# Part V:

# Statistics

Point estimates,
bias & variance,
consistency

# Standard setup in parametric statistics

We assume that data is generated by a particular family of distributions, for example

$$\mathcal{F} = \left\{ N(\mu, \sigma^2) \mid \underbrace{\mu \in \mathbb{R}, \sigma^2 > 0}_{\theta} \right\}.$$

The family $\mathcal{F}$ is called the <u>statistical model</u>.

More generally, $\mathcal{F} = \left\{ f_\theta \mid \theta \in \Theta \right\}$

$f_\theta$ — one particular parameter

$\Theta$ — space of all possible parameters

We are given a sample $X_1, \ldots, X_n \sim f_\theta$ (typically, iid) but the true, underlying $\theta$ is unknown.

# Convention

Parameter space $\Theta$ ("capital theta")

True (unknown) parameter $\theta$ ("lower case theta")

$P_\theta$, $E_\theta$, ... refers to the probability, expectation

under the distribution $f_\theta$

Estimates typically get a "hat": $\hat{\theta}$, $\hat{\mu}$, ...

# Point estimation

The goal of ==point estimation== is to estimate $\theta$.

Def. Given a statistical model $\mathcal{F} = \{ f_\theta \mid \theta \in \Theta \}$, and a sample $X_1, \ldots, X_n \sim F \in \mathcal{F}$. A ==point estimator== $\hat{\theta}_n$ of parameter $\theta$ is a function

$$\hat{\theta}_n := g(X_1, \ldots, X_n)$$

# Bias of an estimator

Def. The **bias** of such an estimator is defined as

$$\text{bias}\left(\hat{\theta}_n\right) := \underbrace{E_\theta\left(\hat{\theta}_n\right)}_{} - \underbrace{\theta}_{\text{true para}}$$

estimate

↳ expectation wrt the distribution $f_\theta$
(the true one!)

**Intuition:** repeat the procedure very often (infinitely often)
and average over the estimate $\hat{\theta}_n$.

An estimate is **unbiased** if its bias is zero.

# Variance and standard error

<u>Def</u>  The ==variance of an estimator== is defined as $\text{Var}_{\theta}(\hat{\theta}_n)$. The corresponding standard deviation is called the ==standard error se==. Typically, se is unknown, but it can be estimated: $\widehat{se}$.

# Example

$$X_1, \ldots, X_n \sim \text{Bernoulli}(p), \quad \text{parameter } p \in [0,1],$$

$$\hat{p}_n := \frac{1}{n} \sum_{i=1}^{n} X_i \quad \text{an estimate of } p.$$

$$E_p(\hat{p}_n) = E_p\left( \frac{1}{n} \sum_{i=1}^{n} X_i \right) = \frac{1}{n} \sum_{i=1}^{n} \underbrace{E_p(X_i)}_{p} = p.$$

Thus, $\hat{p}_n$ is unbiased because

$$E_p(\hat{p}_n) - p = p - p = 0.$$
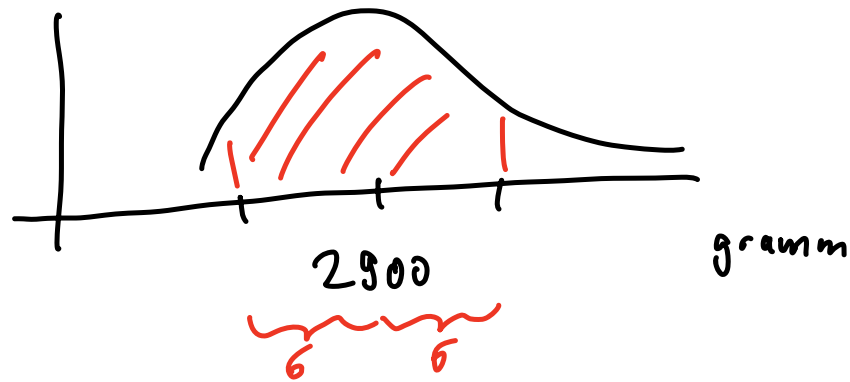
The standard error of this estimate is

$$se = \sqrt{\text{Var}_p(\hat{p}_n)} = \sqrt{\frac{1}{n} \text{Var}_p(X_1)} = \sqrt{\frac{p(1-p)}{n}}$$

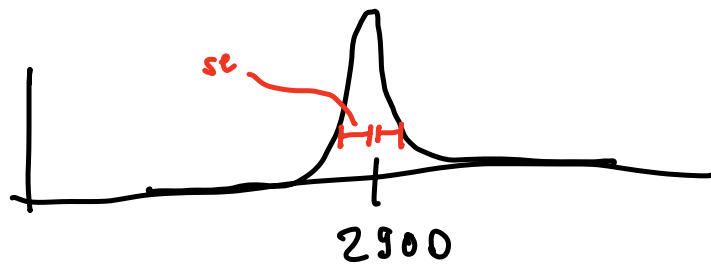We can for example estimate it by

$$\hat{se} = \sqrt{\frac{\hat{p}_n(1-\hat{p}_n)}{n}}.$$

# Example: weight of baby



2900

gramm

Distribution of individual data prs: $\mu = 2900$

$\sigma = 500$



se

2900

distribution of the estimate $\hat{\mu}_u$

# Mean squared error

__Def__    The ==mean squared error $(MSE)$== of an estimate is
the quantity

$$MSE(\hat{\theta}, \theta) = E_\theta \left( ( \hat{\theta}_n - \theta )^2 \right)$$

deterministic

# Bias- Variance decomposition

Theorem : bias-variance-decomposition

$$\underbrace{MSE(\hat{\theta}_n, \theta)}_{\text{how good is our estimate}} = bias^2(\hat{\theta}_n) + Var_\theta(\hat{\theta}_n)$$

**Proof**

$$E_\theta\left((\hat{\Theta}_n - \theta)^2\right) =$$

$$= E_\theta\left((\overbrace{\hat{\theta}_n - E\hat{\theta}_n}^{a} + \overbrace{E\hat{\theta}_n - \theta}^{b})^2\right)$$

$$= E_\theta\left(\overbrace{(\hat{\theta}_n - E\hat{\theta}_n)^2}^{a^2}\right) + 2E_\theta\left(\overbrace{(\hat{\theta}_n - E\hat{\theta}_n)}^{a}\underbrace{\overbrace{(E\hat{\theta}_n - \theta)}^{b}}_{\text{deterministic}}\right) + E_\theta\left(\overbrace{(E\hat{\theta}_n - \theta)^2}^{b^2}\right)$$

$$\underbrace{2(E\hat{\theta}_n - \theta) \cdot \underbrace{E_\theta(\hat{\theta}_n - E\hat{\theta}_n)}_{= E_\theta(\hat{\theta}_n) - E_\theta\hat{\theta}_n = 0}}_{= 0}$$

$$= \underbrace{E_\theta\left((\hat{\theta}_n - E\hat{\theta}_n)^2\right)}_{\text{Var}(\hat{\theta}_n)} \quad + \quad \cancel{E}\Big(\underbrace{(E\hat{\theta}_n - \theta)^2}_{\text{deterministic}}\Big)$$

$$= (E\hat{\theta}_n - \theta)^2$$

$$= \left(\text{bias}(\hat{\theta}_n)\right)^2$$

# Example

$$\mathcal{F} = \{ N(\mu, \sigma^2) \mid \mu \in \mathbb{R}, \sigma > 0 \}$$

Sample: $X_1, \ldots, X_n \sim N(\mu, \sigma^2)$ with unknown $\mu, \sigma^2$, iid

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} X_i \qquad \text{is an unbiased estimate of } \mu.$$

$$\hat{\sigma}_1^2 := \frac{1}{n} \sum_{i=1}^{n} (X_i - \hat{\mu})^2 \qquad \text{"first estimate"}$$

$$\hat{\sigma}_2^2 := \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \hat{\mu})^2 \qquad \text{"second estimate"}$$

$$E(\hat{\sigma}_1^2) = \frac{n-1}{n} \sigma^2 \qquad \text{so} \quad \text{the bias is} \quad \frac{1}{n}\sigma^2$$

$$E(\hat{\sigma}_2^2) = \sigma^2 \qquad \text{unbiased!}$$

$$Var(\hat{\sigma}_1^2) = \frac{2(n-1)\sigma^4}{n^2}$$

$$Var(\hat{\sigma}_2^2) = \frac{2\sigma^4}{n-1}$$

$$MSE(\hat{\sigma}_1^2) = bias^2 + var = \ldots = \left(\frac{2n-1}{n^2}\right)\sigma^4$$

$$MSE(\hat{\sigma}_2^2) = \qquad \ldots \qquad = \frac{2}{n-1}\sigma^4$$

$$\Rightarrow MSE(\hat{\sigma}_1^2) < MSE(\hat{\sigma}_2^2)$$

# Consistent estimator

**Def**    A point estimator $\hat{\theta}_n$ of $\theta$ is <mark>**consistent**</mark>

(strongly consistent) if

$$\hat{\theta}_n \longrightarrow \theta \quad \text{in probability} \quad (a.s.)$$

$$\text{as } n \to \infty$$

**Theorem**    If an estimate satisfies bias $\to 0$ and se $\to 0$

as $n \to \infty$, the the estimate is consistent.

