

POINT ESTIMATION

~~Part of this applies to both the frequentist and Bayesian approach~~

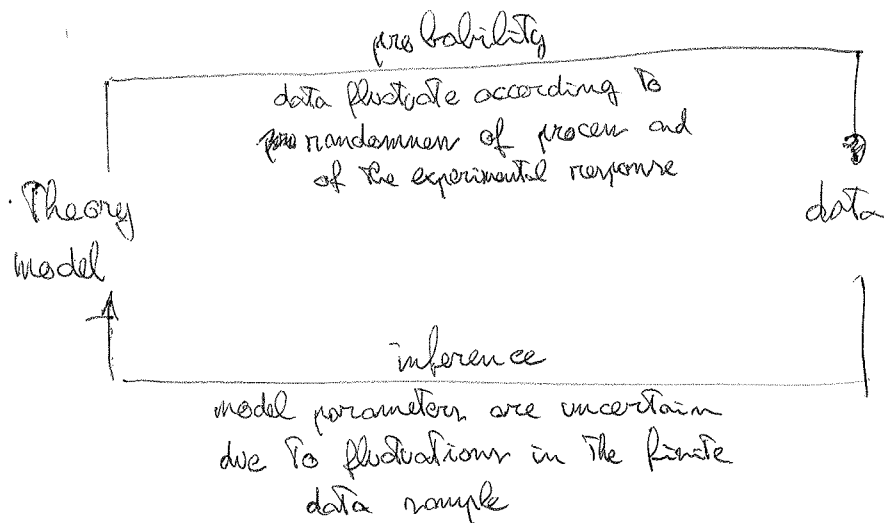
Part of this is frequentist exclusively, even if some methods are easily applicable to a Bayesian approach.

• ~~Estimates and estimators~~

~~An estimate of a parameter is a statistic, i.e. it is a function of the data.~~

• Inference

The inference is the process of determining an estimated value, $\hat{\theta}$ and the corresponding uncertainty of some parameter θ from experimental data.



• Estimator and estimate

The estimate of an unknown parameter is a mathematical procedure to determine the central value of the parameter as a function of the observed data sample.

The function of the data sample that returns the estimate is called "estimator".

Example: If I measure the mass of an object once, the measured value is an estimate of the object mass: ~~mass~~

$$\hat{m}(x) = x$$

Properties of estimator:

- Consistency
- Unbiasedness
- Information content or efficiency
- Robustness

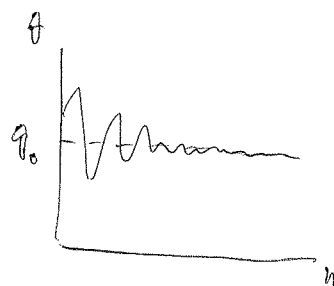
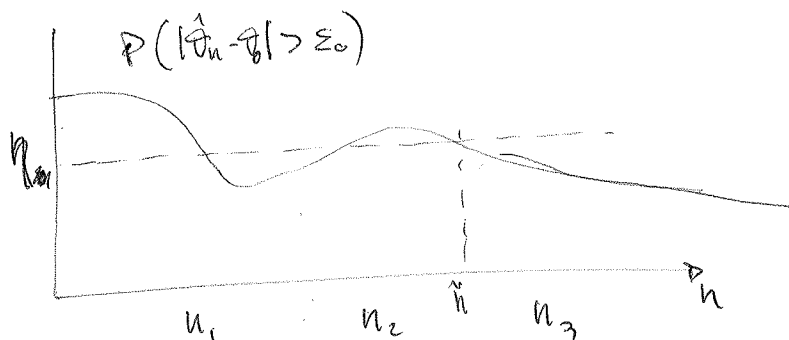
• Consistency

An estimator is consistent if it converges in probability to the true value of the unknown parameter θ_0 .

$$\forall \varepsilon > 0 : \lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta_0| < \varepsilon) = 1$$

In other words, if $\forall \varepsilon > 0$ and $\eta > 0$, There is an \tilde{n} such that

$$P(|\hat{\theta}_n - \theta_0| > \varepsilon) < \eta \quad \forall n > \tilde{n}$$



Notice: \rightarrow The law of large numbers states that the sample mean is a consistent estimator of the parent mean.

\rightarrow Consistency is an asymptotic property. It does not imply that precision is a monotonic function of n .

• Bias

The bias of an estimator is the expected value of the deviation of the parameter estimate from the corresponding true value θ_0 .

