# Statistical Methods in Particle Physics

## Lecture 9

December 10, 2012

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

Winter Semester 2012 / 13

#### Multivariate data analysis methods:

- Fisher discriminant (linear decision boundary)
- Neural networks

Outline

- Kernel density methods
- Support Vector Machines
- Decision trees:
  - Boosting
  - Bagging
- Toolkit for Multivariate Data Analysis: TMVA
  - Framework for "all" MVA-techniques, available in ROOT
- Significance tests

Material from lectures by Helge Voss (May 2009) http://tmva.sourceforge.net/talks.shtml



Nik !

Neyman-Pearson lemma: statistic test

$$t(x) = \frac{P(x|S)}{P(x|B)}$$
 Likelihood ratio

But most of the times we do NOT have access to the "true probability density":

- Try to estimate the PDF<sub>s</sub>(x) and PDF<sub>B</sub>(x) using the Kernel Density Estimators (in limited regions, from training events) → problem coming from the dimensionality!
- Or neglect correlations and use 1-dimensional PDFs
- Or try other approaches to determine the hyperplanes in the feature space, to separate S and B



Construct non-parametric estimators of the pdfs  $f(\vec{x} | H_0), f(\vec{x} | H_1)$ and use these to construct the likelihood ratio

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

*n*-dimensional histogram is a brute force example of this. More clever estimation techniques can get this to work for (somewhat) higher dimension.

See e.g. K. Cranmer, *Kernel Estimation in High Energy Physics*, CPC **136** (2001) 198; hep-ex/0011057; T. Carli and B. Koblitz, *A multi-variate discrimination technique based on range-searching*, NIM A **501** (2003) 576; hep-ex/0211019

## Kernel Density Estimator (KDE)

- Trying to estimate the probability density p(x) in D-dimensional space
- The only thing at our disposal is our "training data"
- Say we want to know p(x) at "this" point "x"
- One expects to find in a volume V around point "x" N\*∫p(x)dx events from a dataset with N events

- Take  $p(\vec{x}) = \frac{1}{Nh^{D}} \sum_{i=1}^{N} K\left(\frac{\vec{x} - \vec{x}_{i}}{h}\right)$
- K kernel density estimator of the probability density
- h bandwidth, smoothing parameter

#### "events" distributed according to p(x)



### **KDE:** Parzen window

we choose as volume the square drawn
 → rectangular kernel function
 Parzen window

$$\mathbf{K} = \sum_{n=1}^{N} k \left( \frac{\mathbf{x} - \mathbf{x}_n}{h} \right), \quad \text{with} \quad k(\mathbf{u}) = \begin{cases} 1, & |u_i| \le \frac{1}{2}, i = 1 \dots D \\ 0, & \text{otherwise} \end{cases}$$

- Determine K from the training data with signal and background mixed together
- Slow: need to sum over N terms
   → take only k-Nearest Neighbours

"events" distributed according to p(x)





Statistical Methods, Lecture 9, December 5,

#### 8

## Kernel Density Estimator

Parzen Window: "rectangular Kernel" → discontinuities at window edges
 → smoother model for p(x) when using smooth Kernel Fuctions: e.g. Gaussian

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{(2\pi h^2)^{1/2}} \exp\left(\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2}\right)^{1/2}$$

- Place a "Gaussian" around each "training data point" and sum up their contributions at arbitrary points "x" → p(x)
- h: "size" of the Kernel → "smoothing parameter"
- there is a large variety of possible Kernel functions









#### **Kernel Density Estimator**



 $\mathbf{p}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^{N} K_{h}(\mathbf{x} - \mathbf{x}_{n})$ 

- : a general probability density estimator using kernel K
- In the size of the Kernel → "smoothing parameter"
- chosen size of the "smoothing-parameter" → more important than kernel function
- h too small: overtraining
- h too large: not sensitive to features in p(x)

for Gaussian kernels, the optimum in terms of MISE (mean integrated squared error) is given by: h<sub>xi</sub>=(4/(3N))<sup>1/5</sup> σ<sub>xi</sub> with σ<sub>xi</sub>=RMS in variable x<sub>i</sub>



■ a drawback of Kernel density estimators: Evaluation for any test events involves ALL TRAINING DATA → typically very time consuming

binary search trees (i.e. Kd-trees) are typically used in kNN methods to speed up searching

#### "Curse of dimensionality"



Filling a D-dimensional histogram to get a mapping of the PDF is typically unfeasable due to lack of Monte Carlo events.