# Confidence sets

# Confidence sets

__Def__    A $(1-\alpha)$ - __confidence interval__ for a parameter

$\theta \in \mathbb{R}$ is an interval $C_n = (a_n, b_n)$ where

$a_n = a(X_1, ..., X_n)$,   $b_n = b(X_1, ..., X_n)$ are functions
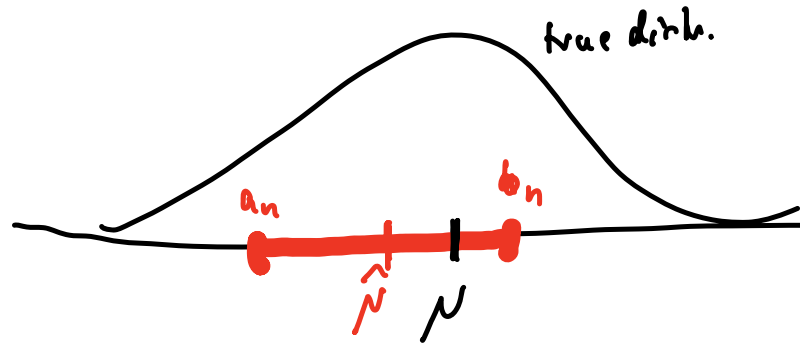
of the sample $X_1, ..., X_n$ such that

$$P_\theta \left( \theta \in C_n \right) \geq 1 - \alpha \quad \text{for all } \theta \in \Theta.$$

true
(unknown)
parameter

deterministic

random
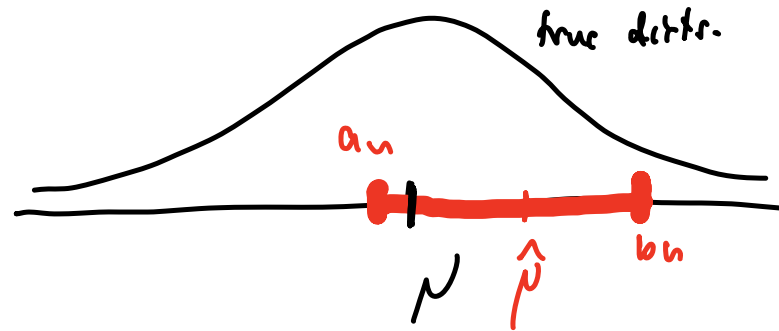
$(1-\alpha)$ is called the coverage of the confidence interval.

# Illustration

First confidence set:

First experiment
$x_1, \ldots, x_n \sim \hat{\mu}$



true distr.

$a_n$   $b_n$

$\hat{\mu}$   $\mu$

Second experiment
$x_1, \ldots, x_n \sim \hat{\mu}$

second confidence set



true distr.

$a_n$

$\mu$   $\hat{\mu}$   $b_n$

... in $(1-\alpha)$ of the repetitions, the true $\mu$ is inside the red interval.

# Example

Coin flips, with $P(X = 1) = p$, $P(X = 0) = 1 - p$,

$p \in [0, 1]$ unknown. Want to estimate it.

$\rightsquigarrow$ observe $X_1, \ldots, X_n \sim f_p$

$\hat{p}_n := \frac{1}{n} \sum_{i=1}^{n} X_i$ . Choose a confidence level $\alpha$,

now want to define $c_n = (a_n, b_n)$. To this end, define

$$\varepsilon_n^2 := \frac{\log(2/\alpha)}{2n} .$$

**Proposition:** $c_n := \left( \hat{p}_n - \varepsilon_n, \hat{p}_n + \varepsilon_n \right)$ is a CI with coverage $1 - \alpha$.
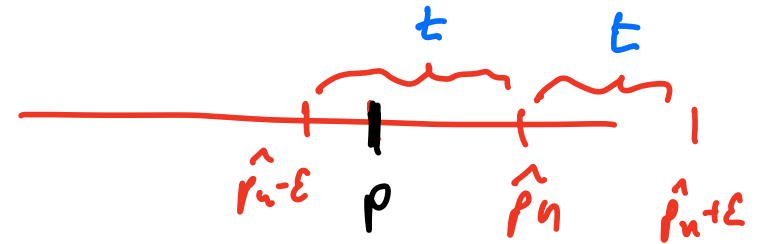
# Proof (example)

**Proof:** By Hoeffding inequality, for any $t$ we have

$$P\left( |\hat{p}_n - p| > t \right) \leq \underbrace{2 \exp(- 2n t^2)}_{\alpha}$$

Set $\alpha := 2\exp(-2nt^2)$

and solve for $t$:

$$\log\left(\frac{\alpha}{2}\right) = -2n t^2 \implies t^2 = \frac{-\log\left(\frac{\alpha}{2}\right)}{2n} = \frac{\log(2/\alpha)}{2n}$$

Choose $\varepsilon_n = t$.

# Maximum likelihood estimator

# Motivating example

$$\mathcal{F} = \{ A \mid A \text{ symmetric, } n \times n \text{ matrix, } a_{ij} \in \{0, 1\} \}$$

adjacency matrices of graphs

Observe $k$ random walks from the graph of length 10.

Goal: reconstruct (estimate) $A$

Idea: among all adjacency matrices $A \in \mathcal{F}$, select the one that has the highest likelihood to have produced the random walks you have observed.

$\leadsto$ Maximum likelihood approach

# Likelihood

More formally: Parametric family $\mathcal{F} = \{ f_\theta \mid \theta \in \Theta \}$,

observe idd points $X_1, \ldots, X_n \sim f_\theta \in \mathcal{F}$.

The likelihood of the data given a parameter $\theta_0$ is

$$P_{\theta_0}(X_1, \ldots, X_n) = P(X_1, \ldots, X_n \mid \theta_0)$$

notation!

$$= \prod_{i=1}^{n} P(X_i \mid \theta_0)$$

# Maximum likelihood

To estimate the true parameter $\theta$, we now select $\hat{\theta}$ such that this likelihood is maximized:

$$\hat{\theta} := \underset{\theta \in \textcircled{H}}{\arg\max} \; P(X_1, ..., X_n \mid \theta) \overset{\text{ind.}}{=} \underset{\theta}{\arg\max} \prod_{i=1}^{n} P(X_i \mid \theta)$$

This is equivalent to the problem

$$\hat{\theta} = \arg\max \log \left( \prod_{i=1}^{n} P(X_i \mid \theta) \right) = \arg\max \sum_{i=1}^{n} \log \underbrace{\underbrace{P(X_i \mid \theta)}_{\in [0,1]}}_{< 0}$$

which is equivalent to minimizing the negative log-likelihood:

$$\hat{\theta} = \underset{\theta}{\arg\min} \sum_{i=1}^{n} \underbrace{- \log P(X_i \mid \theta)}_{> 0}$$

= maximum-likelihood estimator  **MLE**

# Solving max. likelihood problems

Sometimes their optimization problem is easy:

- it might be able to solve it analytically (rare)

- if you are lucky it is convex

- Most typically, it is not convex.

# Example for an analytic solution

Model: $X \sim \text{Poisson}(\lambda)$, this means that

$$P(X=x) = \frac{\lambda^x e^{-\lambda}}{x!}, \quad \text{it has } E(X) = \lambda$$

$$\text{Var}(X) = \lambda.$$
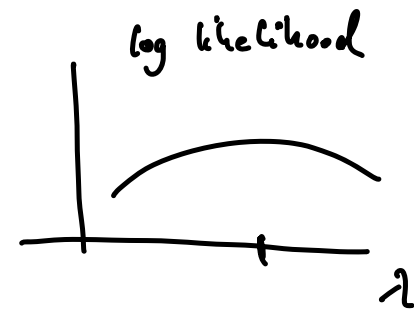
Observe $X_1, \dots, X_n \sim \text{Poisson}(\lambda)$

Want to construct the ML-estimator.

Compute the likelihood:

$$\mathcal{L}(\lambda) = P(X_1, \dots, X_n \mid \lambda) = \prod_{i=1}^{n} \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

# Example (continued)

log likelihood

$$\log(\ldots) = \sum_{i=1}^{n} \log\left(\frac{\lambda^{x_i} e^{-\lambda}}{x_i!}\right)$$

$$= \sum_{i=1}^{n} \left( x_i \log \lambda - \lambda - \log(x_i!) \right)$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad}_{f(\lambda)}$$

Now want to optimize for $\lambda$. Take the derivative (wrt $\lambda$):

$$f'(\lambda) = \sum_{i=1}^{n} \left( \frac{x_i}{\lambda} - 1 \right) = \frac{1}{\lambda}\left(\sum_{i=1}^{n} x_i\right) - n \stackrel{!}{=} 0$$

$$\Rightarrow \quad \underline{\lambda = \frac{1}{n} \sum_{i=1}^{n} x_i}$$

So $\hat{\lambda} := \frac{1}{n} \sum_{i=1}^{n} x_i$ is the ML estimate of $\lambda$.

# MLE properties

From the theory side, MLE often (but not always)
has nice properties:

(1) If the model $\mathcal{F}$ consists of "nice" functions, then
the MLE based on an iid sample is <u>consistent</u>.