Assuming n observations:

$$b_n(\hat{\theta}) = E(\hat{\theta}) - \theta_0 = E(\hat{\theta} - \theta_0)$$

An estimator is unbiased if $\forall n$ and $\forall \theta_0$:

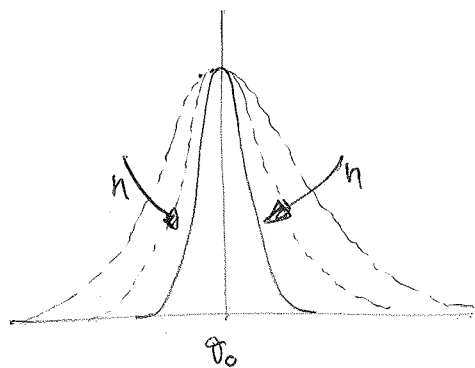
$$b_n(\hat{\theta}) = 0$$

$$E(\hat{\theta}) = \theta_0$$

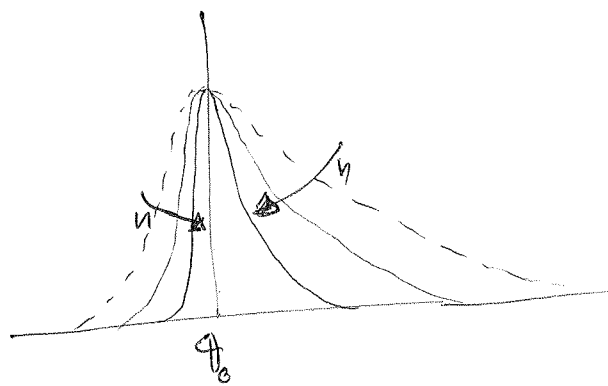
• Bias or consistency

Unbiased

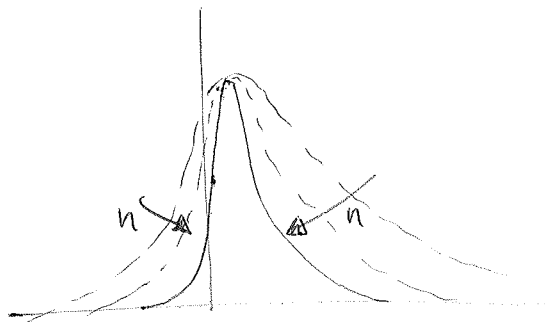
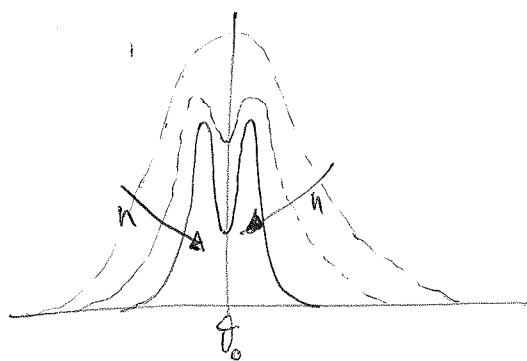
Consistent



Biased



Inconsistent



• Efficiency = minimum variance and Cramer-Rao inequality

Neglecting the bias, the smaller the variance of the estimator, the more certain we are that the estimate is near the true value of the parameter.

One can prove that the variance $V(\hat{\theta})$ of any consistent estimation is subject to a lower bound given by:

$$V(\hat{\theta}) \geq \frac{\left(1 + \frac{\partial b(\theta)}{\partial \theta}\right)^2}{E\left[\left(\frac{\partial \ln L(\hat{\theta}|\theta)}{\partial \theta}\right)^2\right]}$$

→ bias of estimator

→ Fisher information

↓

$V_{CR}(\hat{\theta})$

We define the efficiency of the estimator as: $\varepsilon(\hat{\theta}) = \frac{V_{CR}(\hat{\theta})}{V(\hat{\theta})}$

Any consistent estimator $\hat{\theta}$ has an efficiency which is at most equal to 1.

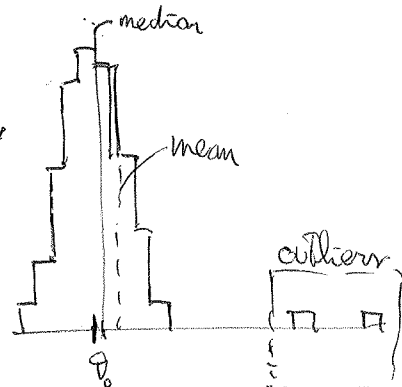
• Robustness

A good estimator should be not ~~too~~ too sensitive if the real distribution has some small deviation from the assumed PDF model.

Entries in the data that introduce visible deviations from the assumed PDF are called outliers.

Example: Suppose we take n measurements of some quantity x , distributed as in the figure.

The median ~~is a~~ estimator is more robust than the mean estimator.



