

Part

Outline

Basic elements

- some vocabulary
- Probability axioms
- some probability distributions
- Two approaches: Frequentist vs. Bayesian
- Hypothesis testing
- Parameter estimation
- Other subjects "nuisance", "spurious", "look elsewhere"

Parts II & III

a quick review of Lecture I

- Probability vs. Statistics
- PDF (and CDF)
 - expectation values, covariance matrix, correlation coefficients
- Frequentist vs. Bayesian
- Some theorems
 - LLN, CLT, Neyman-Pearson lemma, Wilks theorem, etc.

Two interpretations of Probability

- While the classic or frequentist approach can lead to a well-defined probability for a given situation, it is not always usable.
 - \rightarrow In such circumstances one is left with only one option: *Bayesian*.
- When data are scarce → these two approaches can give somewhat different predictions,

but given sufficiently large data sample, they give pretty much the same conclusion. In that case the choice between the two may be regarded arbitrary.

• Perhaps, we may choose one for the main result, and try the other for a cross-check.

Basics

Hypothesis Testing

• In the frequentist approach, we do not, in general, assign probability of a hypothesis itself.

Rather, we compute the probability to accept/reject a hypothesis assuming that it (or some alternative) is true.

In Bayesian, on the other hand, probability of any given hypothesis (*degree of belief*) could be obtained by using the Bayes' theorem:

$$P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H')\pi(H')dH'}$$

which depends on the prior probability $\pi(H)$

Basics

Hypothesis Testing

- A hypothesis *H* specifies the probability for the data (*shown symbolically as x here*),
 often expressed as a function f(x|H)
- The measured data \vec{x} could be anything:
 - * observation of a single particle, a single event, or an entire experiment
 - * uni-/multi-variate, continuous or discrete
- the two kinds:

Basics

- * simple (or "point") hypothesis $-f(\vec{x}|H)$ is completely specified
- * composite hypothesis *H* contains unspecified parameter(s)
- The probability for \vec{x} given *H* is also called the **likelihood** of the hypothesis, written as $L(\vec{x}|H)$

Critical Region - what it is

- Consider e.g. a simple hypothesis H_0 and an alternative H_1
- A (frequentist) test of H_0 :

Basics

Specify a critical region w of the data space Ω such that, assuming H_0 is correct, there is no more than some (small) probability α to observe data in w

Freq. vs. Bayes. Hyp. Testing

 $P(\vec{x} \in w | H_0) \le \alpha$

- α : "size" or "significance level" of the test
- If \vec{x} is observed within w, we reject H_0 with a confidence level 1α

Param. Est. Adv. subjects

data space Ω

critical region w

Basics

Critical Region - how to choose

Param. Est. Adv. subj

• In general, \exists an ∞ number of possible critical regions that give the same significance level α

Freq. vs. Bayes. Hyp. Testing

• Usually, we place the critical region where there is a low probability α for $\vec{x} \in w$ if H_0 is true, but high if the alternative (H_1) is true



Basics

 $t(x_1,\ldots,x_n) = t_{\text{cut ut}}$ where $t(x_1, \ldots, x_n)$ is a scalar test statistic. We can work out the pdfs $g(t|H_0), g(t|H_1), \dots$ The bou for an *n* 2 g(t) I cut an equa accept $H_0 \iff$ reject H_0 Decision boundary is now a 1.5 single 'cut' on *t*, defining $g(t|H_0)$ the critical region. 1 $g(t|H_1)$ where *t* So for an *n*-dimensional 0.5 For the problem we have a PDFs g(corresponding 1-d problem. 0 Decisio 2 3 n 4 'cut' on

 \Rightarrow for a G. Cowan reduced to a 1-dim. problem 5

 $e H_1.$

.....on w of the

data space such that there is no more than some (small) probab

- Rejecting H₀ when it is true is called the Type-I error
 (Q) Given the significance α of the test, what is the maximum probability of Type-I error?
- We might also accept *H*₀ when it is indeed false, and an alternative *H*₁ is true. This is called the **Type-II error**

The probability β of Type-II error:

 $P(\vec{x} \in \Omega - w | H_1) = \beta$

 $1 - \beta$ is called the **power** of the test with respect to H_1

Basics



Basics

exercise on Type-I, II errors

Since $B \to K^* \gamma$ has much higher branching fraction than $B \to \rho \gamma$, the former can be a serious background to the latter. It is crucial to understand the "efficiency" and "fake rate" of K/π identification system of your experiment in this study. The figure below shows the $M_{K\pi}$ invarianbt mass distribution, where one of the pion mass (in $\rho^0 \to \pi^+\pi^-$ decay) is replaced by the Kaon mass, for the $B^0 \to \rho^0 \gamma$ signal candidates (Belle, PRL 2008).



