

Statistical Learning in Physics Analysis

Jens Zimmermann
zimmerm@mppmu.mpg.de



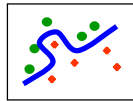
Max-Planck-Institut für Physik, München



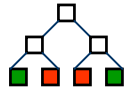
MPI für extraterrestrische Physik, München



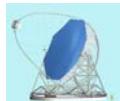
Forschungszentrum Jülich GmbH



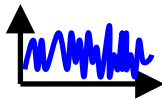
Statistical Learning?



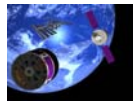
Three Classes of Learning Methods



Applications in Physics Analysis



Training the Learning Methods



Examples

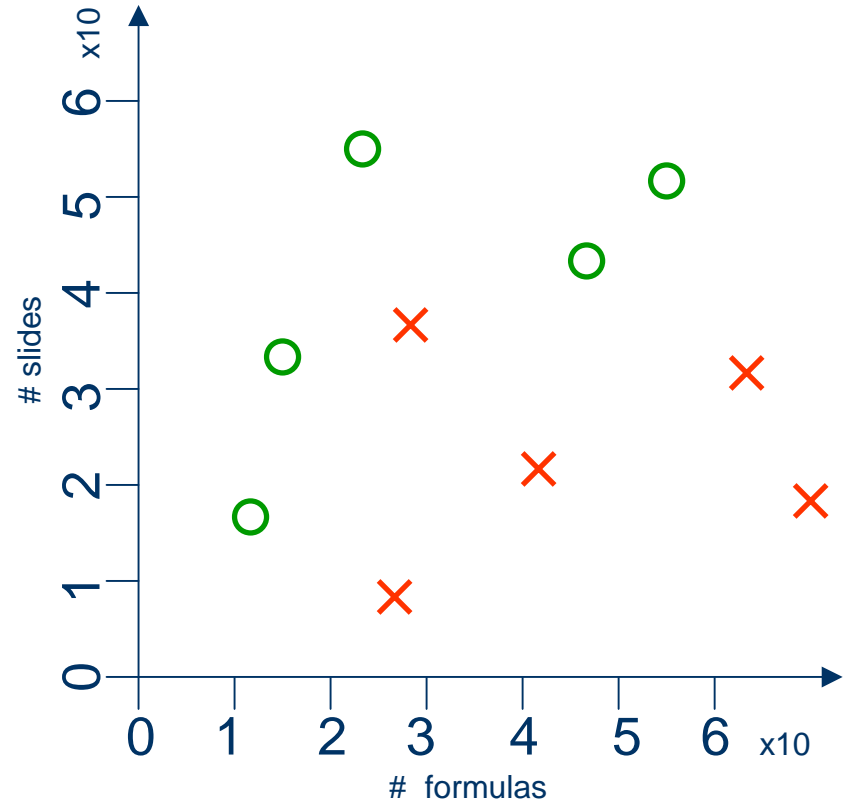


Conclusion

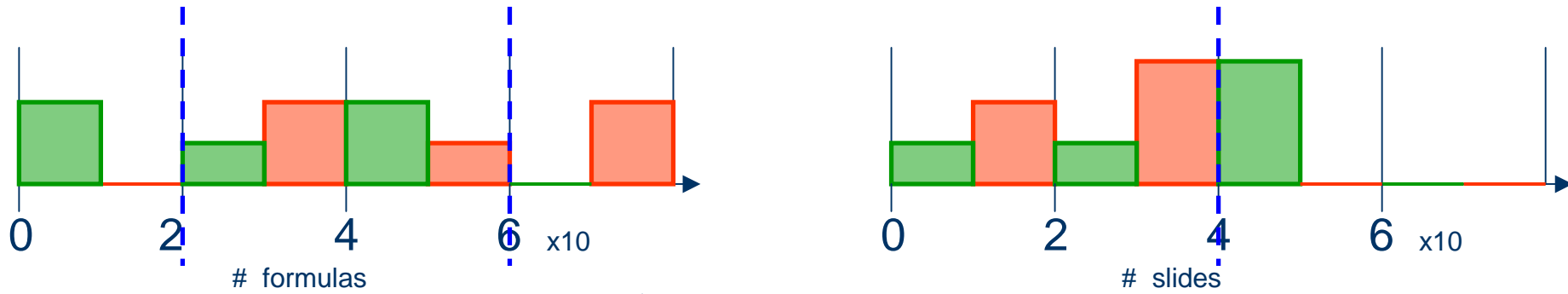
Some Events

Experimentalists ○ Theorists ×

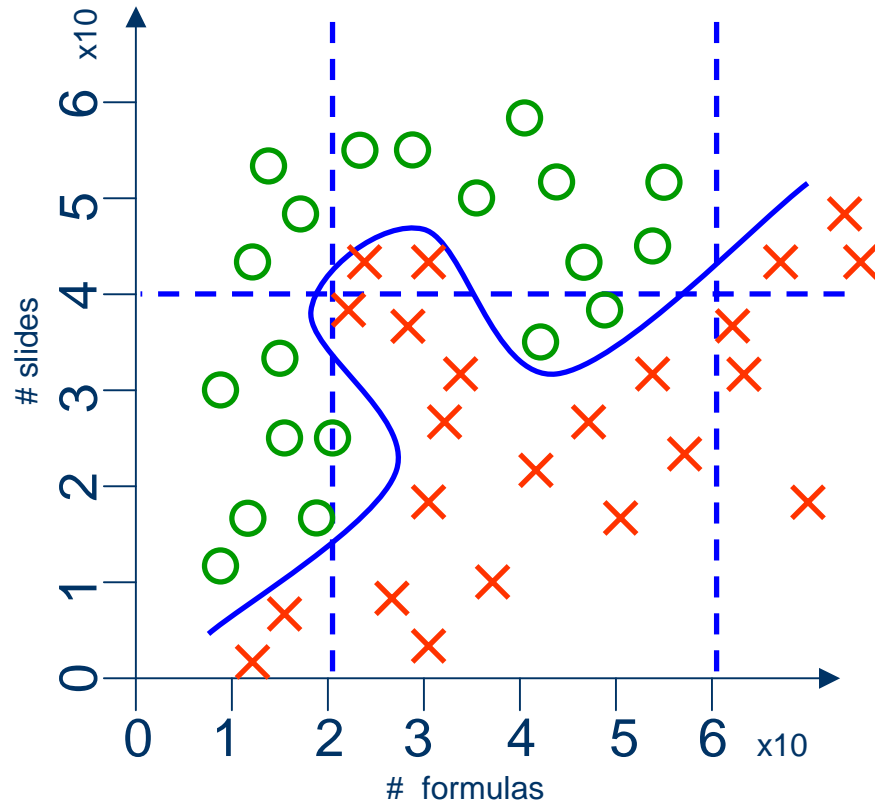
| # formulas | # slides |
|------------|----------|
| 42 | 21 |
| 28 | 8 |
| 71 | 19 |
| 64 | 31 |
| 29 | 36 |
| 15 | 34 |
| 48 | 44 |
| 56 | 51 |
| 25 | 55 |
| 12 | 16 |



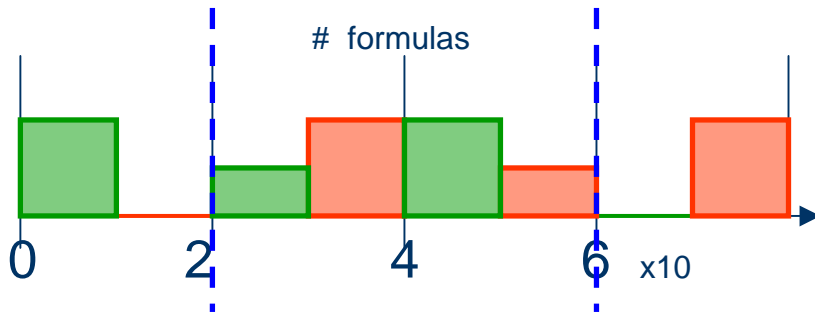
First Analysis



Experimentalists
Theorists



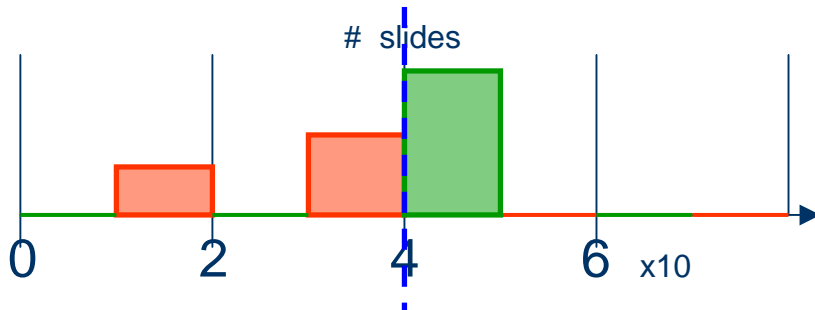
Decision Trees



#formulas < 20 → exp

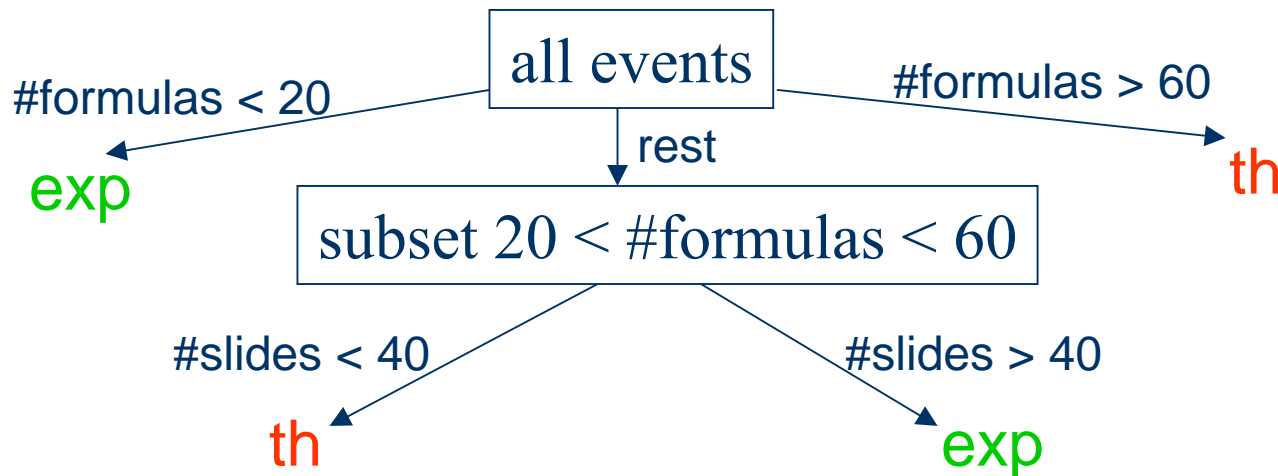
#formulas > 60 → th

20 < #formulas < 60?