Trick: Remove outliers and compute the average on the "clean" sample.

This is called "Trimmed average".

It works, as long as we are confident that we are removing only outliers, and that the outlier ~~population~~ ^{sample} does not overlap too much with the good sample.

• Maximum likelihood estimator

The maximum likelihood estimate of the parameter θ is that value $\hat{\theta}$ for which $L(\vec{x}|\theta)$ is maximal, given the observed data \vec{x} .

Finding the maximum likelihood is ~~also~~ also called finding the "best fit", because it ~~determines~~ determines the value of the parameter that best fits the data.

Suppose we have n measurements ~~of~~ consisting of m values of the data, ~~the sample is~~ so that the total sample is $\vec{x} = \{(x_1^1, \dots, x_m^1), \dots, (x_1^n, \dots, x_m^n)\}$

The likelihood is: $L(\vec{x}|\vec{\theta}) = \prod_{i=1}^n f(x_1^i, \dots, x_m^i|\vec{\theta})$

Since the logarithm is monotonic, finding the $\max(L)$ is equivalent to finding $\max(\ln L)$. The advantage is that we can turn the product into a sum:

$$\ln L = \sum_{i=1}^n \ln f(x_1^i, \dots, x_m^i|\vec{\theta})$$

So we need an algorithm that maximizes $\ln L$, or minimizes $-\ln L$.

• Properties of maximum likelihood:

We need to distinguish between \rightarrow asymptotic properties \rightarrow hold for sufficiently large n
 \rightarrow finite sample properties \rightarrow hold for any n

\rightarrow max \mathcal{L} estimators are consistent: asymptotically, one of the maxima will go arbitrarily close to the true value

\rightarrow max \mathcal{L} estimators might be biased, but $\lim_{n \rightarrow \infty} b(\mathcal{L}) = 0$

\rightarrow max \mathcal{L} estimators are asymptotically normally distributed with minimum variance, and their variance is given by the Cramér-Rao lower bound.

(Therefore, the efficiency of max \mathcal{L} estimators tends asymptotically to 1.)

\rightarrow max \mathcal{L} estimators have asymptotically the lowest variance of any consistent estimator.

\rightarrow max \mathcal{L} estimators are invariant under reparameterization.

If we reparameterize \mathcal{L} , ~~the estimator~~ with $\eta(\theta)$, we find $\hat{\eta} = \eta(\hat{\theta})$

\rightarrow For finite n , \mathcal{L} might have multiple maxima, but we wouldn't know which is the closest to the true one.

\rightarrow However, in my experience this is a rare case that happens only when we have highly (anti)correlated parameters.

• Extended likelihood

Suppose we ~~measure~~ perform n measurements of a random variable x with PDF $f(x|\theta)$.

The number of observations n is itself a random variable, with distribution $P(n|\theta)$.

The likelihood needs to be extended to include $P(n|\theta)$:

$$\mathcal{L} = P(n|\theta) \prod_{i=1}^n f(x_i|\theta)$$

\rightarrow In most cases, $P(n|\theta)$ is a Poisson distribution whose average λ depends on the parameter(s) θ :

$$\mathcal{L} = \frac{e^{-\lambda} \lambda^n}{n!} \prod_{i=1}^n f(x_i|\theta) \quad \text{with } \lambda = \lambda(\theta)$$

→ A common case is when the events can be induced by a signal or background,
 where: signal = new physics that we're interested in
 background = known physics that we wish wouldn't be there to disturb our measurement

The signal and background will have their PDFs f_s and f_b .

~~Then~~ The overall PDF will be,

$$f(n|\vec{\theta}) = \frac{s}{s+b} f_s(n|\vec{\theta}) + \frac{b}{s+b} f_b(n|\vec{\theta})$$

where s and b are the expected number of ~~count~~ signal and background counts.
 Therefore, the expectation value for n will be $\lambda = s+b$.

In the end we have:

$$\begin{aligned} \mathcal{L} &= \frac{(s+b)^n e^{-(s+b)}}{n!} \prod_{i=1}^n \frac{s f_s(n|\vec{\theta}) + b f_b(n|\vec{\theta})}{s+b} \\ &= \frac{e^{-(s+b)}}{n!} \prod_{i=1}^n \left[s f_s(n|\vec{\theta}) + b f_b(n|\vec{\theta}) \right] \end{aligned}$$

Taking the logarithm:

$$\ln \mathcal{L} = -s-b - \ln(n!) + \sum_{i=1}^n \ln \left[s f_s(n|\vec{\theta}) + b f_b(n|\vec{\theta}) \right]$$

↳ Notice that s and b are parameters of the model, in addition to the parameters of f_s and f_b (if any).

