Statistical Methods for LHC Physics



Glen Cowan RHUL Physics www.pp.rhul.ac.uk/~cowan



Physics ColloquiumUniversity of Siegen3 July, 2008

Outline

Quick overview of physics at the Large Hadron Collider (LHC)

New multivariate methods for event selection

Decision trees Support Vector Machines

Some applications of Bayesian methods

Outlook for data analysis at the LHC

Data analysis at the LHC

The LHC experiments are expensive ~ \$10¹⁰ (accelerator and experiments)

the competition is intense (ATLAS vs. CMS) vs. Tevatron

and the stakes are high:



So there is a strong motivation to extract all possible information from the data.

The Large Hadron Collider



Detectors at 4 pp collision points:

- ATLAS CMS general purpose LHCb (b physics)
- ALICE (heavy ion physics)

Counter-rotating proton beams in 27 km circumference ring

pp centre-of-mass energy 14 TeV



The ATLAS detector

2100 physicists37 countries167 universities/labs



Toroid Magnets Solenoid Magnet SCT Tracker Pixel Detector TRT Tracker



25 m diameter
46 m length
7000 tonnes
~10⁸ electronic channels

The Standard Model of particle physics

Matter...



+ gauge bosons...

photon (γ), W[±], Z, gluon (g)

+ relativity + quantum mechanics + symmetries... = Standard Model

25 free parameters (masses, coupling strengths,...).

Includes Higgs boson (not yet seen).

Almost certainly incomplete (e.g. no gravity).

Agrees with all experimental observations so far.

Many candidate extensions to SM (supersymmetry, extra dimensions,...)

A simulated SUSY event in ATLAS



missing transverse energy

Background events



This event from Standard Model ttbar production also has high $p_{\rm T}$ jets and muons, and some missing transverse energy.

 \rightarrow can easily mimic a SUSY event.

LHC data

At LHC, $\sim 10^9$ pp collision events per second, mostly uninteresting

do quick sifting, record ~200 events/sec single event ~ 1 Mbyte 1 "year" $\approx 10^7$ s, 10^{16} pp collisions / year 2×10^9 events recorded / year (~2 Pbyte / year)

For new/rare processes, rates at LHC can be vanishingly small e.g. Higgs bosons detectable per year could be $\sim 10^3$ \rightarrow 'needle in a haystack'

For Standard Model and (many) non-SM processes we can generate simulated data with Monte Carlo programs (including simulation of the detector).