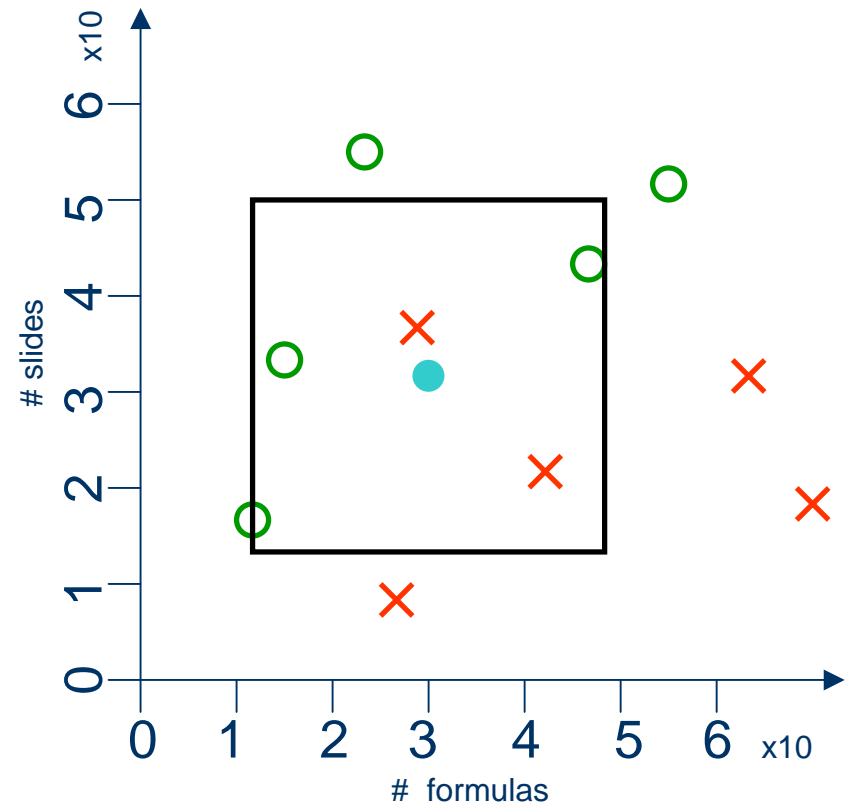
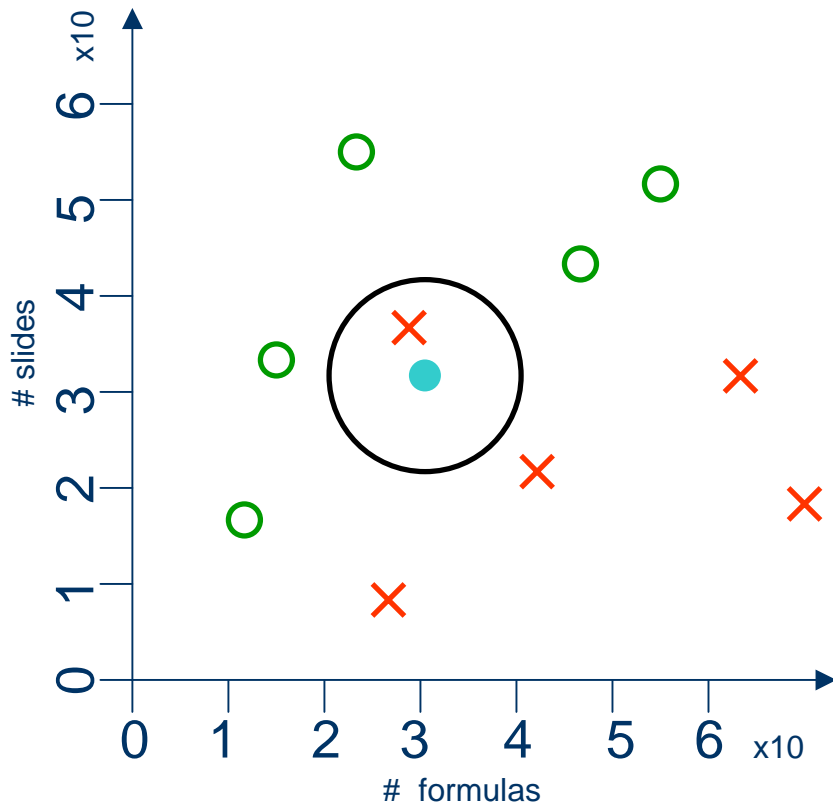
#slides > 40 → exp

#slides < 40 → th



Local Density Estimators

Search for similar events that are already classified and count the members of the two classes.



k-Nearest-Neighbour

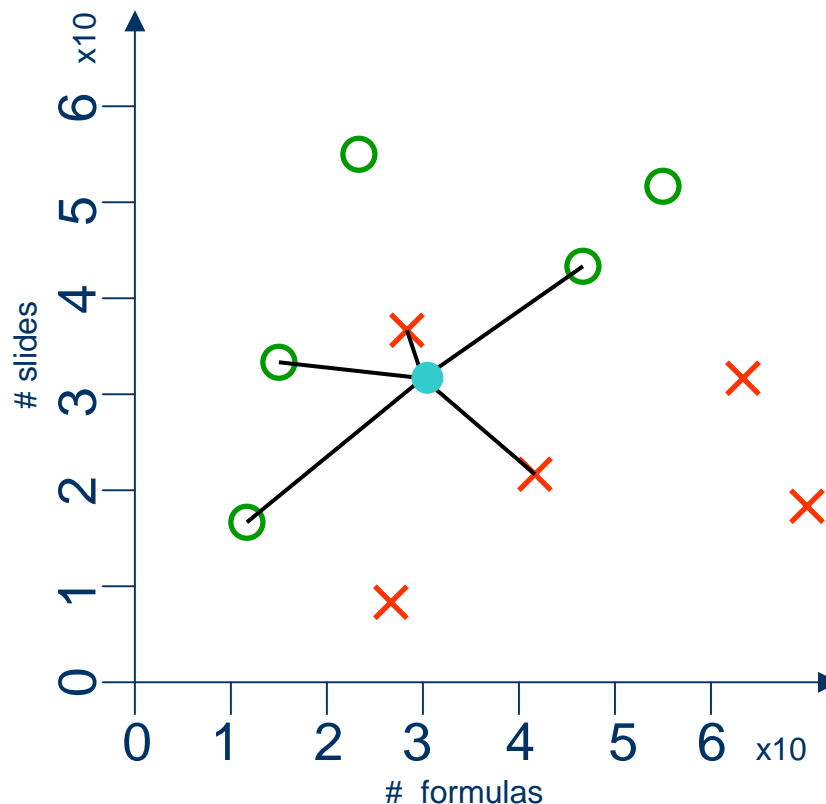
k=1
out= X

k=2
out= ~~X~~

k=3
out= X

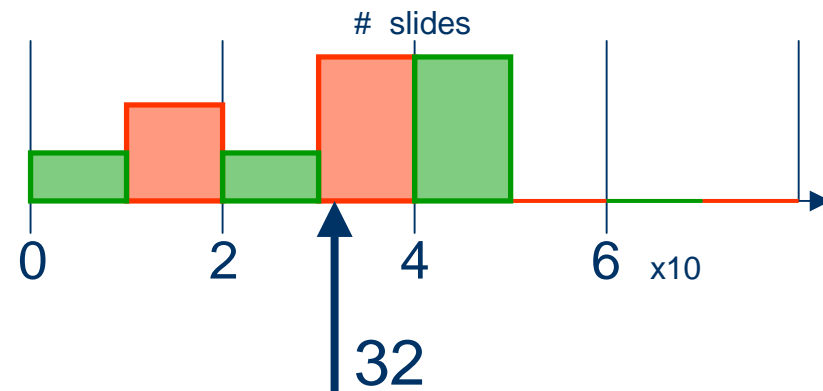
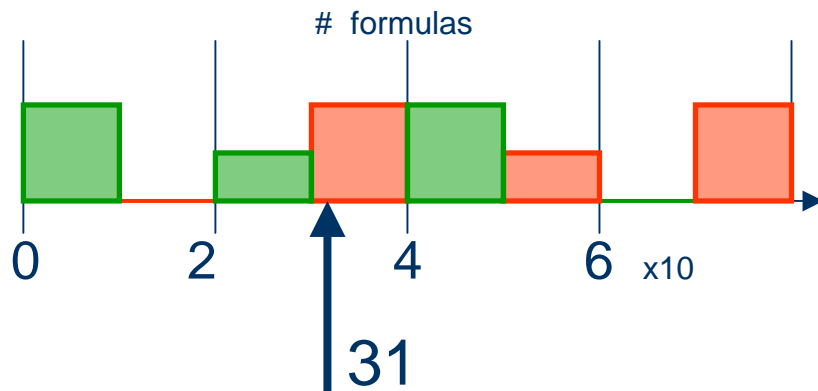
k=4
out= ~~X~~

k=5
out= O

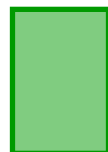


For every evaluation position the distances to each training position need to be determined!

