

HYPOTHESIS TESTING

- Goals: \rightarrow use data to verify or disprove a theory or hypothesis.
 \rightarrow choose between alternative hypotheses

Simple hypothesis = hypothesis which is completely specified

E.g.: Theoretical model and ~~model~~ parameter values

Composite hypothesis = ensemble of more than one simple hypothesis

E.g.: model with free parameters (equivalent to infinite list of hypotheses for all possible values of the parameter)

Goals (more specific wording):

\rightarrow Take H_0 as the null hypothesis (background)

H_1 as the alternative hypothesis (signal + background)

H_0 and H_1 are a complete set: $P(H_0) + P(H_1) = 1$ (Bayesian)

Test of hypothesis = use data to verify/disprove H_0 vs H_1

\rightarrow Take H_0 as a given hypothesis

\bar{H}_0 as all other (unspecified) possible hypotheses

Goodness of fit = use data to verify/disprove H_0 vs \bar{H}_0

• Test statistic

Let \vec{n} be some measured data distributed as:

$f_0(\vec{n}|H_0)$ if H_0 is true

$f_1(\vec{n}|H_1)$ if H_1 is true

Let H_0 and H_1 be a complete set of the alternative hypotheses.

We want to develop a method to determine whether the observed data agree better with H_0 or H_1 .

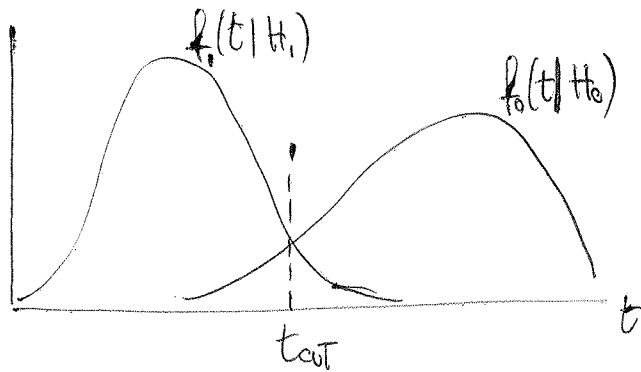
Hyp 1

Instead of using all data \vec{x} , we can use some statistic $t = t(\vec{x})$ and use its PDF to test the hypotheses H_0 and H_1 .

Such statistic is called "Test Statistic".

In general, we will have that $f_0(t|H_0)$ differs from $f_1(t|H_1)$.

Therefore we can set a cut on t , ~~and~~ (before doing the measurement), and "choose" H_0 or H_1 depending on the measured value of t , \hat{t} .



→ Example: Particle identification with scintillating crystal

Suppose we have a scintillating crystal producing light with 2 time constants:

$$\tau_1 = 50 \text{ ns}$$

$$\tau_2 = 200 \text{ ns}$$

Suppose the amplitude ~~of the pulse~~ connected to the second time constant depends on the particle depositing the energy, so that

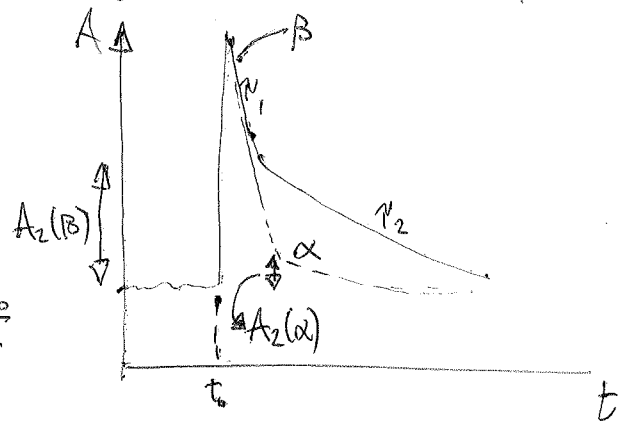
$$A_2(\tau_2, \alpha) \ll A_2(\tau_2, \beta)$$

Instead of using the full pulse shape, we can ~~use~~ fit the pulses with

$$\text{something like: } A(t) = A_1 e^{-\frac{t-t_0}{\tau_1}} + A_2 e^{-\frac{t-t_0}{\tau_2}}$$

$$t = A_2/A_1$$

And use ~~this~~ as a Test Statistic.

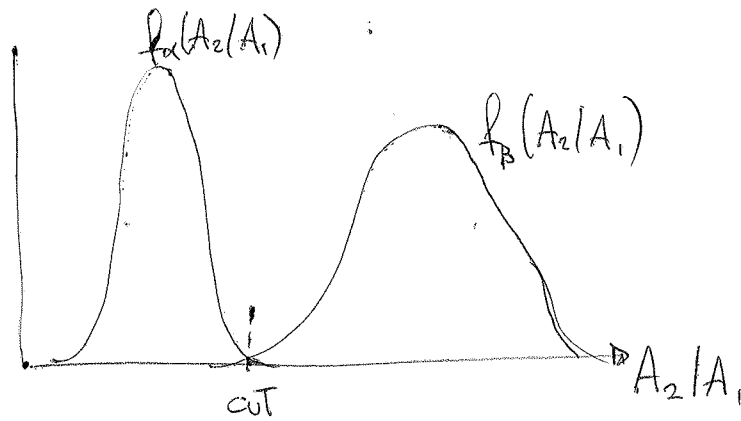


We will proceed as follows:

- 1) Calibrate with an α source to evaluate ~~$f_{\alpha}(\tau_1)$~~ $f_{\alpha}(A_2/A_1)$
- 2) Calibrate with a β source to evaluate ~~$f_{\beta}(\tau_2)$~~ $f_{\beta}(A_2/A_1)$

Fig. 2

3) Decide some cut on A_2/A_1



4) Measure the "physics data" (whatever they are) and use the previous method to distinguish α from β

• Selection, ~~and~~ misidentification and significance

→ Selection efficiency = ~~expected~~ fraction of signal events that are expected to be correctly identified
 $\epsilon_s = 1 - \beta$

→ Misidentification probability = fraction of background events that are expected to be erroneously identified as signal
 $\epsilon_b = \alpha = \text{significance}$

→ Critical region = region where we expect the signal
 w

→ Acceptance region = region where we expect the background
 $W - w$
 $\hat{=}$ region where we accept H_0 or True

In general, the misidentification probability is also called "significance level". When we design a hypothesis test, we need to specify the desired level of significance α , i.e. ~~the amount of fraction of background probability~~ i.e. to which extent we are willing to accept the misidentification of data induced by H_0 with data induced by H_1 :

~~1/11/11~~ $P(t(\frac{\bar{x}}{\sqrt{n}}) \in w \mid H_0) = \alpha$

Hyp 3

Similarly, we can define The "power" of a Test, as The probability of data produced by H_1 to be correctly identified:

$$P(t(\vec{x}) \in w | H_1) = 1 - \beta = \mathcal{E}_s$$

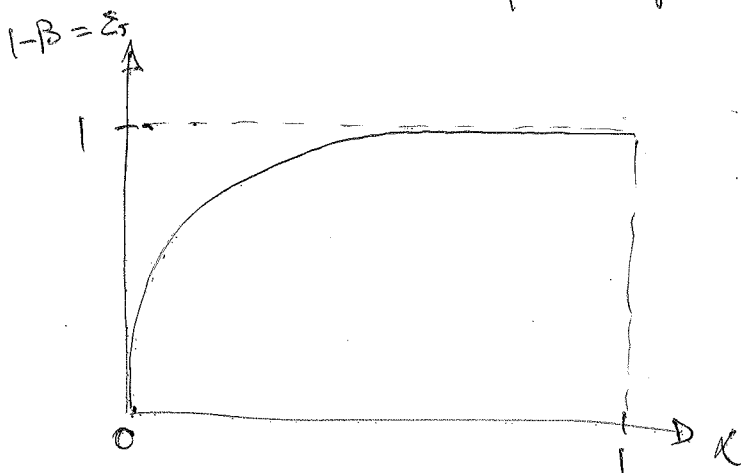
Conversely, β is The probability to misidentify data produced by H_1 as if H_0 were true.

We can define Two Type of errors:

a) Error of The first kind: rejecting H_0 when it is true \Rightarrow prob. = α

b) Error of The second kind: accepting H_0 when it is false \Rightarrow prob. = β

We can plot The so-called "receiver operating characteristic" (ROC) curve:



Clearly, we would like to find The "most powerful" Test of hypothesis!

* Neyman - Pearson Lemma

Finding The most powerful Test is equivalent to finding The best critical region in \mathcal{X} -space.

Suppose we measure The variable \vec{x} with The usual PDFs $f_0(\vec{x} | H_0)$ and $f_1(\vec{x} | H_1)$.

Using The measurement \vec{x} itself as a statistic, we have:

$$\int_w \mathcal{L}_0(\vec{x} | H_0) d\vec{x} = \alpha$$

$$\int_w \mathcal{L}_1(\vec{x} | H_1) d\vec{x} = 1 - \beta$$

where \mathcal{L}_i is f_i ~~mean~~ evaluated at The measured data \vec{x}

Thy. 4

Given a predefined value of α , we want to find the region ~~we~~ w which maximizes $(1-\beta)$.

We can rewrite:

$$1-\beta = \int_w \frac{f_1(\vec{n}|H_1)}{f_0(\vec{n}|H_0)} f_0(\vec{n}|H_0) d\vec{n}$$

$$= E_w \left[\frac{f_1(\vec{n}|H_1)}{f_0(\vec{n}|H_0)} \right]$$

The best critical region w is the one that satisfies:

$$\lambda(\vec{n}) = \frac{f_1(\vec{n}|H_1)}{f_0(\vec{n}|H_0)} \geq k_\alpha$$

with k_α chosen so that the ~~region~~ desired significance is achieved.

This is the Neyman-Pearson lemma.

Notice that: \rightarrow The NP lemma is valid only if the PDFs are known (including the values of their parameters).