A simulated event

X~	
Event listing (summary)	PYTHIA Monte Carlo
I particle/jet KS KF orig p_x p_y p_z E	n > gluino gluino
1 !p+! 21 2212 0 0.000 0.000 7000.000 7000.000 2 !p+! 21 2212 0 0.000 0.000 7000.000 7000.000	$\begin{array}{c} pp \rightarrow gruno-gruno$
3 !9! 21 21 1 0.863 -0.323 1739.862 1739.862 4 !ubar! 21 -2 2 -0.621 -0.163 -777.415 777.415 5 !9! 21 21 3 -2.427 5.486 1487.857 1487.869 6 !9! 21 21 4 -62.910 63.357 -463.274 471.799 7 !~9! 21 1000021 0 314.363 544.843 498.897 979.192 8 !~9! 21 1000021 0 -379.700 -476.000 525.686 980.477 9 !~chi_1-! 21-1000024 7 130.058 112.247 129.860 263.141 10 !sbar! 21 -3 7 259.400 187.468 83.100 330.664 11 !c! 21 4 7 -79.403 242.409 283.026 381.016 12 !~cbi 20! 21 4 7 -79.403 242.409 283.026 381.016	397 pi+ 1 211 209 0.006 0.398 -308.296 308.297 0.140 398 gamma 1 22 211 0.407 0.087-1695.458 1695.458 0.000 399 gamma 1 22 211 0.113 -0.029 -314.822 314.822 0.000 400 (pi0) 11 111 212 0.021 0.122 -103.709 103.709 0.135 401 (pi0) 11 111 212 0.267 -0.058 -94.276 94.276 0.135 402 (pi0) 11 111 212 0.267 -0.052 -144.673 144.674 0.135 403 gamma 1 22 215 -1.581 2.473 3.306 4.421 0.000
12 1 1000023 8 -526.241 -80.371 115.712 585.351 13 19 21 5 8 -51.841 -294.077 389.853 491.098 14 1bbar! 21 -5 8 -0.597 -99.577 21.299 101.944 15 !"chi_10! 21 1000022 9 103.352 81.316 83.457 175.000 16 !s! 21 3 9 5.451 38.374 52.302 65.100 17 !cbar! 21 -4 9 20.839 -7.250 -5.938 22.899 18 "chi_10!</td 21 1000022 12 -136.266 -72.961 53.246 181.914 19 !nu_mu! 21 14 12 -78.263 -24.757 21.719 84.910 20 !nu_mubar! 21 -14 12 -107.801 16.901 38.226 115.620	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
21 gamma 1 22 4 2,636 1,357 0,125 2,967 22 ("chi_1-) 11-1000024 9 129,643 112,440 129,820 262,999 23 ("chi_20) 11 1000023 12 -322,330 -80,817 113,191 382,444 24 "chi_10 1 1000022 15 97,944 77,819 80,917 169,004 25 "chi_10 1 1000022 18 -136,266 -72,961 53,246 181,914 26 nu_mu 1 14 19 -78,263 -24,757 21,719 84,910 27 nu_mubar 1 -14 20 -107,801 16,901 38,226 115,620 28 (Delta++) 11 2224 2 0,222 0,012-2734,287 2734,287 :	41.3 K+1 321 220 $1,380$ $-0,652$ $-0,361$ $1,644$ $0,494$ 414 (pi0)11111 220 $1,078$ $-0,265$ $0,175$ $1,132$ $0,135$ 415 (K_S0)11 310 222 $1,841$ $0,111$ 0.894 $2,109$ $0,498$ 416 K+1 321 223 $0,307$ $0,107$ $0,252$ $0,642$ $0,494$ 417 pi-1 -2112 226 $1,335$ $1,641$ $2,078$ $3,111$ $0,940$ 418 nbar01 -2112 226 $1,335$ $1,641$ $2,078$ $3,111$ $0,940$ 419 (pi0)11111 226 $0,899$ $1,046$ $1,311$ $1,908$ $0,135$ 420 pi+1 211 227 $0,217$ $1,407$ $1,356$ $1,971$ $0,140$ 421 (pi0)11111 227 $1,207$ $2,336$ $2,767$ $3,820$ $0,135$ 422 n01 2112 228 $3,475$ $5,324$ $5,702$ $8,592$ $0,940$ 423 pi-1 -211 228 $1,856$ $2,606$ $2,808$ $4,259$ $0,140$ 424 gamma1 22 229 $-0,012$ $0,247$ $0,421$ $0,489$ $0,000$ 425 gamma1 22 229 $0,025$ $0,034$ $0,009$ $0,043$ $0,000$
•	426 pi+ 1 211 230 2,718 5,229 6,403 8,703 0,140 427 (pi0) 11 111 230 4,109 6,747 7,597 10,961 0,135 428 pi- 1 -211 231 0,551 1,233 1,945 2,372 0,140 429 (pi0) 11 111 231 0,645 1,141 0,922 1,608 0,135 430 gamma 1 22 232 -0,383 1,169 1,208 1,724 0,000 431 gamma 1 22 232 -0,201 0,070 0,060 0,221 0,000

•

Multivariate event selection

Suppose for each event we measure a set of numbers $\vec{\chi} = (\chi_1, \dots, \chi_n)$

 $x_1 = \text{jet } p_T$ $x_2 = \text{missing energy}$ $x_3 = \text{particle i.d. measure, ...}$

 $\vec{\chi}$ follows some *n*-dimensional joint probability density, which depends on the type of event produced, i.e., was it $pp \rightarrow t\bar{t}$, $pp \rightarrow \tilde{g}\tilde{g}$,...



E.g. hypotheses (class labels) $H_0, H_1, ...$ Often simply "signal", "background"

We want to separate (classify) the event types in a way that exploits the information carried in many variables.

Finding an optimal decision boundary

Maybe select events with "cuts":

$$x_i < c_i$$
$$x_j < c_j$$



Or maybe use some other type of decision boundary:



Goal of multivariate analysis is to do this in an "optimal" way.

Statistical Methods for LHC Physics

Test statistics

The decision boundary is a surface in the *n*-dimensional space of input variables, e.g., $y(\vec{x})$ =const.

We can treat the y(x) as a scalar test statistic or discriminating function, and try to define this function so that its distribution has the maximum possible separation between the event types:



Statistical Methods for LHC Physics

Constructing a test statistic

The Neyman-Pearson lemma states: to obtain the highest background rejection for a given signal efficiency (highest power for a given significance level), choose the acceptance region for signal such that

$$\frac{p(\vec{x}|s)}{p(\vec{x}|b)} > c$$

where c is a constant that determines the signal efficiency.

Equivalently, the optimal discriminating function is given by the likelihood ratio: $n(\vec{x}|_0)$

$$y(\vec{x}) = \frac{p(x|s)}{p(\vec{x}|b)}$$

N.B. any monotonic function of this is just as good.

Neyman-Pearson doesn't always help

The problem is that we usually don't have explicit formulae for the pdfs $p(\mathbf{x}|\mathbf{s}), p(\mathbf{x}|\mathbf{b})$, so for a given \mathbf{x} we can't evaluate the likelihood ratio.

Instead we have Monte Carlo models for signal and background processes, so we can produce simulated data:

"training data" events of known type



Naive try: enter each (s,b) event into an *n*-dimensional histogram, use e.g. M bins for each of the *n* dimensions, total of M^n cells.

n is potentially large \rightarrow prohibitively large number of cells to populate, can't generate enough training data.

