

Practical Statistics for Particle Physics Lecture 1

AEPS2018, Quy Nhon, Vietnam

Roger Barlow

The University of Huddersfield

August 2018



Lecture 1: The Basics

1 Probability

- What is it?
- Frequentist Probability
- Conditional Probability and Bayes' Theorem
- Bayesian Probability

2 Probability distributions and their properties

- Expectation Values
- Binomial, Poisson and Gaussian

3 Hypothesis testing

Question: What is Probability?

Typical exam question

Q1 Explain what is meant by the *Probability* P_A of an event A
[1]

Four possible answers

- P_A is number obeying certain mathematical rules.
- P_A is a property of A that determines how often A happens
- For N trials in which A occurs N_A times, P_A is the limit of N_A/N for large N
- P_A is my belief that A will happen, measurable by seeing what odds I will accept in a bet.

Mathematical

Kolmogorov Axioms:



A. N. Kolmogorov

For all $A \subset S$

$$P_A \geq 0$$

$$P_S = 1$$

$$P_{(A \cup B)} = P_A + P_B \text{ if } A \cap B = \phi \text{ and } A, B \subset S$$

From these simple axioms a complete and complicated structure can be erected. E.g. show $P_{\bar{A}} = 1 - P_A$, and show $P_A \leq 1$

But!!!

This says *nothing* about what P_A actually means.

Kolmogorov had frequentist probability in mind, but these axioms apply to any definition.

Classical

or Real probability

Evolved during the 18th-19th century

Developed (Pascal, Laplace and others) to serve the gambling industry.



Two sides to a coin - probability $\frac{1}{2}$ for each face

Likewise 52 cards in a pack, 6 sides to a dice...

Answers questions like 'What is the probability of rolling more than 10 with 2 dice?'



Problem: cannot be applied to continuous variables. Symmetry gives different answers working with θ or $\sin\theta$ or $\cos\theta$. Bertan's paradoxes.

Frequentist

The usual definition taught in schools and undergrad classes

$$P_A = \lim_{N \rightarrow \infty} \frac{N_A}{N}$$

N is the total number of events in the ensemble (or collective)

The probability of a coin landing heads up is $\frac{1}{2}$ because if you toss a coin 1000 times, one side will come down ~ 500 times.

The lifetime of a muon is $2.2\mu\text{s}$ because if you take 1000 muons and wait $2.2\mu\text{s}$, then ~ 368 will remain. The probability that a DM candidate will be found in your detector is $[\textit{insert value}]$ because of 1,000,000 (simulated) DM candidates $[\textit{insert value} \times 1,000,000]$ passed the selection cuts

Important

P_A is not just a property of A , but a joint property of A and the ensemble.

Problems (?) for Frequentist Probability

More than one ensemble

German life insurance companies pay out on 0.4% of 40 year old male clients.

Your friend Hans is 40 today. What is the probability that he will survive to see his 41st birthday?

99.6% is an answer (if he's insured)

But he is also a non-smoker and non-drinker - so maybe 99.8%?

He drives a Harley-Davidson - maybe 99.0%?

All these numbers are acceptable

What is the probability that a K^+ will be recognised by your PID?

Simulating lots of K^+ mesons you can count to get $P = N_{acc}/N_{tot}$

These can be divided by kaon energy, kaon angle, event complexity... and will give different probabilities ... All correct.

There may be no Ensemble

What is the probability that it will rain tomorrow?

There is only one tomorrow. It will either rain or not. P_{rain} is either 0 or 1 and we won't know which until tomorrow gets here

What is the probability that there is a supersymmetric particle with mass below 2 TeV?

There either is or isn't. It is either 0 or 1

Bayes' theorem

Bayes' Theorem applies (and is useful) in **any** probability model

Conditional Probability: $P(A|B)$: probability for A , given that B is true.

Example: $P(\spadesuit A) = \frac{1}{52}$ and $P(\spadesuit A|Black) = \frac{1}{26}$

Theorem

$$P(A|B) = \frac{P(B|A)}{P(B)} \times P(A)$$

Proof.