Shortcoming of nearest-neighbour strategies:

 In higher dimensional classification/regression cases the idea of looking at "training events" in a reasonably small "vicinity" of the space point to be classified becomes difficult:

consider: total phase space volume V=1<sup>D</sup> for a cube of a particular fraction of the volume:

edge length=(fraction of volume)<sup>1/D</sup>



→ 63% of range in each variable necessary → that's not "local" anymore..





$$p(\vec{x}) = \prod_{i=1}^{N} p_i(x_i) \quad \checkmark \quad \begin{array}{l} \text{Holds only if the} \\ \text{components of } x \\ \text{are independent} \end{array}$$

Most importantly, the components of x will generally have non-zero covariances (i.e. they are CORRELATED):

$$V_{ij} = cov[x_i, x_j] = E[x_i x_j] - E[x_i]E[x_j] \neq 0$$



But we can define a set of uncorrelated input variables by a linear transformation, i.e. find the matrix A such that for  $\vec{y} = A \vec{x}$  the covariances are  $cov[x_i, x_i] = 0$ 



For the following suppose that the variables are "decorrelated" in this way for each of  $p(\vec{x}|H_0), p(\vec{x}|H_1)$  separately (since in general their correlations are different).



But even with zero correlation, a multivariate pdf p(x) will in general have non-linearities and thus the decorrelated variables are still not independent



pdf with zero covariance but components still not independent. In fact:

$$p(\mathbf{x}_{2}|\mathbf{x}_{1}) \equiv \frac{p(\mathbf{x}_{1},\mathbf{x}_{2})}{p_{1}(\mathbf{x}_{1})} \neq p_{2}(\mathbf{x}_{2})$$

And therefore:

$$p(x_{1,}x_{2}) \neq p_{1}(x_{1}) p_{2}(x_{2})$$



But if the non-linearities are not too great, it is reasonable to first decorrelate the inputs, and take as our estimator for each pdf:

$$\hat{p}(\vec{x}) = \prod_{i=1}^{n} \hat{p}_{i}(x_{i})$$

So this at least reduces the problem to one of finding estimates of onedimensional pdf's.

The resulting estimated likelihood ratio gives the naïve Bayes classifier (in HEP sometimes called the "likelihood method").

#### **Decision trees**





#### 16

#### Decision trees





overcomes the stability problem



First: in a very simplified way:

• At each step: out of ALL the input variables, find the one for which with a single cut we obtain the best improvement in signal purity

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

where w<sub>i</sub> is the weight of the i'th event

- Iterate until stop criterion reached based on e.g. purity reached, or minimum number of events in a node
- The set of cuts defines the decision boundary

#### **Growing a Decision Tree**



- Why no multiple branches (splits) per node ?
  - Fragments data too quickly; also: multiple splits per node = series of binary node splits
- What about multivariate splits?
  - Time consuming
  - other methods more adapted for such correlatios

#### **Separation Gain**

- What do we mean by "best separation gain"?
- define a measure on how mixed S and B in a node are:
  - Gini-index: (Corrado Gini 1912, typically used to measure income inequality)
    - p (1-p) : p=purity
  - Cross Entropy:
    - -(plnp + (1-p)ln(1-p))
  - Misidentification:

1-max(p,1-p)

 difference in the various indices are small, most commonly used: Gini-index



separation gain: e.g. N<sub>Parent</sub>\*Gini<sub>Parent</sub> – N<sub>left</sub>\*Gini<sub>LeftNode</sub> – N<sub>right</sub>\*Gini<sub>RightNode</sub>

Choose amongst all possible variables and cut values the one that maximised the this.

Helge Voss

#### **Decision Tree Pruning**

One can continue node splitting until all leaf nodes Root are basically pure (using the training sample) node obviously: that's overtraining xi > c1Two possibilities: stop growing earlier  $x_j > c_2$   $x_j < c_2$ generally not a good idea, useless splits might open up subsequent usefull splits S B grow tree to the end and "cut back", nodes that seem statistically dominated: → pruning  $C(T,\alpha) = \sum \sum |y(x) - y(C)| + \alpha N_{\text{leaf nodes}}$ e.g. Cost Complexity pruning: leafs events assign to every sub-tree, T C(T,α) : of T in leaf find subtree T with minmal C(T,α) for given α Loss function prune up to value of a that does not show

regularisaion/ cost parameter

xj > c3 xj < c3

S

xi < c1

В

overtraining in the test sample

#### **Decision Tree Pruning**

"Real life" example of an optimally pruned Decision Tree:





Pruning algorithms are developed and applied on individual trees

optimally pruned single trees are not necessarily optimal in a forest !