Express the following observables in Type-I & Type-II errors. *What are H*₀ & *H*₁*, for each case?*

- $f_{\pi^+ \to K^+}$ = probability of misidentifying a π^+ as a K^+
- $f_{K^+ \to \pi^+}$ = probability of misidentifying a K^+ as a π^+
- ϵ_{K^+} = prob. of identifying a K^+ correctly as a K^+
- ϵ_{π^+} = prob. of identifying a π^+ correctly as a π^+

Definition of a (frequentist) hypothesis test

Defining a multivariate critical region H_1 .

A test of H_0 is defined by specifying a critical region w of the data space such that there is no more than some (small) probable α , assuming H_0 is correct, to observe the data there, i.e.,



Basics

Definition of a (frequentist) hypothesis test

Defining a multivariate childre H1.

A test of H_0 is defined by specifying a critical region w of the data space such that there is no more than some (small) probal some more sophisticated Ways is correct, to observe the data there, i.e.,

linear

Freq.

Basics



(ex) Fisher discriminants, etc.

or nonlinear



(ex) artificial neural net, etc.

algorithms for asimultivariate critical region

A test of H_0 is defined by specifying a critical region w of the data space such that there is no more than some (small) probability of the specific probability of th

Many (old or new) data space such that there is no more than some (small) proba- α , assuming H_0 is correct, to observe the data there, i.e.,

Fisher discriminants

Freq.

$$P(x \in w \mid H_0) \le \alpha$$

- Artificial neural networksquality if data are discrete.
- Boosted decision treesiled the size or
- Kernel density methods

If x is observed in the critical region, reject H_0 .



G. Cowan

Cargese 2012 / Statistics for HEP / Lecture 1

Basics

Software for multivariate analysis

TMVA, Höcker, Stelzer, Tegenfeldt, Voss, Voss, physics/0703039

From tmva.sourceforge.net, also distributed with ROOT Variety of classifiers Good manual

StatPatternRecognition, I. Narsky, physics/0507143

Further info from www.hep.caltech.edu/~narsky/spr.html Also wide variety of methods, many complementary to TMVA Currently appears project no longer to be supported Basics

How to che

• Use Neyman-Pear For a test of size α to obtain the highe choose the critical

everywhere in *w* all where *k* is a consta

• Equivalently, the o

Definition of a (frequentist) hypothesis test

Consider e.g. a simple hypothesis H_0 and alternative H_1 .

A test of H_0 is defined by specifying a critical region w of the data space such that there is no more than some (small) probable α , assuming H_0 is correct, to observe the data there, i.e.,

$$P(x \in w \mid H_0) \le \alpha$$

Need inequality if data are discrete.

 α is called the size or significance level of the test.

If x is observed in the critical region, reject H_0 .



G. Cowan

Cargese 2012 / Statistics for HEP / Lecture 1

 $t(\vec{x}) = P(\vec{x}|H_1) / P(\vec{x}|H_0)$

(Note) Any monotonic function of this leads to the *same test*.

Practical Statistics for Particle Physicists

an application of Neyman-Pearson Lemma



 $P_{i} \equiv P_{i}^{dE/dx} \times P_{i}^{\text{TOF}} \times P_{i}^{\text{Ch}} \quad \text{e.g.} \ (i = \pi \text{ or } K)$ For optimal statistic, construct the likelihood ratio $R_{K/\pi} = P_{K}/P_{\pi}$ (or any ftn. that is monotonic to it) Belle actually used $R_{K/\pi} = P_{K}/(P_{K} + P_{\pi})$ so that $0 \leq R_{K/\pi} \leq 1$