General considerations

In all multivariate analyses we must consider e.g.

Choice of variables to use Functional form of decision boundary (type of classifier) Computational issues Trade-off between sensitivity and complexity Trade-off between statistical and systematic uncertainty

Our choices can depend on goals of the analysis, e.g.,

Event selection for further study Searches for new event types

Decision boundary flexibility

The decision boundary will be defined by some free parameters that we adjust using training data (of known type) to achieve the best separation between the event types.

Goal is to determine the boundary using a finite amount of training data so as to best separate between the event types for an unseen data sample.



Some "standard" multivariate methods

Place cuts on individual variables

Simple, intuitive, in general not optimal

Linear discriminant (e.g. Fisher)

Simple, optimal if the event types are Gaussian distributed with equal covariance, otherwise not optimal.

Probability Density Estimation based methods

Try to estimate $p(\mathbf{x}|\mathbf{s})$, $p(\mathbf{x}|\mathbf{b})$ then use $y(\mathbf{x}) = p(\mathbf{x}|\mathbf{s})/p(\mathbf{x}|\mathbf{b})$.

In principle best, difficult to estimate $p(\mathbf{x})$ for high dimension.

Neural networks

Can produce arbitrary decision boundary (in principle optimal), but can be difficult to train, result non-intuitive.

Decision trees

In a decision tree repeated cuts are made on a single variable until some stop criterion is reached.

The decision as to which variable is used is based on best achieved improvement in signal purity:

$$P = \frac{\sum_{\text{signal}} w_i}{\sum_{\text{signal}} w_i + \sum_{\text{background}} w_i}$$

where w_i is the weight of the *i*th event.

Iterate until stop criterion reached, based e.g. on purity and minimum number of events in a node.



Example by MiniBooNE experiment, B. Roe et al., NIM 543 (2005) 577

Decision trees (2)

The terminal nodes (leaves) are classified as signal or background depending on majority vote (or e.g. signal fraction greater than a specified threshold).

This classifies every point in input-variable space as either signal or background, a decision tree classifier, with the discriminant function

$f(\mathbf{x}) = 1$ if $\mathbf{x} \in \text{signal region}, -1$ otherwise

Decision trees tend to be very sensitive to statistical fluctuations in the training sample.

Methods such as boosting can be used to stabilize the tree.

Boosting

Boosting is a general method of creating a set of classifiers which can be combined to achieve a new classifier that is more stable and has a smaller error than any individual one.

Often applied to decision trees but, can be applied to any classifier.

Suppose we have a training sample T consisting of N events with

 x_1, \dots, x_N event data vectors (each x multivariate) y_1, \dots, y_N true class labels, +1 for signal, -1 for background w_1, \dots, w_N event weights

Now define a rule to create from this an ensemble of training samples T_1, T_2, \dots , derive a classifier from each and average them.

AdaBoost

A successful boosting algorithm is AdaBoost (Freund & Schapire, 1997). First initialize the training sample T_1 using the original

 $x_{1},...,x_{N} \quad \text{event data vectors}$ $y_{1},...,y_{N} \quad \text{true class labels (+1 or -1)}$ $w_{1}^{(1)},...,w_{N}^{(1)} \quad \text{event weights}$ with the weights equal and normalized such that $\sum_{i=1}^{N} w_{i}^{(1)} = 1.$

Train the classifier $f_1(\mathbf{x})$ (e.g. a decision tree) using the weights $\mathbf{w}^{(1)}$ so as to minimize the classification error rate,

$$\varepsilon_1 = \sum_{i=1}^N w_i^{(1)} I(y_i f_1(x_i) \le 0)$$
 ,

where I(X) = 1 if X is true and is zero otherwise.

Statistical Methods for LHC Physics

Updating the event weights (AdaBoost)

Assign a score to the *k*th classifier based on its error rate:

$$\alpha_k = \ln \frac{1 - \varepsilon_k}{\varepsilon_k}$$

Define the training sample for step k+1 from that of k by updating the event weights according to

$$W_{i}^{(k+1)} = W_{i}^{(k)} \frac{e^{-\alpha_{k}f_{k}(\mathbf{x}_{i})y_{i}/2}}{Z_{k}}$$
Normalize so that
 $i = \text{event index}$
 $k = \text{training sample index}$

$$\sum_{i} W_{i}^{(k+1)} = 1$$
Iterate K times, final classifier is
 $f(\mathbf{x}) = \sum_{k=1}^{K} \alpha_{k} f_{k}(\mathbf{x}, T_{k})$

BDT example from MiniBooNE

~200 input variables for each event (v interaction producing e, μ or π). Each individual tree is relatively weak, with a misclassification error rate ~ 0.4 - 0.45



B. Roe et al., NIM 543 (2005) 577

Monitoring overtraining

From MiniBooNE example

black = background red = signal



Boosted decision tree summary

Advantage of boosted decision tree is it can handle a large number of inputs. Those that provide little/no separation are rarely used as tree splitters are effectively ignored.