(2) If $\mathcal{F}$ consists of "nice" functions, the MLE estimate $\hat{\theta}_{MLE}$
is <u>asymptotically normal</u>:

$$\frac{\hat{\theta}_{MLE} - \theta}{se} \xrightarrow{\text{in distr.}} N(0,1) \qquad \text{and}$$
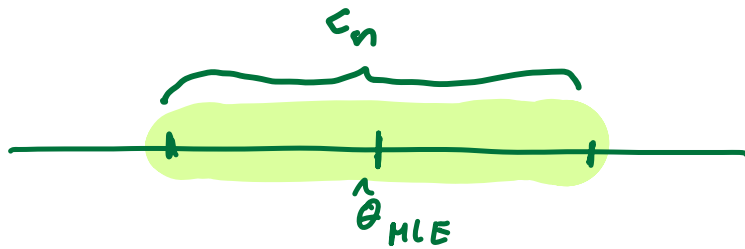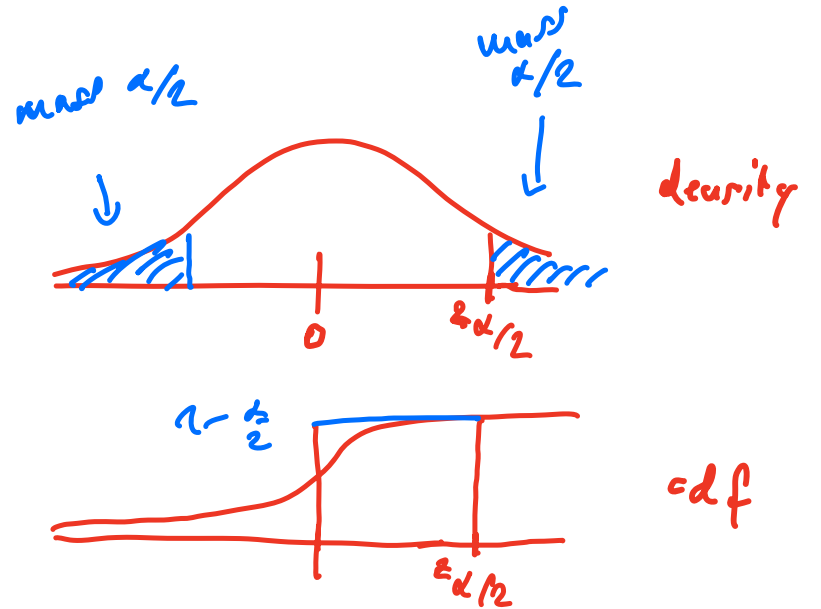
$$\frac{\hat{\theta}_{MLE} - \theta}{\hat{se}} \xrightarrow{\text{in distr.}} N(0,1)$$

(§) This can be used to construct (approximate) confidence intervals:

$$c_n := \left( \hat{\theta}_{MLE} - \underbrace{z_{\alpha/2} \cdot \hat{se}}_{-\varepsilon}, \quad \hat{\theta}_{MLE} + \underbrace{z_{\alpha/2} \cdot \hat{se}}_{+\varepsilon} \right)$$

where $z_{\alpha/2} := \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$

cdf of $N(0,1)$

mass $\alpha/2$     mass $\alpha/2$

density

$0$   $z_{\alpha/2}$

$c_n$

$\hat{\theta}_{MLE}$

$1 - \frac{\alpha}{2}$

cdf

$z_{\alpha/2}$

$c_n$ is an "approximate CI" in the sense that

$$P_\theta \left( \theta \in c_n \right) \longrightarrow 1 - \alpha \quad as \quad n \to \infty.$$

# Sufficiency & Identifiability

# Sufficiency

Intuition: given sample $X_1, \ldots, X_n \sim f_\theta \in \mathcal{F}$

we typically count the (large) sample to a statistic

$$T(X_1, \ldots, X_n) \qquad \text{( in the extreme case, one number). ?}$$

Question: can we recover the true parameter $\theta$ from this statistic

If yes one would like to call the statistics $T(X_1, \ldots, X_n)$

"sufficient".

# Sufficiency

Which properties would we need to assert sufficiency?

- when we observe two samples $X_1, \ldots, X_n$ and $X_1', \ldots, X_n'$, and $T(X_1, \ldots, X_n) = T(X_1', \ldots, X')$, then we would infer the same $\Theta$.

- when we know $T(X_1, \ldots, X_n)$, then we would need some way to recover the likelihood of the data.

Formal definition is technical, skipped.
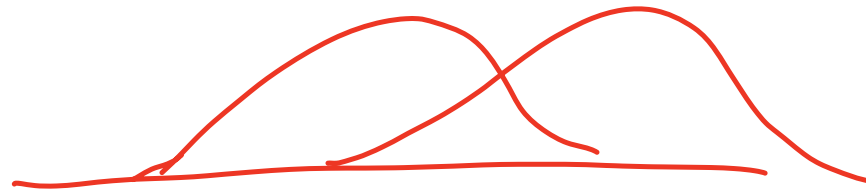
# Identifiability

Sometimes families of distributions can be described in different ways with different sets of parameters.

Def A $\underline{parameter\ \theta}$ for a family $\mathcal{F} = \{f_\theta \mid \theta \in \Theta\}$ is $\underline{identifiable}$ if distinct values of $\theta$ correspond to distinct pdfs in $\mathcal{F}$;

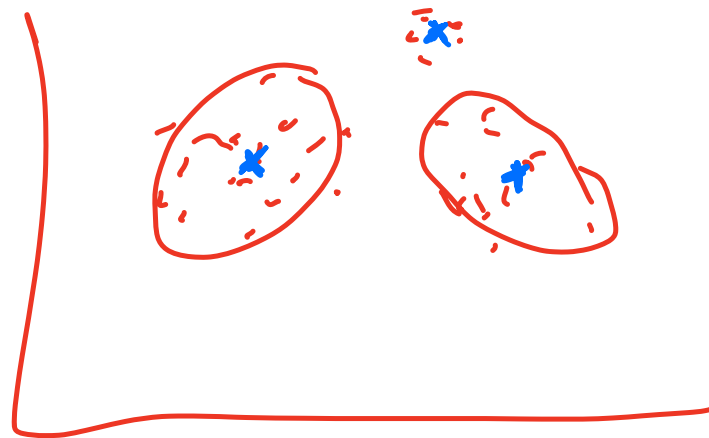$$\theta \neq \theta' \implies f_\theta \neq f_{\theta'}$$

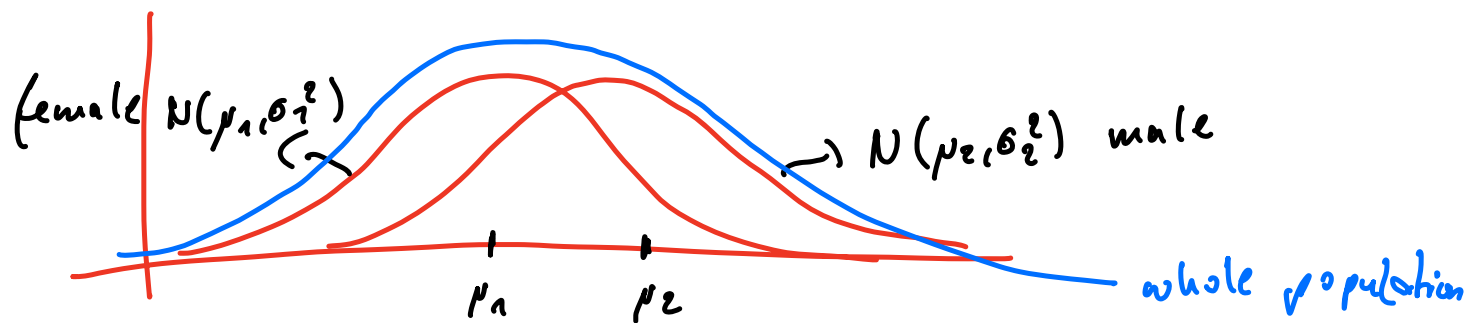# Example (identifiability)

Example:  Mixture distributions

$$\mathcal{F} = \left\{ \sum \alpha_i \, N(\mu_i, \sigma_i^2) \right\} \qquad \text{with} \quad \sum \alpha_i = 1$$



Example 1-d

Example 2-d

female $N(\mu_1, \sigma_1^2)$      $N(\mu_2, \sigma_2^2)$ male

whole population

$\mu_1$    $\mu_2$
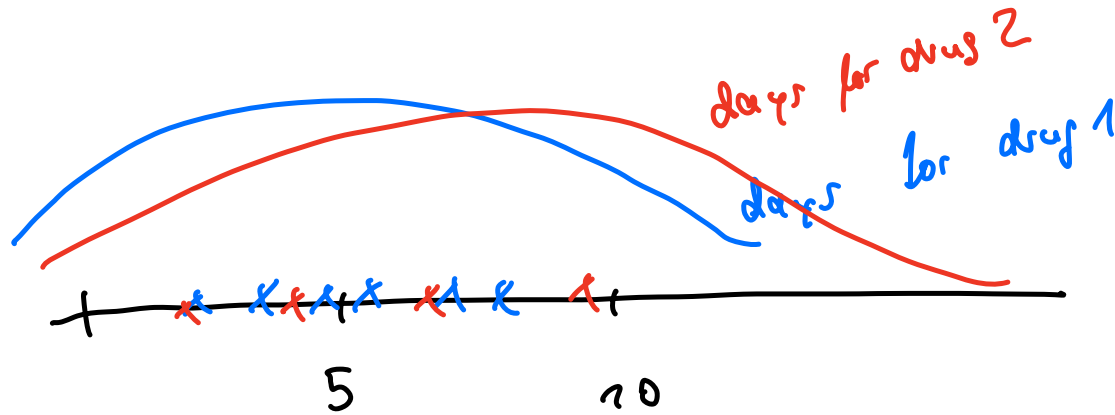
You observe samples from the whole population.