~~Otherwise, this~~

\rightarrow The NP lemma provides ~~an upper~~ the most-powerful test, ~~but~~ if we don't know the parameter values, ~~then~~ the power of any test will be \leq than that of NP.

Practical instructions (assuming parameter values are known):

- 1) Evaluate $f_0(\vec{n}|H_0)$ and $f_1(\vec{n}|H_1)$
- 2) Evaluate $\lambda(\vec{n})$ and find region w
- 3) Do your measurement, obtaining data \vec{n} .
- 4) If $\lambda(\vec{n}) > k_\alpha \Rightarrow H_1$ is considered True
- If $\lambda(\vec{n}) \leq k_\alpha \Rightarrow H_0$ is considered True

• Projective likelihood ratio test: Wilks' Theorem

Suppose we have two "nested" hypotheses H_0 and H_1 , so that the parameter values of H_0 are a special case of H_1 (e.g. signal = 0).

~~We can define the~~

We can divide the total parameter space Ω in ~~two parts~~ Ω_0 and Ω_1 .

$$\begin{array}{ll} H_0: \vec{\theta} \in \Omega_0 \subset \Omega & H_0: \vec{\theta} \in \Omega_0 \subset \Omega \\ H_1: \vec{\theta} \in \Omega_1 = \Omega - \Omega_0 & H_1: \vec{\theta} \in \Omega_1 = \Omega - \Omega_0 \end{array}$$

We can define the maximum \mathcal{L} ratio:

$$\lambda = \frac{\max_{\vec{\theta} \in \Omega_0} \mathcal{L}_0(\vec{\kappa} | H_0)}{\max_{\vec{\theta} \in \Omega_1} \mathcal{L}_1(\vec{\kappa} | H_1)} \Rightarrow -2 \ln \lambda = \chi^2$$

Wilks Theorem: If H_0 is true and for $n \rightarrow \infty$,

$-2 \ln \lambda$ has a χ^2 distribution with $\text{NDOF} = \dim(\Omega_1) - \dim(\Omega_0)$

What does it mean?

→ Example: Deviation of parameter μ from predicted value

$H_0: \mu = \mu_0 = \text{value predicted by Theory}$

$H_1: \mu \neq \mu_0$

$$\chi^2(\mu_0) = -2 \ln \frac{\max_{\vec{\theta}} \prod_{i=1}^n \mathcal{L}(\vec{\kappa}_i | \mu_0, \vec{\theta})}{\max_{\mu, \vec{\theta}} \prod_{i=1}^n \mathcal{L}(\vec{\kappa}_i | \mu, \vec{\theta})}$$

$\left. \begin{array}{l} \text{Y maximized over } \vec{\theta} \\ \text{for a fixed value of } \mu = \mu_0 \end{array} \right\}$

$\left. \begin{array}{l} \text{max } \mathcal{L} \text{ over entire} \\ \text{parameter space} \end{array} \right\}$

• Discoveries and upper limits

- Suppose we are searching for a new physics process. We make a measurement and we need to quote ~~the~~ a result. How do we decide whether the data tell us that there is new physics?

→ Frequentist approach: measure the "significance", i.e. the probability that a background statistical fluctuation produces a fake signal at least as intense as the measured one.

→ Bayesian approach: quantify the ~~degree~~ posterior degree of belief on the hypotheses H_0 and H_1 .

• P-value

To claim a discovery, we need to determine that the data are sufficiently inconsistent with the H_0 -only hypothesis H_0 .

⇒ We can use a test statistic t to measure such inconsistency!

p-value = probability p that the test statistic t ~~measures~~ assumes a value greater or equal to the measured value \hat{t} due to an overfluctuation of the background.

↳ The p-value has a uniform distribution in $[0, 1]$ if H_0 is true

↳ The p-value tends to have small values if H_1 is true

→ Example: Event counting experiment

Take the number of observed events n as a test statistic.

p-value = probability to measure $\geq n$ events under the H_0 hypothesis.

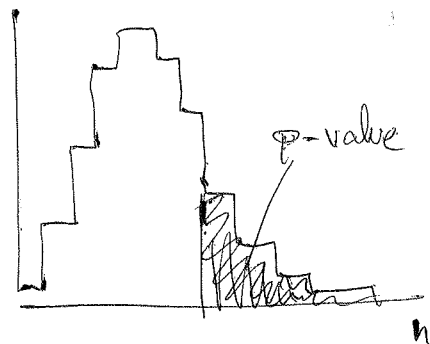
→ Example: p-value for Poisson Counting

H_0 : Poisson counting with $\lambda_b = 4.5$, $\lambda_s = 0$

H_1 : ~~known~~ $\lambda_b = 4.5$, $\lambda_s > 0$

Measurement $\rightarrow n = 8$

$P(n|\lambda=4.5)$



$$p\text{-value} = P(n \geq 8 | H_0) = \sum_{n=8}^{\infty} \text{Poisson}(n, 4.5) = \dots = 0.087$$

• Significance level

Instead of quoting a p-value, we normally quote the number of STDs that correspond to an area equal to the p-value under the right tail of a standard normal distribution:

$$\phi = \int_z^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz = 1 - G(z) = G(-z) = \frac{1}{2} \left[1 - \text{erf}\left(\frac{z}{\sqrt{2}}\right) \right]$$

z is called "significance level".