Maximum Likelihood / Naive Bayes



$$p_{Th} = \frac{2}{5} \times \frac{3}{5} = 0.24$$



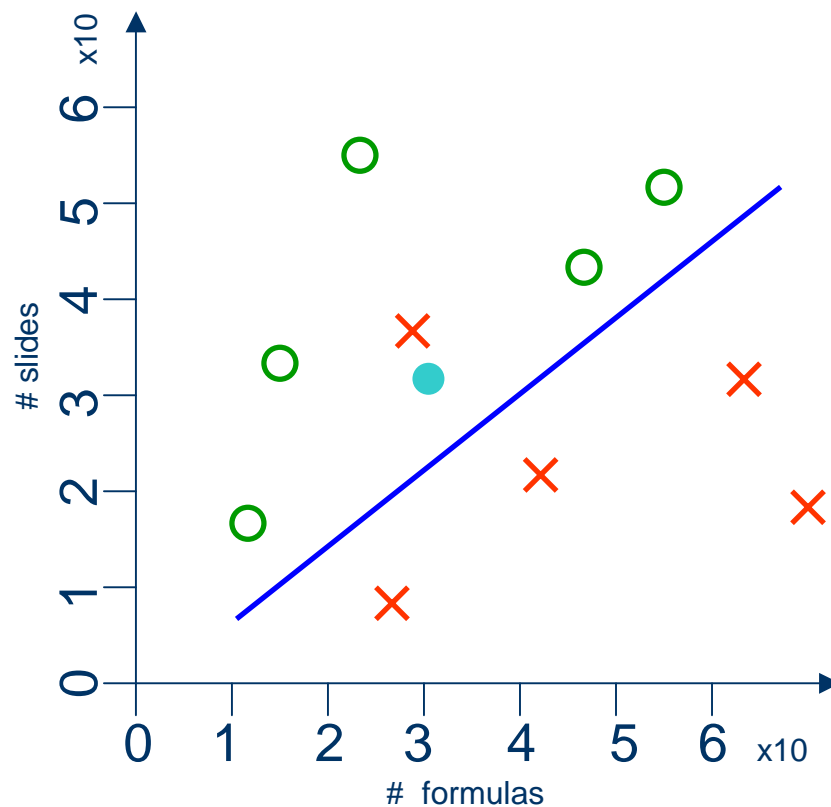
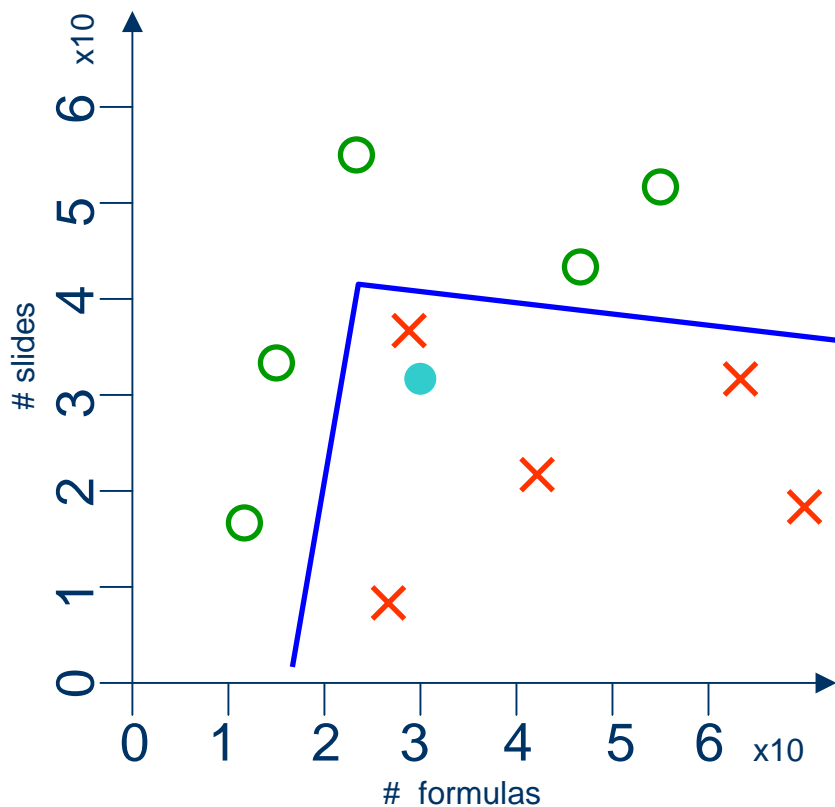
$$p_{Exp} = \frac{1}{5} \times \frac{1}{5} = 0.04$$

out=×

Correlation gets lost completely by projection!

Methods Based on Linear Separation

Divide the input space into regions separated by one or more hyperplanes.
Extrapolation is done!



Linear Discriminant Analysis

$$y_k = \gamma_0 + \gamma_1 x_{k,1} + \dots + \gamma_d x_{k,d} + \sqrt{v} \xi_k \quad 1 \leq k \leq n \quad A = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \dots & x_{n,d} \end{pmatrix}$$

Fisher
1930

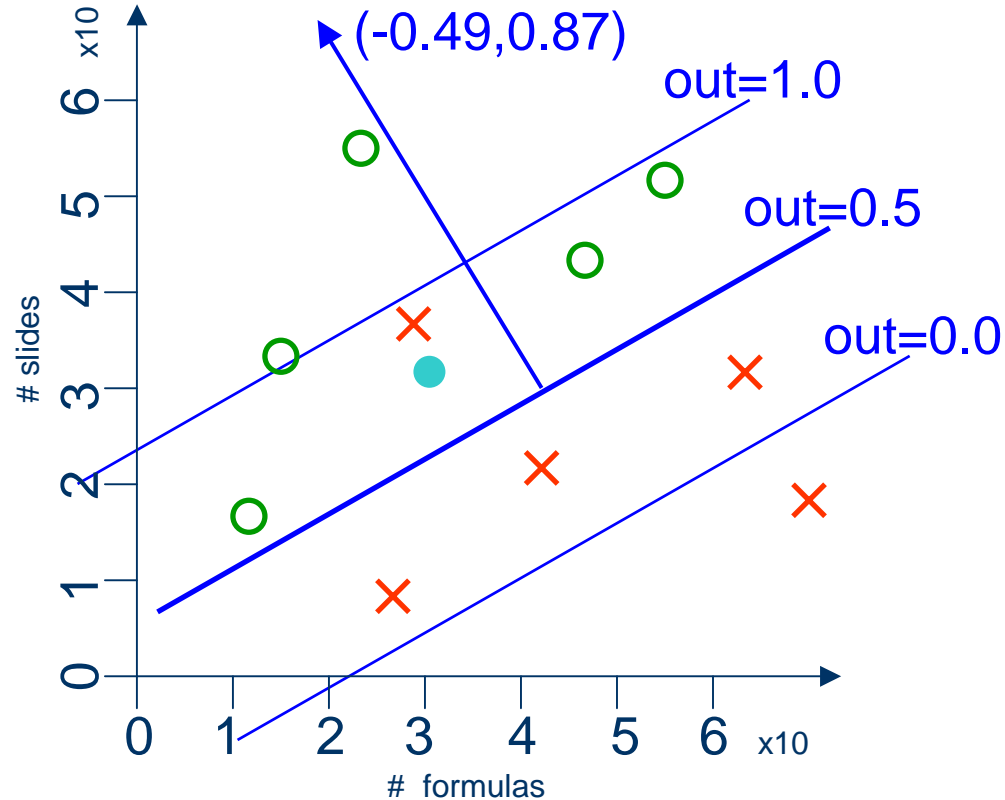


$$Y = A\gamma$$

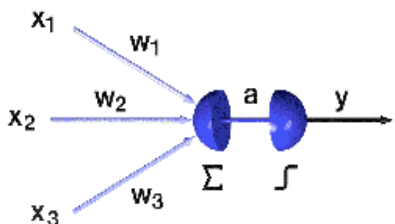
$$\hat{\gamma} = (A^T A)^{-1} A^T Y$$

$$\text{out} = 0.3 - 0.012 \times \text{formulas} + 0.021 \times \text{slides}$$

Only one separating hyperplane is usually not enough!
Can we combine two or more?

