Here  $\alpha$  is equivalent to the number of heads in the trials and  $\beta$  corresponds to the number of tails. Here a few visualizations of beta distributions for different  $\alpha$  and  $\beta$ .



→ As more tosses are done and heads continue to come in larger proportion the peak narrows which increases our confidence in the fairness of the coin.

### Posterior belief distribution

So let's put it all together

$$\begin{aligned}
 P(\theta | z, N) &= \frac{P(z, N | \theta) P(\theta)}{P(z, N)} \\
 &= \frac{\theta^z (1-\theta)^{N-z} \cdot \overset{\alpha}{\underset{\sim}{G}}(1-\theta)^{\beta-1}}{B(z, \beta) \cdot P(z, N)} \\
 &= \frac{\theta^{z+\alpha-1} (1-\theta)^{N-z+\beta-1}}{B(z+\alpha, N-z+\beta)} \\
 \Rightarrow P(\theta | z+\alpha, N-z+\beta)
 \end{aligned}$$

So this means by knowing the mean and standard deviation (as we can derive  $\alpha$  and  $\beta$  from them) of our belief about  $\theta$  or the parameter  $\theta$  and by observing the number of heads in  $N$  flips, we can update our belief about  $\theta$ .

For example you could suppose the coin is biased and has a mean of 0.6 with  $\sigma = 0.1$  then  $\alpha = 13.8$  and  $\beta = 9.2$

which means our distribution will be biased on the right.  
And suppose you then observed 80 heads ( $z = 80$ ) in 100 flips ( $N = 100$ )

So our posterior probability would look like this

$$P(\theta | z+\alpha, N-z+\beta) = P(\theta | 93.8, 29.2)$$

### Bayes' factor

The Bayes' factor is the equivalent of the p-value in the Bayesian framework.

It is defined as the ratio of the likelihood of one particular hypothesis to the likelihood of another hypothesis.

Typically, it is used to find the ratio of the likelihood of an alternative hypothesis to a null hypothesis.

For example, if the Bayes Factor is 5 then it means the alternative hypothesis is 5 times as likely as the null hypothesis given the data.

Similar to p-values we can use thresholds to decide when we should reject a null hypothesis (e.g. 10)

A Bayes' factor of 1 means that there is no evidence.

Some statisticians believe that the Bayes factor offers an advantage over p-values because it allows you to quantify the evidence for and against two competing hypotheses.

Credibility Intervals (analogous to confidence intervals in frequentist statistics)

They allow us to describe the unknown true parameter not as a fixed value but with a probability distribution.

Therefore we can make probabilistic statements about the parameter falling within that range.

As the Bayesian inference returns a distribution of possible outcomes, the credibility interval is the range containing a particular percentage of probable values.

For instance, the 95% credible interval is simply the central portion of the posterior distribution that contains 95% of the values.

It allows us to say something like "given the observed data, the outcome has 95% probability of falling within this range"