Y. Kwon (Yonsei Univ.)

Practical Statistics for Particle Physicists

Oct. 12–25, 2016



Kyle Cranmer (NYU)

CERN School HEP, Romania, Sept. 2011

Y. Kwon (Yonsei Univ.)

Practical Statistics for Particle Physicists



(Quiz) With Neyman **Definition lefth for we misy**) Have **THE** is test way to optimize the critical region ("Cut"). Then why should we bother with multivariate analyses such as artificial neural network action that there is no more than some (small) probability

 H_0 is correct, to observe the data there, i.e.,

 $P(x \in w \mid H_0) \leq \alpha$ data space Ω ity if data are size or evel of the test. ed in the $|H_0\rangle$ Ansele The modeling of P(x H) may $|H_0\rangle = P(\downarrow$ not be perfect, if the correlation w are not taken properly into account. G. Cowan

are not taken properly into accour This will become more serious for higher dimensions of **x**.



Y. Kwon (Yonsei Univ.)

Practical Statistics for Particle Physicists

Oct. 12-25, 2016



By User:Repapetilto @ Wikipedia & User:Chen-Pan Liao @ Wikipedia - File:P value.png, CC BY-SA 3.0, https:// commons.wikimedia.org/w/index.php?curid=36661887

In short, *p*-value is the 'size' of a test against a given hypothesis.

Y. Kwon (Yonsei Univ.)

Practical Statistics for Particle Physicists



Basics Freq. vs. Bayes. Hyp. Testing

Param. Est.

Adv. subjects

Remember?

Gaussian (Normal) distribution



Table 36.1: Area of the tails α outside $\pm \delta$ from the mean of a Gaussian distribution.

Significance and the *p*-value

Often we quote the significance Z, for a given *p*-value

• Z = the number of standard dev. that a Gaussian random variable would fluctuate in one direction to give the same *p*-value



 $p = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \, dx = 1 - \Phi(Z) \qquad 1 - \text{TMath::Freq}$

 $Z = \Phi^{-1}(1-p)$ TMath::NormQuantile

(Ex) Z = 5 (a "5-sigma effect") $\Leftrightarrow p = 2.9 \times 10^{-7}$

p-value example: a fair coin?

We toss a coin N = 20 times and get n = 17 heads. Test whether this coin is 'fair' or not.

Hypothesis H_0 : the coin is fair ($\mu = 50\%$ chance for head)

Example: significance of a signal

Assume both
$$n_s$$
, n_b are Poisson.
 $P(n;s,b) = \frac{(s+b)^n}{n!} e^{-(s+b)}$

Suppose b = 0.5 (assume precise), and we observe $n_{obs} = 5$. Can we claim evidence for a signal excess? Give *p*-value for the null hypothesis s = 0.

1983 Korean Baseball Champion HaiTai Tigers' Lineup



			BA	OBA	SLG	HR	RBI	SB
1	김일권	CF	.275	.345	.364	6	26	48
2	서정환	SS	.257	.320	.339	3	34	13
3	김성한	1B	.327	.401	.448	7	40	13
4	김봉연	DH	.280	.371	.552	22	59	2
5	김종모	LF	.350	.404	.524	11	44	7
6	김준환	BF	.248	.308	.362	10	43	11
7	김무종	С	.262	.313	.453	12	60	2
8	양승호	3B	.236	.292	.309	2	11	3
9	차영화	2B	.266	.308	.323	1	23	16

(observation) Six out of 9 starting hitters have family name 'Kim'.

(fact) According to census, ~20% of all Koreans have family name 'Kim'.

(Hypothesis to test) The manager of 1983 Tigers (himself a 'Kim') has a bias toward players with family name 'Kim'.