$$0.5 \, N(\mu_1, \sigma_1^2) + 0.5 \, N(\mu_2, \sigma_2^2)$$

It is impossible without further knowledge to identify the original distribution parameters just from observing the full distribution (you don't know who was female and who not)

# Hypothesis testing

# Motivation

Example: Two drugs $D_1, D_2$, we measure number of days to recovery for both drugs: $x_1, \ldots, x_n$ treated with $D_1$
$x_1', \ldots, x_u'$
$D_2$



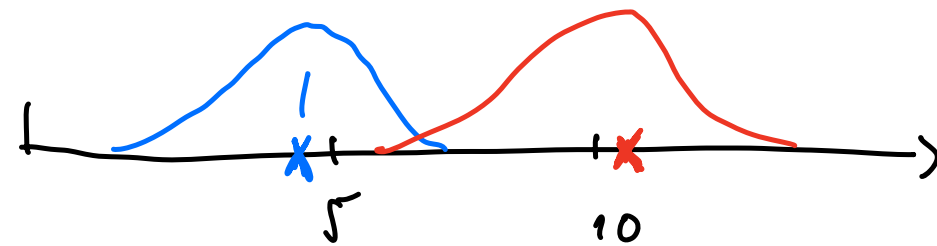days for drug 2

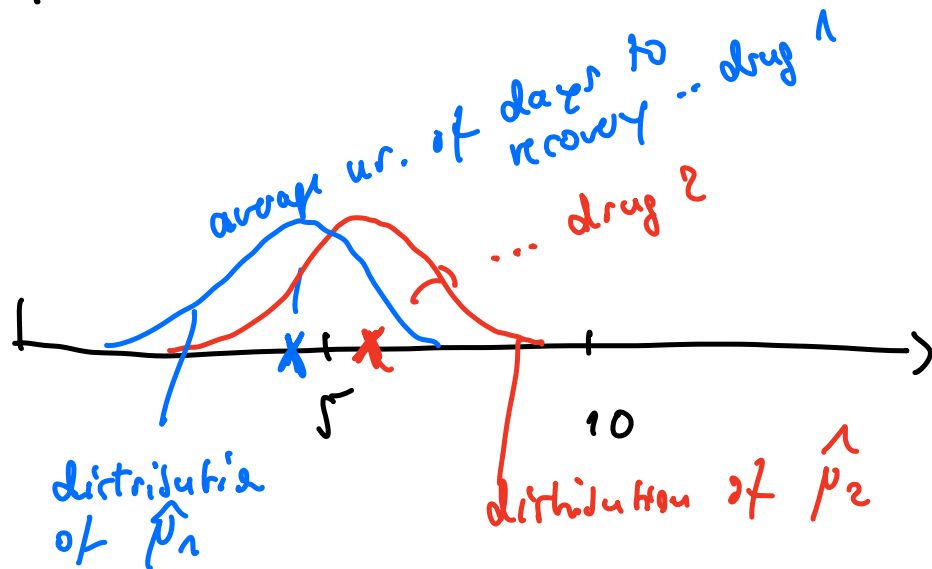days for drug 1

5     10

Question: is Drug 1 better than Drug 2 ?

# General idea

... will be: consider the distributions of the estimates $\hat{\mu}_1$, $\hat{\mu}_2$.

If they are "far apart", we would say that they are different...



average nr. of days to recovery .. drug 1

... drug 2

distribution of $\hat{\mu}_1$

distribution of $\hat{\mu}_2$

But how do we know what "far apart" is?

## Example

Want to test whether a coin is fair.

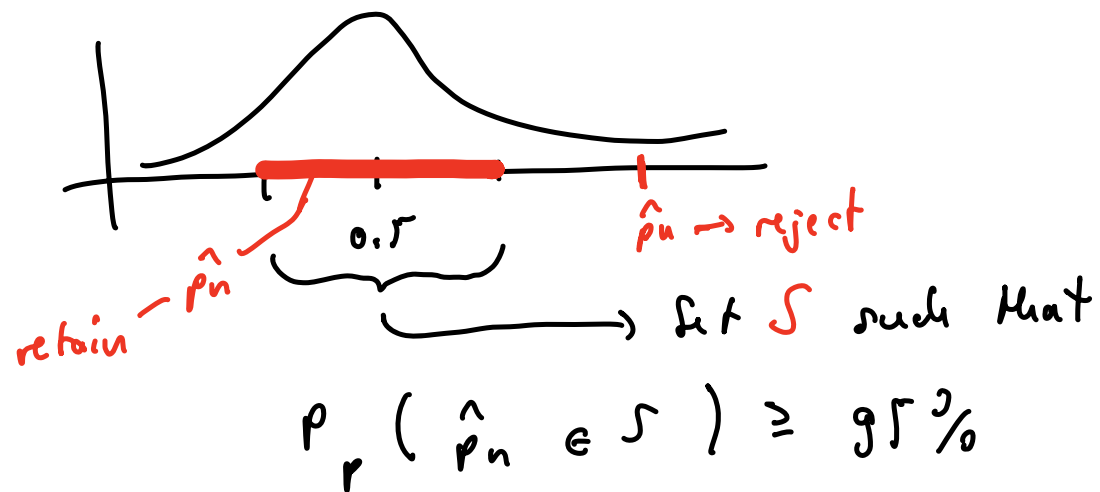**Null hypothesis:** $H_0$ : coin is fair

**Alternative hypothesis:** $H_1$ : coin is unfair

Sample many coin flips and estimate $\hat{p}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$.

We want to **reject $H_0$** if $\hat{p}_n$ is "far away" from $0.5$.

Question: "far away"?

Look at the distribution of $\hat{p}$ under the null hypothesis:



retain — $\hat{p}_n$

$0.5$

$\hat{p}_n \rightarrow$ reject

$\longrightarrow$ set $S$ such that

$$P_r \left( \hat{p}_n \in S \right) \geq 95\%$$

# More formal setup

Statistical model $\mathcal{F} = \{ f_\theta \mid \theta \in \Theta \}$. Assume that

$$\Theta_0 \subset \Theta , \quad \Theta_1 \subset \Theta , \quad \Theta_0 \cap \Theta_1 = \emptyset .$$

Want to test

$$\underbrace{H_0 : \theta \in \Theta_0}_{\text{null hyp.}} \qquad \text{against} \qquad \underbrace{H_1 : \theta \in \Theta_1}_{\text{alternative hyp.}} .$$

Sample data from the unknown $f_\theta$, compute a **test statistic**
$T(x_1, \ldots, x_n)$. Now we construct a **rejection region $R_n$**

such that
$$T(x_1, \ldots, x_n) \in R_n \implies \text{reject } H_0$$
$$T(x_1, \ldots, x_n) \in R_n \implies \text{retain } H_0$$

Typical hypotheses are of the form

- $H_0 : \quad \theta = \theta_0 \qquad$ vs $\qquad H_1 : \quad \theta \neq \theta_0$

- $H_0 : \quad \theta < \theta_0 \qquad$ vs $\qquad H_1 : \quad \theta \geq \theta_0$

Two types of error can occur:

|  | Test retains $H_0$ | Test rejects $H_0$ |
|---|---|---|
| $H_0$ true | :) | Type I error |
| $H_1$ true | Type II error | :) |

# Power of a test, $\beta$

**Def**   The *power function* of a test with rejection region $R$

is the function

$$\beta(\theta) := P_\theta(T(x) \in R)$$

- If $\theta \in \Theta_0$ then $T(x)$ should not end up in $R$.
for such $\theta$, $\beta(\theta) = P(\text{Type I error})$.
Ideally, $\beta(\theta)$ should be small.

- If $\theta \in \Theta_1$ then we hope that $T(x) \in R$. So

$$\beta(\theta) = 1 - P(\text{Type II error}).$$

Ideally, $\beta(\theta)$ is large.

# Level of a test, $\alpha$

<u>Def</u> We say that a ==test is of level $\alpha$ if==

$$\sup_{\theta \in \Theta_0} \beta(\theta) \leq \alpha$$

Intuition: worst case guarantee no matter which $\theta \in \Theta_0$ we pick, the type I error is not larger than $\alpha$.

( Intuition to remember: ==$\alpha \; \hat{=} \;$ type I error== )

# Standard approach for testing

Standard procedure: We fix the level $\alpha$ of a test in advance, for example $0.05$ or $0.01$.

Then we can also look at the type $I$ error. For example among several tests of level $\alpha$, you might now choose the one that has the smallest type-$II$-error.

Notation used often in literature:

$$\alpha = P(\text{type } I \text{ error}) \qquad \alpha = \text{level of the test}$$

$$\beta = P(\text{type } II \text{ error}) \qquad 1-\beta = \text{power of a test}$$

Remark: the power of a test is typically evaluated when we test against a concrete hypothesis $\theta_1 \in \Theta_1$. We say "the power of the test against alternative $\theta_1 \in \Theta_1$."

# Uniformly most powerful test

**Def** Let $\mathcal{T}$ be a set of tests of level $\alpha$ for testing

$$H_0: \theta \in \Theta_0 \quad vs \quad H_1: \theta \notin \Theta_0.$$

A test in $\mathcal{T}$ with power function $\beta(\theta)$ is

uniformly most powerful (UMP) if

$$\beta(\theta) \geq \beta'(\theta) \quad \text{for all } \theta \in \Theta_0^C$$

and for all $\beta'$ that are power functions
for other tests in $\mathcal{T}$.

