

## Modern Methods of Statistical Data Analysis

From parameter estimation to deep learning — A guided tour of probability

## Lecture 4

# Introduction to Parameter Estimation & the Method of Maximum Likelihood

P.-D. Dr. Roger Wolf

roger.wolf@kit.edu

Dr. Jan Kieseler

jan.kieseler@kit.edu

Dr. Pablo Goldenzweig

pablo.goldenzweig@kit.edu

Dr. Slavomira Stefkova

slavomira.stefkova@kit.edu

#### Today

• What goes into such fits?

- First part of lecture: general concepts of parameter estimation
- Second part: the method of maximal likelihood



#### i.i.d. random variables

- Consider *n* variables  $X_1, X_2, \ldots, X_n$ 
  - $X_1, X_2, \ldots, X_n$  are independent and identically distributed (i.i.d.) if
    - $X_1, X_2, \ldots, X_n$  are **independent**, and
    - all have the **same PMF** (if discrete) or **PDF** (if continuous)
      - $E[X_i] = \mu$  for i = 1, ..., n
      - $Var[X_i] = \sigma^2$  for i = 1, ..., n

**Quick check:** Are  $X_1, X_2, \ldots, X_n$  i.i.d. with the following distributions?

- 1.  $X_i \sim \text{Exp}(\tau), X_i \text{ independent } \checkmark$
- 2.  $X_i \sim \text{Exp}(\tau_i), X_i$  independent (unless  $\tau_i$  equal)

3.  $X_i \sim \text{Exp}(\tau), X_1 = X_2 = \cdots = X_n$  dependent!  $(x_1 = x_2 = \cdots = x_n)$ 

4.  $X_i \sim Bin(n_i, p), X_i$  independent (unless  $n_i$  equal)

Review

#### Working with the CLT

Let 
$$X_1, X_2, \ldots, X_n$$
 be i.i.d., where  $E[X_i] = \mu$  and  $Var[X_i] = \sigma^2$ .

As 
$$n \to \infty$$
:  

$$\sum_{i=1}^{n} X_{i} \sim \mathcal{N}(n\mu, n\sigma^{2})$$
Sum of i.i.d. RVs
$$\frac{1}{n} \sum_{i=1}^{n} X_{i} \sim \mathcal{N}(\mu, \frac{\sigma^{2}}{n})$$
Average of i.i.d. RVs
(sample mean)
$$\frac{1}{n} \sum_{i=1}^{n} X_{i} \sim \mathcal{N}(\mu, \frac{\sigma^{2}}{n})$$
Interpret: As we increase *n*
(the size of our sample):
• The variance of our sample
mean  $\sigma^{2}/n$  decreases

• The probability that our sample mean  $\bar{X}$  is close to the true mean  $\mu$  increases

#### CLT: Key take home message

No matter what the distribution of the population is, the distribution of mean samples from the population will always be Normally distributed



i.e., No matter what the distribution of the sample is, if you sample batches of data from that distribution and take the mean of each batch, <u>the mean values from</u> <u>those batches will be Normally distributed</u>

Review

#### The CLT and real life

- Central Limit Theorem:
  - Sample mean  $\bar{X} \sim \mathcal{N}(\mu, \sigma^2/n)$
  - If we know  $\mu$  and  $\sigma^2$ , we can compute probabilities on the sample mean  $\bar{X}$  of a given sample size n

#### In real life:

- Yes, the CLT still holds...
- But we often don't know  $\mu$  or  $\sigma^2$  of our original distribution
- However, we can collect data (a sample of size *n*)
- Question: How can we estimate the values  $\mu$  or  $\sigma^2$  from our sample?
  - Answer: Parameter estimation!

Review

#### Parameter estimation: General concepts (i)



Since it is assumed that the observations are all independent and that each  $X_i$  are described by the same PDF f(x)

#### What if we don't know f(x)?

Central problem of statistics: use x to learn properties of f(x) The RVs we have been discussing are parametric models:

Distribution = model + parameter  $\theta$ 

For each of the distributions below, what are the parameters  $\theta$  ?

