# INFORMATION

- ## ~~Fisher~~ Information

When we perform an experiment, we typically collect a huge amount of data that
we need to clean and reduce in order to make a statement on whatever quantity
we are interested on.

Examples → In CUORE, we have ~ 200 TB of raw data, but our publications
    just report the result on the halflife of an isotope.
    → CMS or ATLAS have PB of data, but just measured the
        Higgs mass and cross section.

We need to define a method to select the useful information.
But first we need to define the requirements for what we call information:

→ The information should increase with the number of observations

→ The information should be conditional on what we want to learn from
    the experiment.
    Data which are irrelevant to the hypothesis under test should contain
        no information.
    → The greater the information, the better should be the precision of the experiment.


- ## Likelihood

Let's take a real random variable $\vec{x}$ with PDF $f(\vec{x}|\vec{\theta})$,
    where $\vec{\theta}$ is a set of real parameters.
The set of allowed values of $\vec{x}$ is $\Omega_{\theta}$, which might depend on $\vec{\theta}$.

Suppose we make a set of $n$ observations of $\vec{x}$: $\vec{x}_1, \dots, \vec{x}_n$

The joint PDF of $\vec{x}$ is:  $P(\vec{x}|\vec{\theta}) = P(\vec{x}_1, \dots, \vec{x}_n|\vec{\theta}) = \prod_{i=1}^{n} f(\vec{x}_i|\vec{\theta})$

Since the values $\vec{x}_i$ are fixed (they are measured!), $P$ is no longer a PDF,
but only a function of $\vec{\theta}$, and we denote it as $\mathcal{L}$:

$$\boxed{\mathcal{L}(\vec{\theta}) = \mathcal{L}(\vec{x}|\vec{\theta}) = \prod_{i=1}^{n} f(\vec{x}_i|\vec{\theta})}$$

- Statistic

A statistic is any new random variable $t = t(\vec{x}_1, ..., \vec{x}_n)$.

For example, the average $\langle \vec{x} \rangle$ is a statistic.

- Fisher information

Assume that  1) $\Omega_\theta$ is independent of $\vec{\theta}$

           2) $\mathcal{L}(\vec{x}|\vec{\theta})$ is regular enough so that the operator

$$\frac{\partial^2}{\partial\theta_i \partial\theta_j} \quad \text{and} \quad \int d\vec{x} \quad \text{commute.}$$

The Fisher information given by an observation $n$ about the parameter $\theta$ is defined as:

$$I_n(\theta) = E\left[\left(\frac{\partial \ln \mathcal{L}(x|\theta)}{\partial\theta}\right)^2\right]$$

$$= \int_\Omega \left(\frac{\partial \ln \mathcal{L}(x|\theta)}{\partial\theta}\right)^2 \mathcal{L}(x|\theta)\, dx$$

If $\vec{\theta}$ has $k$ dimensions, $I_n(\theta)$ is a $k \times k$ matrix:

$$\left[I_n(\vec{\theta})\right]_{ij} = \int_\Omega \frac{\partial \ln \mathcal{L}}{\partial\theta_i} \cdot \frac{\partial \ln \mathcal{L}}{\partial\theta_j} \cdot \mathcal{L} \cdot dx$$

Equivalently, one can prove that:

$$\left[I_n(\vec{\theta})\right]_{ij} = -E\left[\frac{\partial^2}{\partial\theta_i \partial\theta_j} \ln \mathcal{L}\right]$$

The Fisher information is additive:  $I_N(\theta) = N\, I_1(\theta)$

● Sufficiency

A statistic $t = t(\vec{x})$ is sufficient for $\theta$ if the conditional density function of $\vec{x}$ given $t$, $f(\vec{x}|t)$ is independent of $\theta$.

If $t$ is a sufficient statistic, any strictly monotonic function of $t$ is also a sufficient statistic.

$\Rightarrow$ There is as much information about $\theta$ in $T$ as there is in the original data $\vec{x}$.

$\Rightarrow$ No other function of the data can give any further information about $\theta$.

Example: The set $t = \vec{x}$ is sufficient, since it carries all the initial information. However, it provides no data reduction, so it is useless.

If $t(\vec{x})$ is a sufficient statistic for $\theta$, the likelihood factorises as:
$$\mathscr{L}(\vec{x}|\vec{\theta}) = g(t, \vec{\theta}) h(\vec{x}) \qquad \text{and viceversa}$$

where: $h(\vec{x})$ does not depend on $\vec{\theta}$

$g(t, \vec{\theta}) \propto A(t|\theta)$, the conditional probability density for $t$ given $\theta$.

Therefore $A(t|\theta) = \int \mathscr{L}(\vec{x}|\vec{\theta}) d\vec{x}$

In general, for any statistic $t$:
$$I_t(\vec{\theta}) \leq I_x(\vec{\theta})$$
with the equality if and only if $t$ is a sufficient statistic.

In other words, the information provided by a sufficient statistic is the same as that of the original sample $\vec{x}$.