Easy to deal with inputs of mixed types (real, integer, categorical...).

If a tree has only a few leaves it is easy to visualize (but rarely use only a single tree).

There are a number of boosting algorithms, which differ primarily in the rule for updating the weights (ɛ-Boost, LogitBoost,...)

Other ways of combining weaker classifiers: Bagging (Boostrap-Aggregating), generates the ensemble of classifiers by random sampling with replacement from the full training sample.

Support Vector Machines

Support Vector Machines (SVMs) are an example of a kernel-based classifier, which exploits a nonlinear mapping of the input variables onto a higher dimensional feature space.

The SVM finds a linear decision boundary in the higher dimensional space.

But thanks to the "kernel trick" one does not every have to write down explicitly the feature space transformation.

Some references for kernel methods and SVMs:

The books mentioned on Monday

C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition,

research.microsoft.com/~cburges/papers/SVMTutorial.pdf

N. Cristianini and J.Shawe-Taylor. An Introduction to Support Vector Machines and other kernel-based learning methods. Cambridge University Press, 2000. The TMVA manual (!)

Linear SVMs

Consider a training data set consisting of

 x_1, \dots, x_N event data vectors y_1, \dots, y_N true class labels (+1 or -1)

Suppose the classes can be separated by a hyperplane defined by a normal vector *w* and scalar offset *b* (the "bias"): $y(x) = x \cdot w + b$

 $\begin{aligned} \mathbf{x}_{i} \cdot \mathbf{w} + b \ge +1 & \text{for all } y_{i} = +1 \\ \mathbf{x}_{i} \cdot \mathbf{w} + b \le -1 & \text{for all } y_{i} = -1 \\ \text{or equivalently} \\ y_{i}(\mathbf{x}_{i} \cdot \mathbf{w} + b) - 1 \ge 0 & \text{for all } i \end{aligned}$

Margin and support vectors

The distance between the hyperplanes defined by y(x) = +1 and y(x) = -1 is called the margin, which is:



If the training data are perfectly separated then this means there are no points inside the margin.

Suppose there are points on the margin (this is equivalent to defining the scale of w). These points are called support vectors.

Linear SVM classifier

We can define the classifier using

$$f(\mathbf{x}) = \operatorname{sign}(\mathbf{x} \cdot \mathbf{w} + b)$$

which is +1 for points on one side of the hyperplane and -1 on the other. The best classifier should have a large margin, so to maximize

margin =
$$\frac{2}{\|\boldsymbol{w}\|}$$

we can minimize $\|\boldsymbol{w}\|^2$ subject to the constraints

$$y_i(\mathbf{x_i} \cdot \mathbf{w} + b) - 1 \ge 0$$
 for all i

Lagrangian formulation

This constrained minimization problem can be reformulated using a Lagrangian

$$L = \frac{1}{2} \| \mathbf{w} \|^{2} - \sum_{i=1}^{N} \alpha_{i} (y_{i} (\mathbf{x}_{i} \cdot \mathbf{w} + b) - 1)$$

positive Lagrange multipliers α_{i}

We need to minimize L with respect to w and b and maximize with respect to α_i .

There is an α_i for every training point. Those that lie on the margin (the support vectors) have $\alpha_i > 0$, all others have $\alpha_i = 0$. The solution can be written (sum only contains

$$\mathbf{w} = \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i}$$

(sum only contains support vectors)

Statistical Methods for LHC Physics

Dual formulation

The classifier function is thus

$$f(\mathbf{x}) = \operatorname{sign}(\mathbf{x} \cdot \mathbf{w} + b) = \operatorname{sign}\left(\sum_{i} \alpha_{i} y_{i} \mathbf{x} \cdot \mathbf{x}_{i} + b\right)$$

It can be shown that one finds the same solution a by minimizing the dual Lagrangian

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

So this means that both the classifier function and the Lagrangian only involve dot products of vectors in the input variable space.

Nonseparable data

If the training data points cannot be separated by a hyperplane, one can redefine the constraints by adding slack variables ξ_i :

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) + \xi_i - 1 \ge 0$$
 with $\xi_i \ge 0$ for all i

Thus the training point x_i is allowed to be up to a distance ξ_i on the wrong side of the boundary, and $\xi_i = 0$ at or on the right side of the boundary.

For an error to occur we have $\xi_i > 1$, so

 $\sum \xi_i$



is an upper bound on the number of training errors.

Cost function for nonseparable case

To limit the magnitudes of the ξ_i we can define the error function that we minimize to determine *w* to be

$$E(\boldsymbol{w}) = \frac{1}{2} \|\boldsymbol{w}\|^2 + C \left(\sum_i \xi_i\right)^k$$

where *C* is a cost parameter we must choose that limits the amount of misclassification. It turns out that for k=1 or 2 this is a quadratic programming problem and furthermore for k=1 it corresponds to minimizing the same dual Lagrangian

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

where the constraints on the α_i become $0 \le \alpha_i \le C$.