By convention, we claim: \rightarrow evidence if $z \geq 3$

\rightarrow observation/discovery if $z \geq 5$

• Significance for Poissonian counting

H_0 : expected b events (b known)

H_1 : expected $b+s$ events ($s \geq 0$ unknown)

The likelihood is: ~~$L(n|s,b)$~~ $L(n|s,b) = \frac{(s+b)^n}{n!} e^{-(s+b)}$

To compute the significance, we need to compare the measured number of events n with the expected background b under the H_0 hypothesis ($s=0$).

→ If b is large, we can approximate the \mathcal{L} with a Gaussian with $\mu=b$ and $\sigma=\sqrt{b}$.

An excess $n-b=s$ must be compared with \sqrt{b} .

The significance will be: $z = \frac{n-b}{\sqrt{b}} = \frac{s}{\sqrt{b}}$

→ If b ~~then~~ is large and has some large uncertainty σ_b ,

The significance will be: $z = \frac{n-b}{\sqrt{b+\sigma_b^2}}$

→ If b is small, one can prove that the significance is:

$$z = \sqrt{2 \left[(s+b) \ln \left(1 + \frac{s}{b} \right) - s \right]}$$

• Significance with likelihood ratio

Take again two nested hypotheses H_0 and H_1 , with $H_0 = H_1(\beta=0)$
 ↳ signal strength

We can define the Test Statistic:

$$\lambda(s, \vec{\theta}) = \frac{\mathcal{L}_{s+b}(\vec{n} | s, \vec{\theta})}{\mathcal{L}_b(\vec{n} | \vec{\theta})}$$

→ Notice that we inverted numerator and denominator w.r.t. Wilks Theorem.

A minimum of $-2 \ln \lambda$ at $\beta = \hat{\beta}$ indicates the possible presence of a signal with strength $\hat{\beta}$.

According to Wilks Theorem, $2 \ln \lambda$ follows a χ^2 distrib with 1 DOF.

An approximate estimate of the significance is: $z = \sqrt{2 \ln \lambda(\hat{\beta})}$

→ This is a local significance that can be used if we have a "perfect" prior knowledge of the other parameters $\vec{\theta}$.

→ If we estimate $\vec{\theta}$ from the data, we need to consider the "look elsewhere effect".

• Significance with Toy-PC

A more general approach can be ~~not~~ obtained with Toy-PC.

- 1) Generate many Toy-PC datasets with no signal ($s=0$), obtaining an approximate distribution of $-2 \ln \lambda$
- 2) The p-value is the probability that λ is \leq the observed value $\hat{\lambda}$:

$$p = P(\lambda(\theta) \leq \hat{\lambda})$$

This is equal to the fraction of Toy-PC for which $\lambda(\theta) \leq \hat{\lambda}$.

• Bayesian method for hypothesis testing

~~In case~~

Suppose we have two hypotheses H_0 and H_1 representing a complete set:

$$P(H_0 | \vec{n}) + P(H_1 | \vec{n}) = 1$$

We can apply the Bayes Theorem to the two hypotheses:

$$P(H_i | \vec{n}) = \frac{P(\vec{n} | H_i) \pi(H_i)}{P(\vec{n})}$$

$$\text{Where: } P(\vec{n}) = P(\vec{n} | H_0) \pi(H_0) + P(\vec{n} | H_1) \pi(H_1)$$

$$\text{and } P(\vec{n} | H_0) = \int_{\Omega_\theta} \mathcal{L}_0(\vec{n} | \vec{\theta}) \pi(\vec{\theta}) d\vec{\theta} = \text{"evidence" of } H_0$$

$$P(\vec{n} | H_1) = \int_{\Omega_{\theta,s}} \mathcal{L}_1(\vec{n} | s, \vec{\theta}) \pi(\vec{\theta}) \pi(s) ds d\vec{\theta} \quad s = \text{signal parameter}$$

Again, we need to assign some prior both to the parameters $s, \vec{\theta}$, and to the models H_0 and H_1 .

A general choice is: $\pi(H_0) = \pi(H_1) = 0.5$

At this point we set a threshold on $P(H_0)$ in order to claim the "evidence" or discovery.

• Bayes Factor / Ratio

If H_0 and H_1 are not a complete set of hypotheses, we can't compute $P(H_1 | \vec{n})$ ~~or~~ $P(\vec{n})$, and therefore $P(H_1 | \vec{n})$.

However, we can compute the ratio:

$$\frac{P(H_1 | \vec{n})}{P(H_0 | \vec{n})} = \frac{\cancel{P(\vec{n} | H_1)} \pi(H_1)}{\cancel{P(\vec{n} | H_0)} \pi(H_0)}$$

↓ Posterior odds
 ↓ Bayes factor
 ↓ prior odds

If $\pi(H_0) = \pi(H_1)$, The Posterior odds are identical to the Bayes factor.