The probability that A and B are both true can be written in two ways

$$P(A|B) \times P(B) = P(A \& B) = P(B|A) \times P(A)$$

Throw away middle term and divide by $P(B)$



Bayes' theorem

Examples

Example

$$P(\spadesuit A | Black) = \frac{P(Black | \spadesuit A)}{P(Black)} P(\spadesuit A) = \frac{1}{\frac{1}{2}} \times \frac{1}{52} = \frac{1}{26}$$

Example

Example: In a beam which is 90% π , 10% K , kaons have 95% probability of giving no Cherenkov signal; pions have 5% probability of giving none. What is the probability that a particle that gave no signal is a K ?

$$P(K | no\ signal) = \frac{P(no\ signal | K)}{P(no\ signal)} \times P(K) = \frac{0.95}{0.95 \times 0.1 + 0.05 \times 0.9} \times 0.1 = 0.68$$

This uses the (often handy) breakdown:

$$P(B) = P(B|A) \times P(A) + P(B|\bar{A}) \times \overline{P(A)}$$

Bayesian Probability

Probability expresses your belief in A . 1 represents certainty, 0 represents total disbelief

Intermediate values can be calibrated by asking whether you would prefer to bet on A , or on a white ball being drawn from an urn containing a mix of white and black balls.

This avoids the limitations of frequentist probability - coins, dice, kaons, rain tomorrow, existence of SUSY can all have probabilities.



Bayesian Probability and Bayes Theorem

Re-write Bayes' theorem as

$$P(\text{Theory}|\text{Data}) = \frac{P(\text{Data}|\text{Theory})}{P(\text{Data})} \times P(\text{Theory})$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

Works sensibly

Data predicted by theory boosts belief - moderated by probability it could happen anyway

Can be chained.

Posterior from first experiment can be prior for second experiment. And so on. (Order doesn't matter)

From Prior Probability to Prior Distribution

Suppose theory contains parameter a : (mass, coupling, decay rate...)

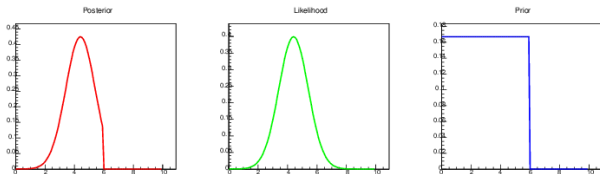
Prior probability distribution $P_0(a)$

$\int_{a_1}^{a_2} P_0(a) da$ is your prior belief that a lies between a_1 and a_2

$\int_{-\infty}^{\infty} P_0(a) da = 1$ (or: your prior belief that the theory is correct)

Generalise the number $P(\text{data}|\text{theory})$ to the function $L(x|a)$

Bayes' Theorem given data x the posterior is : $P_1(a) \propto L(x|a)P_0(a)$



If range of a infinite, $P_0(a)$ may be vanishingly small ('improper prior'). Not a problem. Just normalise $P_1(a)$

Shortcomings of Bayesian Probability

Subjective Probability

Your $P_0(a)$ and my $P_0(a)$ may be different. How can we compare results?

What is the right prior?

Is the wrong question.

'Principle of ignorance' - take $P(a)$ constant (uniform distribution). But then not constant in a^2 or \sqrt{a} or $\ln a$, which are equally valid parameters.

Jeffreys' Objective Priors

Choose a flat prior in a transformed variable a' for which the Fisher information, $-\left\langle \frac{\partial^2 L(x;a)}{\partial a^2} \right\rangle$ is flat. Not universally adopted for various reasons.