Nonlinear SVM

So far we have only reformulated a way to determine a linear classifier, which we know is useful only in limited circumstances.

But the important extension to nonlinear classifiers comes from first transforming the input variables to feature space:

$$\vec{\varphi}(\mathbf{x}) = (\varphi_1(\mathbf{x}), \dots, \varphi_m(\mathbf{x}))$$

These will behave just as our new "input variables". Everything about the mathematical formulation of the SVM will look the same as before except with $\phi(x)$ appearing in the place of x.

Only dot products

Recall the SVM problem was formulated entirely in terms of dot products of the input variables, e.g., the classifier is

$$f(\mathbf{x}) = \operatorname{sign}\left(\sum_{i} \alpha_{i} y_{i} \mathbf{x} \cdot \mathbf{x}_{i} + b\right)$$

so in the feature space this becomes

$$f(\mathbf{x}) = \operatorname{sign}\left(\sum_{i} \alpha_{i} y_{i} \vec{\varphi}(\mathbf{x}) \cdot \vec{\varphi}(\mathbf{x}_{i}) + b\right)$$
The Kernel trick

How do the dot products help? It turns on that a broad class of kernel functions can be written in the form:

$$K(\mathbf{x}, \mathbf{x}') = \vec{\varphi}(\mathbf{x}) \cdot \vec{\varphi}(\mathbf{x}')$$

Functions having this property must satisfy Mercer's condition

$$\int K(\mathbf{x}, \mathbf{x}')g(\mathbf{x})g(\mathbf{x}')d\mathbf{x}d\mathbf{x}' \ge 0$$

for any function g where $\int g^2(x) dx$ is finite.

So we don't even need to find explicitly the feature space transformation $\phi(x)$, we only need a kernel.

Finding kernels

There are a number of techniques for finding kernels, e.g., constructing new ones from known ones according to certain rules (cf. Bishop Ch 6).

Frequently used kernels to construct classifiers are e.g.

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + \theta)^p$$
 polynomial

$$K(\mathbf{x}, \mathbf{x'}) = \exp\left(\frac{-\|\mathbf{x}-\mathbf{x'}\|^2}{2\sigma^2}\right)$$
 Gaussian

$$K(\mathbf{x}, \mathbf{x'}) = \tanh(\kappa(\mathbf{x} \cdot \mathbf{x'}) + \theta)$$
 sigmoidal

Using an SVM

To use an SVM the user must as a minimum choose

a kernel function (e.g. Gaussian)

any free parameters in the kernel (e.g. the σ of the Gaussian) the cost parameter *C* (plays role of regularization parameter)

The training is relatively straightforward because, in contrast to neural networks, the function to be minimized has a single global minimum.

Furthermore evaluating the classifier only requires that one retain and sum over the support vectors, a relatively small number of points.

SVM in HEP

SVMs are very popular in the Machine Learning community but have yet to find wide application in HEP. Here is an early example from a CDF top quark anlaysis (A. Vaiciulis, contribution to PHYSTAT02).



Multivariate analysis discussion

For all methods, need to check:

Sensitivity to statistically unimportant variables (best to drop those that don't provide discrimination);

Level of smoothness in decision boundary (sensitivity to over-training)

Given the test variable, next step is e.g., select *n* events and estimate a cross section of signal: $\hat{\sigma}_s = (n-b)/\varepsilon_s L$

Now need to estimate systematic error...

If e.g. training (MC) data \neq Nature, test variable is not optimal, but not necessarily biased.

But our estimates of background *b* and efficiencies would then be biased if based on MC. (True also for 'simple cuts'.)

Multivariate analysis discussion (2)

But in a cut-based analysis it may be easier to avoid regions where untested features of MC are strongly influencing the decision boundary.

Look at control samples to test joint distributions of inputs.

Try to estimate backgrounds directly from the data (sidebands).

The purpose of the statistical test is often to select objects for further study and then measure their properties.

Need to avoid input variables that are correlated with the properties of the selected objects that you want to study. (Not always easy; correlations may be poorly known.)

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

Comparing multivariate methods (TMVA)



Choose the best one!

Bayesian vs. frequentist methods

Two schools of statistics use different interpretations of probability:

I. Relative frequency (frequentist statistics):

$$P(A) = \lim_{n \to \infty} rac{\operatorname{times outcome is} A}{n}$$

II. Subjective probability (Bayesian statistics):

P(A) = degree of belief that A is true

In particle physics frequency interpretation most used, but subjective probability can be more natural for non-repeatable phenomena: systematic uncertainties, probability that Higgs boson exists... Frequentist Statistics – general philosophy In frequentist statistics, probabilities are associated only with the data, i.e., outcomes of repeatable observations.