One can then set some thresholds on the Bayes factor (or on the posterior odds) to claim "evidence" and "discovery".

→ Example: Evidence (Bayes factor)

• Numerical ~~costs~~ and practical considerations

When running a Bayesian analysis, we might need to face three different problems involving 3 different algorithms:

- 1) Finding the global mode of posterior → Minimizer algorithm
- 2) ~~Finding~~ Interval estimation → MCMC
- 3) Computing "significance" (doing model testing) (or Bayes factor) → n-dimensional integration of full posterior PDF

At the moment, there is no algorithm that does all 3 of them at the same time. Moreover, MCMC and integrators are inefficient.

⇒ If you have an idea for an algorithm that can do all 3 things with a high efficiency (no discarded points) and that can work for $\dim > 50$, please let me know, because I want to work with you!

Hyp. II

• ~~Goodness of fit~~

• Sensitivity

Before doing a measurement, or when planning an experiment, we can ask ourselves what is the expected probability to make a discovery, or what's the expected 90% limit we will place.

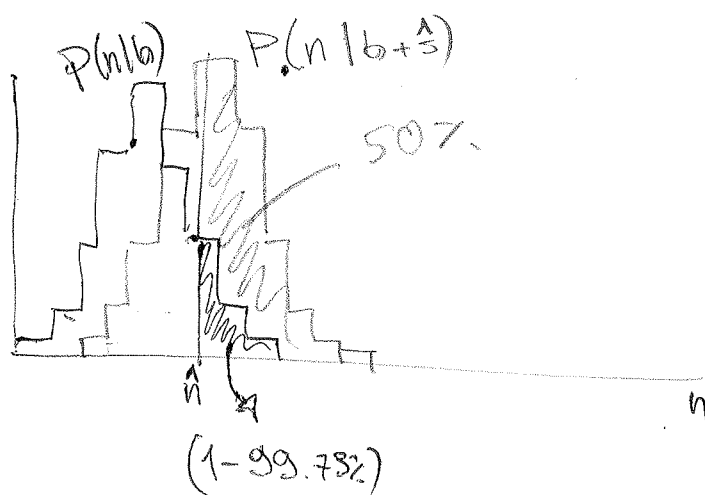
→ Discovery ~~probability~~ sensitivity = expected number of signal events (or signal strength) for which an experiment has a 50% chance to observe an excess over the background at 3 σ or 50% CL significance.

→ Exclusion sensitivity = expected number of signal events (or signal strength) that an experiment has 50% chance to exclude at 90% CL

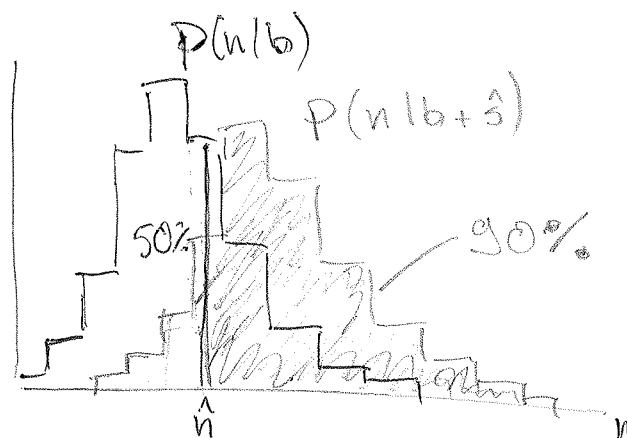
→ For a counting experiment with known expected background b , the two sensitivities can be computed by solving these equations:

$$\text{Discovery} = \begin{cases} P(n \leq \hat{n} | b) \geq 99.73\% \quad (\text{for } 3\sigma) \\ P(n \geq \hat{n} | b + \hat{s}) \geq 50\% \end{cases}$$