↳ Notice that the term $\ln(n!)$ does not depend on any parameter, so it can be omitted from the maximization.

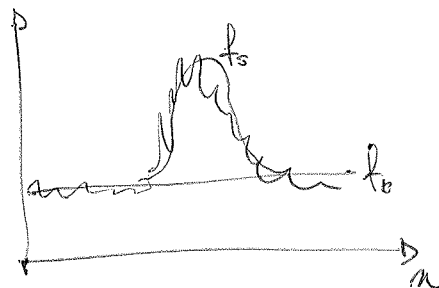
↳ Further simplifications might be possible, depending on the actual form of f_s and f_b .

Example: Extended likelihood fit of Gaussian peak over flat background

$$f_s = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(n-\mu)^2}{2\sigma^2}\right)$$

$$f_b = \frac{1}{\Delta n} \quad \Delta n = \text{fit range (fixed)}$$

Parameters: $\vec{\theta} = \{s, b, \mu, \sigma\}$



• Binned likelihood

Suppose the number of measurement n is so large that computing $\sum_{i=1}^n \ln f(x_i; \theta)$ would take too long.

We can simplify the problem (from the computational point of view) by binning the data n in $m \ll n$ bins.

If the data are ~~multidimensional~~ vectors \vec{x} , we can do multidimensional bins.

The likelihood will be a multinomial:

$$\mathcal{L} = n! \prod_{\text{bin } i=1}^m \frac{p_i^{k_i}}{k_i!}$$

where: i is the bin index (not the event index!)

k_i is the number of events in bin i
 p_i is the expectation value for the number of counts in bin i

$$p_i = \int_{\Delta x_i} f(x|\theta) dx \quad \Delta x_i = \text{bin width}$$

The likelihood will be a multinomial, times an extended term for the total number of measurements n :

$$\mathcal{L} = P(n|\theta) \left[\prod_{i=1}^m \frac{p_i^{k_i}}{k_i!} \right] n! \quad \text{where: } i \text{ is the bin index (not the event index!)}$$

k_i is the number of events in bin i

p_i is the probability associated to bin i

→ If $P(n|\theta)$ is a Poisson distribution with expectation value λ :

$$\lambda = \sum_{i=1}^m \lambda_i$$

λ_i = expected number of events in bin i

$$p_i = \frac{\lambda_i}{\lambda}$$

$$= \lambda \cdot \int_{\Delta x_i} f(x|\theta) dx \quad \Delta x_i = \text{width of bin } i$$

$$\mathcal{L} = \frac{e^{-\lambda} \lambda^n}{n!} n! \prod_{i=1}^m \left(\frac{\lambda_i}{\lambda} \right)^{k_i} \frac{1}{k_i!} = e^{-\lambda} \lambda^n \prod_{i=1}^m \frac{\lambda_i^{k_i}}{\lambda^{k_i} \cdot k_i!}$$

$$\mathcal{L} = \prod_{i=1}^m \frac{e^{-\lambda_i} \lambda_i^{k_i}}{k_i!}$$

⇒ The binned version of an extended \mathcal{L} with a Poisson ~~distribution~~ for n , in the product of a Poisson term for each bin!

Notice: This derivation is different than the one reported in literature, e.g. in the Cowan.

[Point 7]

• ~~Minimizer~~ algorithm

→ ~~Minimizer~~ algorithm available in ~~also~~ pretty much any programming language

• Least-squares method

Sometimes this is also referred as χ^2 method, although ~~this is not the case~~.
The squared residuals behave as a χ^2 only in specific circumstances.

Consider a set of measurements with ^{Gaussian} ~~symmetric~~ uncertainties $y_i \pm \sigma_i$, $i=1, \dots, n$.
Assume each y_i corresponds to a perfectly known x_i .

This would be for example the case of a calibration function, where x is the literature energy value of γ lines, and y is the amplitude measured by our detector.

Assume the relation between y and x is: $y = y(x, \vec{\theta})$

Notice: ~~It is just~~ $y(x, \vec{\theta})$ is just a function, not a PDF.