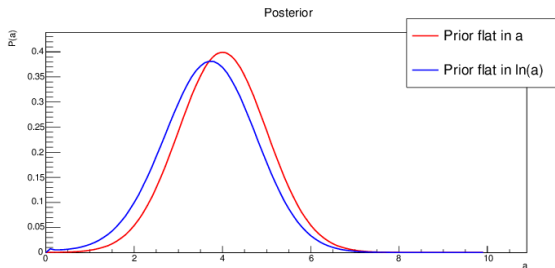
With lots of data, $P_1(a)$ decouples from $P_0(a)$. But not with little data..

Right thing to do: try several forms of prior and examine spread of results ('robustness under choice of prior')

Just an example

Measure $a = 4.0 \pm 1.0$.

Likelihood is Gaussian (coming up!)



Taking a prior uniform in a gives a posterior with a mean of 4.0 and a standard deviation of 1.0 (red curve)

Taking a prior uniform in $\ln a$ shifts the posterior significantly.

Exercise

Try this for yourself with various values

Heretical idea - maybe classical probability still has a place?

Quantum Mechanics gives probabilities.

If P_A is not 'real' - either because it depends on an arbitrary ensemble, or because it is a subjective belief - then it looks like there is nothing 'real' in the universe.

The state of a coin - or an electron spin - having probability $\frac{1}{2}$ makes sense.

The lifetime of a muon - i.e. probability per unit time that it will decay - seems to be a well-defined quantity, a property of the muon and independent of any ensemble, or any Bayesian belief.

The probability a muon will produce a signal in your muon detector seems like a well-defined quantity, if you specify the 4 momentum and the state of the detector ...

Of course the inverse probability "What is the probability that a muon signal in my detector comes from a real muon, not background" is not intrinsically defined.

Perhaps classical probability has a place in physics - but not in interpreting results.

Do not mention this to a statistician or they will think you're crazy

Integer Values

Numbers of positive tracks, numbers of identified muons, numbers of events..

Generically call this r . Probabilities $P(r)$

Real-number Values

Energies, angles, invariant masses...

Generically call this x . Probability Density Functions $P(x)$.

$P(x)$ has dimensions of $[x]^{-1}$. $\int_{x_1}^{x_2} P(x) dx$ or $P(x) dx$ are probabilities

Sometimes also use cumulative $C(x) = \int_{-\infty}^x P(x') dx'$

Mean, Standard deviation, and expectation values

From $P(r)$ or $P(x)$ can form the **Expectation Value**

$$\langle f \rangle = \sum_r f(r)P(r) \quad \text{or} \quad \langle f \rangle = \int f(x)P(x) dx$$

Sometimes written $E(f)$

In particular the **mean** $\mu = \langle r \rangle = \sum_r rP(r)$ or $\langle x \rangle = \int xP(x) dx$
and higher **moments** $\mu_k = \langle r^k \rangle = \sum_r r^k P(r)$ or $\langle x^k \rangle = \int x^k P(x) dx$
and **central moments**

$$\mu'_k = \langle (r - \mu)^k \rangle = \sum_r (r - \mu)^k P(r) \quad \text{or} \quad \langle (x - \mu)^k \rangle = \int (x - \mu)^k P(x) dx$$

The Variance and Standard Deviation

$$\mu'_2 = V = \sum_r (r - \mu)^2 P(r) = \langle r^2 \rangle - \langle r \rangle^2$$
$$\text{or } \int (x - \mu)^2 P(x) dx = \langle x^2 \rangle - \langle x \rangle^2$$

The **standard deviation** is the square root of the variance $\sigma = \sqrt{V}$

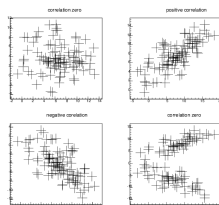
Statisticians usually use variance. Physicists usually use standard deviation
Skew is $\langle (x - \langle x \rangle)^3 \rangle / \sigma^3$ and Kurtosis is $\langle (x - \langle x \rangle)^4 \rangle / \sigma^4 - 3$

Covariance and Correlation

2-dimensional data (x, y)

Form $\langle x \rangle, \langle y \rangle, \sigma_x$ etc

Also other quantities



Covariance

$$\text{Cov}(x, y) = \langle (x - \mu_x)(y - \mu_y) \rangle = \langle xy \rangle - \langle x \rangle \langle y \rangle$$

Correlation

$$\rho = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

ρ lies between 1 (complete correlation) and -1 (complete anticorrelation).