**Remark:** In practice it is often impossible to find an UMP test.

Neyman - Pearson - lemma
&
likelihood ratio tests

# Neyman - Pearson

__Theorem__  Suppose we test $H_0 : \theta = \theta_0$ against $H_1 : \theta = \theta_1$.

Consider

$$T = \frac{\mathcal{L}(\theta_1)}{\mathcal{L}(\theta_0)} = \frac{\prod\limits_{i=1}^{n} f(x_i \mid \theta_1)}{\prod\limits_{i=1}^{n} f(x_i \mid \theta_0)} \Bigg\} \quad \text{==likelihood ratio.==}$$

Assume we reject $H_0$ if $T > k$ (for some $k$).

If we choose $k$ such that $P_{\theta_0}(T > k) = \alpha$,

then this is the ==most powerful level-$\alpha$-test.==

## More general likelihood-ratio-test:

Parameter space $\Theta$, $\Theta_0 \subset \Theta$, $\Theta_1 = \Theta_0^c$. Then we

consider the test statistic

$$\tilde{T} = \frac{\sup\limits_{\theta \in \Theta_0} \mathcal{L}(\theta)}{\sup\limits_{\theta \in \Theta_1} \mathcal{L}(\theta)} \qquad \text{or even simpler} \qquad T = \frac{\sup\limits_{\theta \in \Theta_0} \mathcal{L}(\theta)}{\sup\limits_{\theta \in \Theta} \mathcal{L}(\theta)}$$

and we determine a parameter $\lambda$ such that the rejection region

is of the form $R = \{T \leq \lambda\}$.

In practice the difficulties are
- compute the suprema (in practice)
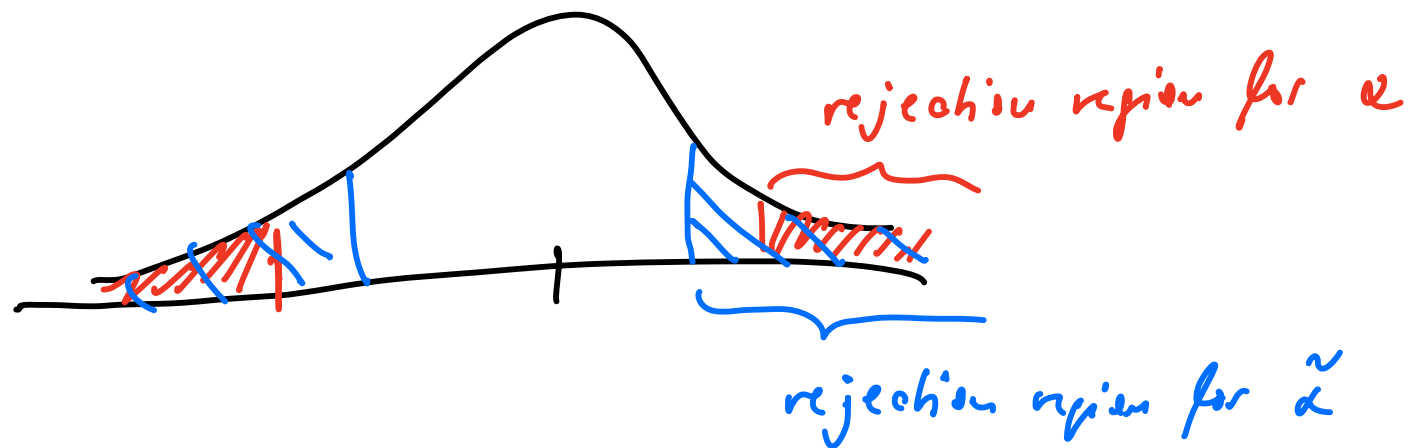- fix $R$, fix $\lambda$ (in theory)

p - values

# p - values

Consider a test at level $\alpha$, and denote its
rejection region as $R_\alpha$.

Recall: $\alpha = P(\text{Type} - \text{I} - \text{error})$.

The smaller $\alpha$, the more difficult does it get to reject $H_0$.

(we often even have that $\alpha < \tilde{\alpha} \Rightarrow R_\alpha \subset R_{\tilde{\alpha}}$)



rejection region for $\alpha$

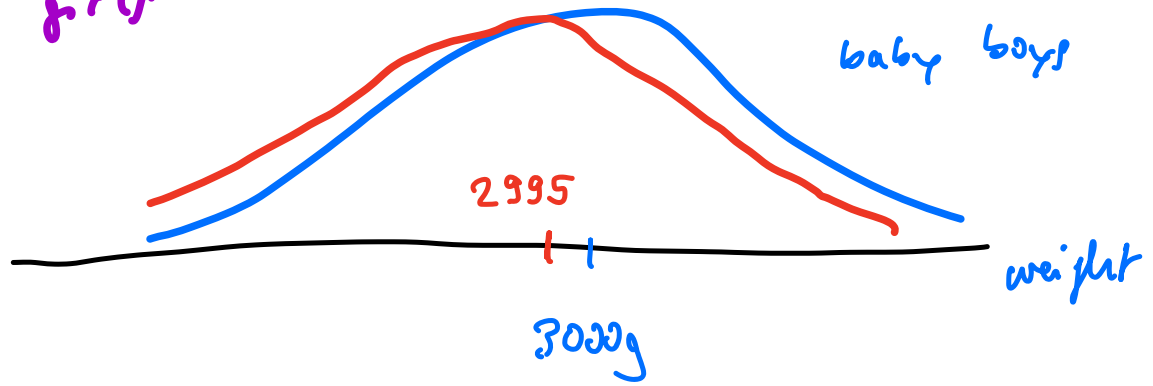rejection region for $\tilde{\alpha}$

# p-value

**Def** The ==p-value== is defined as

$$p = \inf \{ \alpha \mid T(x_1, \dots, x_n) \in R_\alpha \}$$

i.e. the smallest $\alpha$ for which the level-$\alpha$-test would reject the null hypothesis.

Intuition: smaller p-values are "better", more evidence for rejecting the null (less error).

Example baby boys and girls

baby boys

2995

3000g

weight

sample distribution

Sample many baby girls, many baby boys

$\hat{\mu}_g$ , $\hat{se}_g$

mean weight $\hat{\mu}_b$ , & $\hat{se}_b$

distribution of
test statistics

2995                3000

for a large test will find a statistically
significant difference.   ~ small $p$

Multiple testing

# Motivation

Example: gene expression data

| | patients with cancer $(n = 20)$ | | | | | Control group $(n = 20)$ | | | |
|---|---|---|---|---|---|---|---|---|---|
| gene 1 | | | | | | | | | |
| gene 2 | | | | | | | | | |
| ... | | | | | | | | | |
| gene 17 | 0.5 | 0.2 | 0.9 | 0.8 | 0-5 | 0.01 | 0.05 | 0.1 | 0.02 |
| ... | | | | | | | | | |
| gene 105 | | | | | | | | | |
| gene 1000 | | | | | | | | | |

$=: m$

$\alpha \%$ of the test will "ring a bell"

Assume we run, for each gene, a test of level $\alpha$

$P( \text{ test } i \text{ makes type-I-error}) = 5\%$. Now we have $m$ tests.

$P\left(\text{at least one of the tests makes a type-I-error}\right) =$

$= P\left(t_1 \text{ makes error } \underline{or} \ t_2 \text{ error } \underline{or} \ .. \quad \underline{or} \ t_m \text{ makes error}\right)$

$= 1 - P\left(\text{no error in } t_1 \ \underline{and} \text{ no error in } t_2 \text{ and } ...\right) = \overset{\text{assume independence}}{\Longleftarrow}$

$= 1 - \prod_{i=1}^{m} P\left(\text{no error in } t_i\right) = \underline{1 - (0.95)^m} \xrightarrow[m \to \infty]{} 1$

$m = 1 \quad \Rightarrow \quad * = 0.05$

$m = 10 \quad \Rightarrow \quad * = 0.40$

$m = 50 \quad \Rightarrow \quad * = 0.92$

Many "wrong" tests!

# Family-wise error rate (FWER)

Definition: Consider a family of $m$ tests. The family-wise error rate (FWER) is the probability that at least one type-I-error occurs in the family:

$$
\begin{aligned}
FWER = P( & t_1 \text{ makes type-I-error } \underline{or} \\
& b_2 \quad \cdots \\
& \cdots \\
& b_m \text{ makes type-I-error } ).
\end{aligned}
$$

# Bonferroni correction

Assume we run $m$ tests, and we want to achieve a FWER $\alpha$ (e.g. $\alpha = 0.05$). Then we run the individual tests with level $\frac{\alpha}{m} =: \alpha_{single}$. Then:

$$FWER = P(\text{at least one type-I-error}) =$$

$$= P(t_1 \text{ error } \underline{or} \ t_2 \ldots) \leq$$

$$\leq \sum_{i=1}^{m} \underbrace{P(t_i \text{ makes error})}_{\alpha_{single}} = m \cdot \alpha_{single} = m \cdot \frac{\alpha}{m} = \alpha.$$

# Bonferoni, discussion

Bonferoni controls the FWER.

Advantage: simple, correct

Disadvantage: too conservative, low power (high type-II-error)
the test barely discovers anything!

# False discovery rate, FDR

**Def** Assume we have a family of $m$ tests. We call

$$E\left( \frac{\#\text{ false rejections}}{\#\text{ all rejections}} \right) =: \quad FDR$$

the false discovery rate.