If the measurements y_i are distributed around the curve $y(x, \vec{\theta})$ according to a Gaussian distribution with STD σ_i , the likelihood will be a product of Gaussian PDFs:

$$\mathcal{L}(\vec{y} | \vec{\theta}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi} \sigma_i} \exp \left(-\frac{(y_i - y(x_i, \vec{\theta}))^2}{2 \sigma_i^2} \right)$$

Maximizing \mathcal{L} is equivalent to minimize $-2 \ln \mathcal{L}$:

$$-2 \ln \mathcal{L} = \underbrace{\sum_{i=1}^n \frac{(y_i - y(x_i, \vec{\theta}))^2}{\sigma_i^2}}_{\text{This is a sum of standard normal variables, so it is a } \chi^2 \text{ dist!}} + 2 \underbrace{\sum_{i=1}^n \ln(2\pi \sigma_i^2)}_{\text{does not depend on } \vec{\theta}, \text{ so we can drop it}}$$

This is a sum of standard normal variables, so it is a χ^2 dist!

$$\chi^2(\vec{\theta}) = \sum_{i=1}^n \frac{(y_i - y(x_i, \vec{\theta}))^2}{\sigma_i^2}$$

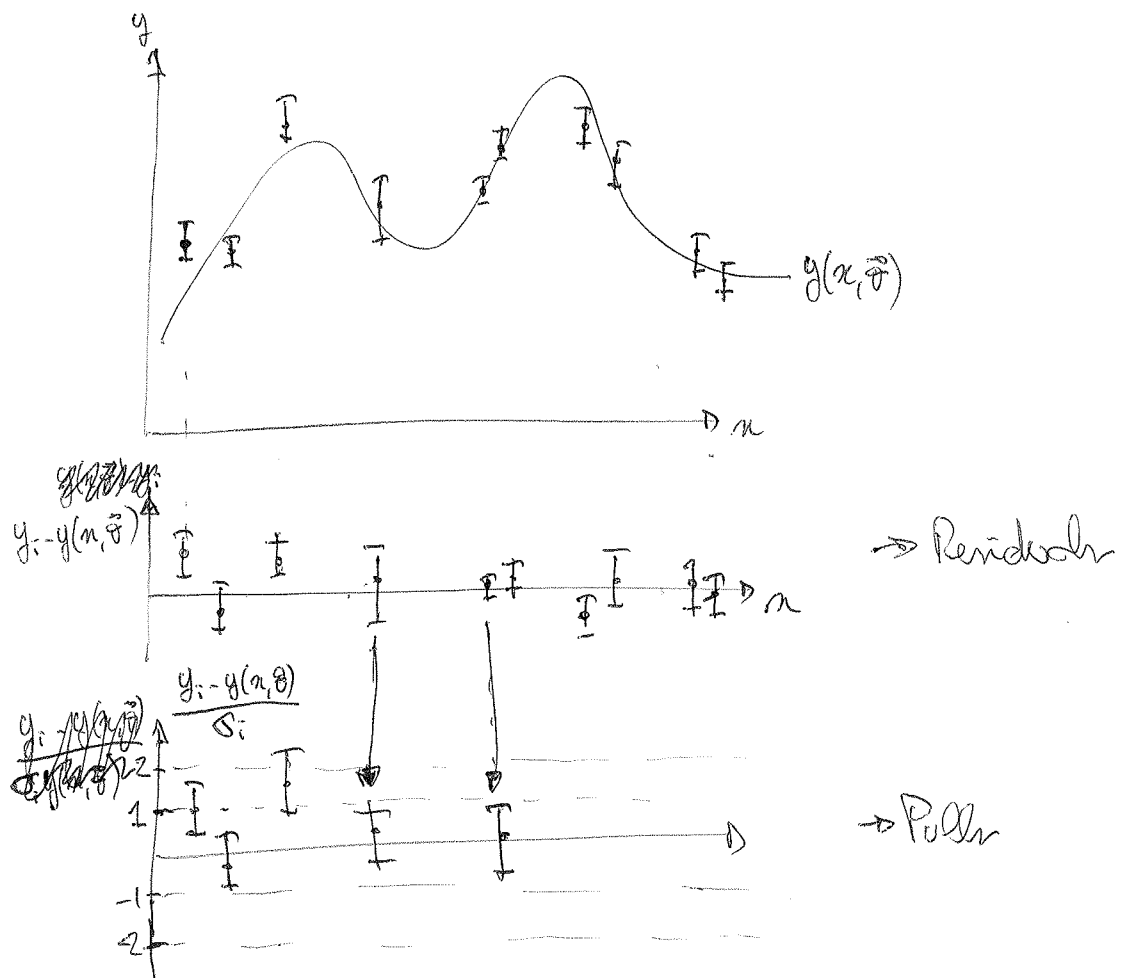
Point 8

Notice That: \rightarrow If the parameters $\vec{\theta}$ are known, this is truly a χ^2 distro.

\rightarrow If $\vec{\theta}$ is not known, That combination of $\vec{\theta}$ values that minimizes $\sum \frac{(y_i - y(\cdot))^2}{\sigma_i^2}$ ~~is~~ is a χ^2 distribution only ~~if~~ if all measurements y_i are uncorrelated.

