Statistical Methods in Particle Physics

Lecture 8 December 3, 2012

Silvia Masciocchi, GSI Darmstadt s.masciocchi@gsi.de

Winter Semester 2012 / 13





Today we talk of **STATISTICAL TESTS**

 \rightarrow a way to find criteria to select candidate events (or particles, ...) for further analysis (e.g. signal vs background)

The goal of test statistics is to make a statement about how well the observed data stand in agreement with given predicted probabilities, i.e. with *hypotheses*



A *hypothesis H* specifies the probability for the data i.e., the outcome of the observation, here symbolically "x" We can write:

$\mathbf{x} \sim \mathbf{f}(\mathbf{x} | \mathbf{H})$

x can be uni-/multivariate, continuous or discrete x could represent for example the observation of a single particle, a single event, or an entire "experiment"

Possible values of x form the **sample space S** (or "data space")

The probability for x given H is also called the **likelihood of the** hypothesis, written L(x|H)

Use the ALICE Time Projection Chamber to identify the particle species: electron, muon, pion, kaon, proton, deuteron

"x" = particle momentum (p), specific energy loss in TPC (dE/dx) (and more)



Example: I want to select electrons (hypothesis H_1) from all other particles (hypothesis H_0)

In Bayesian approach: Can add prior hypotheses on the relative particle abundances (e.g. you see that pions are many more!)



Goal is to make some statement based on the observed data x, as to the validity of the possible hypotheses.

A test of hypothesis H_0 is defined by specifying a **critical region** W (also called **rejection** region) of the data space S, such that there is no more than some (small) probability α , assuming H_0 is correct, to observe the data there:

 $\mathsf{P}(\mathsf{x} \in \mathsf{W} | \mathsf{H}_0) \leq \alpha$

If x is observed there, reject H_0 .

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

The complementary region is called acceptance region.



Looking for electrons: null hypothesis H₀ is to be a hadron



s.masciocchi@gsi.de

Statistical Methods, Lecture 8, December 3, 2012



Event by event, we need to decide whether to take it as signal or as background





Rejecting the hypothesis H_0 when it is true is a Type-I error. The maximum probability for this is the **size of the test**:

$$\mathsf{P}(\mathsf{x} \in \mathsf{W} | \mathsf{H}_0) \leq \alpha$$

But we might also accept H_0 when it is false and an alternative H_1 is true. This is called Type-II error, and occurs with probability:

$$\mathsf{P}(\mathsf{x} \in \mathsf{S} - \mathsf{W} | \mathsf{H}_1) = \beta$$

One minus this is called the power of the test with respect to the alternative hypothesis H_1 :

Power = $1 - \beta$







We have a data sample with two kinds of events, corresponding to hypotheses H_0 (background) and H_1 (signal).

We want to select those of type H_1 .

Each event is a point in \vec{x} space (n dimensions).

What 'decision boundary' should we use to accept/reject events as belonging to event types H_0 or H_1 ?

One possibility is to select events with several 'cuts': e.g.

$$x_i < C_i$$
$$x_j < C_j$$





But we can also use some other sort of decision boundary !!



How can we formalize this to choose the boundary in an 'optimal' way?

In addition:

- \vec{x} is the result of the measurements, n can be large
- \vec{x} follows some joint pdf in an n-dimensional space

Usually it is awkward to work with multidimensions!

At first we try to construct a test statistic of lower dimension (e.g. scalar):

- Compactify the data
- Try not to loose the ability to discriminate between hypotheses



Multivariate analysis (MVA)

- Map the n-dimensional space of the observable variables ("feature" space of our measurements) to one dimensional output
 - $\mathbb{IR}^n \rightarrow \mathbb{IR}$

 $(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow \mathbf{t}(\mathbf{\vec{x}})$

Test statistic

- There are model classes for this
 - Various types: linear, non-linear, flexible, less flexible
- We can use previous knowledge, "known" or "previously solved" cases
- The resulting class (description) should have good generalization properties

Often associated with the term of "machine learning"





The decision boundary can be defined by an equation of the form:

 $t(x_1, ..., x_n) = constant = t_{cut}$

where $t(x_1, ..., x_n)$ is a scalar test statistic

We can work out the pdf's:

 $g(t|H_0), g(t|H_1)$

Decision boundary is now a single 'cut' on t, which divides the space into the critical (rejection region) and the acceptance region.

This defines a TEST: if the data fall in the critical region, we reject H_0



The probability to reject background hypothesis for a background event (background efficiency) is:



Suppose only one type of background b.