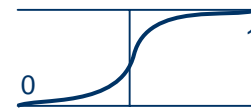


Neural Networks

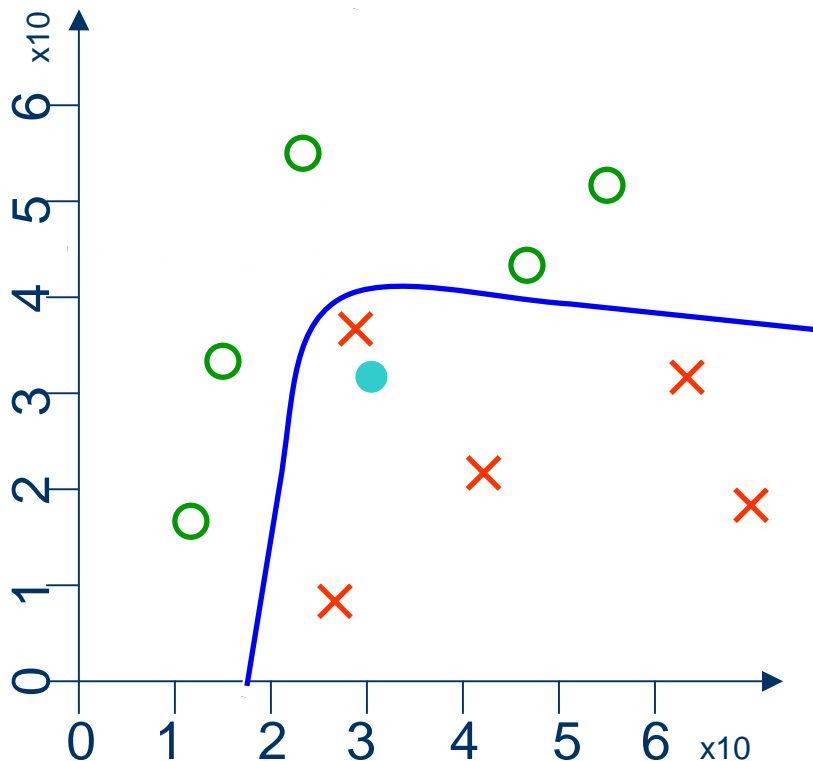
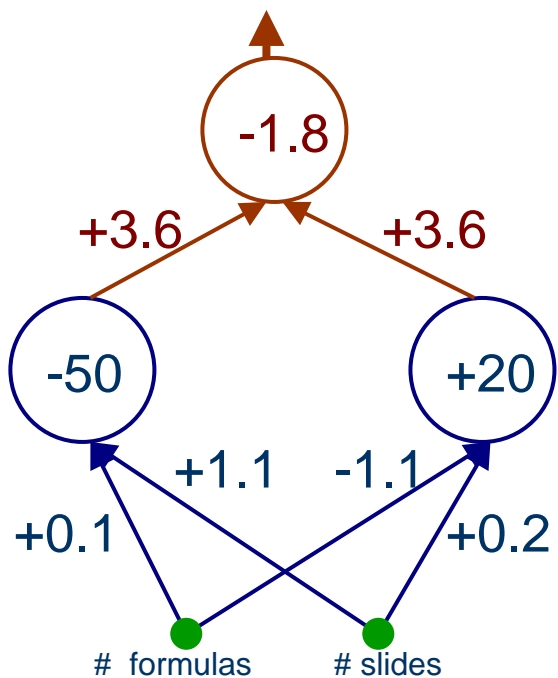


$$y = \sigma\left(\sum w_i x_i + s\right)$$

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$



Train NN with two hidden neurons (gradient descent): $E = \frac{1}{N} \sum_{i=1}^N (y_i - out(x_i))^2$



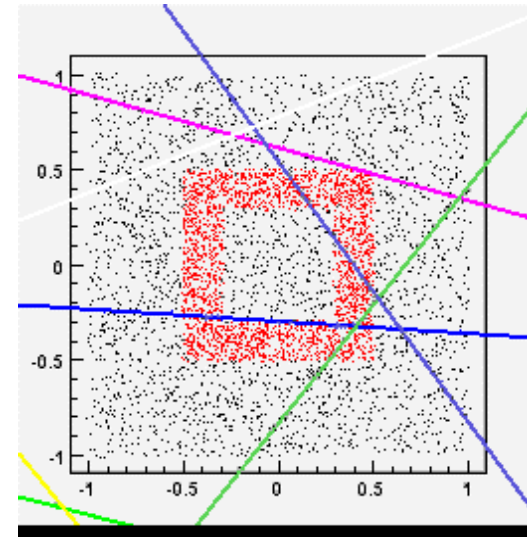
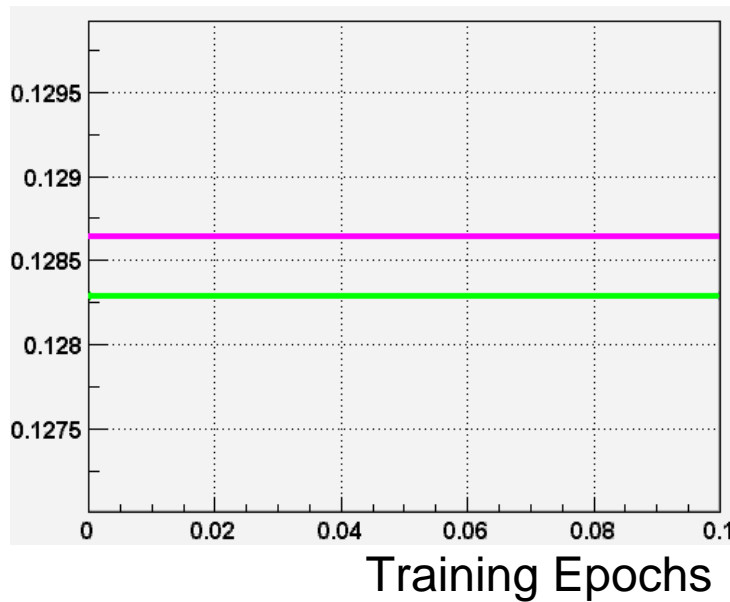
NN Training