\rightarrow If The individual measurements y_i are not normally distributed, or if they are correlated, or if $y(x, \vec{\theta})$ does not perfectly describe the data, The least-square method can still be used ~~for~~ for estimating parameters, but it will not behave as a χ^2 .

\hookrightarrow This can actually be used to test the goodness of fit!



• ~~Least squares for binned histograms~~

• Linear regression

Suppose you have a set of measurements y_i with Gaussian uncertainties σ_i .

Suppose $y(x, \vec{\theta})$ is of the form $y = ax + b$ with $\vec{\theta} = (a, b)$.

The χ^2 is:
$$\chi^2(a, b) = \sum_{i=1}^n \frac{(y_i - ax_i - b)^2}{\sigma_i^2}$$

For convenience, let's introduce weights:
$$w_i = \frac{1/\sigma_i^2}{\sum_{i=1}^n 1/\sigma_i^2}$$

So that:
$$\sum w_i = 1$$

We can minimize χ^2 analytically by imposing:
$$\begin{cases} \frac{\partial \chi^2}{\partial a} = 0 \\ \frac{\partial \chi^2}{\partial b} = 0 \end{cases}$$

This corresponds to:
$$\begin{cases} 2 \sum w_i (y_i - ax_i - b)x_i = 0 \\ 2 \sum w_i (y_i - ax_i - b) = 0 \end{cases}$$

~~$$\sum w_i y_i = \sum w_i$$~~

$$\begin{cases} \sum w_i x_i y_i = a \sum w_i x_i^2 + b \sum w_i x_i \\ \sum w_i y_i = a \sum w_i x_i + b \end{cases}$$

Which can be easily solved for a and b , obtaining:

$$\hat{a} = \bar{y} - b \bar{x}$$

$$\hat{b} = \frac{\text{cov}(x, y)}{V[x]}$$

→ A similar approach can be used (sometimes) ~~to solve~~ for nonlinear functions.

→ Otherwise we need numerical methods.

- Numerical and historical considerations

Often times, we deal with histograms rather than scatter plots.

In the past, the computational power was not enough to numerically maximize \mathcal{L} , ~~on the other hand~~, especially when dealing with factorials.

On the other hand, it is much easier to minimize the sum of squares.

This has led to the ~~use of the~~ application of the least-squares method to the fit of histograms.

- Least squares method for fitting histograms

Assume we have a set of measurements x_i following a PDF $f(x_i | \vec{\theta})$.

~~Assume we bin the spectrum~~ Assume we bin the data, and that every bin has a number of entries k_i large enough so that the corresponding Poisson distribution can be approximated by a Gaussian with ~~variance given~~ by the ~~strength~~ expected number of counts in that bin.

with mean given by the expected number of counts in that bin:

$\lambda =$ expected number of total counts

$$\lambda_i = \int_{\text{bin } i} \lambda f(x_i | \vec{\theta}) dx$$

The likelihood reads:
$$\mathcal{L} = \prod_{\text{bin } i=1}^m \frac{1}{\sqrt{2\pi} \sigma_i} \exp\left(-\frac{(k_i - \lambda_i(\vec{\theta}))^2}{2 \sigma_i^2}\right)$$

Taking the logarithm:
$$-2 \ln \mathcal{L} = \sum_{i=1}^m \ln(2\pi \sigma_i^2) + \sum_{i=1}^m \frac{(k_i - \lambda_i(\vec{\theta}))^2}{\sigma_i^2}$$

The problem is, we don't know σ_i .

We can make some considerations: ~~is is~~

$\rightarrow \sigma_i^2$ is the variance of the number of counts (measured or expected) in bin i . If this is large, it is equal to the number of counts.

But the logarithm of a large number is small and varies very slowly, so we can neglect the term $\sum \ln(2\pi \sigma_i^2)$.

At this point we have:
$$\chi^2 = \sum_{i=1}^n \frac{(k_i - \lambda_i(\vec{\theta}))^2}{\sigma_i^2}$$

Now we can make two approximations/assumptions:

→ Neyman's χ^2 : $\sigma_i^2 = k_i = \text{measured number of counts in bin } i$

$$\chi_N^2 = \sum_{i=1}^m \frac{(k_i - \lambda_i(\vec{\theta}))^2}{k_i}$$

↳ Standard for many pre-defined fit methods,
e.g. in ROOT/TMinuit

↳ Fails miserably if one bin is empty

→ Pearson's χ^2 : $\sigma_i^2 = \lambda_i = \text{expected number of counts in bin } i$

$$\chi_P^2 = \sum_{i=1}^m \frac{(k_i - \lambda_i(\vec{\theta}))^2}{\lambda_i(\vec{\theta})}$$

↳ Available/optional in some pre-defined fit methods

↳ Works also for empty bins

↳ Also called "modified least squares method"

• Limitations of least-squares fit for binned histograms.