<b>1</b> . Ber( <i>p</i> )	$\theta = p$	
2. Poi $(\lambda)$	$\theta = \lambda$	In the real world, we
3. Uni $(\alpha, \beta)$	$\boldsymbol{\theta} = (\alpha, \beta)$	don't know the <i>true</i> parameters, but we do
<b>4.</b> $\mathcal{N}(\mu, \sigma^2)$	$\boldsymbol{\theta} = (\mu, \sigma^2)$	get to observe data
<b>5</b> . $Y = mX + b$	$\boldsymbol{\theta} = (m, b)$	

In general: consider a random variable X distributed according to a PDF  $f(x; \theta)$ 

- The functional form of  $f(x; \theta)$  is known
- The value of at least one parameter  $\theta$  (or parameters  $\theta = (\theta_1, \dots, \theta_m)$ ) is *not known*

Goal: construct a function of observed x to estimate  $\theta$ 

#### Key definitions: Statistic & Estimator

<u>def</u> statistic  $\equiv$  A function of the observed measurements which contains no unknown parameters.

def estimator  $\equiv$  A statistic used to estimate some property of a PDF (e.g.,  $\mu$ ,  $\sigma$ , or other parameter).



Modern Methods of Data Analysis

#### Estimating $\hat{\theta}$ from given data *x*: parameter fitting

#### Since $\hat{\theta}$ is a function of RVs $\Rightarrow$ it is *itself* a RV

i.e. if the entire experiment  $x = (x_1, \dots, x_n)$  is repeated,  $\hat{\theta}$  would change and be described according to PDF  $g(\hat{\theta}; \theta)$ 

**Sampling distribution** (*The PDF of a statistic*) Depends on <u>true value</u>

#### Expectation value of Estimator $\hat{\theta}$

#### **General:** a(x) is a function of RVs distributed according to f(x)

$$E[a(x)] = \int_{-\infty}^{\infty} ag(a)da \stackrel{!}{=} \int_{-\infty}^{\infty} a(x)f(x)dx \qquad \text{Eqn. 1.44, Cowan}$$

Since  

$$g(a)da = \int_{dS} f(x)dx \xrightarrow{}_{\text{multiply}} ag(a)da = a \int_{dS} f(x)dx \xrightarrow{}_{\text{Integrate}} \int_{-\infty}^{\infty} ag(a)da = \int_{-\infty}^{\infty} af(x)dx$$
Thus (replacing a with  $\hat{\theta}$ )  

$$E[\hat{\theta}(\mathbf{x})] = \int \hat{\theta}g(\hat{\theta};\theta)d\hat{\theta} = \int \cdots \int \hat{\theta}(\mathbf{x})f(x_{1};\theta) \cdots f(x_{n};\theta)dx_{1} \cdots dx_{n}$$

Interpret: This is the expected mean of  $\hat{\theta}$  from an infinite # of similar experiments, each with a sample of size *n* 

Modern Methods of Data Analysis

#### Bias & mean squared error (MSE)



**Interpret:** sum of squares of statistical and systematic uncertainties

 $\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ 

sample mean

(not population mean)

**Suppose:** sample of size *n* of RV *X* with values  $x_1, x_2, \ldots, x_n$ 

**<u>Assume</u>**: X is distributed according to some PDF f(x)f(x)

**<u>Need</u>**: a function of the  $x_i$  to be an estimator for the expectation

value (population mean) of x,  $\mu$ .

**Important property given by the weak law of large #'s:** If the variance of x exists, then  $\overline{x}$  is a **consistent** estimator for the population mean  $\mu$ 

i.e., for  $n \to \infty$ ,  $\bar{x}$  converges to  $\mu$ 

**Expectation value** of estimator  $\bar{x}$ 

For alue  

$$E[\bar{x}] = E\left[\frac{1}{n}\sum_{i=1}^{n}x_i\right] = \frac{1}{n}\sum_{i=1}^{n}E[x_i] = \frac{1}{n}\sum_{i=1}^{n}\mu = \mu$$

Modern Methods of Data Analysis

## Estimators for variance



#### Proof of <u>unbiased</u> estimators for variance

1. Show that the sample variance  $s^2$  and the statistic  $S^2$  are unbiased estimators of the population variance  $\sigma^2$  with

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2}$$
 and  $S^{2} = \frac{1}{n} \sum_{i=1}^{n} (x_{i} - \mu)^{2}$ . (1)

with  $\bar{x}$  denoting the sample mean and  $\mu$  being the population mean. I.e. if one does not know the population mean and needs to estimate it from the sample via  $\bar{x}$ , one has to use  $s^2$  to obtain an unbiased estimator for the population variance.