Benjamini/Hochberg: Controlling FDR

# Benjamini/Hochberg (1998) approach:

- Fix FDR $\alpha$ in advance.

- Run the $m$ individual tests and evaluate their $p$-values.

- Sort $p$-values increasingly: $p_{(1)} \leq p_{(2)} \leq p_{(3)} \leq \ldots \leq p_{(m)}$

- Define thresholds $l_i := i \cdot \frac{\alpha}{m}$

- Find the largest index $i_0$ such that $p_{(i_0)} \leq l_{i_0}$.

  (below the red line)

- Reject the hypotheses for $i = 1, \ldots, i_0$, retain all the others.

# Theorem (Benjamini-Hochberg)

Theorem :     If the Benjamini-Hochberg procedure is applied (and the tests are independent), then regardless of how many null hypotheses are true and regardless of the distribution of $p$-values when the null is false, we obtain     FDR $\leq \alpha$.

(Remark: similar approach also works without independence assumption, many modifications exist.)

# Intuition

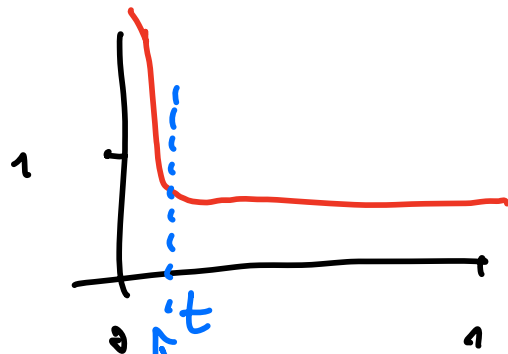- Under the null hypothesis, the p-values always have a uniform distribution on $[0, 1]$.



density of p-values under $H_0$

density of p-values under $H_1$

ideally it should look like this

If we have some $H_0$ and some $H_1$ being true over here the density would maybe look like that:



here we have (hopefully) many of the $H_1$ s but we also have some $H_0$ s.

Goal: set threshold $t$ such that FDR satisfies what we want.

Integral of the pink area: Expected number of p-values corresponding to $H_1$ that are below $t_1$

Integral of blue area:       ...       $H_0$

$H_0$

$1$

$H_1$

$t_1$          $t_2$

By moving $t$ from $0$ to $1$ we control the FDR:

For $t_1$, the FDR is small

$t_2$                              large

# General Remarks

- BH tends to have more power than Bonferroni

- BH controls FDR, not FWER (overall type-I-error)!

- BH works best in sparse regime where only few tests reject the null

- BH gives guarantees on FDR, but in general does not minimize it.

- When all the $H_0$ are true, BH $\approx$ Bonferroni.

# Non-parametric tests

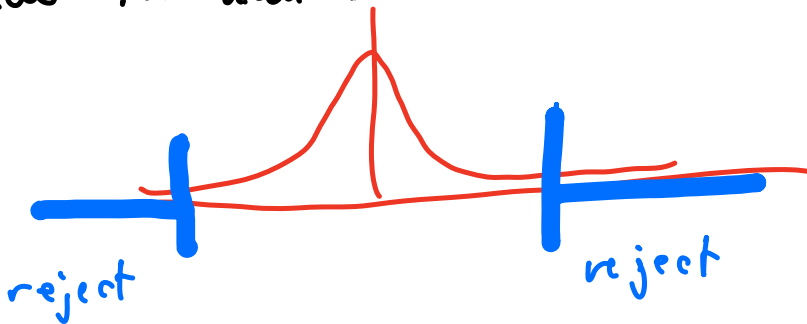# Non-parametric tests

Standard (parametric scenario):

- Statistical model $\mathcal{F} = \{ f_\theta \mid \theta \in \Theta \}$



distribution of the samples

- Observe data, compute a test statistics $T_n$, for example the mean $\bar{X}$

- Need to know the distribution of the test statistics $T_n$ under the null distribution:



reject                    reject

distribution of $T_n$ under the null hyp.

# Goodness-of-fit test (gof)

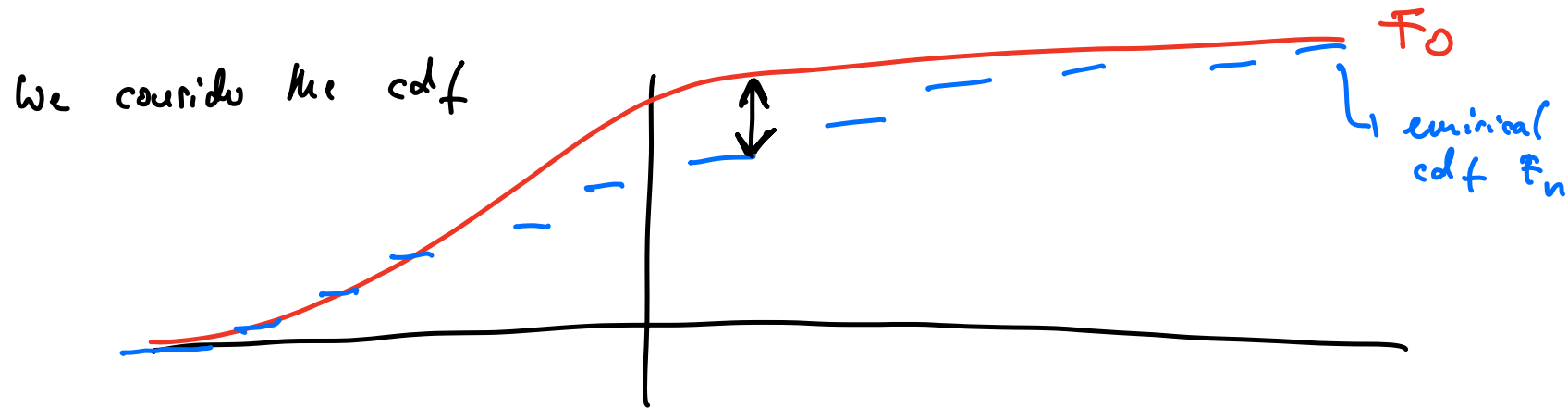==Goodness-of-fit tests:==  Goal is to test whether a data set comes from a particular distribution $F_0$

$H_0: \quad F(x) = F_0(x)$

         ↑ true distribution that generated the data

$H_1: \quad F(x) \neq F_0(x)$

# Kolmogorov-Smirnov test for gof

We consider the cdf



$F_0$

emirical cdf $F_n$

$F_0$ = cdf of the given distribution

$F_n$ = cdf of the data

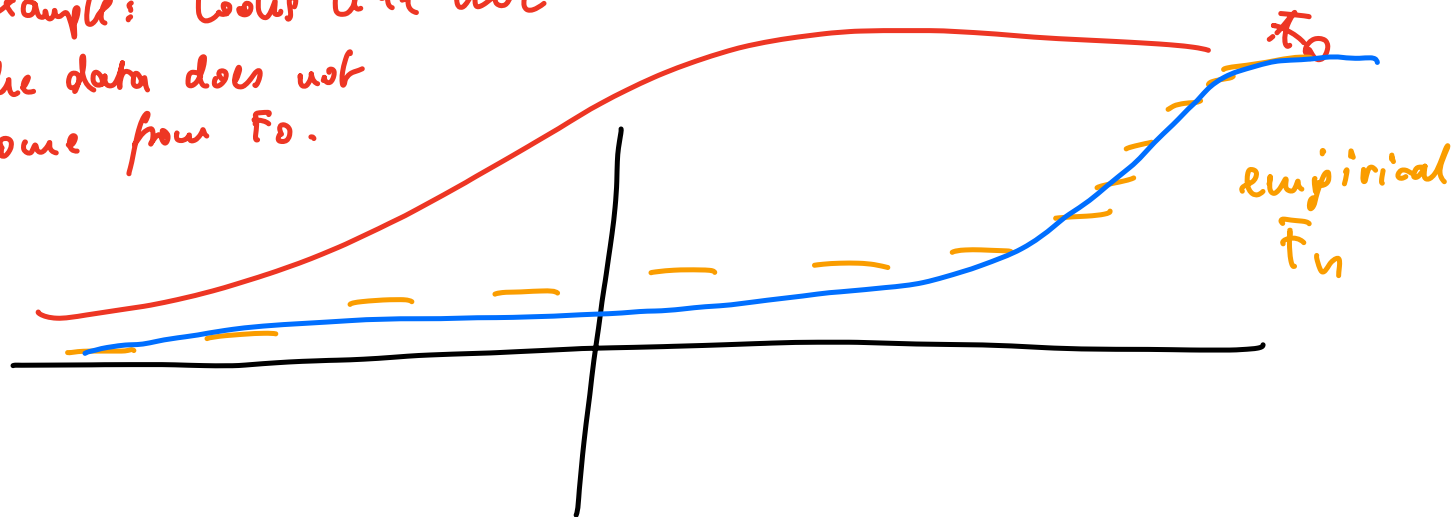$$D_n := \sup_{x \in \mathbb{R}} \left| F_n(x) - F_0(x) \right|$$

By the Glivenko-Cantelli theorem we know that under the null hypothesis, $F_n \rightarrow F_0$ uniformly. a.s.

It is possible to compute the distribution of $D_n$, and it is independent of $F_0$, it just depends on $n$.

From this we can compute rejection thresholds. and design a test.



Example: looks like here the data does not come from $F_0$.