$\rho = 0$ if x and y are independent.

Covariance and Correlation (continued)

Many Dimensions ($x_1, x_2, x_3 \dots x_n$)

Covariance matrix $\mathbf{V}_{ij} = \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle$

Correlation matrix $\rho_{ij} = \frac{\mathbf{V}_{ij}}{\sigma_i \sigma_j}$

Diagonal of \mathbf{V} is σ_i^2

Diagonal of ρ is 1.

The Binomial Distribution

Binomial: Number of successes in N trials, each with probability p of success

$$P(r; p, N) = \frac{N!}{r!(N-r)!} p^r q^{1-r} \quad (q \equiv 1 - p)$$

Binomial distributions
for

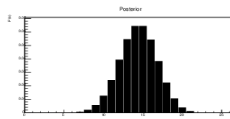
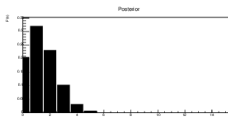
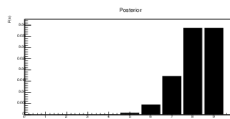
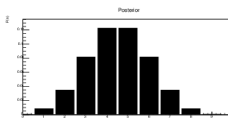
(1) $N = 10, p = 0.6$

(2) $N = 10, p = 0.9$

(3) $N = 15, p = 0.1$

(4) $N = 25, p = 0.6$

Mean $\mu = Np$, Variance $V = Npq$, Standard Deviation $\sigma = \sqrt{Npq}$

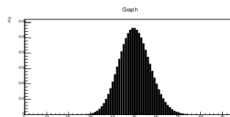
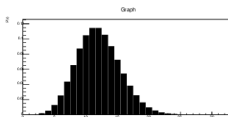
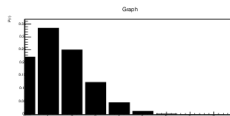
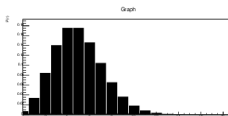


The Poisson Distribution

Number of events occurring at random rate λ

$$P(r; \lambda) = e^{-\lambda} \frac{\lambda^r}{r!}$$

Limit of binomial as $N \rightarrow \infty$, $p \rightarrow 0$ with $np = \lambda = \text{constant}$



Poisson distributions for

- (1) $\lambda = 5$
- (2) $\lambda = 1.5$
- (3) $\lambda = 12$
- (4) $\lambda = 50$

Mean $\mu = \lambda$, Variance $V = \lambda$, Standard Deviation $\sigma = \sqrt{\lambda} = \sqrt{\mu}$

Meet this **a lot** as it applies to event counts - on their own or in histogram bins

Pop Quiz

You need to know the efficiency of your PID system for positrons

Find 1000 data events where 2 tracks have a combined mass of 3.1 GeV (J/ψ) and negative track is identified as an e^- . ('Tag-and-probe' technique)

In 900 events the e^+ is also identified. In 100 events it is not. Efficiency is 90%

What about the error?

Colleague A says $\sqrt{900} = 30$ so efficiency is $90.0 \pm 3.0\%$

Colleague B says $\sqrt{100} = 10$ so efficiency is $90.0 \pm 1.0\%$

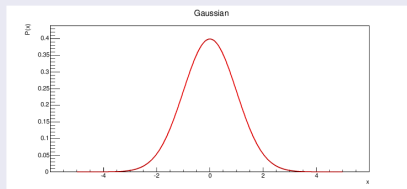
Which is right?

The Gaussian

The Formula

$$P(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The Curve



Only 1 Gaussian curve, as μ and σ are just location and scale parameters

Properties

Mean is μ and standard deviation σ . Skew and kurtosis are 0.