→ For bins with few events, k_i are not approximately normally distributed.

~~Therefore, for~~

→ In presence of empty bins, the χ_N^2 is usually modified so that empty bins are ignored. This leads to a significant bias.

→ The χ_P^2 also has some (smaller) bias in the presence of empty bins.

⇒ Do not use χ^2 fits on histograms!

• Point estimation in practice

So, how do we minimize the χ^2 or maximize L ?

Every programming language has a tool (or many tools) for performing χ^2 fit or L fit with Poisson ~~statistical~~ models.

The issue ~~isn't~~ arises when we want to tweak the function that we want to minimize for precision or economy (speed) reasons.

Possible tools are:

→ ROOT (C++ and Python) → ~~Default fit~~

root.cern.ch

→ Available χ^2 fit for scatter plots
(called TGraphErrors)

→ Default χ^2_N for histograms

→ χ^2_P fit with option "P"

→ L fit with option "L"

→ Custom L possible, but complicated

→ BAT (C++ and ROOT) → It's actually a framework for Bayesian fit, but it can be easily forced to perform frequentist ones.

github.com/bat/bat

→ Nice feature: The user must implement the L .

Then one can use `model → FindModel()` to call the underlying ROOT fitter (Minuit) using the custom L .

→ numpy + scipy.optimize (Python) → easy to define custom L which is ~~able~~ forced to "minimize"

→ Julia → minimizer available e.g. from "Optim" ~~library~~ package

→ can minimize custom L

• Exercises / examples:

- Unbinned vs Binned fit (extended vs limited likelihood)
- Scatter Plot
- Scatter plot with overestimated uncertainty
- Scatter plot with unaccounted systematic
- Likelihood vs χ^2 for Poisson distributed counts
- Efficiency curve fit with Binomial, Poisson and Gaussian distribution
- ~~Radioactive decay fit (only N , $N + T_{1/2}$)~~

• Detection efficiency and weights

Whenever we operate a detector, we have an efficiency $\varepsilon < 1$ of detecting the events due to several reasons:

- a) Geometry: not all parts of the detector are sensitive (active)
- b) Particle type: we might be sensitive only to some particles
- c) Timing: we might have dead times, pile-up times, ...
- d) Noise: the signal induced by a particle might be happening in a period affected by high noise
- e) Energy/amplitude: The efficiency has some roll-up curve as a function of the signal amplitude/energy.

We have two options: 1) Correct the data. This ~~has~~ introduces the problem that we need to account for the ~~is~~ imputed correction ~~is~~ by assigning a corresponding uncertainty to the data points.

2) Include the efficiency in the fit model.

This is the exact method, and is preferable whenever possible.

• Efficiency-corrected likelihood

Suppose each "true" event has an efficiency $\varepsilon(\vec{\pi})$ to be detected.

ε depends on the parameter $\vec{\theta}$ directly, if we fit some efficiency curve, or indirectly, i.e. through $\vec{\pi}$.

The PDF of detected events $\vec{\pi}$ will be: $\varepsilon(\vec{\pi}|\vec{\theta}) f(\vec{\pi}|\vec{\theta})$

The likelihood will be: ~~Wrong~~ $\mathcal{L} = \frac{e^{-\lambda} \lambda^n}{n!} \prod_i \varepsilon(\vec{\pi}_i|\vec{\theta}) f(\vec{\pi}_i|\vec{\theta})$

→ Example: Looking for peak near threshold

Suppose we have an X-ray detector with the following properties:

Energy resolution $\sigma = 1$ keV

Trigger Efficiency: $\varepsilon(E, t, \sigma_t) = \frac{P}{2} \left[1 + \operatorname{erf} \left(\frac{E - t}{\sigma_t} \right) \right]$ with $P = 0.9$
 $t = 2$ keV
 $\sigma_t = 1$ keV

Suppose we have an exponential background:

$$f_b(E|\lambda) = \lambda \exp(-\lambda E) \quad \text{with } \lambda = 3 \text{ keV}^{-1}$$

Suppose we look for an X-ray signal peak at 5 keV:

$$f_s(E|\mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left[-\frac{(E - \mu)^2}{2\sigma^2} \right] \quad \text{with } \mu = 5 \text{ keV}, \sigma = 1 \text{ keV}$$