Review proofs at home: ILIAS: /Reading material / L04 / UnbiasedEstimators

#### **Estimator for covariance**

 Similarly one can show that this expression is an unbiased estimator for the covariance of two random variables x and y:

$$\widehat{V}_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y}) = \frac{n}{n-1} (\overline{x}\overline{y} - \overline{x}\overline{y})$$

• This can be normalized by the square-root of the estimators of the sample variance  $s_x$  and  $s_y$ 

$$r = \frac{\widehat{V}_{xy}}{s_x s_y} = \frac{\sum_{i=1}^n (x_i - \overline{x})(y_i - \overline{y})}{\left(\sum_{j=1}^n (x_j - \overline{x})^2 \cdot \sum_{k=1}^n (y_k - \overline{y})^2\right)^{1/2}}$$
$$= \frac{\overline{xy} - \overline{x} \overline{y}}{\sqrt{(\overline{x^2} - \overline{x}^2)(\overline{y^2} - \overline{y}^2)}}.$$

# Variance of estimators

Now were talking about the variance of the estimators we just defined

• Variance of sample mean  $\bar{x}$ 

This was step b) in the proof on slide 15

$$V[\overline{x}] = E[\overline{x}^{2}] - (E[\overline{x}])^{2} = E\left[\left(\frac{1}{n}\sum_{i=1}^{n}x_{i}\right)\left(\frac{1}{n}\sum_{j=1}^{n}x_{j}\right)\right] - \mu^{2}$$
$$= \frac{1}{n^{2}}\sum_{i,j=1}^{n}E[x_{i}x_{j}] - \mu^{2}$$
$$= \frac{1}{n^{2}}\left[(n^{2} - n)\mu^{2} + n(\mu^{2} + \sigma^{2})\right] - \mu^{2} = \frac{\sigma^{2}}{n},$$

This expresses the wellknown result that the standard deviation of the mean of  $\boldsymbol{n}$  measurements of  $\boldsymbol{x}$  is equal to the standard deviation of f(x) itself divided by  $\sqrt{\boldsymbol{n}}$ .

where  $\sigma^2$  is the variance of f(x) and we have used  $E[x_i x_j] = \mu^2$  for  $i \neq j$  and  $E[x_i^2] = \mu^2 + \sigma^2$ 

• In a similar way one can show that the variance of  $s^2$  is

$$V[s^2] = \frac{1}{n} \left( \mu_4 - \frac{n-3}{n-1} \mu_2^2 \right), \qquad \qquad \int_{-\infty}^{\infty} (x-\mu)^n f(x) dx = \mu_n,$$

#### Expectation value and variance of correlation

 The expectation value and variance of the correlation coefficient r depend on higher moments of the joint PDF *f*(*x*, *y*)

$$E[r] = \rho - \frac{\rho(1-\rho^2)}{2n} + O(n^{-2})$$
$$V[r] = \frac{1}{n} (1-\rho^2)^2 + O(n^{-2}).$$

# Answer Time: Quiz 3

#### Monte Carlo method recap

How would you write an algorithm to integrate a multi-dimensional function (of dimensionality N) using the 'acceptance-rejection' method? Address in particular: what ingredients you need and sketch out the algorithm explicitly.

1) Generate N uniform random numbers in intervals of  $[x_{\min}^{i}, x_{\max}^{i}]$ 

• Denote these as  $u = (u_1, u_2, \dots, u_N)$ 

2) Evaluate the function we wish to integrate:  $f(u_1, u_2, ..., u_N)$ 

3) Generate another random number  $\nu$  between 0 and  $f_{\rm max}$ 

• Reject the event if  $\nu > f(u_1, u_2, \dots, u_N)$ ; otherwise accept

4) Repeat 1-3

The integral is then given by total number of accepted events divided by the total number of tried events

# Take 5



# Likelihood (function)

#### $PDF \rightarrow Likelihood$ (Interpretation)

Suppose a measurement of the mass m of an elementary particle yields the value  $m_0$ , and it is known that the measuring apparatus yields values normally distributed about the unknown true mass  $m_t$ , with a known rms deviation  $\sigma_m$ 