8 hidden neurons = 8 separating lines

Test-Error



Train-Error



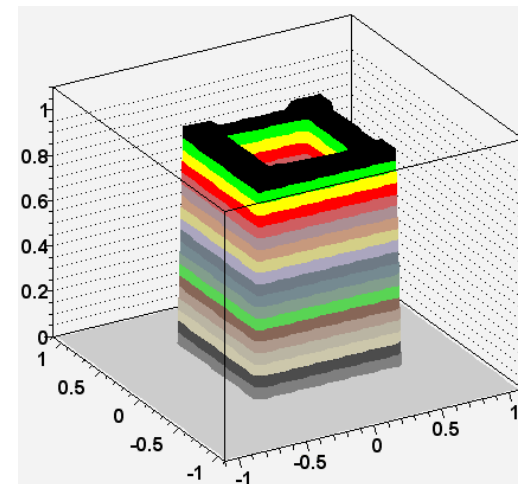
signal



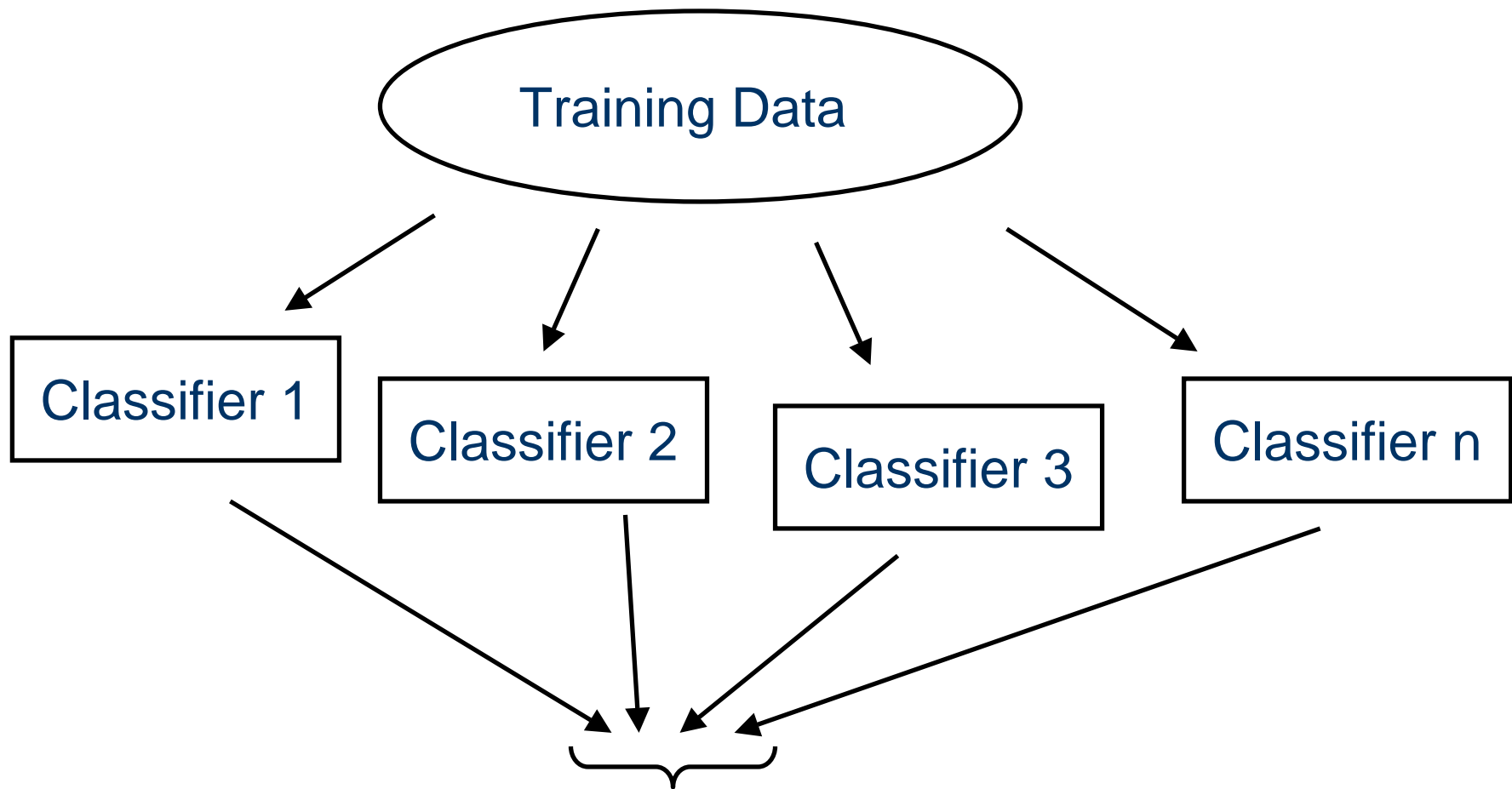
background



$$E = \frac{1}{N} \sum_{i=1}^N (y_i - out(x_i))^2$$



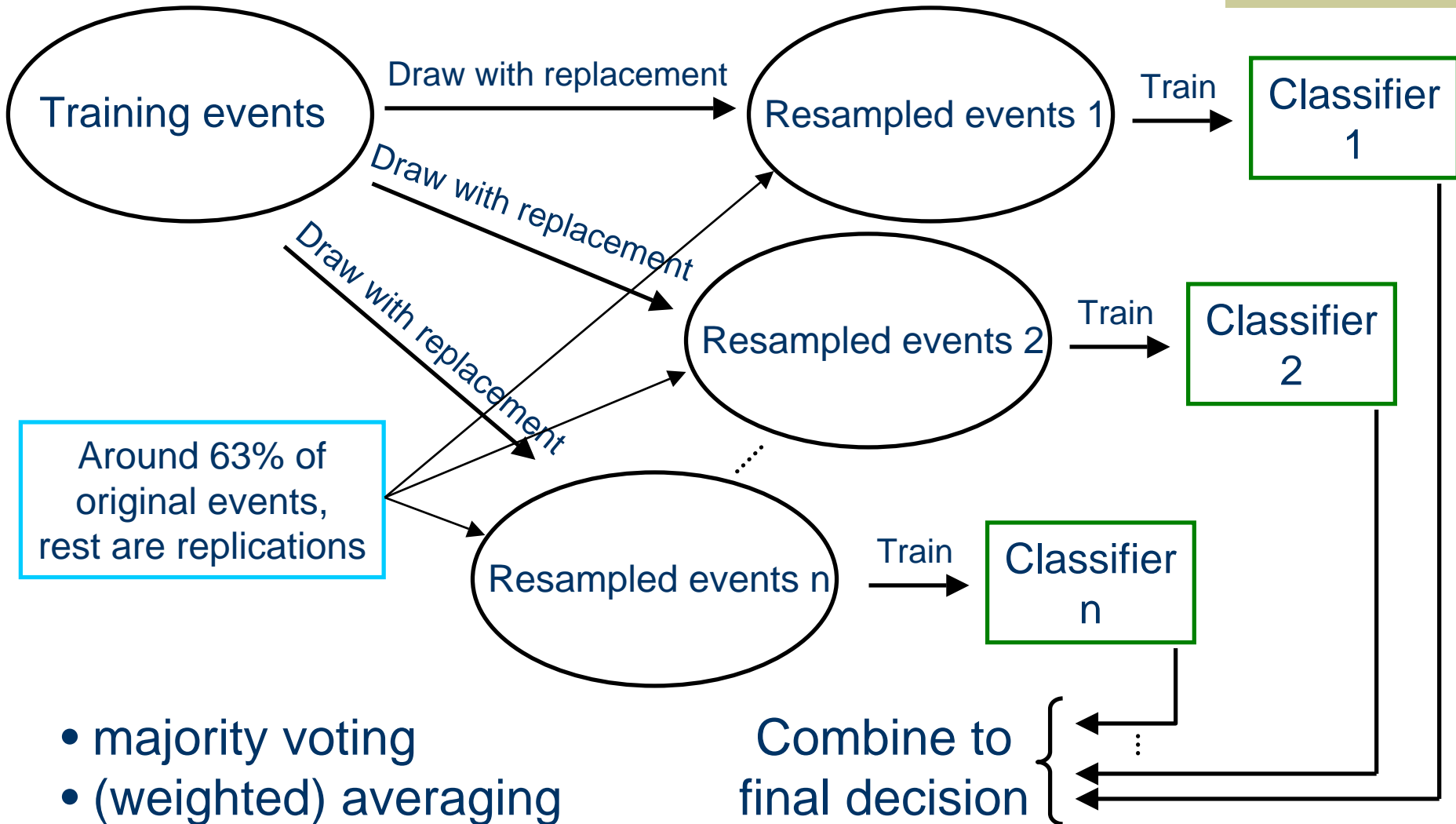
Meta-Learning Strategies



Combine different classifications
to one final decision

Bagging – Procedure

Bootstrap aggregating



- majority voting
- (weighted) averaging

Aim is to create strong classifiers which are as independent as possible.

Random Forests

Basis:

Decision Tree (CART)
without pruning

Modification:

At each node of the tree:
Search only through a
randomly selected subset
of all features

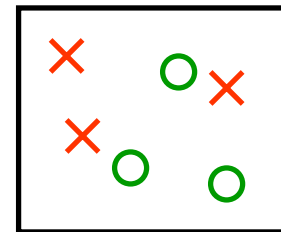
Use Bagging on this classifier

Tree, Randomness, Combination

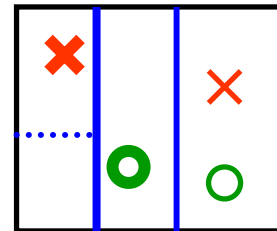
→ RF



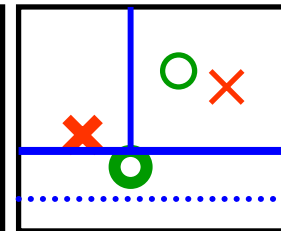
Training:



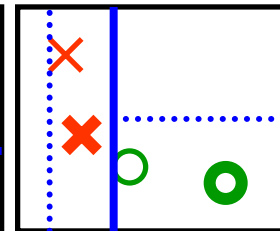
Create
3 trees



1 - 2,1

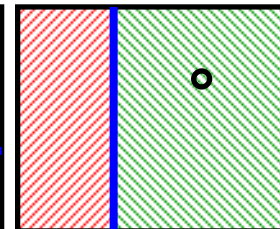
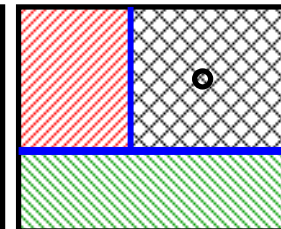
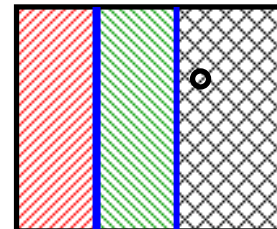


2 - 2,1

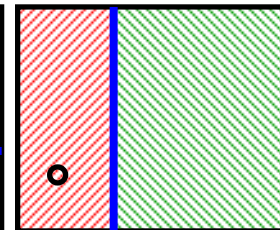
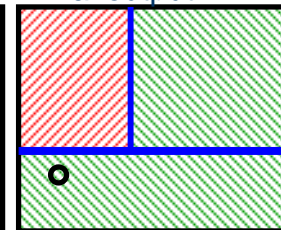
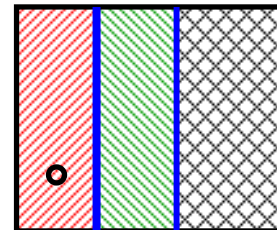


1 - 1,2

Testing/Evaluation:



final output =



final output =

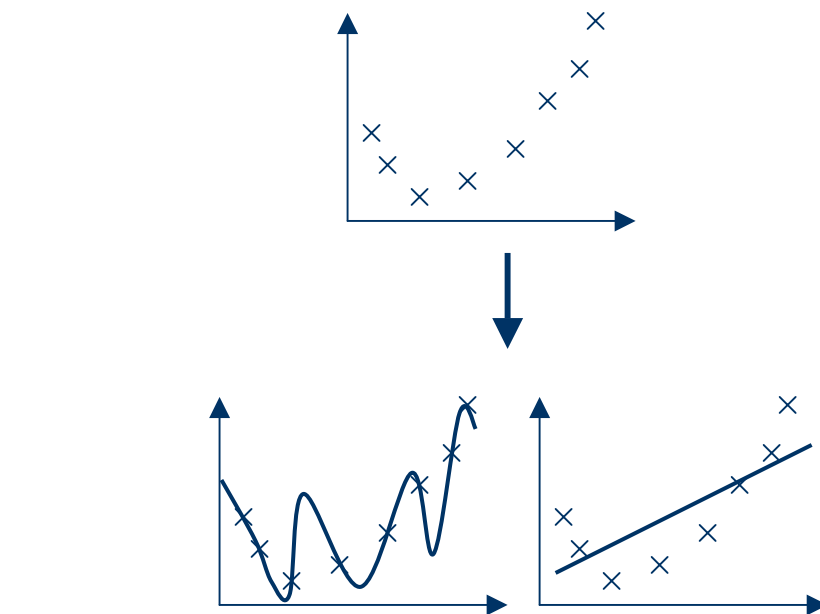
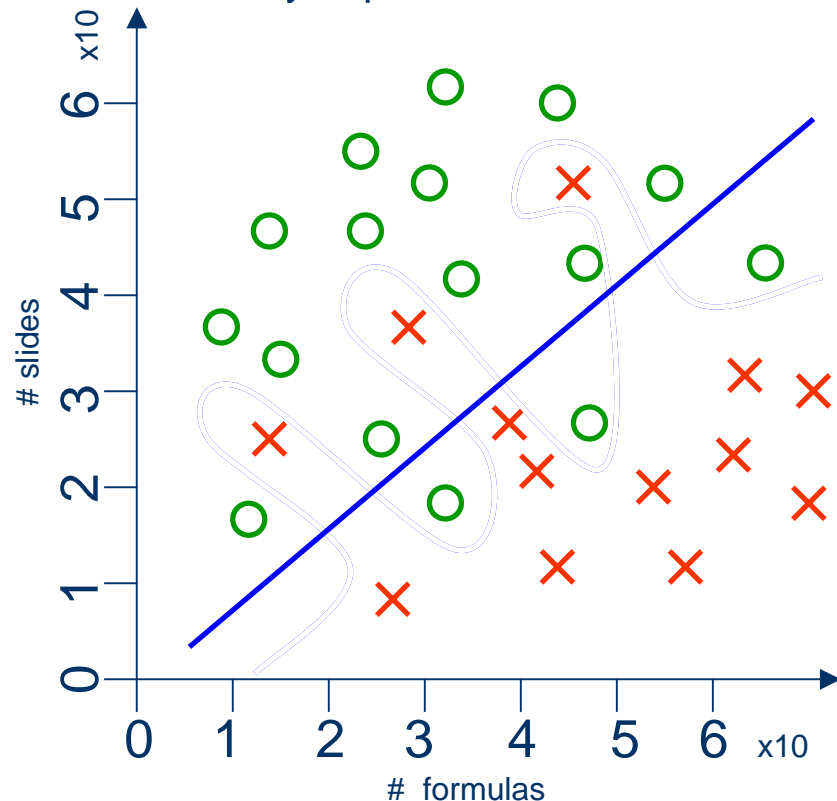
Training of Statistical Learning Methods

Statistical Learning Method:

From N examples (\vec{x}_i, y_i) infer a rule $\vec{x} \rightarrow out(\vec{x})$

Important: Generalisation vs. Overtraining

Easily separable but with noise?