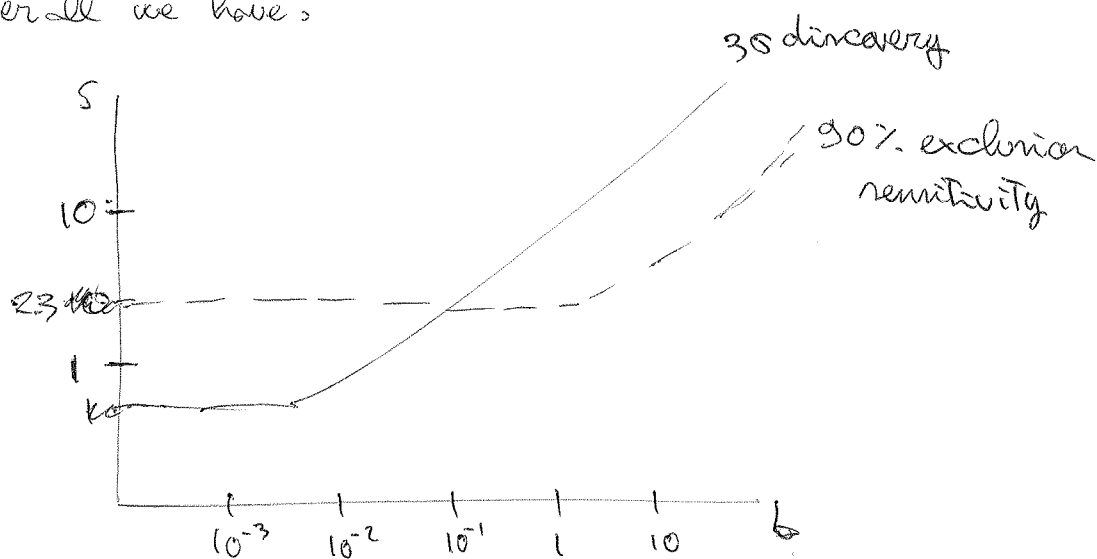
- ↳ Find $\min(\hat{n})$ that satisfies eq. 1
- ↳ Find $\hat{s} = \min(s)$ that satisfies eq. 2
- ↳ \hat{s} is the discovery sensitivity



Exclusion =
$$\begin{cases} P(n \leq \hat{n} | b) \geq 50\% \\ P(n \geq \hat{n} | b+s) \geq 90\% \end{cases}$$



Overall we have:



GOODNESS OF FIT

• Address The question: How well do The data agree with The functional form predicted by a hypothesis H_0 ?

→ Here we're testing only one null hypothesis H_0 , or one hypothesis H_0 against The infinite and unspecified set of alternatives T_0 to H_0 .

↳ The Theoretical basis is less robust than in hypothesis testing.

Nevertheless we can successfully perform The test, using a frequentist approach.

• Procedure (general for all methods):

1) Choose/construct a Test statistic $t(\vec{n})$ which is sensitive to The level of agreement between The data and The hypothesis H_0 .

Let's assume That we choose t so That larger values of t indicate a worse agreement.

~~2) Compute The probability p That, assuming H_0 to be true, the~~

2) Assume H_0 true and compute The probability p That, repeating The measurement many times, we get a value of t greater or equal than The actually measured one. This is The p-value.

↳ Small p = bad agreement between data and H_0 = bad FIT

↳ p-value calculation could depend or not on The PDF of t .

• EX.: Test statistic for Poisson distrib

Assume we measure a discrete variable n with Poisson PDF with $\lambda = 17.3$, and obtain $n = 12$.

What is The p-value, or The level of compatibility with The hypothesis of n belonging to a Poisson with $\lambda = 17.3$?

~~for $t(\vec{n}) = (n - \lambda)$ take $t = |n - \lambda|$~~

$$P = \sum_{n \geq 5.3} \frac{e^{-\lambda} \lambda^n}{n!} = \sum_{n=1}^{12} \frac{e^{-\lambda} \lambda^n}{n!} + \sum_{n=23}^{\infty} \frac{e^{-\lambda} \lambda^n}{n!} = 0.229$$

$t \geq 5.3$

[Hyp. 14]

• Distribution-free Test

A goodness of fit Test is distribution-free if the distribution of t is known independently of H_0 .

↳ Also the p -value is independent of H_0

↳ We can compute p for any H_0 , and compare it to tabulated data that were calculated once for all!

↳ Eventually p might depend on the number of events, number of bins in a histogram, or number of constraints in a fit.

• Distribution-free Tests for histograms

Suppose we measure n times a ~~scalar~~ variable x with PDF $f(x)$.

Then we have n values of the Test statistic t , with PDF $f(t)$.

Assume n is a Poisson variable.

If we bin the n values of t , we get a histogram where each bin follows a Poisson statistic.

↳ We have lost the dependence on $f(t)$

• Pearson's χ^2 Test for histograms

Let's assume the number of entries in each bin n_i is large enough that we can approximate the corresponding Poisson to a Gaussian.

Then we can use as a statistic:

$$\chi^2 = \sum_{i=1}^m \frac{(n_i - \lambda_i)^2}{V[\lambda_i]} \quad \text{where } m = \# \text{ of bins}$$