Our likelihood will be:

$$\mathcal{L} = \frac{e^{-\lambda} \lambda^n}{n!} \prod_{i=1}^n \left[\frac{s}{s+b} \varepsilon(E_i|t, \sigma_t) f_s(E_i|\mu, \sigma) + \frac{b}{s+b} f_b(E_i|\lambda) \right] \quad \text{with } v = s+b$$

$$\mathcal{L} = \frac{e^{-v} v^n}{n!} \prod_{i=1}^n \frac{s}{s+b} \varepsilon(E_i|t, \sigma_t) f_s(E_i|\mu, \sigma) + \frac{b}{s+b} f_b(E_i|\lambda) \varepsilon(E_i|t, \sigma_t) \quad \text{with } v = s+b$$

$\int \left[\frac{s}{s+b} \varepsilon f_s + \frac{b}{s+b} \varepsilon f_b \right] dE$
↓ does not depend on E_i

$$\mathcal{L} \propto \frac{e^{-v} v^n}{n!} \prod \left[\frac{s}{s+b} \varepsilon f_s + \frac{b}{s+b} \varepsilon f_b \right]$$

Point 15

• Simultaneous Fit

Suppose we have two sets of data \vec{n} and \vec{y} that ~~depend~~ have some parameters in common:

$$f(\vec{n}|\vec{\theta}, \vec{v}) \quad f(\vec{y}|\vec{\theta}, \vec{v})$$

We can run a combined simultaneous fit simply by multiplying the two likelihoods:

$$\mathcal{L}(\vec{n}, \vec{y}|\vec{\theta}, \vec{v}, \vec{v}) = \mathcal{L}(\vec{n}|\vec{\theta}, \vec{v}) \cdot \mathcal{L}(\vec{y}|\vec{\theta}, \vec{v})$$

↓

Example: Combined efficiency and signal + bkg fit

Suppose we have a set of data \vec{E} corresponding to the number of triggered events from ~~where~~ injected events.

We inject 1000 events with a pulse generator into a detector for 20 energy values between ~~100~~ 1 and 20 keV, and count how many survive.

Assume the efficiency follows the curve:

$$\varepsilon(E) = \frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{E - \mu_E}{\sigma_E} \right) \right] \quad \text{with} \quad \begin{array}{l} p = 0.9 \\ \mu_E = 3 \text{ keV} \\ \sigma = 1 \text{ keV} \end{array} \quad \text{for Fit parameters}$$

Suppose then we measure some data \vec{E} distributed as:

a) A Gaussian signal: $f_s(\vec{E}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(E-\mu)^2}{2\sigma^2}\right)$ with $\begin{array}{l} \mu = 5 \text{ keV} \\ \sigma = 1.2 \text{ keV} \end{array} \quad \text{Fixed}$

b) A flat background: $f_b(\vec{E}) = \frac{1}{\Delta E}$ with $\begin{array}{l} \Delta E = E_{\max} - E_{\min} \\ E_{\min} = 1 \text{ keV} \\ E_{\max} = 20 \text{ keV} \end{array}$

~~Suppose the true~~

The likelihood for the efficiency measurement is:

$$\mathcal{L}(\vec{k} | \mu_e, \sigma_e, \rho) = \prod_{i=1}^{n_k} \frac{n!}{k_i! (n-k_i)!} \rho^{k_i} (1-\rho)^{n-k_i} \text{ with } n =$$

$$\mathcal{L}(\vec{k} | \rho, \mu_e, \sigma_e) = \prod_{i=1}^{n_k} \frac{n!}{k_i! (n-k_i)!} \Sigma(E_i)^{k_i} (1 - \Sigma(E_i))^{n-k_i} \text{ with } n = 1000 \text{ injected events}$$

$E_i = 1, 2, \dots, 20 \text{ keV}$

$n_k = 20 \text{ energy points}$

The likelihood for the physics data \vec{E} is:

$$\mathcal{L}(\vec{E} | s, b, \rho) = \frac{e^{-\lambda} \lambda^{n_E}}{n_E!} \prod_{i=1}^{n_E} \left[\frac{\rho s}{\lambda} f_s(E_i) + \frac{\rho b}{\lambda} f_b(E_i) \right]$$

$$\text{with } \lambda = \rho(s+b) \sum_{i=1}^{n_E} \Sigma(E_i)$$

$n_E = \text{number of detected events}$

$s = \text{expected value for signal events}$

$b = \text{expected value for background events}$

The ~~total~~ combined likelihood is:

$$\mathcal{L}(\vec{k}, \vec{E} | \rho, \mu_e, \sigma_e, s, b) = \mathcal{L}(\vec{k} | \rho, \mu_e, \sigma_e) \cdot \mathcal{L}(\vec{E} | s, b, \rho)$$