Too high degree of polynomial results in interpolation
but too low degree means bad approximation

Applications in Physics Analysis

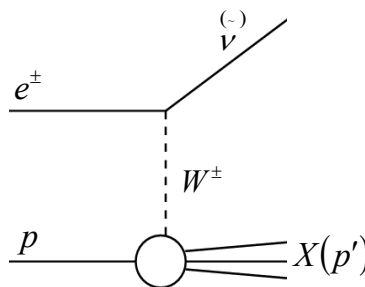
Classification
Online
„Trigger“



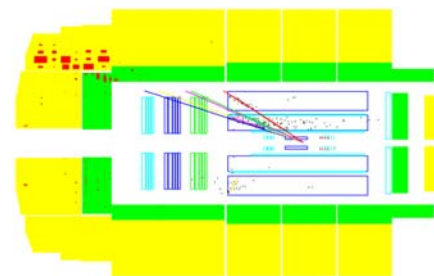
H1



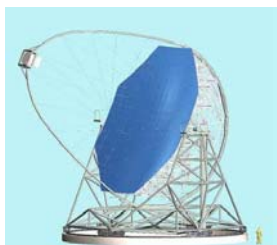
L2NN



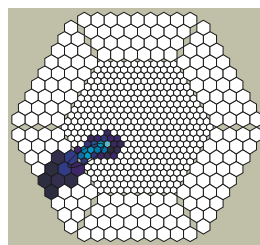
Charged Current



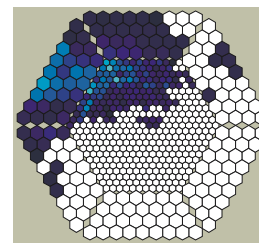
Classification
Offline
„Purification“



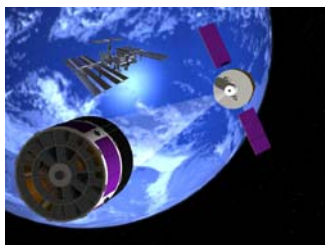
MAGIC



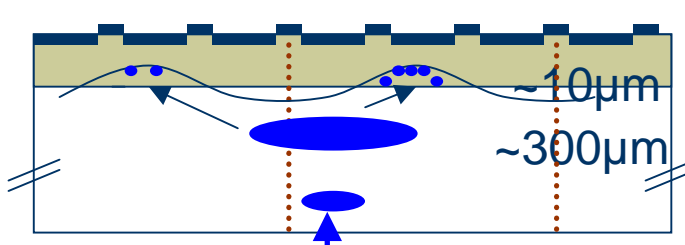
Gamma vs. Hadron



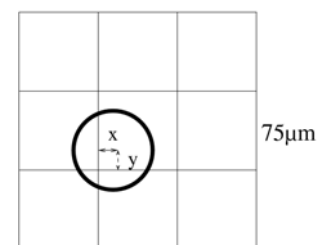
Regression



XEUS

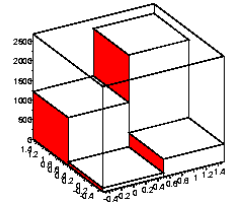
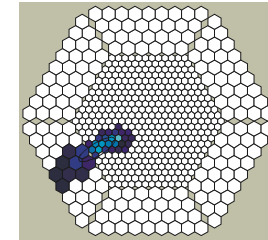
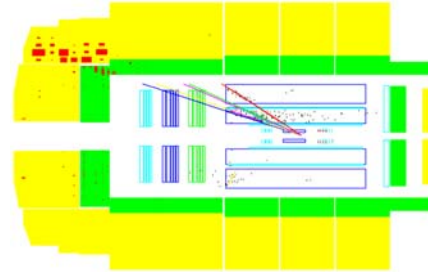


X-ray CCD

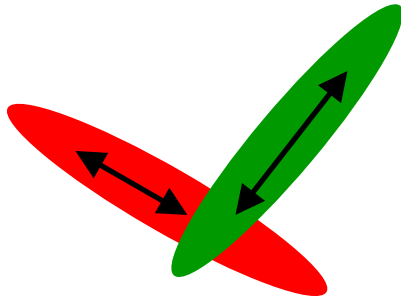


Features

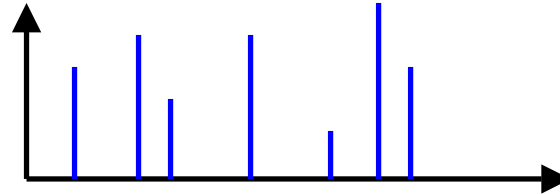
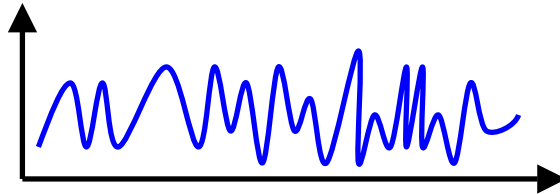
Choose raw quantities



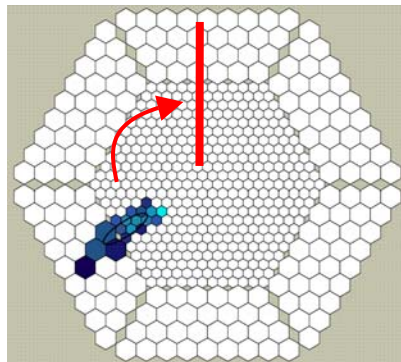
PCA



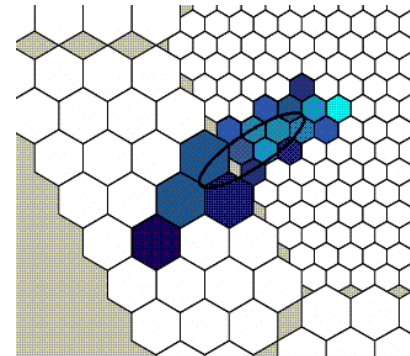
FFT



Symmetrie

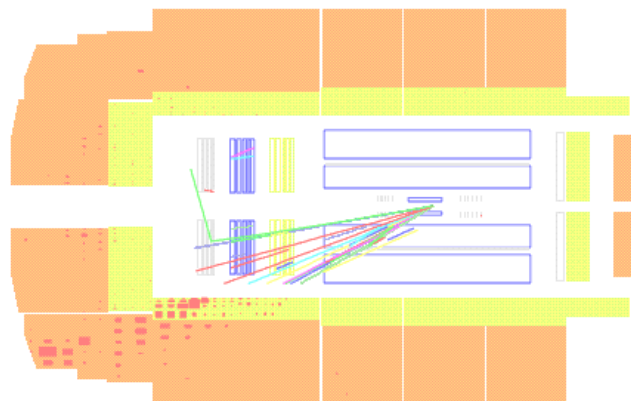
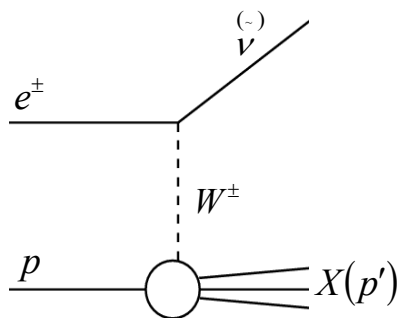