Probability = limiting frequency

Probabilities such as

P (Higgs boson exists), *P* (0.117 < α_s < 0.121),

etc. are either 0 or 1, but we don't know which. The tools of frequentist statistics tell us what to expect, under the assumption of certain probabilities, about hypothetical repeated observations.

The preferred theories (models, hypotheses, ...) are those for which our observations would be considered 'usual'.

Bayesian Statistics – general philosophy

In Bayesian statistics, interpretation of probability extended to degree of belief (subjective probability). Use this for hypotheses:

probability of the data assuming hypothesis H (the likelihood) prior probability, i.e., before seeing the data $P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H) dH}$ posterior probability, i.e., after seeing the data over all possible hypotheses

Bayesian methods can provide more natural treatment of non-repeatable phenomena:

systematic uncertainties, probability that Higgs boson exists,...

No golden rule for priors ("if-then" character of Bayes' thm.)

Statistical vs. systematic errors Statistical errors:

How much would the result fluctuate upon repetition of the measurement?

Implies some set of assumptions to define probability of outcome of the measurement.

Systematic errors:

What is the uncertainty in my result due to uncertainty in my assumptions, e.g.,

model (theoretical) uncertainty; modelling of measurement apparatus.

Usually taken to mean the sources of error do not vary upon repetition of the measurement. Often result from uncertain value of, e.g., calibration constants, efficiencies, etc.

Systematic errors and nuisance parameters

Response of measurement apparatus is never modelled perfectly:



Model can be made to approximate better the truth by including more free parameters.

systematic uncertainty \leftrightarrow nuisance parameters

Example: fitting a straight line

Data: $(x_i, y_i, \sigma_i), i = 1, ..., n$.

Model: measured y_i independent, Gaussian: $y_i \sim N(\mu(x_i), \sigma_i^2)$

 $\mu(x;\theta_0,\theta_1) = \theta_0 + \theta_1 x ,$ 1.8 data • fit -1.6 assume x_i and σ_i known. 1.4 yGoal: estimate θ_{a} 1.2 (don't care about θ_i). 1 0.8 1.5 0.5 1 0 2

x

Frequentist approach

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

$$\chi^2(\theta_0, \theta_1) = -2\ln L(\theta_0, \theta_1) + \text{const} = \sum_{i=1}^n \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}.$$

Standard deviations from tangent lines to contour

 $\chi^2 = \chi^2_{\rm min} + 1 \; .$

Correlation between $\hat{\theta}_0, \hat{\theta}_1$ causes errors to increase.



Frequentist case with a measurement t_1 of θ_1

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

The information on θ_1

improves accuracy of $\hat{\theta}_0$.



Bayesian method

We need to associate prior probabilities with θ_0 and θ_1 , e.g.,

$$\pi(\theta_0, \theta_1) = \pi_0(\theta_0) \pi_1(\theta_1) \quad \text{reflects 'prior ignorance', in any} \\ \pi_0(\theta_0) = \text{const.} \quad \text{case much broader than } L(\theta_0) \\ \pi_1(\theta_1) = \frac{1}{\sqrt{2\pi\sigma_{t_1}}} e^{-(\theta_1 - t_1)^2/2\sigma_{t_1}^2} \quad \leftarrow \text{based on previous} \\ \text{measurement} \end{cases}$$

Putting this into Bayes' theorem gives:

$$p(\theta_0, \theta_1 | \vec{y}) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} e^{-(y_i - \mu(x_i; \theta_0, \theta_1))^2 / 2\sigma_i^2} \pi_0 \frac{1}{\sqrt{2\pi\sigma_{t_1}}} e^{-(\theta_1 - t_1)^2 / 2\sigma_{t_1}^2}$$

$$posterior \propto likelihood \times prior$$

Bayesian method (continued)

We then integrate (marginalize) $p(\theta_0, \theta_1 | x)$ to find $p(\theta_0 | x)$:

$$p(\theta_0|x) = \int p(\theta_0, \theta_1|x) d\theta_1$$

In this example we can do the integral (rare). We find

$$p(\theta_0|x) = \frac{1}{\sqrt{2\pi}\sigma_{\theta_0}} e^{-(\theta_0 - \hat{\theta}_0)^2 / 2\sigma_{\theta_0}^2} \text{ with}$$
$$\hat{\theta}_0 = \text{ same as ML estimator}$$
$$\sigma_{\theta_0} = \sigma_{\hat{\theta}_0} \text{ (same as before)}$$

Ability to marginalize over nuisance parameters is an important feature of Bayesian statistics.

Digression: marginalization with MCMC Bayesian computations involve integrals like

$$p(\theta_0|x) = \int p(\theta_0, \theta_1|x) d\theta_1$$

often high dimensionality and impossible in closed form, also impossible with 'normal' acceptance-rejection Monte Carlo.

Markov Chain Monte Carlo (MCMC) has revolutionized Bayesian computation.

Google for 'MCMC', 'Metropolis', 'Bayesian computation', ...

MCMC generates correlated sequence of random numbers: cannot use for many applications, e.g., detector MC; effective stat. error greater than \sqrt{n} .