The terminal nodes (or 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 the input-variable space (or feature space) as either signal or background

 $\rightarrow$  A decision tree classifier, with discriminant function:

$$f(\vec{x}) = \begin{cases} 1 & \text{if } \vec{x} \in \text{signal region} \\ -1 & \text{otherwise} \end{cases}$$

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 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

Suppose we have a training sample T consisting of N events with:  $x_1, ..., x_N$  event data vectors (each x multivariate)  $y_1, ..., y_N$  true class labels (+1 for signal, -1 for background)  $w_1, ..., w_N$  event weights

Now define a rule to create from this an ensemble of training samples T1, T2, ...

Derive a classifier from each and average them

#### Boosting



Helge Voss

12

#### Adaptive Boosting (AdaBoost)



AdaBoost re-weights events misclassified by previous classifier by:

$$\frac{1 - f_{err}}{f_{err}} \quad \text{with :}$$

$$f_{err} = \frac{\text{misclassified events}}{\text{all events}}$$

 AdaBoost weights the classifiers also using the error rate of the individual classifier according to:

$$y(x) = \sum_{i}^{N_{\text{Classifier}}} \log\left(\frac{1 - f_{\text{err}}^{(i)}}{f_{\text{err}}^{(i)}}\right) C^{(i)}(x)$$

#### **Bagging and Randomised Trees**

other classifier combinations:

- Bagging:
  - combine trees grown from "bootstrap" samples
  - (i.e re-sample training data with replacement)

Randomised Trees: (Random Forest: trademark L.Breiman, A.Cutler)

- combine trees grown with:
  - random subsets of the training data only
  - consider at each node only a random subsets of variables for the split
  - NO Pruning!

These combined classifiers work surprisingly well, are very stable and almost perfect "out of the box" classifiers

Helge Voss

#### AdaBoost: A simple demonstration

var(i) > x

в

 $var(i) \le x$ 

S

The example: (somewhat artificial...but nice for demonstration) :

- Data file with three "bumps"
- Weak classifier (i.e. one single simple "cut" ↔ decision tree stumps )



Two reasonable cuts: a) Var0 > 0.5 → ε<sub>signal</sub>=66% ε<sub>bkg</sub> ≈ 0% misclassified events in total 16.5% or b) Var0 < -0.5 → ε<sub>signal</sub>=33% ε<sub>bkg</sub> ≈ 0% misclassified events in total 33%

#### the training of a single decision tree stump will find "cut a)"

#### AdaBoost: A simple demonstration

The first "tree", choosing cut a) will give an error fraction: err = 0.165

→ before building the next "tree": weight wrong classified training events by (1-err/err)) ≈ 5

the next "tree" sees essentially the following data sample:



Helge Voss

**BDT** response

16

#### AdaBoost vs other Combined Classifiers

Sometimes people present "boosting" as nothing else then just "smearing" in order to make the Decision Trees more stable w.r.t statistical fluctuations in the training.

→clever "boosting" however can do more, than for example:

- Random Forests

Bagging

as in this case, pure statistical fluctuations are not enough to enhance the 2<sup>nd</sup> peak sufficiently

however: a "fully grown decision tree" is much more than a "weak classifier"

→ "stabilization" aspect is more important



Surprisingly: Often using smaller trees (weaker classifiers) in AdaBoost and other clever boosting algorithms (i.e. gradient boost) seems to give overall significantly better performance !

Helge Voss

#### **Boosting at Work**





- 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 and are effectively ignored.
- Easy to deal with inputs of mixed types (real, integer, categorical, ...)
- It a tree had only a few leaves it is easy to visualize (but rarely we use only a single tree)
- There are a number of boosting algorithms, which differ primarily in the rule for updating the weights (ε-Boost, LogitBoost, ...)
- Other ways of combining the weaker classifiers: Bagging (Bootstrap-Aggregating) generates the ensemble of classifiers by random sampling with replacement from the full training sample.