The probability density for obtaining the value mgiven the true mass  $m_t$ 

By writing  $\mathscr{L}$  instead of Pwe draw attention to the fact that we are considering its behavior for different values of  $m_t$ given the particular measured datum  $m = m_0$ 



 $\mathscr{L}$  is the probability, under the assumption of the theory  $(m_t)$ , to observe the data which were actually observed  $(m_0)$ 

ILIAS /Reading material /L04 /WhyIsntEveryPhysicistABayesian\_RCousins.pdf https://www.astro.princeton.edu/~strauss/AST303/bayesian\_paper.pdf

#### **Recall Bayes' Theorem terminology**

Where's  $\mathscr{L}$ ?



Review

# posterior $P(F \mid E) = \frac{P(E \mid F) P(F)}{P(E)}$ P(E)normalization constant

#### Likelihood function

Under the assumption of the hypothesis  $f(x; \theta)$ , including the value of  $\theta$ 

The probability that 
$$x_i$$
 in  $[x_i, x_i + dx_i]$  for all  $i = \prod_{i=1}^n f(x_i; \theta) dx_i$ 

Since the  $dx_i$  do not depend on the parameter  $\theta$ 

$$\mathscr{L}(\theta) = \prod_{i=1}^{n} f(x_i; \theta) \longleftarrow \mathsf{PDF of } x$$

parameter(s) we wish to estimate

This is the *joint PDF* for the  $x_i$ , although it is treated as a function of the parameter  $\theta$  (i.e., it's really  $\mathscr{L}(x_i; \theta)$ ). The  $x_i$ , on the other hand, are treated as **fixed** (i.e., <u>the</u> <u>experiment is over</u>). Maximum likelihood estimation is just a systematic way of searching for the parameter values of our *chosen distribution* that maximize the probability of observing the data that we observe

*i.e., we have the data but want to learn about the model, specifically the model's parameters.* 

#### Maximum likelihood estimators

 $\frac{\text{def}}{\text{maximum likelihood estimators for the parameters}} \equiv \max_{\substack{\text{maximize the}\\\text{likelihood function}}}$ 

$$\frac{\partial \mathscr{L}}{\partial \theta_i} = 0, \quad i = 1, \dots, m$$

Often not so easy...

#### Solution: Take the logarithm

- $\rightarrow$  Monotonically increasing (param. val. which maximizes  $\mathscr{L}$ , maximizes  $\log \mathscr{L}$ )
- $\rightarrow$  exponentials in *f* are converted into simple factors

 $\longrightarrow \prod \mathrel{\Rightarrow} \sum$ 

$$\log \mathscr{L}(\theta) = \log \left( \prod_{i=1}^{n} f(x_i; \theta) \right) = \sum_{i=1}^{n} \log f(x_i; \theta)$$

#### Good vs. bad ML estimates

- A sample of 50 observations (shown as tick marks) of a Gaussian random variable.
  - Left: mean and width parameters that maximize the ML (solid shows estimated, dashed shows true PDF)
  - Right: The PDF evaluated with parameters far from the true values, giving lower likelihoods





**Given**: The proper decay times of unstable particles of a certain type have been **measured** for *n* decays, yielding values  $t_i, \ldots t_n$ 

**Choose**: The exponential PDF with mean  $\tau$  as a **hypothesis** for the distribution of *t* 

**Task: Estimate** the value of the parameter  $\tau$ 

$$\log \mathscr{L}(\tau) = \sum_{i=1}^{n} \log f(t_i; \tau) = \sum_{i=1}^{n} \left( \log \frac{1}{\tau} - \frac{t_i}{\tau} \right)$$

$$\underset{\log \mathscr{L} \text{ w.r.t. } \tau}{\text{Maximize}} \longrightarrow \frac{\partial \log \mathscr{L}(\tau)}{\partial \tau} = 0$$

$$\underset{\text{estimator } \hat{\tau}}{\text{Gives the ML}} \longrightarrow \hat{\tau} = \frac{1}{n} \sum_{i=1}^{n} t_i \quad \checkmark$$

$$E[\hat{\tau}(t_1, t_2, \dots, t_n)] = \frac{1}{n} \sum_{i=1}^n \tau = \tau$$

 $f(t;\tau) = \frac{1}{\tau}e^{-t/\tau}$ 

 $\hat{ au}$  is an unbiased estimator for au

Read up on Gaussian example at home (pg74, Cowan)

Exponential distribution

So far: Method to determine heta via ML and get  $\hat{ heta}$  🧹

**<u>Now</u>**: What is the variance of  $\hat{\theta}$  ?

i.e., if we **repeat our experiment** a large number of times, how widely spread out would  $\hat{\theta}$  be?