Overall fractions of signal and background events are π_s and π_b (prior probabilities).

Suppose we select signal events with t>t $_{cut}$. What is the PURITY of the selected sample?

PURITY means the probability to be signal given that the event was accepted. Using Bayes' theorem we find:

$$\mathsf{P}(\mathbf{s}|\mathbf{t} > \mathbf{t}_{\mathsf{cut}}) = \frac{\mathsf{P}(\mathbf{t} > \mathbf{t}_{\mathsf{cut}}|\mathbf{s})\pi_{\mathsf{s}}}{\mathsf{P}(\mathbf{t} > \mathbf{t}_{\mathsf{cut}}|\mathbf{s})\pi_{\mathsf{s}} + \mathsf{P}(\mathbf{t} > \mathbf{t}_{\mathsf{cut}}|\mathbf{b})\pi_{\mathsf{b}}} = \frac{\epsilon_{\mathsf{s}}\pi_{\mathsf{s}}}{\epsilon_{\mathsf{s}}\pi_{\mathsf{s}} + \epsilon_{\mathsf{b}}\pi_{\mathsf{b}}}$$

 \rightarrow the purity depends on the prior probabilities as well as on the signal and background efficiencies !!



How can we choose a test's critical region in an "optimal way"?

Neyman-Pearson lemma states:

To get the highest power for a given significance level (or highest purity for a given efficiency) in a test of H_0 (background) versus H_1 (signal), the critical region should have:

$$\frac{P(x|H_1)}{P(x|H_0)} > c$$

inside the region, and \leq c outside, where c is a constant which determines the power.

Equivalently, optimal scalar test statistics is:

$$\mathbf{t}(\mathbf{x}) = \frac{\mathbf{P}(\mathbf{x}|\mathbf{H}_1)}{\mathbf{P}(\mathbf{x}|\mathbf{H}_0)}$$

Likelihood ratio

Neyman-Pearson lemma





y(x) is the discriminating function given by your estimator (i.e. the likelihood ratio)
 varying y(x)>"cut" moves the working point (efficiency and purity) along the ROC curve

- where to choose your working point? → need to know prior probabilities (abundances)
 - measurement of signal cross section:
 discovery of a signal (typically: S<<B):
 - precision measurement:
 - trigger selection:

maximum of S/ $\sqrt{(S+B)}$ or equiv. $\sqrt{(\epsilon \cdot p)}$ maximum of S/ $\sqrt{(B)}$ high purity (p) high efficiency (ϵ)



Usually we do **NOT** have explicit formulae for the pdfs $P(x|H_0), P(x|H_1)$

What we usually have are Monte Carlo models for signal and background processes, so we can produce simulated data, and enter each event into an n-dimensional histogram.

But then we need M bins for each of the n dimensions

\rightarrow total of Mⁿ cells !!

If n is large, then we end up with a prohibitively large number of cells to populate with Monte Carlo data !!!

Compromise solution:

Make Ansatz for form of the test statistic t(x) with fewer parameters; determine them (e.g. using MC) to give best discrimination between signal and background!





Distinguish between 2 processes:

 $H_0: e^+e^- \rightarrow WW \rightarrow hadrons (usually 4 jets)$

 $H_1: e^+e^- \rightarrow q\overline{q} \rightarrow hadrons (usually 2 jets)$

For each event we measure \vec{X} (n. of hadrons, their momenta, jets, missing energy, angles between jets, etc etc)

According to Neyman-Pearson, to select WW's we should cut on

$$\mathbf{t}(\vec{\mathbf{x}}) = \frac{\mathbf{f}(\vec{\mathbf{x}}|\mathbf{H}_0)}{\mathbf{f}(\vec{\mathbf{x}}|\mathbf{H}_1)}$$

But we do not know entirely these pdf's !!! Partly help with MC, partly simplify / transform the description of t



$$\mathbf{t}(\mathbf{\vec{x}}) = \sum_{i=1}^{H} \mathbf{a}_{i} \mathbf{x}_{i} = \mathbf{\vec{a}}^{\mathsf{T}} \mathbf{\vec{x}}$$

Choose the parameters $a_1, ..., a_n$ so that the pdf's g(t|s), g(t|b)have maximum SEPARATION:

We want large distance between the mean values and small widths



Fisher: maximize
$$J(\vec{a}) = \frac{(\tau_s - \tau_b)^2}{\Sigma_s^2 + \Sigma_b^2}$$