# MEASUREMENT THEORY

In general, whenever we perform a measurement, we need to convey the result in a clear and synthetic way. Often times our result is a number (or a set of numbers) that will/should be used by others in the future, so we need to minimize the possible ambiguity on the underlying meaning of the quantity we quote.

Suppose ~~that~~ we collect some data $\vec{n}$ distributed with a PDF $f(\vec{n}|\vec{\theta})$, and want to make a statement on ~~some of the par~~ one parameter $\theta$ (out of the vector $\vec{\theta}$).

We can ask the following questions:

→ Based on the measured data $\vec{n}$, what is the single value $\hat{\theta}$ that is closest to the true (unknown) value of $\theta$?

⇒ Point estimation

→ Based on the measured data $\vec{n}$, what is the range of values that is most likely to include the true (unknown) value of $\theta$?

⇒ Interval estimation

~~Based on the~~ m

→ Is our model $f(\vec{n}|\vec{\theta})$ good enough to describe the measured data?

⇒ Goodness of fit

→ In the case we want to test the existence of new physics, e.g. the presence of a ~~any~~ new signal over a known background, are the measured data described better by the background-only or by the signal + background model?

⇒ Hypothesis Testing

- Addressed question vs required method

Each of the 4 questions listed above requires the use of dedicated statistical and computational methods.
Understanding the relation between addressed questions and required methods is fundamental, and will save you a lot of time in the future (trust me)!

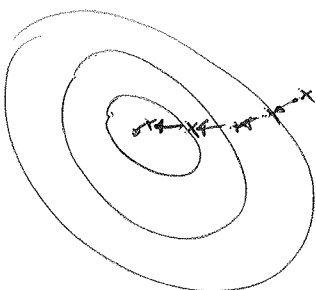~~Point estimation → Maximize like Find parameter values $\hat{\theta}$ that maximize the likelihood (or minimize the $x^2$)~~

Point estimation → Find the parameter values $\hat{\theta}$ that maximize $\mathcal{L}$ → Minimizer algorithms (gradient follower)

Interval estimation → Study the Tails of $\mathcal{L}$ or of the posterior $P(\vec{\theta}|\vec{n})$ → Study all possible combinations of $\vec{\theta}$ giving $\vec{n}$
  → Toy MC
  → Map the posterior $P(\vec{\theta}|\vec{n})$
    → Markov Chain MC
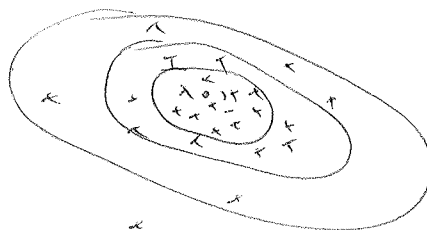
Goodness of Fit → Quantify the probability of a random fluctuation to give a worse fit → ~~tests~~ Analytical method (e.g. $x^2$) plus lookup Tables
  → Toy MC

Hypothesis Testing → Compare the signal strength with the probability of a similarly strong background fluctuation → Toy MC + some method to compare validity of alternative hypotheses.

Minimizer:



Mapper:

So far, we've used a very vague language on purpose. To be more specific, we need to choose either the frequentist or the Bayesian approach, and specify the questions addressed by each of them.

- Frequentist approach

~~Point Estimation~~

Assumptions: The true value of the parameter $\theta$ is fixed but unknown. We cannot associate a PDF to $\theta$, but just to the data $x$.

Point estimation: ~~What~~ Based on the measured data, what's our best "estimate" for the fixed unknown parameter? What's the estimate that is closer to the true value?

Interval estimation: Based on the measured data, what interval contains the true value with a predefined amount of probability (e.g. 68%)?
⌐→ ~~this has to be true also if we repeat~~
 → If we repeat the measurement 100 times, we will have 100 different intervals, ~~and~~ the true value will be contained in them 68 times

Goodness of fit: ~~Is~~ ~~Does~~ Does my model provide a suitable description of the data, or is there any indication that it should be modified somehow?

Hypothesis Testing: Based on the data, which among ~~the~~ two (or more) alternative hypotheses is true?
⌐ ~~What is the probability that~~
 └→ Assuming H0 is true, what is the probability that the data will fake H1 (and viceversa)?

• Bayesian approach

In the Bayesian approach, the probability is interpreted as a "degree of belief"
and can be therefore applied to a wider range of ~~you~~ elements, including:

- random variables
- (true) parameters of a model
- hypotheses

Point estimation: based on the measured data, what is the most probable value for
the parameter $\theta$?

Interval estimation: based on the measured data, what is the interval ~~that that~~
~~contains the of the~~ of the PDF of $\theta$, $f(\theta)$, that
contains a given amount of probability (e.g. 68%)?

Goodness of fit: This question makes no sense in the Bayesian approach, because
we cannot ~~compare~~ compare one hypothesis with N unknown ones.

Hypothesis Testing: based on the data, what is the ratio of the probabilities
of hypotheses $H_0$ and $H_1$?

Item 4