Suppose hypothesis *H* predicts pdf  $f(\vec{x}|H)$  for a set of observations  $\vec{x} = (x_{1,}...,x_{n})$ 

We observe a single point in this space:  $\vec{x}_{obs}$ 

What can we say about the validity of *H* in light of the data?

Decide what part of the data space represents less compatibility with *H* than does the point  $\vec{x}_{obs}$  (Not unique!)







Express 'goodness-of-fit' by giving the p-value for H:

p = probability, under assumption of H, to observe data with equal or lesser compatibility with H relative to the data we got

NOTE! This is NOT the probability that H is true!

In frequentist statistics we don't talk about P(H) (unless H represents a repeatable observation).

In Bayesian statistics we do. Use Bayes' theorem to obtain

$$\mathsf{P}(\mathsf{H}|\mathbf{\vec{x}}) = \frac{\mathsf{P}(\mathbf{\vec{x}}|\mathsf{H}) \pi(\mathsf{H})}{\int \mathsf{P}(\mathbf{\vec{x}}|\mathsf{H}) \pi(\mathsf{H}) \, \mathsf{d}\mathsf{H}}$$

where  $\pi(H)$  is the prior probability for H.

For now stick with the frequentist approach.



#### Testing whether a coin is "fair"

Probability to observe n heads in N coin tosses is binomial:

$$P(n;p,N) = \frac{N!}{n! (N-n)!} p^n (1-p)^{N-n}$$

Hypothesis H: the coin is fair (p=0.5)

Suppose we toss the coin N=20 times and get n=17 heads.

Region of data space with equal or lesser compatibility with H relative to n=17 is: n=17, 18, 19, 20, 0, 1, 2, 3 Adding up the probabilities for these values gives: P(n=0, 1, 2, 3, 17, 18, 19, or 20) = 0.0026

i.e. p=0.0026 is the probability of obtaining such a bizarre result (or more so) "by chance", under the assumption of H



Suppose we observe n events. These can consist of:

- n<sub>b</sub> events from known processes (background)
- n<sub>s</sub> events from a new process (signal)

If  $n_s$ ,  $n_b$  are Poisson random variables with means s, b, then  $n=n_s+n_b$  is also Poisson, with mean s+b

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

Suppose b=0.5, and we observe nobs=5. Should we claim evidence for a new discovery?



Suppose we observe n events. These can consist of:

- n<sub>b</sub> events from known processes (background)
- n<sub>s</sub> events from a new process (signal)

If  $n_s$ ,  $n_b$  are Poisson random variables with means s, b, then  $n=n_s+n_b$  is also Poisson, with mean s+b

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

Suppose b=0.5, and we observe nobs=5. Should we claim evidence for a new discovery?

```
Give p-value for hypothesis s=0:
p-value = P( n \ge 5; b=0.5, s=0)
= 1.7 x 10<sup>-4</sup> \neq P(s=0) !!
```

Often define significance Z as the number of standard deviations that a Gaussian variable would fluctuate in one direction to give the same p-value



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





## The significance of a peak



In the two bins with the peak, 11 entries found with b = 3.2The p-value for the s=0 hypothesis is:

P( n 
$$\ge$$
 11; b=3.2, s=0) = 5.0 x 10<sup>-4</sup>



But ... did we know where to look for the peak?

→ give P(n≥11) in any 2 adjacent bins

Is the observed width consistent with the expected x resolution?

 $\rightarrow$  take x window several times the expected resolution How many bins x distributions have we looked at?

→ look at a thousand of them, you'll find a  $10^{-3}$  effect Did we adjust the cuts to "enhance" the peak?

 $\rightarrow$  freeze cuts, repeat analysis with new data How about the bins to the sides of the peak ... (too low!) Should we publish ??



HEP folklore is to claim discovery when  $p= 2.9 \times 10^{-7}$ , corresponding to a significance Z=5.

This is very subjective and really should depend on the prior probability of the phenomenon in question, e.g.

Phenomenon	Reasonable p-value for discovery
D <sup>o</sup> D <sup>o</sup> mixing	~0.05
Higgs	~10 <sup>-7</sup> (?)
Life on Mars	~ <b>10</b> <sup>-10</sup>
Astrology	~10 <sup>-20</sup>

One should also consider the degree to which the data are compatible with the new phenomenon, not only the level of disagreement with the null-hypothesis: p-value is only the first step !!!