Example: comparison of hypotheses

Given a set of data resulting from the measurement of some observable

 $\Omega = \{-1.0, -0.9, -0.7, -0.1, 0.0, 0.1, 0.2, 0.5, 0.6, 1.0\}$

where the total number of data N = 10, determine which of the following models is a better description of the data:

- H_0 : the data are distributed according to a Standard Gaussian ($\mu =$ 0, $\sigma = 1$);
- H_1 : the data are uniformly distributed over |-1, +1|

For a test statistic, we may use

$$R = \frac{P(H_0|\text{data})}{P(H_1|\text{data})} = \frac{P(\text{data}|H_0)P(H_0)}{P(\text{data}|H_1)P(H_1)} \qquad R \ge 1 ??$$

Example: comparison of hypotheses

Given a set of data resulting from the measurement of some observable

 $\Omega = \{-1.0, -0.9, -0.7, -0.1, 0.0, 0.1, 0.2, 0.5, 0.6, 1.0\}$



taken from A. Bevan's book

Basics Freq. vs. Bayes. Hyp. Testing Param. Est. Adv. subjects Control of a Chequentist Hypothesis Con

A test of H_0 is defined by specifying a critical region w of the

- In general, we cannot find a single critical region that gives the maximum power for all possible alternatives (no "uniformly most powerful" test)
- In HEP, we often try to construct a test of the Standard Model as H₀ (or sometimes called "background only")
 such that we have a well specified *false discovery rate* α (=prob. to reject H₀ when it is true),
 and high power w.r.t. some interesting alternative H₁, e.g. SUSY, Z', etc.
- But, there is no such thing as a *model-independent* test.
 Any statistical test will inevitably have high power w.r.t. some alternatives and less for others

Basics

Freq. vs. Bayes. Hyp. Testing

but

m

U

t→Wb

Intervals

BR(t-Wb)= T(t-Wa

 $= \frac{|V_{eb}|^2}{|V_{eb}|^2 + |V_{eb}|^2}$ $= \frac{|V_{eb}|^2}{|V_{eb}|^2 + |V_{eb}|^2}$ $\cong \frac{(6.9745)^2}{(0.0074)^2 + (0.044)^2 \cdot (0.7145)^2}$

Adv. subjects

Intervals

from 'Big Bang Theory'

Definition of a (frequentist) hypothesis test

Consider a α o simple hypothesis U and alternative U

2016 Review of Particle Physics.

Please use this CITATION: C. Patrignani et al. (Particle Data Group), Chin. Phys. C, 40, 1000

 $t \rightarrow Wb$

• $\Gamma(t \rightarrow Wb)/\Gamma(t \rightarrow Wq (q = b, s, d))$

OUR AVERAGE assumes that the systematic uncertainties are uncorrelated.

VALUE	DOCUMENT ID			COMMENT				
0.957 ± 0.034	OUR AVERAGE Error includes scale factor of 1.5.							
0.87 ± 0.07	1 AALTONEN	2014G (CDF	$\ell\ell + E_T + \geq 2j$ (0,1,				
$1.014 \pm 0.003 \pm 0.032$	2 KHACHATRYAN	2014E	CMS	ℓ ℓ + $\not\!$				
0.94 ±0.09	³ AALTONEN	2013G (CDF	$\ell + E_T + \geq$ 3jets (\geq				
0.90 ± 0.04	4 ABAZOV	2011X I	D0					

Measurement with errors

Let's say we are reporting a single measurement

$$x = a \pm b$$

Frequentist interpretation

Repeating the measurement many times under identical conditions ("ensemble"), the estimated interval will vary each time. In 68.3% of those results, the true value of x will lie within the interval.

Result of each measurement is a sampling from a Gaussian distribution G(μ,σ)

- We may not know µ
- We have some idea about σ -- experimental sensitivity

when $\mu \pm \sigma$ is not enough...

If the PDF of the estimator is not Gaussian, or if there are physical boundaries on the possible values of the parameter,

one usually quotes an interval given a confidence level.

Basics

a Bayesian procedure for intervals

$$1 - \alpha = \int_{\theta_{\text{lo}}}^{\theta_{\text{up}}} p(\theta | \boldsymbol{x}) \, d\theta$$

If the physical value is non-negative, one may choose a prior:

$$\pi(s) = \begin{cases} 0 & s < 0\\ 1 & s \ge 0 \end{cases}$$

Likelihood for *s*, given *b*, is

$$P(n|s) = \frac{(s+b)^n}{n!}e^{-(s+b)}$$

If what we seek is of a very low (or no) signal, interval \rightarrow UL Then, $\int_{a}^{s_{up}} P(n|s) \pi(s) ds$

$$1 - \alpha = \int_{-\infty}^{\log p} p(s|n) ds = \frac{\int_{-\infty}^{\infty} P(n|s) \pi(s) ds}{\int_{-\infty}^{\infty} P(n|s) \pi(s) ds}$$

$$F_{\chi^2}^{-1}: \text{ inverse of the CDF} \quad \Rightarrow s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1} [1 - \alpha; 2(n+1)] - b_{39}$$