Fit

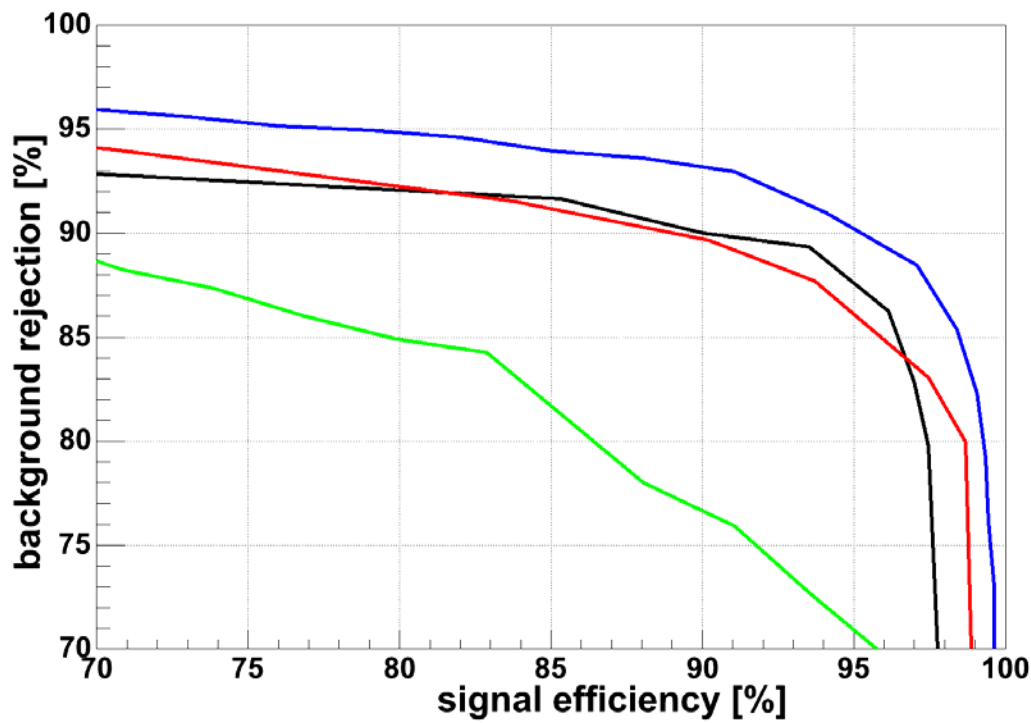
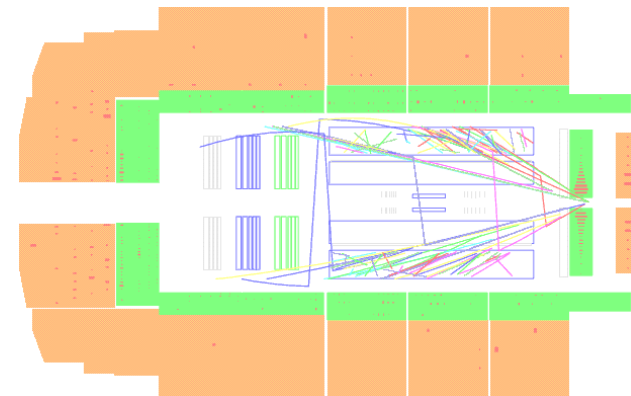


H1 - Triggering Charged Current Interactions

Charged Current Interaction (signal)



Beam Gas Event (background)



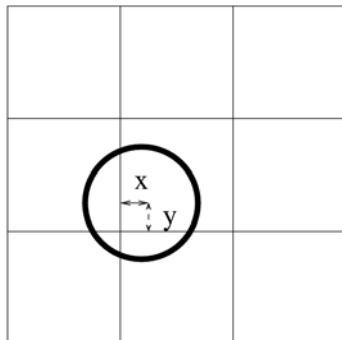
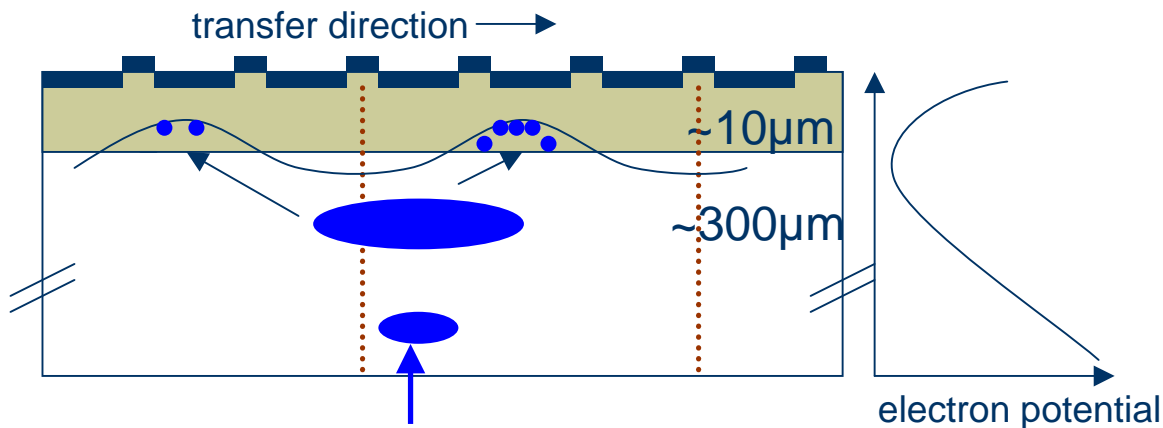
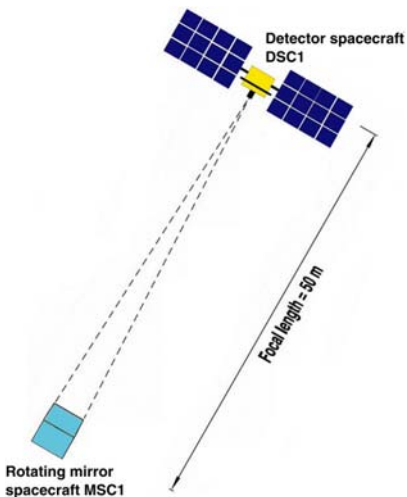
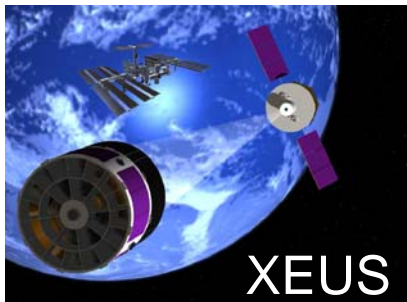
Neural Network

Decision Tree C4.5

k-Nearest-Neighbour-Search

Maximum Likelihood

Regression of the Incident Position of X-ray Photons



$$x_{COM} = \frac{\sum x_i c_i}{\sum c_i}$$

$$x_{CCOM} = x_{COM} + \Delta(x_{COM})$$

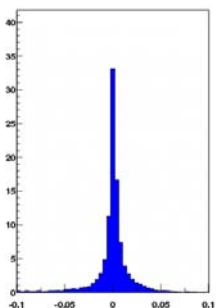
σ of reconstruction in μm :

Neural Networks 3.6

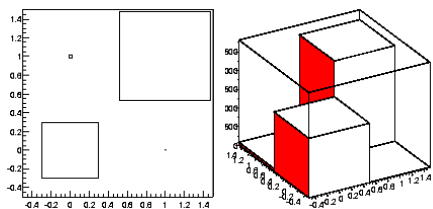
k-Nearest-Neighbour 3.7

classical methods

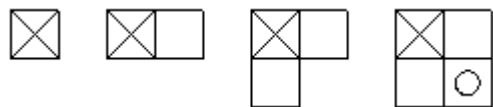
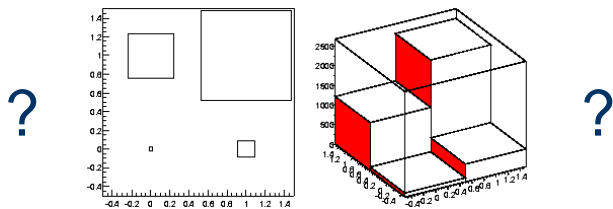
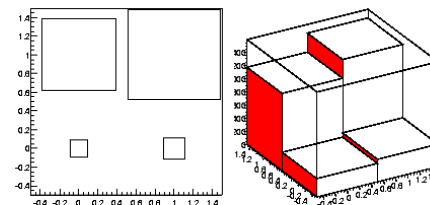
ETA 3.9 CCOM 4.0