$F_0$

empirical $\overline{F}_n$

# two sample test

$X_1, \ldots, X_n \sim F_1$ a first sample

distributed according to $F_1$,

$Y_1, \ldots, Y_m \sim F_2$ a second sample distributed acc. to $F_2$.
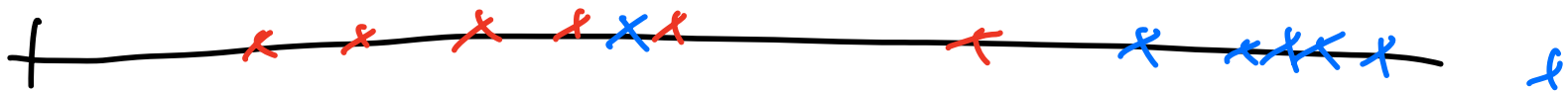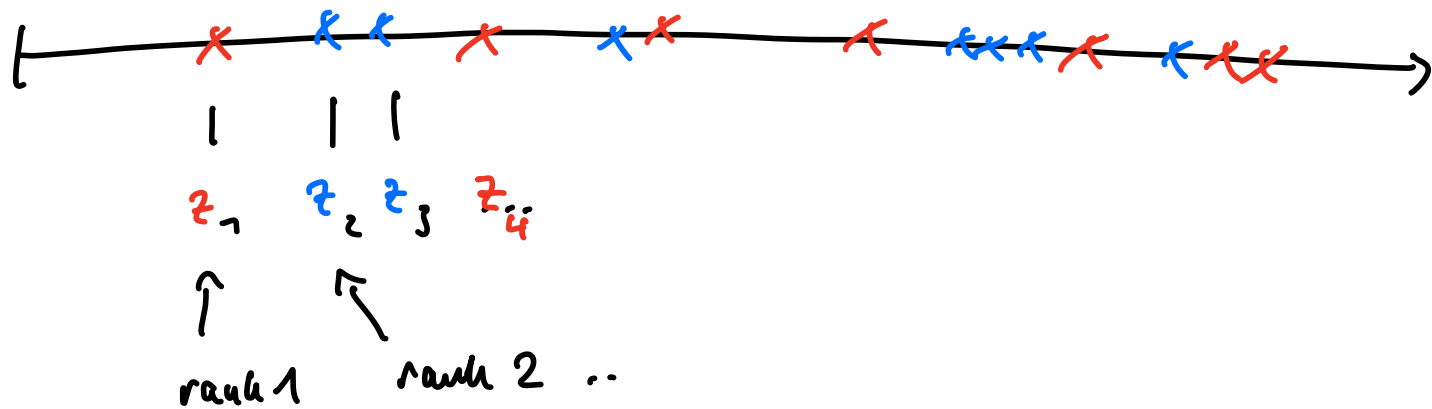
Question: $F_1 = F_2$ ?

$H_0: \quad F_1 = F_2 \qquad , \quad H_1: \quad F_1 \neq F_2$

# Wilcoxon - Mann-Whitney test (based on ranks)

Test:
- "Pool the sample": $\underline{x_1, \ldots, x_n,}$ $\underline{y_1, \ldots, y_m}$ $\in \mathbb{R}$

- Sort the pooled sample in increasing order and retrieve the rank of all points $\leadsto$ rank($x_i$)
  rank($y_i$)
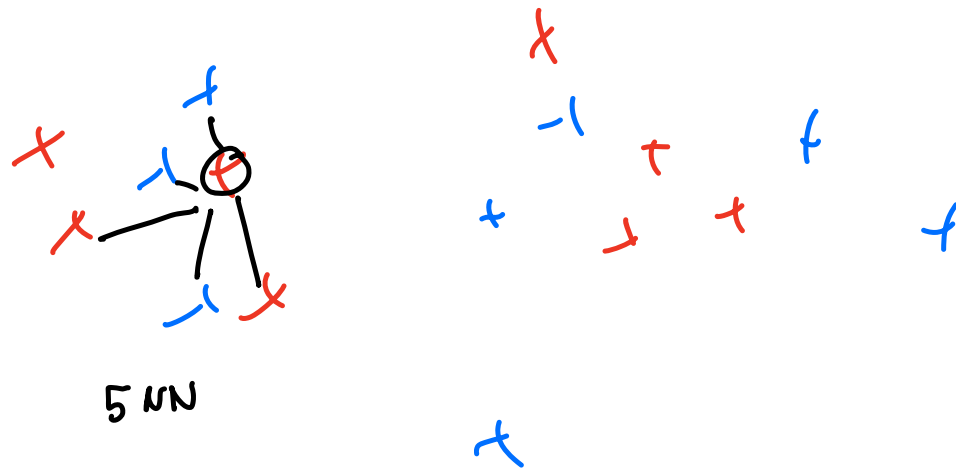


rank 1    rank 2 ..

- Compute the rank sums for both groups:

red group: $\quad W_{red} = \sum_{i \in \text{red population}} \text{rank}(x_i)$

$$W_{blue} = \sum_{i \in \text{blue pop.}} \text{rank}(y_i)$$

- If $|W_{red} - W_{blue}|$ is small, we retain $H_0$, if large, reject $H_0$.

# Extension to a multivariate setting using k nearest neighbors
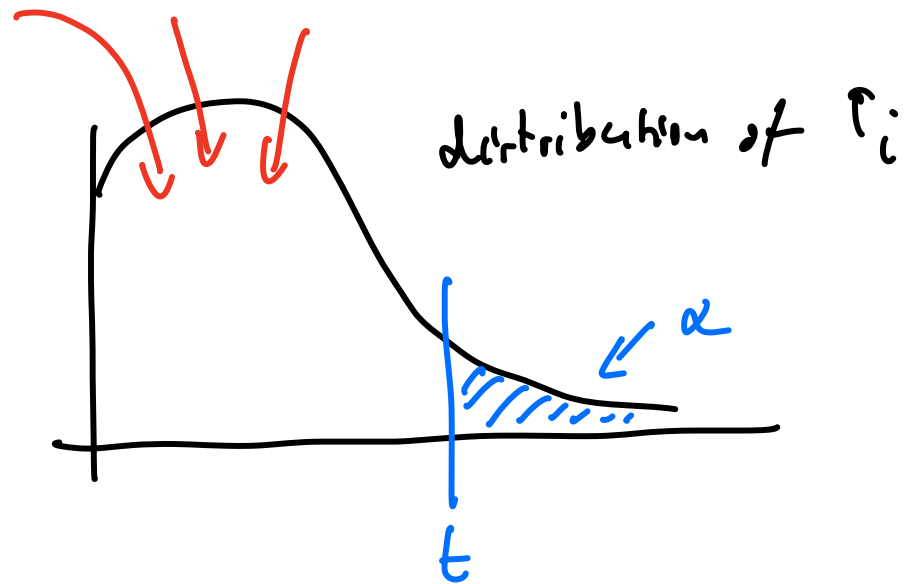
- Two samples, we pool them:



5 NN

- For each point, we look at the colors of the
  k nearest neighbors:

- Under the null hypothesis we expect that the
  number of red neighbors ≈ number of blue neighbors.

# Permutation (randomization) tests

- Sample $X_1, \ldots, X_n$  <u>group A</u>  $\rightsquigarrow$ mean $\bar{X}$

  $Y_1, \ldots, Y_n$  <u>group B</u>  mean $\bar{Y}$  $\rightarrow$ ?

- Compute observed statistics $T_{observed} = $ mean(red) $-$ mean(blue)

- Pool the sample

- For $k = 1, \ldots, 10^3$: shuffle the group membership ("colors")

- Compute the difference $T_k = $ mean(red) $-$ mean(blue)

~) $T_1, T_2, \ldots, T_{1000}$



distribution of $T_i$

- Find $\alpha$-quantile to determine rejection threshold.

- Check whether the observed $T_{observed}$ on the true data is $\leq t$.

# Bootstrap tests

# Motivation

Motivation: $X_1, \ldots, X_n \sim F$, no knowledge on $F$ want to estimate a parameter $\theta = t(F)$. You purvah an estimate $\hat{\theta}$ based on $X_1 \ldots, X_n$, want to know how reliable $\hat{\theta}$ is.

The first thing to look at is the standard error se.

- If we have assumptions on $F$, we can analytically compute the distribution of $\hat{\theta}$, the $\hat{se}$, ... (this is rare!)

distribution of $\hat{\theta}$

$\hat{\theta}$

(1)

- We could also try to obtain many samples

$$X_1^{(1)}, \dots, X_n^{(1)}$$

$$X_1^{(2)} \dots X_n^{(2)}$$

$$\vdots$$

$$X_1^{(m)} \dots X_n^{(m)}$$

and then estimate the distribution of $\hat{\theta}$:

Problem: need too many samples.

(2)

$\hat{\theta}$

Then we could maybe build a hist on this.

# Idea of the bootstrap

- Given the sample $X_1, \ldots, X_n \rightsquigarrow$ estimate $\hat{\Theta}_{orig}$

- Draw a subsample of $X_1 \ldots X_n$, compute $\hat{\Theta}^*$, repeat very often

(3)



of $\Theta^*$

histogram based on resampled data

"Hope: histogram of $\hat{\Theta}^*$ (3) is "close" to histogram of $\hat{\Theta}$ (2), which is close to (1)

Example: estimate the standard error of an estimate $\hat{\Theta}$

# Algorithm in pseudocode

Input: $X_1, \ldots, X_n$ — number of original sample points

For $b = 1, \ldots, B$ — number of bootstrap replications

- Sample $X_1^*, \ldots, X_n^*$ uniformly with replacement from $X_1, \ldots, X_n$

- Estimate the parameter $\hat{\theta}_b^*$

gives us $\hat{\theta}_1^*, \ldots, \hat{\theta}_B^*$

Estimate the standard error $\hat{se}$ of the original estimate $\hat{\theta}$ by the standard dev. of the bootstrap replicates:

$$\hat{se}_B := \left( \frac{1}{B-1} \sum_{b=1}^{B} \left( \hat{\theta}_b^* - \left( \underbrace{\frac{1}{B} \sum_{i=1}^{B} \hat{\theta}_i^*}_{\text{mean of replicates}} \right) \right)^2 \right)^{1/2}$$

Does it always work?