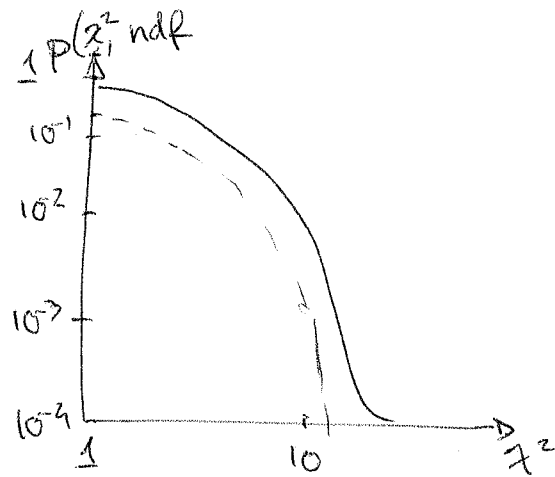
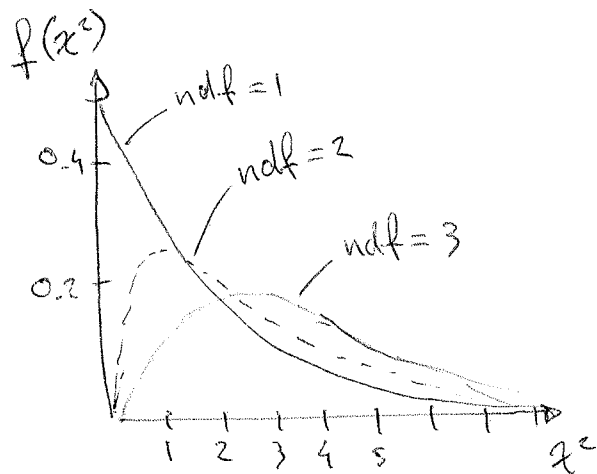
$\lambda_i = \text{expectation value for bin } i$
(depends on $f(x^2)$)

$V[\lambda_i] = \text{variance for } \lambda_i$

The PDF of χ^2 is:

$$f(\chi^2, ndf) = \frac{1}{2^{\frac{ndf}{2}} \Gamma(\frac{ndf}{2})} (\chi^2)^{\frac{ndf}{2}-1} e^{-\frac{\chi^2}{2}} \rightarrow \begin{matrix} \text{mean} = ndf \\ V[\chi^2] = 2 \text{ ndf} \end{matrix}$$

The p-value is: $p = P(\chi^2, n) = \int_{\chi^2}^{+\infty} f(z, n) dz$



If we fit $f(n|\vec{\theta})$ on the data, and $\dim(\vec{\theta}) = d$, the minimize value of χ^2 will follow a $\chi^2(n-d)$ distribution.

~~If the data~~

Possible situations:

- χ^2 Too small (Too good fit) → errors have been overestimated, or the data have been selected
↳ usually if $p < 0.05$
- χ^2 Too large → Hypothesis H_0 is wrong, or there are unaccounted correlations in the data

Reduced $\chi^2 = \frac{\chi^2}{ndf}$ → gives also some info on agreement between data and H_0
↳ not as informative as p-value

→ What if n is too small to use the χ^2 probability?
↳ Use a MC approach!

Wald - Wolfowitz run-Test

Notice: The Pearson's χ^2 Test does not take into account the sign of the deviations.

The following two cases would give exactly the same χ^2 :

+	+	+	+	+
<hr/>				
+	+	+	+	+

+	+	+	+	+
<hr/>				
+	+	+	+	+

Let's define as "run" each region ~~with~~ of measurements with residuals of the same sign.

→ The number of runs r is binomial

Denoting with n_+ = number of measurements with positive residuals
 n_- = " " " " " negative "

Number of possible combinations: $\frac{n!}{n_+! n_-!}$

Expected number of runs: $E[r] = 1 + \frac{2 n_+ n_-}{n}$

Variance: $V[r] = \frac{2 n_+ n_- (2 n_+ n_- - n)}{n^2 (n - 1)}$

With $n \geq 20$, r can be approximated by a Gaussian, so we have:

$$p = \frac{r - E[r]}{\sqrt{V[r]}}$$

• Combining Tests

Suppose that we run the χ^2 and run-test on some data, obtaining the p-values \hat{p}_1 and \hat{p}_2 . How do we combine them?

Let's assume that: 1) p_1 and p_2 are uniformly distributed in $[0, 1]$

2) ~~\hat{p}_1 and \hat{p}_2 are independent~~

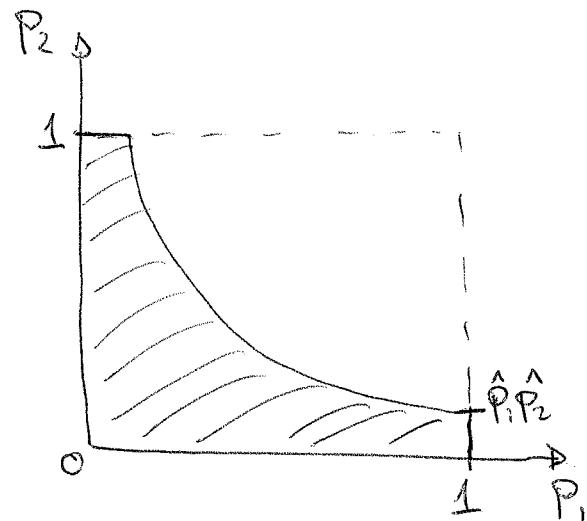
2) t_1 and t_2 are independent: $P(t_1, t_2 | H_0) = P(t_1 | H_0) P(t_2 | H_0)$

This might be hard to prove, but holds for the χ^2 and run-test, because the first does not use the information of the sign, and the second uses only the sign.