The Central Limit Theorem

Why the Gaussian is so important

If the variable X is the sum of N variables $x_1, x_2 \dots x_N$ then

- 1 Means add: $\langle X \rangle = \langle x_1 \rangle + \langle x_2 \rangle + \dots + \langle x_N \rangle$
- 2 Variances add: $V_X = V_1 + V_2 + \dots + V_N$
- 3 If the variables x_i are independent and identically distributed (i.i.d.) then $P(X)$ tends to a Gaussian for large N

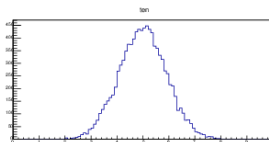
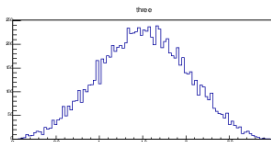
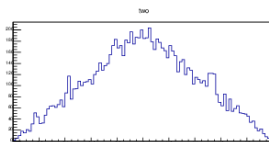
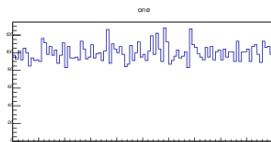
(1) is obvious

(2) is pretty obvious, and means that standard deviations add in quadrature, and that the standard deviation of an average falls like $\frac{1}{\sqrt{N}}$

(3) applies **whatever** the form of the original $p(x)$

Demonstration

Take a uniform distribution from 0 to 1. It is flat. Add two such numbers and the distribution is triangular, between 0 and 2.



With 3 numbers, it gets curved. With 10 numbers it looks pretty Gaussian

Proof

Introduce the **Characteristic Function** $\langle e^{ikx} \rangle = \int e^{ikx} P(x) dx = \tilde{P}(k)$

Expand the exponential as a series

$$\langle e^{ikx} \rangle = \langle 1 + ikx + \frac{(ikx)^2}{2!} + \frac{(ikx)^3}{3!} \dots \rangle = 1 + ik \langle x \rangle + (ik)^2 \frac{\langle x^2 \rangle}{2!} + (ik)^3 \frac{\langle x^3 \rangle}{3!} \dots$$

Take logarithm and use expansion $\ln(1+z) = z - \frac{z^2}{2} + \frac{z^3}{3} \dots$

this gives power series in (ik) , where coefficient $\frac{\kappa_r}{r!}$ of $(ik)^r$ is made up of expectation values of x of total power r

$$\kappa_1 = \langle x \rangle, \kappa_2 = \langle x^2 \rangle - \langle x \rangle^2, \kappa_3 = \langle x^3 \rangle - 3 \langle x^2 \rangle \langle x \rangle + 2 \langle x \rangle^3$$

... These are called the **Semi-invariant cumulants of Thièle**. Under a change of scale α , $\kappa_r \rightarrow \alpha^r \kappa_r$. Under a change in location only κ_1 changes.

If X is the sum of i.i.d. random variables: $x_1 + x_2 + x_3 \dots$ then $\tilde{P}(X)$ is the convolution of $P(x)$ with itself N times

The FT of convolution is the product of the individual FTs

The logarithm of a product is the sum of the logarithms

So $\tilde{P}(X)$ has cumulants $K_r = N \kappa_r$

To make graphs commensurate, need to scale X axis by standard deviation, which grows like \sqrt{N} . Cumulants of scaled graph $K'_r = N^{1-r/2} \kappa_r$

As $N \rightarrow \infty$ these vanish for $r > 2$. Leaving a quadratic.

If the log is a quadratic, the exponential is a Gaussian. So $\tilde{P}(X)$ is Gaussian.

The FT of a Gaussian is a Gaussian. QED.

Gaussian or Normal?

Statisticians call it the 'Normal' distribution. Physicists don't. But be prepared.

Even if the distributions are not identical, the CLT tends to apply, unless one (or two) dominates.

Most 'errors' fit this, being compounded of many different sources.

Hypothesis Testing: What is it?