Hypotheses: k = 0,1Measurement: \vec{x} i, j = 1,, n (components) Means and variances for the x_i:

$$(\mu_{k})_{i} = \int x_{i} f(\vec{x} | H_{k}) d\vec{x}$$
$$(V_{k})_{ij} = \int (x - \mu_{k})_{i} (x - \mu_{k})_{j} f(\vec{x} | H_{k}) d\vec{x}$$

In terms of mean and variance of $t(\vec{x})$ this becomes:

$$\tau_{\mathbf{k}} = \int \mathbf{t}(\mathbf{\vec{x}}) \, \mathbf{f}(\mathbf{\vec{x}} | \mathbf{H}_{\mathbf{k}}) \, \mathbf{d}\mathbf{\vec{x}} = \mathbf{\vec{a}}^{\mathsf{T}} \mathbf{\vec{\mu}}_{\mathbf{k}}$$
$$\Sigma_{\mathbf{k}}^{2} = \int (\mathbf{t}(\mathbf{\vec{x}}) - \tau_{\mathbf{k}})^{2} \, \mathbf{f}(\mathbf{\vec{x}} | \mathbf{H}_{\mathbf{k}}) \, \mathbf{d}\mathbf{\vec{x}} = \mathbf{\vec{a}}^{\mathsf{T}} \, \mathbf{V}_{\mathbf{k}} \, \mathbf{\vec{a}}$$

Reasoning to come to Fisher's statement





The numerator of $J(\vec{a})$ is:

$$(\tau_0 - \tau_1)^2 = \sum_{i,j=1}^n a_i a_j (\mu_0 - \mu_1)_i (\mu_0 - \mu_1)_j$$
$$= \sum_{i,j=1}^n a_i a_j B_{ij} = \vec{a}^T B \vec{a}$$

The denominator:

$$\Sigma_0^2 + \Sigma_1^2 = \sum_{i,j=1}^n a_i a_j (V_0 + V_1)_{ij} = \vec{a}^T W \vec{a}$$

Maximize $J(\vec{a}) = \frac{\vec{a}^T B \vec{a}}{\vec{a}^T W \vec{a}} = \frac{\text{separation between classes}}{\text{separation within classes}}$



Setting: $\frac{\partial J}{\partial a_i} = 0$ gives Fisher's linear discriminant function:

$$t(\vec{x}) = \vec{a}^{T}\vec{x}$$
, with $\vec{a} \propto W^{-1}(\vec{\mu_{0}} - \vec{\mu_{1}})$



Corresponds to a linear decision boundary



Fisher linear discriminant analysis determines a canonical direction for which the data is most separated when projected on a line in this direction. The solid gray line shows the canonical direction.



Another illustration - 2





The squares are projected points on a line inclined at the angle θ with respect to the origin. When θ is adjusted so the projected points are aligned with the gray line, the points are maximally separated in the sense that the ratio of between-classes variances to within-classes variance is maximized.

We obtain equivalent separation between hypotheses if we multiply the a_i by a common scale factor and add an arbitrary offset a_0 :

$$t(\vec{x}) = a_0 + \sum_{i=1}^n a_i x_i$$

Thus we can fix the mean values under the hypotheses H_0 and H_1 to arbitrary values as 0 and 1.

Then maximizing

$$J(\vec{a}) = \frac{(\tau_{s} - \tau_{b})^{2}}{\Sigma_{s}^{2} + \Sigma_{b}^{2}}$$

Is equivalent to minimizing

$$\Sigma_0^2 + \Sigma_1^2 = E_0 [(t - \tau_0)^2] + E_1 [(t - \tau_1)^2]$$

A type of least squares principle !!!



Many new (and some old) methods:

- Fisher discriminant *(linear decision boundary)*
- Neural networks
- Kernel density methods
- Support Vector Machines
- Decision trees:
 - Boosting
 - Bagging
- Toolkit for Multivariate Data Analysis: TMVA
 - Framework for "all" MVA-techniques, available in ROOT

Linear decision boundaries



A linear decision boundary is only optimal when both classes follow multivariate Gaussians with equal covariances and different means





For other cases, a linear boundary is almost useless

We can try to find a transformation $x_{1,}..., x_n \rightarrow \phi_1(\vec{x}), ..., \phi_m(\vec{x})$ So that the transformed "feature space" variables can be separated better by a linear boundary:



The optimal decision boundary may not be a hyperplane \rightarrow non linear test statistic !!