# Consistency result for bootstrap

**Theorem** ( Consistency of the estimate of the standard errors)

- Assume that $X_1, \ldots, X_n \sim F$, iid, and

$$E\left(\|X_1\|^2\right) < \infty .$$

- Let $\hat{\theta}_n = g(X_1, \ldots, X_n)$ be the parameter that we estimate. Assume that $g$ is continuously differentiable in a neighborhood of $\mu = E X_1$, with a non-zero gradient. Then the bootstrap estimate of the standard error is strongly consistent.

# Example where it goes wrong

$X_1, \ldots, X_n \sim \text{Uniform}[0, \theta]$, where $\theta \in [0, 1]$, unknown.

Want to estimate $\theta$. The ML estimate of $\theta$ is simply the largest number we observe:
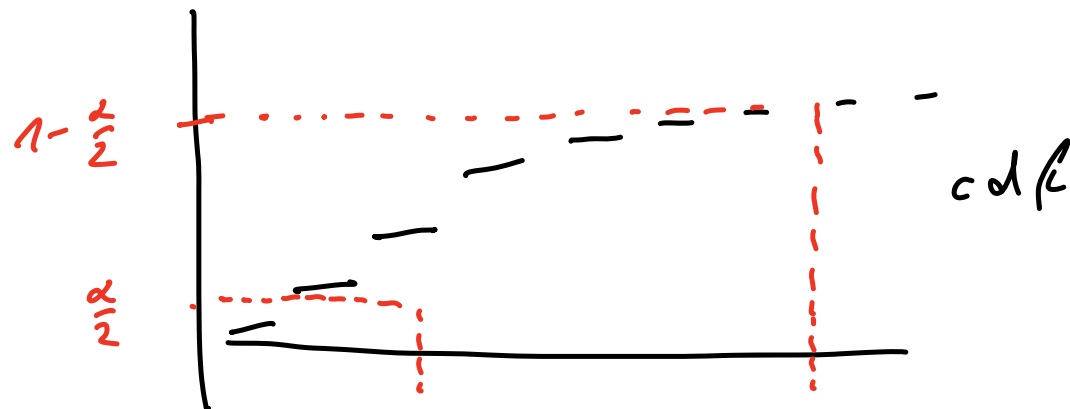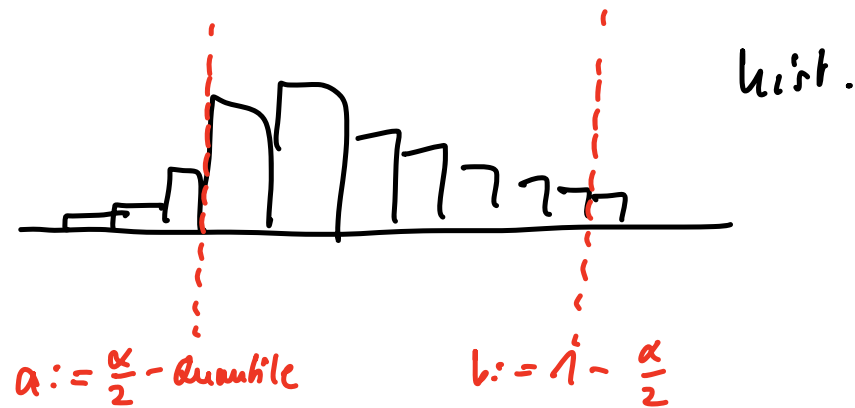
$$\hat{\theta} = \max_{i=1\ldots n} X_i.$$

Estimating the $\hat{se}$ by bootstrap is going to fail.

Estimating tails or extreme values by bootstrap is problematic.

# Confidence sets by bootstrap

Bootstrap - percentile - method:

- Given sample $x_1, \ldots, x_n$, estimate $\hat{\theta}$

- Generate bootstrap replicates $\hat{\theta}_1^*, \ldots, \hat{\theta}_B^*$
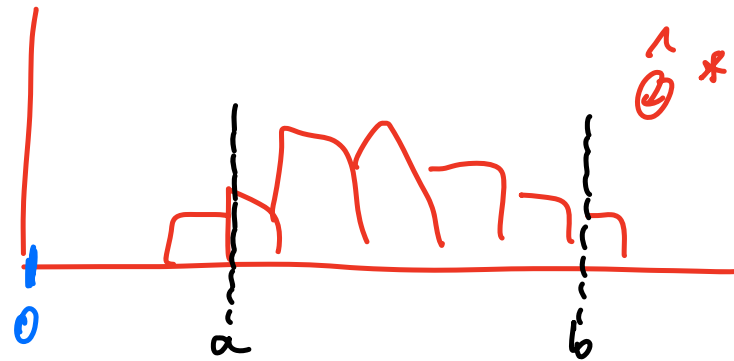
- look at the histogram of the $\hat{\theta}_b^*$



hist.

$a := \frac{\alpha}{2} - $ Quantile          $b := 1 - \frac{\alpha}{2}$

$1 - \frac{\alpha}{2}$

$\frac{\alpha}{2}$

cdf

- $CI = [a, b]$

  It has coverage $1 - \alpha$ because

  $$P_\theta \left( \hat{\theta} \in CI \right) \geq 1 - \alpha$$

  <span style="color:red">(approximately, because $n, \theta$ finite)</span>

Subsequently you can construct bootstrap

tests in the obvious way



$H_0: \hat{\theta} = \theta \quad vs \quad H_1: \hat{\theta} \neq \theta$

Bayesian statistics

# Frequentist vs. Bayesian statistics

## Frequentist statistics:

- Probability = limiting frequency

- parameters $\theta$ are constants, we cannot assign probabilities to them

- statistics behaves well when repeated often

## Bayesian statistics

- probability = degree of belief

- parameters do have probabilities

- have a prior belief about the world, update it based on observed data.

# Bayesian statistics: the model

Assume a statistical model $\{f_\theta \mid \theta \in \Theta\}$, as in frequentist approach.

It encodes our prior assumptions on the data-generating process in general.
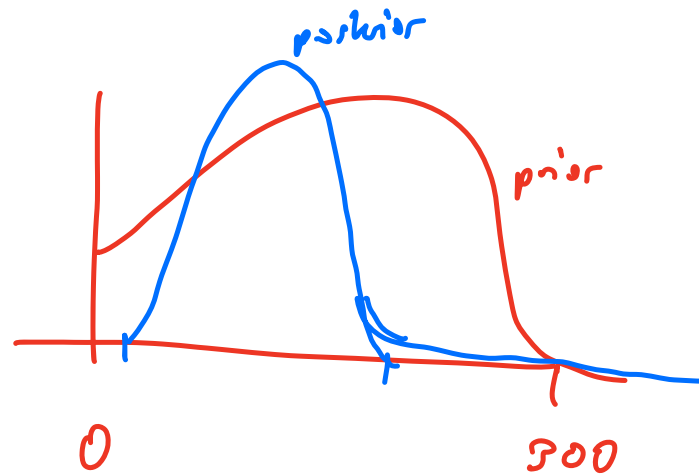
$\Theta$ unknown, want to estimate it.

# Bayesian approach: prior distribution

We assume that we have a prior belief about the parameters $\Theta$:

$f(\Theta)$ prior distribution

the parameters, not the data

# Bayesian statistics: the likelihood

Observe data $X_1, \ldots, X_n$ iid from some of the $f_\theta$ ($\theta$ unknown).

We call $f(X | \theta)$ the likelihood of the data given the parameter $\theta$

density → $f$

data → $X$

parameter → $\theta$

(In frequentist world, we could now use MLE to select the para that maximizes the likelihood)

# Bayesian statistics: posterior

Now we update our belief & we compute the

posterior using Bayes rule:     $f(\theta \mid X_1 \ldots X_n)$

$$\underbrace{f(\theta \mid X_1 \ldots X_n)}_{\text{posterior}} = \frac{\overbrace{f(X_1 \ldots X_n \mid \theta)}^{\text{likelihood}} \cdot \overbrace{f(\theta)}^{\text{prior}}}{\underbrace{\int f(X_1 \ldots X_n \mid \theta)\, f(\theta)\, d\theta}_{\substack{\text{normalizing constant} \\ \text{(does not depend on } \theta \text{ any more)}}}}$$

The posterior is a distribution.

# Statistics derived from posterior

- Now you can make statements based on the posterior.

  - If you want to return one "best guess" for $\theta$,
    you could use
    - max of posterior (MAP)

      - mean of posterior

- You can construct confidence interval:

  find $a, b$ such that

  $$P(\theta \in [a, b]) = 95\%.$$

# Discussion

**Advantages:**

- easy to interpret

- natural way to incorporate prior knowledge

**Disadvantages**

- analytic solutions are rare, typically you have to solve computationally hard problems

- need to choose a prior