3 methods:

- 1. Analytical calculation of  $V[\hat{\theta}]$
- 2. MC method
- 3. RCF bound / graphical technique

Continue using the exponential function as an example since <u>applicable to all 3 methods</u>

$$f(t;\tau) = \frac{1}{\tau}e^{-t/\tau}$$

For a **very limited** number of cases, one can compute the variances of the ML estimators **analytically** 

 $n \infty$ 

$$V[\hat{\tau}] = E[\hat{\tau}^2] - (E[\hat{\tau}])^2 \qquad \text{Recall:} \quad E[x] = \int_{-\infty}^{\infty} xf(x)ds \quad (1.39, \text{Cowan})$$

$$= \int \dots \int \left(\frac{1}{n} \sum_{i=1}^n t_i\right)^2 \frac{1}{\tau} e^{-t_1/\tau} \dots \frac{1}{\tau} e^{-t_n/\tau} dt_1 \dots dt_n$$

$$-\left(\int \dots \int \left(\frac{1}{n} \sum_{i=1}^n t_i\right) \frac{1}{\tau} e^{-t_1/\tau} \dots \frac{1}{\tau} e^{-t_n/\tau} dt_1 \dots dt_n\right)^2$$

$$= \frac{\tau^2}{n} \qquad \text{No surprise here: Recall we derived this on slide 17}$$

#### How would we report this?

# $\hat{\tau} \pm \hat{\sigma}_{\hat{\tau}}$

Estimate from one experiment

Estimated standard deviation one expects  $\hat{\tau}$  to vary by if experiment is repeated many times with the same number of measurements per experiment

#### Variance of ML estimators: MC method (i)

Often too difficult to compute the variances analytically Use the MC method to investigate the distribution of the ML estimates

- 1. Simulate a large # of experiments
- 2. Compute the ML estimates each time
- $\rightarrow$

"true" parameter = the estimated value from the real experiment

3. Look at how the resulting values are distributed

Unbiased estimator for the variance of a PDF

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2}$$

Compute *s* for the ML estimates obtained from the MC experiments and give this as the statistical error of the parameter estimated from the real measurement

#### Variance of ML estimators: MC method (ii)



#### Variance of ML estimators: MC method (iii)



But what if we want the decay constant instead of the mean lifetime?

$$\lambda = \frac{1}{\tau}$$

In general: transformational invariance

 $\partial L/\partial \theta = 0$  implies  $\partial L/\partial a = 0$  at  $a = a(\theta)$ unless  $\partial a/\partial \theta = 0$ .

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial \theta} = 0.$$

#### But with

$$\hat{\lambda} = 1/\hat{\tau} = n/\sum_{i=1}^{n} t_i.$$

$$E[\hat{\lambda}] = \lambda \; \frac{n}{n-1}$$

...biased for small n

#### Variance of ML estimators: RCF bound (i)

#### Often too difficult to compute the variances analytically and a MC study involves a significant effort

Proof: ILIAS /Reading material / L04 /RCF\_proof.pdf

Use the Rao-Cremer-Frechet (RCF) inequality (aka information inequality) to compute the lower bound on an estimator's variance



## For the **exponential** distribution with mean $\tau$

Recall we chose this distribution since it can be solved analytically



 $\partial b/\partial \tau = 0, \ b = 0$  (slide 29)



In the case of equality (minimum variance), the estimator is said to be efficient

#### **Key points:**

- If efficient estimators exist for a given problem, the ML method will find them
- 2) ML estimators are always efficient in the large sample limit (exept when the extent of the sample space depends on the estimated parameter)

...but what if you <u>can't</u> compute the RCF bound analytically? One can estimate  $V^{-1}$  by evaluating the second derivative with the measured data and the ML estimates  $\hat{\theta}$ :

 $\left(\hat{V^{-1}}\right)_{ij} = \frac{\partial^2 \log \mathscr{L}}{\partial \theta_i \partial \theta_j} \bigg|_{\theta = \hat{\theta}}$ 

For a single parameter  $\theta$  this reduces to:



The routines **MIGRAD** and **HESSE** in **MINUIT** determine numerically the matrix of second derivatives of  $\underline{\log \mathcal{L}}$  using finite differences, evaluate it at the ML estimates, and invert to find the covariance matrix.