The probability that $p_1 p_2 \leq \hat{p}_1 \hat{p}_2$ is:

$$P(p_1 p_2 \leq \hat{p}_1 \hat{p}_2) = \int_0^1 dp_1 \int_0^{\hat{p}_1 \hat{p}_2 / p_1} dp_2$$

$$= \hat{p}_1 \hat{p}_2 [1 - \ln(\hat{p}_1 \hat{p}_2)]$$



$$\Rightarrow P(p_1 p_2 \leq \hat{p}_1 \hat{p}_2) > \hat{p}_1 \hat{p}_2$$

→ With n Tests, we can ~~compute~~ prove that:

$$p = -2 \ln \prod_{i=1}^n \hat{p}_i$$

is distributed as a χ^2 with $2n$ def.

- χ^2 Test for unbinned data

Suppose we measure n ~~times~~ n times and fit it with $f(n|\vec{\theta})$

We can still run a χ^2 Test by binning the data n in m bins:

$$\chi^2 = 2 \sum_{\substack{i=1 \\ n_i \neq 0}}^m n_i \ln \frac{n_i}{\hat{\lambda}_i}$$

Multinomial core

where n_i = number of events in bin i
 $\hat{\lambda}_i$ = expectation value for n_i obtained from fit

For Poisson distributed data:

$$\chi^2 = 2 \sum_{\substack{i=1 \\ n_i \neq 0}}^m n_i \ln \frac{n_i}{\hat{\lambda}_i} + \hat{\lambda}_i - n_i \rightarrow \text{In large-sample limit, it follows a } \chi^2 \text{ with } \text{dof} (m-d)$$

- Test using max- \mathcal{L} estimate

Suppose we use \mathcal{L}_{\max} as a Test Statistic, and compare the measured value $\hat{\mathcal{L}}_{\max}$ to the set of \mathcal{L}_{\max} from Toy-TC experiments, where we set the value of $\vec{\theta}$ to their expected true values.

We have that: ~~1)~~ \mathcal{L}_{\max} distributions are not well separated under different hypotheses

↳ Do not use \mathcal{L}_{\max} as a Test-Statistic for GOF.

• Kolmogorov - Smirnov Test

Take a set of n measurements $x_i = x_1, \dots, x_n$ ordered in increasing values of x .

The discrete cumulative distribution ~~is~~ is:

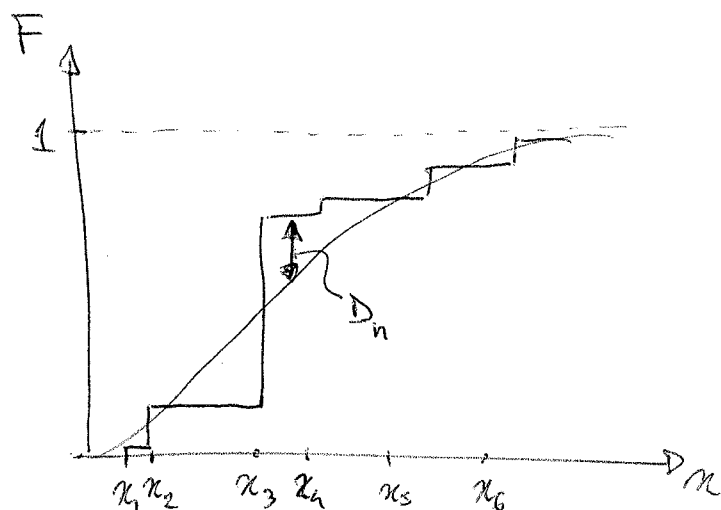
$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \vartheta(x - x_i) \quad \text{with } \vartheta(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

This can be compared with The cumulative PDF of $f(x)$:

$$F(x) = \int_{-\infty}^x f(y) dy$$

We can take as a Test - statistic:

$$D_n = \max |F_n(x) - F(x)|$$



For large n , D_n converges to 0 in probability.

One can prove that: 1) The distribution of $K = \sqrt{n} D_n$ does not depend on $f(x)$

2) The probability that $K \leq k$ with Kolmogorov is

The Kolmogorov distribution:

$$P(K \leq k) = \frac{\sqrt{2\pi}}{k} \sum_{i=1}^{\infty} e^{-(2i-1)^2 \pi^2 / 8k^2}$$

If the parameters $\vec{\theta}$ of $f(x|\vec{\theta})$ are obtained from a fit, one cannot use the Kolmogorov distribution but has to empirically obtain the distribution of K from log-PLC.

→ The KS Test can be used also to compare two measurements, and to test if they are produced from the same PDF!

• Smirnov - Gromer - Von Mises Test

Use as test statistics: $W^2 = \int_{-\infty}^{+\infty} [F_n(n) - F(n)]^2 dF(n) \quad f(n) dn$

↳ Instead of using the single point where the difference is largest, we use the integral of the squared difference

• Anderson - Darling Test:

$$A^2 = n \int_{-\infty}^{+\infty} \frac{(F_n(n) - F(n))^2}{F(n)(1 - F(n))} dF(n) \quad \rightarrow \text{put more weight on the tails of the distribution}$$