Making choices

Is this track a pion or a kaon?

Is this event signal or background?

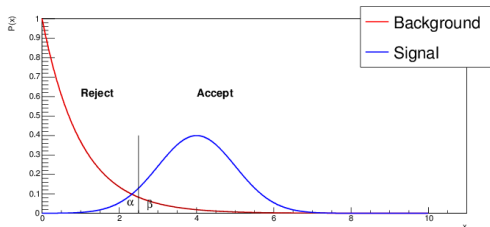
Is the detector performance degrading with time?

Does the data agree with the Standard Model prediction or not?

Type I and Type II errors

Suppose you measure some parameter x which is related to what you are trying to measure.

(May well be output from a neural network or other ML system)



Imposing a cut as shown:

Lose fraction α of signal. ('Type I error'). α is the **significance**

Admit fraction β of background. ('Type II error'). $1 - \beta$ is the **power**

Where should I put the cut?

Strategy for the cut depends on three things - hypothesis testing only covers one of them

Performance

α and β as functions of the cut value

Prior signal to noise ratio

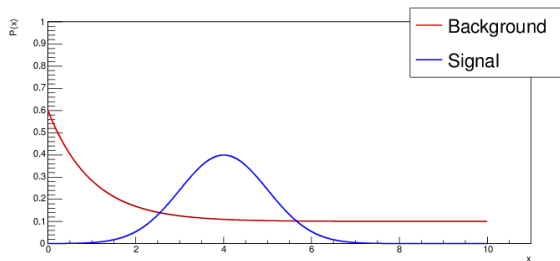
These plots are normalised to 1. Red curve is (probably) MUCH bigger.

Penalties

You have a trade-off between efficiency and purity: what are they worth?
In medical decisions, type I errors are **much** worse than type II.

The Neymann-Pearson Lemma

Suppose S and B curves are more complicated - or x is multidimensional?

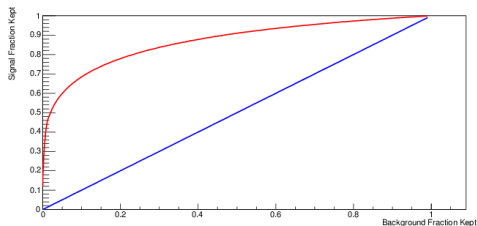


N & P say: include regions of greatest $\frac{S(x)}{B(x)}$ (ratio of likelihoods)

For a given α , this gives the smallest β ('Most powerful at a given significance')

Proof: if you remove a small region from 'accept' to 'reject' it has to be replaced by equivalent which (by construction) brings more background. However complicated, problem reduces to a single monotonic variable $\frac{S}{B}$

Efficiency, Purity, and ROC plots



Performance shown by ROC ('Receiver Operating Characteristic') plots. Plot fraction of background accepted (β) against fraction of signal retained ($1 - \alpha$).

Effect of increasing cut goes from very loose at top right (all data accepted) to very tight at bottom left (all data rejected).

Diagonal line corresponds to no discrimination - curves identical

The further the actual line bulges away from that, the better

Warning: 'Background Efficacy', 'Contamination', 'Purity' are used ambiguously

The Null Hypothesis

To show an effect is present

- Eating broccoli makes you smart
- Facebook advertising increases sales
- A new drug increases patient survival rates
- The data shows Beyond-the-Standard-Model physics

You have to try (your best) and fail to disprove the opposite: The **Null Hypothesis** H_0

- Broccoli lovers and broccoli loathers have the same IQ
- Sales are independent of the Facebook advertising budget
- The survival rates for old and new treatments are the same
- The Standard Model (functions or Monte-Carlo) describe the data

If the null hypothesis is not tenable, you've proved your point

α - the 'significance' - is the probability that the null hypothesis will be wrongly rejected, and you'll claim an effect where there isn't any.

There is a minefield of difficulties. Correlation v. Causation. Multiple trials and self-censorship... More on this in Part 3.