#### Intuition

#### Why is this how we calculate the standard errors?

The easy way to think about this is to recognize that the **curvature of the** *likelihood function* tells us how certain we are about our estimate of our parameters. The **more curved** the likelihood function, the **more certainty** we have that we have estimated the right parameter. <u>The second derivative of the</u> *likelihood function is a measure of the likelihood function's curvature* - this is why it provides our estimate of the uncertainty with which we have estimated our parameters.

If the curvature is small, then the likelihood surface is flat around its maximum value (the MLE). If the curvature is large and thus the variance is small, the likelihood is strongly curved at the maximum.

Key ingredient: Taylor expand log likelihood function around maximum:

$$\log \mathscr{L}(\theta) = \log \mathscr{L}(\hat{\theta}) + \left[\frac{\partial \log \mathscr{L}}{\partial \theta}\right]_{\theta=\hat{\theta}} (\theta-\hat{\theta}) + \frac{1}{2!} \left[\frac{\partial^2 \log \mathscr{L}}{\partial \theta^2}\right]_{\theta=\hat{\theta}} (\theta-\hat{\theta})^2 + \dots$$

$$= 0$$

$$at \ maximum$$

$$\hat{\sigma^2}_{\theta} = \left(-1/\frac{\partial^2 \log \mathscr{L}}{\partial \theta^2}\right)\Big|_{\theta=\hat{\theta}} RCF$$

$$RCF$$

$$\log \mathscr{L}(\theta) = \log \mathscr{L}_{max} - \frac{(\theta-\hat{\theta})^2}{2\hat{\sigma^2}_{\hat{\theta}}} \text{ or } \left[\log \mathscr{L}(\hat{\theta}\pm\hat{\sigma}_{\hat{\theta}}) = \log \mathscr{L}_{max} - \frac{1}{2}$$

$$\text{With known maximum, can determine estimator of variance by solving what value of  $\sigma_{\theta}$  gives a likelihood value of  $\log \mathscr{L}_{max} - 1/2$$$

#### Variance of ML estimators: graphical technique (ii)

#### Return to our previous example using the exponential distribution

-52.5  $\log L(\tau)$ Reading off from the curve  $\hat{\tau} - \Delta \hat{\tau}$ •  $\Delta \hat{\tau} = 0.137$ •  $\Delta \hat{\tau}_{+} = 0.165$ -53 Both reasonably close and we find -53.5 •  $\hat{\sigma}_{\hat{\tau}} \approx \Delta \hat{\tau}_{-} \approx \Delta \hat{\tau}_{+} \approx 0.15$  We will later make a -54 reinterpretation of the interval 0.8 1  $[\tau - \sigma_{\tau}, \tau + \sigma_{\tau}]$  as an approximation of the 68.3% central confidence interval



(slides 31 & 34)

http://www.pp.rhul.ac.uk/~cowan/sda

# Take a second



### Another ML example: 2 parameters (i)

- Consider a particle reaction where each scattering event is characterized by a certain scattering angle  $\Theta$  (or equivalently  $x = \cos \Theta$ )
  - Suppose now a theory predicts that this follows an angular distribution given by

$$f(x; \alpha, \beta) = \frac{1 + \alpha x + \beta x^2}{2 + 2\beta/3}.$$
e.g.  $\alpha = 0, \beta = 1$  is the LO QED expectation for  $e^+e^- \rightarrow \mu^+\mu^-$ 

• To make this slightly more complicated, let's assume we have a finite detector acceptance, such that  $x_{min} < x < x_{max}$ :

$$f(x;\alpha,\beta) = \frac{1+\alpha x+\beta x^2}{(x_{\max}-x_{\min})+\frac{\alpha}{2}(x_{\max}^2-x_{\min}^2)+\frac{\beta}{3}(x_{\max}^3-x_{\min}^3)}.$$