Freq. vs. Bayes. Hyp. Testing Intervals Adv. subjects (Ex) UL on Poisson parameter

- Consider again the case of observing $n \sim \text{Poisson}(s + b)$. Suppose b = 4.5 and $n_{\text{obs}} = 5$. Find upper limit on *s* at 95% CL.
- Relevant alternative is s = 0, resulting in critical region at low n.
- The *p*-value of hypothesized *s* is $P(n \le n_{obs}; s, b)$. Therefore, the upper limit s_{up} at $CL = 1 - \alpha$ is obtained from

Basics

Frequentist "confidence intervals"

on repeated measurements

Remember frequentist approach is always about repeated measurements!

"confidence interval"

= intervals constructed to include the true value of the parameter with a probability \geq (a specified value)

Frequentist "confidence intervals"

Consider a pdf f(x;\theta) $P(x_1 < x < x_2;\theta) = 1 - \alpha = \int_{x_1}^{x_2} f(x;\theta) dx$

- *x* : outcome of an experiment
- $\boldsymbol{\theta}$: unknown parameter for which we set the interval



for Frequentist UL, the 90% (or whatever) integration is done above the UL



$\frac{1}{90\%}$			n=0	M_{LOTA} $(< q_{OT})$ n = 0		for Frequent whatever) in high you high high high high high high high hig			tist UL, the 90% (or integration is done above the UL) $\overline{}$		
$1 - \alpha = 90\%$		$1 - \alpha = $	$1 - \alpha = 95\%$		$1 - \alpha = 90\%$						
n	$\mu_{ m lo}$	$\mu_{ m up}$	$\mu_{ m lo}$	$\mu_{ m up}$		\overline{n}	μ_1	μ_2	μ_1	μ_2	
0		2.30	—	3.00		0	0.00	2.44	0.00	3.09	
1	0.105	3.89	0.051	4.74		1	0.11	4.36	0.05	5.14	
2	0.532	5.32	0.355	6.30		2	0.53	5.91	0.36	6.72	
3	1.10	6.68	0.818	7.75		3	1.10	7.42	0.82	8.25	
4	1.74	7.99	1.37	9.15		4	1.47	8.60	1.37	9.76	
4	1.74	7.99	1.37	9.15		4	1.47	8.60	1.37	9.76	3

Table 38.3: naive frequentist interval

Table 38.4: Feldman-Cousins interval
Phys. Rev. D57, 3873 (1998)"unified approach"44

Confidence interval from inversion of a test

- For confidence intervals for a parameter θ, define a test of size α for the hypothesized value θ (repeat this for all θ)
 - If the observed data falls in the critical region, reject the value θ .
 - The values that are *not rejected* constitutes a **confidence interval** for μ at confidence level $CL = 1 \alpha$.
- By construction the confidence interval will contain the true value of θ with probability $\geq 1 \alpha$.
 - * The interval depends on the choice of the test (critical region).
 - * If the test is formulated in terms of a *p*-value, p_{θ} , then the confidence interval represents those values of θ for which $p_{\theta} > \alpha$.
 - * To find the end points of the interval, set $p_{\theta} = \alpha$ and solve for θ .

Coincidence of frequentist and Bayesian intervals

If the expected background is zero, the Bayesian upper limit (for a Poisson RV) becomes equal to the limit determined by frequentist approach.

$$s_{\rm up} = \frac{1}{2} F_{\chi^2}^{-1} \left[p, 2(n+1) \right] - b$$
$$= \frac{1}{2} F_{\chi^2}^{-1} \left(1 - \alpha; 2(n+1) \right)$$

For more details, you may read e.g. a statistics review in PDG. <u>http://pdg.lbl.gov/2015/reviews/rpp2015-rev-statistics.pdf</u> Basics

Parameter Estimation

Basics of parameter estimation

• The parameters of a PDF are constants characterizing its shape, e.g.

$$f(x;\theta) = \frac{1}{\theta}e^{-x/\theta}$$

where θ is the parameter, while x is the random variable.