Many methods of multivariate statistical methods:

- Neural networks
- Support vector machines
- Kernel density methods
- Decision trees
- TMVA





If we want to go to the "arbitrary" non-linear decision boundaries, t(x) needs to be constructed in "any" non-linear fashion

$$t(\vec{x}) = \sum_{i}^{M} (w_{i} h_{i}(\vec{x}))$$

- Think of $h_i(x)$ as a set of "basis" functions
- If h(x) is sufficiently general (i.e. non linear), a linear combination of "enough" basis functions (M) should allow to describe any possible discriminating function t(x)

We take the $h_i(x)$ to be such that:

Where:

- a₀: threshold
- A(x): activation function





We take the $h_i(x)$ to be such that:

$$\mathbf{t}(\vec{\mathbf{x}}) = \sum_{i}^{M} \mathbf{w}_{0i} \mathbf{A}(\mathbf{w}_{i0} + \sum_{j=1}^{n} \mathbf{w}_{ij} \mathbf{x}_{j})$$

 $t(\vec{x}) \quad \text{is} \quad$

- A linear combination of
 - non-linear functions of
 - linear combination of
 - the input data









- Nodes in hidden layer represent the "activation functions" whose arguments are linear combinations of input variables → non linear response to the input
- The output is a linear combination of the output of the activation functions at the internal nodes
- Input to the layers from preceding nodes only \rightarrow feed forward network (no backward loops)
- It is straightforward to extend this to "several" input layers

Multilayer perceptron (MLP)





Nodes \rightarrow neurons

Links (weights) \rightarrow synapses

→ Neural network: try to simulate reactions of a brain to certain stimulus (input data)

Use training events to adjust the weights such that:

- $t(x) \rightarrow 0$ for background events
- $t(x) \rightarrow 1$ for signal events



t(x) is a very "wiggly" function with many local minima. A global overall fit in the many parameters is possible but not the most efficient method to train neutral networks ...



Use smarter methods instead of a global overall fit in the many parameters:

- Back propagation: learn from experience, gradually adjust your perception to match reality
- Online learning: learn event by event and not only at the end of your life from the entire experience

- Start with random weights
- Adjust weights in each step a bit, in the direction of the steepest descent of the loss function
- Training is repeated n times over the whole data sample: HOW OFTEN??

NOTE: for online learning, the training events should be mixed randomly, otherwise you first steer in a wrong direction from which it is afterward hard to get out again !!

Overtraining



Very careful not to OVERDO with the training !!



NN: cross validation



- Many (all) classifiers have tuning parameters that need to be controlled against overtraining:
 - Number of training cycles, number of nodes (neural net)
 - Smoothing parameters

• ...

- The more free parameters a classifiers has to adjust internally → more prone to overtraining
- More training data \rightarrow better training results
- Divide the data set into "training" and "test" samples (reduces the training data)

Train	Train	Train	Train	Test

What is the best network architecture?



 Theoretically a single hidden layer is enough for any problem, provided one allows for sufficient number of nodes.

(K.Weierstrass theorem)

• "Relatively little is known concerning advantages and disadvantages of using a single hidden layer with many nodes over many hidden layers with fewer nodes. The mathematics and approximation theory of the MLP model with more than one hidden layer is not very well understood"

...."nonetheless there seems to be reason to conjecture that the two hidden layer model may be significantly more promising than the single hidden layer model"

A.Pinkus, "Approximation theory of the MLP model with neural networks", Acta Numerica (1999), pp. 143-195









In order to identify individual particles (PID track by track) compare the energy deposit (dE/dx) and signal temporal shape with results from test beam !

 \rightarrow in MC the energy deposit in the TRD chambers is not reproduced well enough, cannot use MC for comparison or training

 \rightarrow use test beams where clean beams of electrons or pions of well defined energy hit the chambers





We use different methods to evaluate the PID, with increasing amount of information used:

 \rightarrow 1-dimensional likelihood: only the total charge is compared







We use different methods to evaluate the PID, with increasing amount of information used:

 \rightarrow 2-dimensional likelihood:select 2 regions and compare each







Multivariate data analysis methods:

- Fisher discriminant (linear decision boundary)
- Neural networks
- Kernel density methods
- Support Vector Machines
- Decision trees:
 - Boosting
 - Bagging
- Toolkit for Multivariate Data Analysis: TMVA
 - Framework for "all" MVA-techniques, available in ROOT