Basic idea: sample multidimensional $\vec{\theta}$, look, e.g., only at distribution of parameters of interest.

Example: posterior pdf from MCMC Sample the posterior pdf from previous example with MCMC:



Summarize pdf of parameter of interest with, e.g., mean, median, standard deviation, etc.

Although numerical values of answer here same as in frequentist case, interpretation is different (sometimes unimportant?)

Bayesian method with vague prior

Suppose we don't have a previous measurement of θ_1 but rather some vague information, e.g., a theorist tells us:

 $\theta_1 \ge 0$ (essentially certain);

 θ_1 should have order of magnitude less than 0.1 'or so'.

Under pressure, the theorist sketches the following prior: $\pi_1(\theta_1) = \frac{1}{\tau} e^{-\theta_1/\tau}, \quad \theta_1 \ge 0, \quad \tau = 0.1.$

From this we will obtain posterior probabilities for θ_0 (next slide).

We do not need to get the theorist to 'commit' to this prior; final result has 'if-then' character.

Sensitivity to prior

Vary $\pi(\theta)$ to explore how extreme your prior beliefs would have to be to justify various conclusions (sensitivity analysis).

Try exponential with different mean values...

Try different functional forms...



Outlook for Bayesian methods in HEP Bayesian methods allow (indeed require) prior information about the parameters being fitted.

This type of prior information can be difficult to incorporate into a frequentist analysis

This will be particularly relevant when estimating uncertainties on predictions of LHC observables that may stem from theoretical uncertainties, parton densities based on inconsistent data, etc.

Prior ignorance is not well defined. If that's what you've got, don't expect Bayesian methods to provide a unique solution. Try a reasonable variation of priors -- if that yields large variations in the posterior, you don't have much information coming in from the data.

You do not have to be exclusively a Bayesian or a Frequentist Use the right tool for the right job

Outlook for data analysis at the LHC

Recent developments from Machine Learning provide some new tools for event classification with a number of advantages over those methods commonly used in HEP, e.g.,

> Boosted decision trees Support Vector Machines

Bayesian methods can allow for a more natural treatment of non-repeatable phenomena, e.g., model uncertainties. MCMC can be used to marginalize posterior probabilities

Software for these methods now much more easily available, expect rapid development as the LHC begins to produce real results.

Quotes I like

"Alles sollte so einfach wie möglich sein, aber nicht einfacher."

– A. Einstein

"If you believe in something you don't understand, you suffer,..." – Stevie Wonder

Extra slides

Probability – quick review

Frequentist (*A* = outcome of repeatable observation):

$$P(A) = \lim_{n \to \infty} \frac{\text{outcome is } A}{n}$$

Subjective (*A* = hypothesis):

$$P(A) =$$
degree of belief that A is true

Conditional probability:
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Bayes' theorem:
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)}$$

Bayes optimal analysis

From Bayes' theorem we can compute the posterior odds:



which is proportional to the likelihood ratio.

So placing a cut on the likelihood ratio is equivalent to ensuring a minimum posterior odds ratio for the selected sample.

Classification viewed as a statistical test

Probability to reject H_0 if it is true (type I error): $\alpha = \int_{R_1} f(y|H_0) dy$ $\alpha =$ significance level, size of test, false discovery rate

Probability to accept H_0 if H_1 is true (type II error): $\beta = \int_{R_0} f(y|H_1) dy$ $1 - \beta$ = power of test with respect to H_1

Equivalently if e.g. $H_0 =$ background, $H_1 =$ signal, use efficiencies:

$$\varepsilon_{s} = \int_{R_{1}} f(y|H_{1}) dy = 1 - \beta = \text{power}$$
 $\varepsilon_{b} = \int_{R_{0}} f(y|H_{0}) dy = 1 - \alpha$

Purity / misclassification rate

Consider the probability that an event assigned to a certain category is classified correctly (i.e., the purity).

Use Bayes' theorem:



posterior probability

N.B. purity depends on the prior probabilities for an event to be signal or background (~s, b cross sections).

Imperfect pdf estimation

What if the approximation we use (e.g., parametric form, assumption of variable independence, etc.) to estimate $p(\mathbf{x})$ is wrong?

If we use poor estimates to construct the test variable

$$y(\vec{x}) = \frac{\hat{p}(\vec{x}|H_0)}{\hat{p}(\vec{x}|H_1)}$$

then the discrimination between the event classes will not be optimal.

But can this cause us e.g. to make a false discovery?

Even if the estimate of p(x) used in the discriminating variable are imperfect, this will not affect the accuracy of the distributions $f(y|H_0)$, $f(y|H_1)$; this only depends on the reliability of the training data.



Discovery = number of events found in search region incompatible with background-only hypothesis. Maximize the probability of this happening by setting y_{cut} for maximum s/\sqrt{b} (roughly true).