Suppose we have a sample of observed values, *x*.
 We want to find some function of the data to *estimate* the parameter(s): *θ*(*x*).
 Often *θ* is called an estimator.

Basics

- Properties of estimators
- If we were to repeat the entire measurement, the set of estimates would follow a PDF:



- We want small (or zero) bias (\Rightarrow syst. error): $b = E[\hat{\theta}] \theta$ $b = E[\hat{\theta}]$
- and we want a small variance (\Rightarrow stat. error): $V[\hat{\theta}]$

Y. Kwon (Yonsei Univ.)

Basics

Practical Statistics for Particle Physicists



he

)a

Bias vs. Consistency



Basics

Practical Statistics for Particle Physicists

Likelihood fusider for simple hypothesis H_0 and alternative H_1 . A test of H_0 is defined by specifying a critical region w of the

- Suppose the entire result of an experiment (*set of measurements*) is a collection of numbers \vec{x} , and suppose the joint PDF for the data \vec{x} is a function depending on a set of parameters $\vec{\theta}$: $f(\vec{x}; \vec{\theta})$
- Evaluate this function with the measured data \vec{x} , regarding this as a function of $\vec{\theta}$ only. This is the **likelihood function**.

$$L(\vec{\theta}) = f(\vec{x}; \vec{\theta}) \ (\vec{x}, \text{fixed})$$

critical region, reject H_0 .

G. Cowan

Cargese 2012 / Statistics for HEP / Lecture 1

Basics

critical region w

Ja

Ω

The likelihood function for i.i.d. data

Freq. vs. Bayes. Hyp. Testing

i.i.d. = independent and identically distributed

Consider *n* independent observations of {x : x₁, · · · , x_n}, where x follows f(x, θ).
 The joint PDF for the whole data sample is:

$$f(x_1, \cdots, x_n; \vec{\theta}) = \prod_{i=1}^n f(x_i; \vec{\theta})$$

• In this case, the likelihood function is

$$L(\vec{\theta}) = \prod_{i=1}^{n} f(x_i; \vec{\theta}) \quad (x_i \text{ constant})$$

Param. Est.

So we define the max. likelihood (ML) estimator(s) to be the parameter value(s) for which the L becomes maximum.

Y. Kwon (Yonsei Univ.)

Basics

Practical Statistics for Particle Physicists

Oct. 12–25, 2016

52

1 1

the ba

Ω

Basics

Freq. vs. Bayes. Hyp. Testing Param. Est. Adv. subjects

ML estimator example: fitting to a straight line

- Suppose we have a set of data:
 (x_i, y_i, σ_i), i = 1, · · · , n.
- Modeling: y_i are independent and follow y_i ~ G(μ(x_i), σ_i) (G: Gaussian) where μ(x_i) are modelled as μ(x; θ₀, θ₁) = θ₀ + θ₁x

Assume x_i and σ_i are known.

Goal: to estimate θ₀
 Here, let's suppose we don't care about θ₁ (an example of a *nuisance parameter*)



Basics

Freq. vs. Bayes. Hyp. Testing Param. Est. Adv. subjects

ML fit with Gaussian data

• In this example, the *y_i* are assumed independent, so that likelihood function is a product of Gaussians:

$$L(\theta_0, \theta_1) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right]$$

• Then maximizing *L* is equivalent to minimizing

$$\chi^{2}(\theta_{0},\theta_{1}) = -2\ln L(\theta_{0},\theta_{1}) + C = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i};\theta_{0},\theta_{1}))^{2}}{\sigma_{i}^{2}}$$

i.e., for Gaussian data, ML fitting is the same as the method of least squares

VVIIK'S TREOFEM

ML fit or Least-square fit?

- Solution Consider we have a random variable $x \in [0, 3]$, and a distribution f(x).
- In a series of measurements, we obtained
 - 9 events in [0,1), 10 events in [1,2), and 8 events in [2,3]
 - We have a model of uniform f(x), and would like to estimate the mean value of $\int f(x) dx$ for each histogram bin.
- Run a thought-experiment, comparing
 - maximum likelihood method, and least-square method
 - Do they give the same result?