The p-value is a function of the data, and is thus itself a random variable with a given distribution. Suppose the p-value of H is found from a test statistic t(x) as:

$$p_{H} = \int_{t}^{\infty} f(t'|H) dt'$$

The pdf of  $p_{H}$  under assumption of H is:

$$g(p_{H}|H) = \frac{f(t|H)}{\left|\partial p_{H}/\partial t\right|} = \frac{f(t|H)}{f(t|H)} = 1 \qquad (0 \le p_{H} \le 1)$$

 $g(p_H|H')$ 

In general for continuous data, under assumption of H,  $p_{H} \sim$  uniform in [0,1]

And is concentrated toward zero for some (broad) class of alternatives

 $p_H$ 



The probability to find the *p*-value of  $H_0$ ,  $p_0$ , less than  $\alpha$  is

$$\mathsf{P}(\mathsf{p}_0 \leq \alpha | \mathsf{H}_0) = \alpha$$

We started by defining critical region in the original data space ( $\mathbf{x}$ ), then reformulated this in terms of a scalar test statistic  $t(\mathbf{x})$ .

We can take this one step further and define the critical region of a test of  $H_0$  with size  $\alpha$  as the set of data space where  $p_0 \le \alpha$ .

Formally the *p*-value relates only to  $H_0$ , but the resulting test will have a given power with respect to a given alternative  $H_1$ .

Test statistic for comparing observed data (n<sub>i</sub> independent) to predicted mean values

$$\vec{n} = (n_1, \dots, n_N)$$
$$\vec{v} = (v_1, \dots, v_N)$$

$$\chi^2 = \sum_{i=1}^{N} \frac{(n_i - v_i)^2}{\sigma_i^2}$$
, where  $\sigma_i^2 = V[n_i]$  Pearson's  $\chi^2$  statistic

 $\chi^2$  = sum of squares of the deviations of the I'th measurement from the I'th prediction, using  $\sigma_i$  as the 'yardstick' for comparison

For  $n_i \sim Poisson(v_i)$  we have  $V[n_i] = v_i$ , so this becomes:

$$\chi^{2} = \sum_{i=1}^{N} \frac{(n_{i} - v_{i})^{2}}{v_{i}}$$

If  $n_i$  are Gaussian with mean  $v_i$  and standard deviation  $\sigma_i$ , namely  $n_i \sim N(v_i, \sigma_i^2)$ , then Pearson's X<sup>2</sup> will follow the X<sup>2</sup> pdf (here for X<sup>2</sup> =z)

$$f_{X^2}(z;N) = \frac{1}{2^{N/2}\Gamma(N/2)} z^{N/2-1} e^{-z/2}$$

If the n<sub>i</sub> are Poisson with  $v_i >> 1$  (in practice OK for  $v_i > 5$ ) then the Poisson distribution becomes Gaussian and therefore Pearson's X<sup>2</sup> statistic here as well follows the X<sup>2</sup> pdf.

The  $X^2$  value obtained from the data then gives the p-value:

$$p = \int_{X^2}^{\infty} f_{X^2}(z;N) dz$$





Recall that for the chi-square pdf for N degrees of freedom  $X^2$ 

$$\mathsf{E}[\mathsf{z}] = \mathsf{N}, \quad \mathsf{V}[\mathsf{z}] = 2\mathsf{N}$$

This makes sense; if the hypothesized  $v_i$  are right, the rms deviation of  $n_i$  from  $v_i$  is  $\sigma_i$ , so each term in the sum contributes ~ 1.

One often sees  $\chi^2/N$  reported as a measure of goodness-of-fit. BUT! Better to give the  $\chi^2$  and N separately. Consider e.g.:

$$X^2$$
 = 15, N=10 → p-value = 0.13  
 $X^2$  = 150, N=100 → p-value = 9.0 x 10<sup>-4</sup>

i.e. for N large, even a X<sup>2</sup> per dof only a bit greater than one can imply a small p-value, i.e., poor goodness-of-fit

## Summary of p-value



- Introduction to significance tests:
  - p-value expresses the level of agreement between data and hypothesis
  - p-value is NOT the probability of the hypothesis!
- P-value can be used to define a critical region, i.e. region of data space where  $p < \alpha$
- Widely used X<sup>2</sup> test:
  - Statistic = sum of (data prediction)2 / variance
  - Often  $X^2$  chi-square pdf  $\rightarrow$  use to get p-value
  - Otherwise may need to use MC