Statistical Methods for LHC Physics

Controlling false discovery

So for a reliable discovery what we depend on is an accurate estimate of the expected number of background events, and this accuracy only depends on the quality of the training data; works for any function y(x).

But we do not blindly rely on the MC model for the training data for background; we need to test it by comparing to real data in control samples where no signal is expected.

The ability to perform these tests will depend on on the complexity of the analysis methods.

MCMC basics: Metropolis-Hastings algorithm Goal: given an *n*-dimensional pdf $p(\vec{\theta})$, generate a sequence of points $\vec{\theta}_1, \vec{\theta}_2, \vec{\theta}_3, \dots$ Proposal density $q(\vec{\theta}; \vec{\theta}_0)$ 1) Start at some point $\dot{\theta_0}$ 2) Generate $\vec{\theta} \sim q(\vec{\theta}; \vec{\theta}_0)$ e.g. Gaussian centred about θ_0 3) Form Hastings test ratio $\alpha = \min \left[1, \frac{p(\vec{\theta})q(\vec{\theta}_0; \vec{\theta})}{p(\vec{\theta}_0)q(\vec{\theta}_0; \vec{\theta}_0)} \right]$

- 4) Generate $u \sim \text{Uniform}[0, 1]$
- 5) If $u \le \alpha$, $\vec{\theta}_1 = \vec{\theta}$, \leftarrow move to proposed point

else
$$\vec{\theta}_1 = \vec{\theta}_0$$
 \leftarrow old point repeated

Iterate (\mathbf{b})

Statistical Methods for LHC Physics

Metropolis-Hastings (continued)

This rule produces a *correlated* sequence of points (note how each new point depends on the previous one).

For our purposes this correlation is not fatal, but statistical errors larger than naive \sqrt{n} .

The proposal density can be (almost) anything, but choose so as to minimize autocorrelation. Often take proposal density symmetric: $q(\vec{\theta}; \vec{\theta}_0) = q(\vec{\theta}_0; \vec{\theta})$ Test ratio is (*Metropolis*-Hastings): $\alpha = \min \left[1, \frac{p(\vec{\theta})}{p(\vec{\theta}_0)}\right]$

I.e. if the proposed step is to a point of higher $p(\vec{\theta})$, take it; if not, only take the step with probability $p(\vec{\theta})/p(\vec{\theta}_0)$. If proposed step rejected, hop in place.

Metropolis-Hastings caveats

Actually one can only prove that the sequence of points follows the desired pdf in the limit where it runs forever.

There may be a "burn-in" period where the sequence does not initially follow $p(\vec{\theta})$.

Unfortunately there are few useful theorems to tell us when the sequence has converged.

Look at trace plots, autocorrelation.

Check result with different proposal density.
A more general fit (symbolic) Given measurements: $y_i \pm \sigma_i^{\text{stat}} \pm \sigma_i^{\text{sys}}$, i = 1, ..., n, and (usually) covariances: V_{ij}^{stat} , V_{ij}^{sys} . Predicted value: $\mu(x_i; \theta)$, expectation value $E[y_i] = \mu(x_i; \theta) + b_i$ control variable parameters bias

Often take: $V_{ij} = V_{ij}^{\text{stat}} + V_{ij}^{\text{sys}}$

Minimize $\chi^2(\theta) = (\vec{y} - \vec{\mu}(\theta))^T V^{-1} (\vec{y} - \vec{\mu}(\theta))$

Equivalent to maximizing $L(\theta) \approx e^{-\chi^2/2}$, i.e., least squares same as maximum likelihood using a Gaussian likelihood function.



To get desired probability for θ , integrate (marginalize) over b:

$$p(\theta|\vec{y}) = \int p(\theta, \vec{b}|\vec{y}) d\vec{b}$$

→ Posterior is Gaussian with mode same as least squares estimator, σ_{θ} same as from $\chi^2 = \chi^2_{\min} + 1$. (Back where we started!) The error on the error

Some systematic errors are well determined

Error from finite Monte Carlo sample

Some are less obvious

Do analysis in *n* 'equally valid' ways and extract systematic error from 'spread' in results.

Some are educated guesses

Guess possible size of missing terms in perturbation series; vary renormalization scale $(\mu/2 < Q < 2\mu)$

Can we incorporate the 'error on the error'?

(cf. G. D'Agostini 1999; Dose & von der Linden 1999)

A prior for bias $\pi_b(b)$ with longer tails



Statistical Methods for LHC Physics



with mode and standard deviation: $\sigma_s = 0.5$: $\hat{\mu} = 1.000 \pm 0.072$

Simple test with inconsistent data

Case #2: there is an outlier

Posterior $p(\mu|y)$:



\rightarrow Bayesian fit less sensitive to outlier.

 \rightarrow Error now connected to goodness-of-fit.

Goodness-of-fit vs. size of error

In LS fit, value of minimized χ^2 does not affect size of error on fitted parameter.

In Bayesian analysis with non-Gaussian prior for systematics, a high χ^2 corresponds to a larger error (and vice versa).