Bayesian likelihood function

• Suppose our *L*-function contains two parameters θ_0 and θ_1 , where we have some knoweldege about the prior probability on θ_1 from previous measurements:

$$\pi(\theta_0, \theta_1) = \pi_0(\theta_0)\pi_1(\theta_1)$$

$$\pi_0(\theta_0) = \text{const.}$$

$$\pi_1(\theta_1) = \frac{1}{\sqrt{2\pi}\sigma_p} e^{-(\theta_1 - \theta_p)^2/2\sigma_p^2}$$

• Putting this into the Bayes' theorem gives the posterior probability:

$$p(\theta_0, \theta_1 | \vec{x}) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} e^{-(y_i - \mu(x_i; \theta_0, \theta_1))^2 / 2\sigma_i^2} \pi_0 \frac{1}{\sqrt{2\pi}\sigma_p} e^{-(\theta_1 - \theta_p)^2 / 2\sigma_p^2}$$

• Then, $p(\theta_0|\vec{x}) = \int p(\theta_0, \theta_1|\vec{x}) \ d\theta_1$

with alternative priors

Suppose we don't have a previous measurement of θ₁ but rather a theorist saying that θ₁ should be > 0 and not too much greater than, say, 0.1 or so. In that case, we may try modeling the prior for θ₁ as something like

$$\pi_1(heta_1) = rac{1}{ au} e^{- heta_1/ au}, \; heta_1 \geq 0, \; au = 0.1$$

• From this we obtain (numerically) the posterior PDF for θ_0



• This plot summarizes all knowledge about θ_0 .

Exercises

1. Setting limits:

(a) The parameter S is measured to be -1.1 ± 0.4 . What is the Bayesian 90% CL upper limit on S given that S is physically bound within the interval [-1, +1]?



(b) Determine the 90% CL upper limit on the signal yield μ given a background expectation of one event (ignore systematic uncertainties) and an observed yield of one event for a rare decay search.

(Example) a T2K result

PRL 107, 041801 (2011)



T2K observed 6 candidate events of $v_{\mu} \rightarrow v_{e}$ while a background of 1.5±03 events is expected.

- How significant is this signal?
- How to include the systematic uncertainty in the analysis?
- What is the relevant 'limit' from this result?

(Ex) "Bayesian classifier"

- 3. The ratio of energy and momentum (E/p) for a charged particle is used to identify electrons against other particle types. Assume that it takes a value within the interval [0.0, 1.2].
 - (a) Assuming that the production of all three types of charged particles are equal, define a Bayesian classifier based on a Gaussian distributions with $\mu = 1.0$ for an electron, 0.38 for π^{\pm} and 0.08 for μ^{\pm} . For simplicity, let's assume all three PDFs have a width of $\sigma = 0.2$. Classify the tracks with the following measured values of E/p: 0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9 and 1.0.
 - (b) Repeat the above problem, but with the assumption that the expected fraction of charged particles in the experiment is 80% for π^{\pm} , 10% for μ^{\pm} and 10% for e^{\pm} .

(Ex) "Bayesian classifier"

- 3. The ratio of energy and momentum (E/p) for a charged particle is used to identify electrons against other particle types. Assume that it takes a value within the interval [0.0, 1.2].
 - (a) Assuming that the production of all three types of charged particles are equal, define a Bayesian classifier based on a Gaussian distributions with $\mu = 1.0$ for an electron, 0.38 for π^{\pm} and 0.08 for μ^{\pm} . For simplicity, let's assume all three PDFs have a width of $\sigma = 0.2$. Classify the tracks with the following measured values of E/p: 0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9 and 1.0.
 - (b) Repeat the above problem, but with the assumption that the expected fraction of charged particles in the experiment is 80% for π^{\pm} , 10% for μ^{\pm} and 10% for e^{\pm} .
- Given some data that can be tested against a set of classifications given by hypotheses H_i .
- For each event ω_j in the data set we compute $P(\omega_j|H_i)$ for all *i*.
- Find *i* where H_i gives the largest value of P over all *i*'s, and we classify the event ω_j as belonging to the category j.