## Another ML example: 2 parameters (ii)

- MC experiment with 2000 events according to
  - $\alpha = 0.5, \beta = 0.5$

• 
$$x_{min} = -0.95, x_{max} = 0.95$$

 Numerically maximizing the loglikelihood function we find



$$\hat{\alpha} = 0.508 \pm 0.052,$$

$$\hat{\beta} = 0.47 \pm 0.11, \quad \text{from second} \\ \text{derivatives} \\ \text{from ML maximum} \quad \widehat{\cos}[\hat{\alpha}, \hat{\beta}] = 0.0026 \quad r = 0.46.$$

## Another ML example: 2 parameters (iii)

- Can validate this result by using MC techniques
  - Let's produce 500 similar experiments, all with 2000 events with the true values for  $\alpha = 0.5$ ,  $\beta = 0.5$
  - Sample means, standard deviations, covariance and correlation





Sample variance, covariance, correlation: Good agreement with RCF bound

#### Another ML example: 2 parameters (iv)

Can also draw likelihood contour in 2D

• Contour of 
$$\log \mathscr{L} = \log \mathscr{L}_{\max} -$$



 Large sample limit contour given by

$$\frac{1}{1-\rho^2}\left[\left(\frac{\alpha-\hat{\alpha}}{\sigma_{\hat{\alpha}}}\right)^2+\left(\frac{\beta-\hat{\beta}}{\sigma_{\hat{\beta}}}\right)^2-2\rho\left(\frac{\alpha-\hat{\alpha}}{\sigma_{\hat{\alpha}}}\right)\left(\frac{\beta-\hat{\beta}}{\sigma_{\hat{\beta}}}\right)\right]=1.$$

• I.e. ellipse with center at  $(\hat{\alpha}, \hat{\beta})$ and angle  $\phi$  with respect to the  $\alpha$ -axis

$$\tan 2\phi = \frac{2\rho\sigma_{\hat{\alpha}}\sigma_{\hat{\beta}}}{\sigma_{\hat{\alpha}}^2 - \sigma_{\hat{\beta}}^2}.$$

#### For next time

- Required reading
  - Cowan textbook: chapters 5 & 6.1-6.8
  - Reading material / L04 /
    - VarianceOfMLEstimators
    - UnbiasedEstimators
    - RCF\_proof
- Extra reading for fun: /Reading material / L04 /
  - WhyIsntEveryPhysicistABayesian\_RCousins

Listed as 'fun' but you should **really** read it!

#### Next time

- Extended maximum likelihood
- Maximum likelihood with binned data
- Testing goodness-of-fit
- $\chi^2$  method

# Quiz Time: 4th Round

2. Let us assume you have a simple PDF of the form

$$f(x;\lambda) = 1 + \lambda \left(x - 0.5\right), \qquad (2)$$

with a sample space S spanning the interval [0, 1] such that  $\int_S f(x; \lambda) dx = \int_0^1 f(x; \lambda) dx = 1$ .

- a) First sketch the PDF for  $\lambda = 1, 0, -1$ .
- b) Five measurements were done giving x = (0.89, 0.03, 0.50, 0.31, 0.49). Calculate the log-likelihood function for three different values of  $\lambda = 1, -0.5, -1$  by hand.
- c) The log-likelihood function is in good approximation a parabolic function, i.e. can be described by a polynomial of second order as a function of the tested value  $\lambda$ . Calculate the coefficients a, b, c of log  $L(\lambda) = a\lambda^2 + b\lambda + c$ . For what value of  $\lambda$  is log  $L(\lambda)$  maximal?
- d) Sketch the log-likelihood function. Using the graphical method, i.e.

$$\log L(\hat{\lambda} \pm \hat{\sigma}_{\hat{\lambda}}) = \log L_{\max} - \frac{1}{2}, \qquad (3)$$

determine the uncertainty  $\hat{\sigma}_{\hat{\lambda}}$  of the estimated parameter  $\hat{\lambda}$  (with  $\hat{\lambda}$  denoting the value where  $\log L(\hat{\lambda}) = \log L_{\max}$  and  $L_{\max}$  is the maximal likelihood value).

#### Bibliography

- Part of the material presented in this lecture is adapted from the following sources. See the active links (when available) for a complete reference
  - **Probability for CS** (Stanford): slides 3-4
  - Statistical Data Analysis textbook by G. Cowan (U. London): all figures with white background