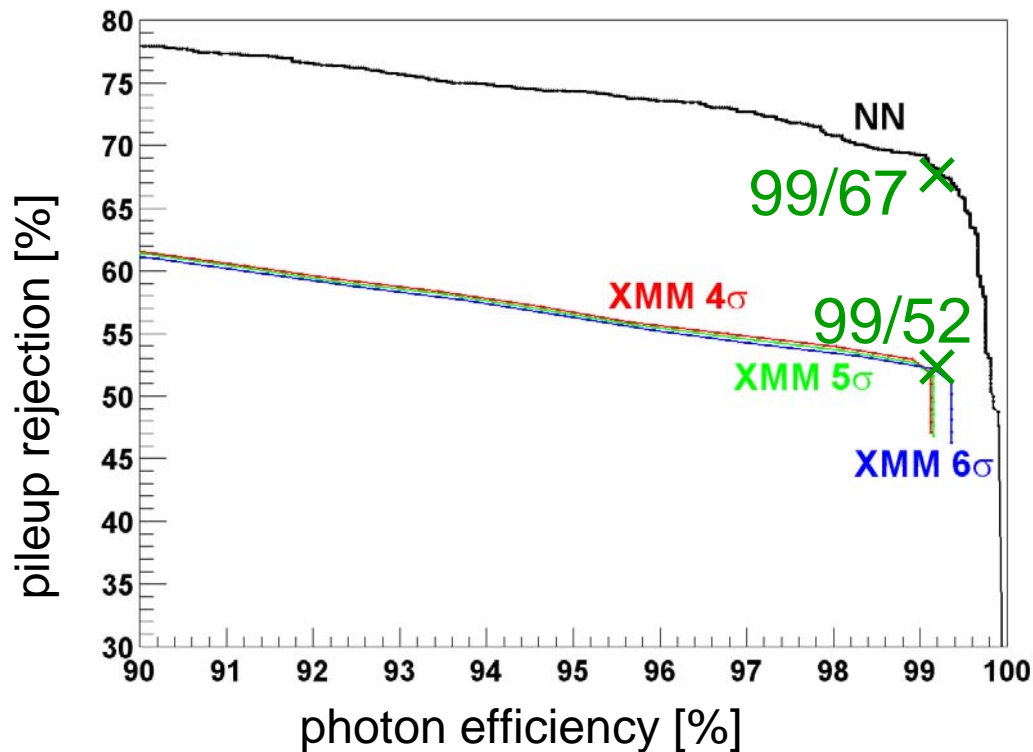
Pileup Recognition – Setup



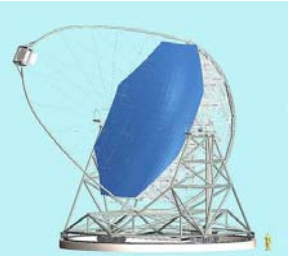
Pileup vs. Single photon



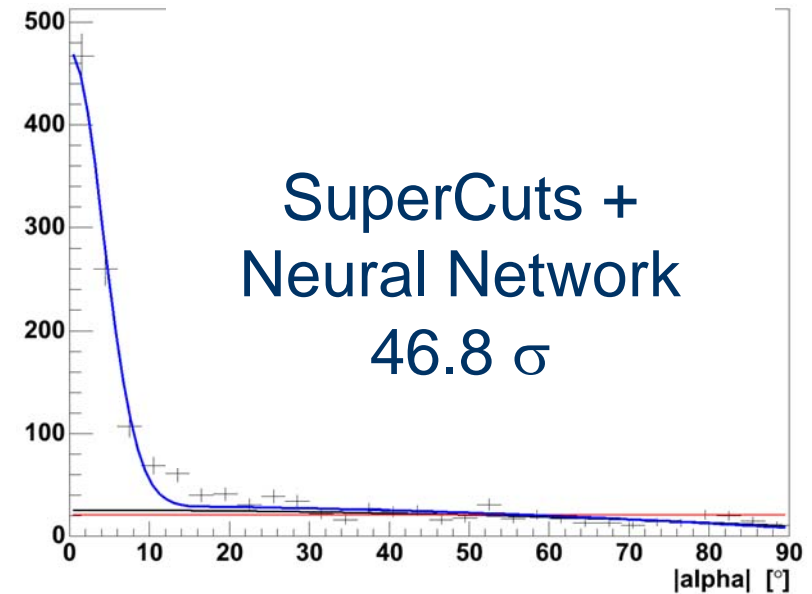
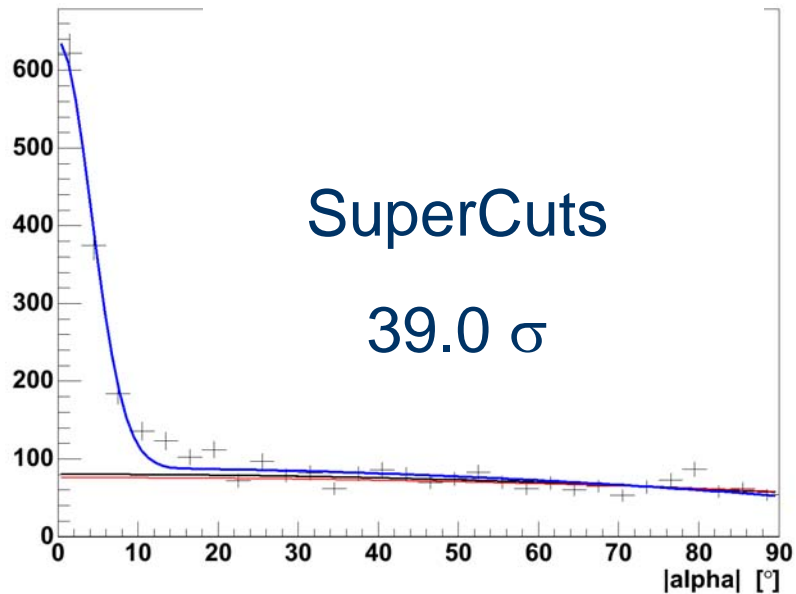
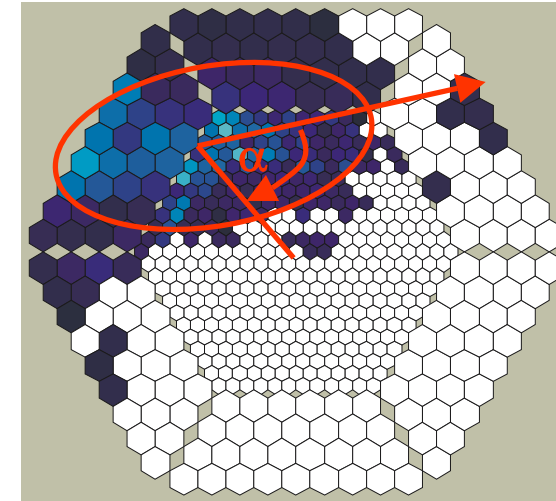
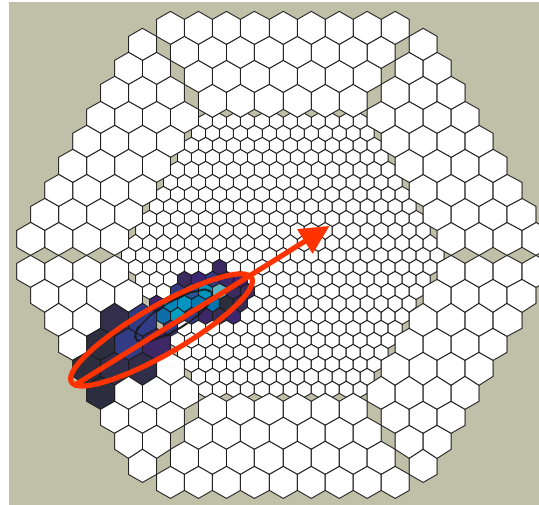
classical algorithm
„XMM“



MAGIC - Gamma/Hadron Separation



Observation
Mkn421
22.04.04



Summary

- Three classes of statistical learning methods

- Decision Trees (Bagging)

C4.5

CART

Random Forest

- Local Density Estimators

Maximum Likelihood

k-Nearest-Neighbour

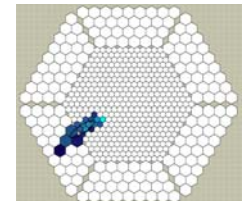
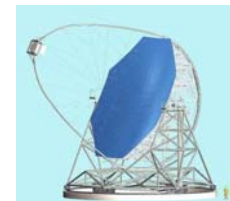
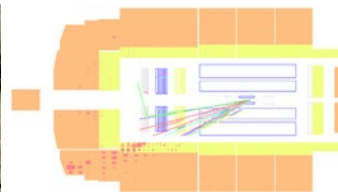
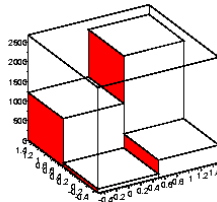
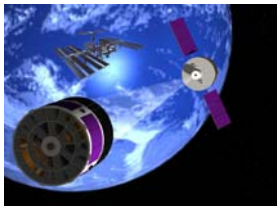
- Linear Separation

Linear Discriminant Analysis

Neural Networks

Support Vector Machines

- Many applications in current experiments and analysis



- Compared to classical methods usually at least small improvements

Conclusion

Up To You

Recognize possible applications

Intelligent Trigger? Purification?

Control Training and Evaluation

Overtraining? Systematic Uncertainties?

Compare classical and different
statistical learning methods

Time? Performance? Convenience?

Win efficiency / significance



Theory of Communication: Minimum Description Length Principle

Hypothesis H and Data D

$$P(H | D) \cdot P(D) = P(H \cap D) = P(D | H) \cdot P(H)$$

$$P(H | D) = \frac{P(H)P(D | H)}{P(D)}$$

Bayes



18th century

Our hypothesis should have the maximum probability given the data:

maximize $P(H | D)$

maximize $P(H)P(D | H)$

maximize $\log P(H) + \log P(D | H)$



Shannon 1948

$$\begin{array}{l} \leftarrow \\ \downarrow \end{array} I(X) \propto -\log P(X)$$

→ MDLP minimize $I(H) + I(D | H)$



Rissanen
1990

Statistical Learning Theory: Structural Risk Minimization

We have N training events with input x_i and correct output y_i

$$R_{emp}(\alpha) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i, \alpha))^2 \quad R(\alpha) = \int (y - f(x, \alpha))^2 dP_{X,Y}$$

empirical risk actual risk

Relationship is uniform convergence:

$$P\left(\sup_{\alpha} |R(\alpha) - R_{emp}(\alpha)| > \varepsilon\right) \xrightarrow{N \rightarrow \infty} 0$$

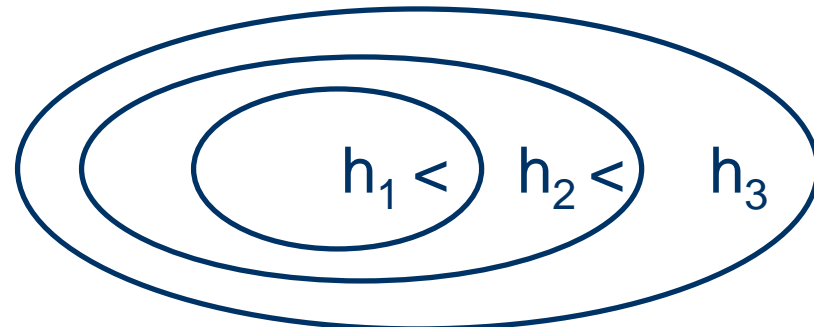
Upper bound for the actual risk (Vapnik):



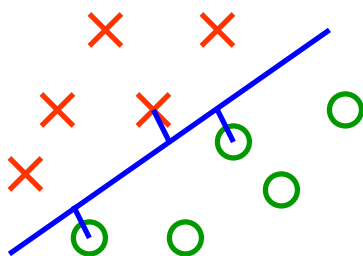
1996

$$R(\alpha) \leq R_{emp}(\alpha) + \varepsilon(N, p, h) \quad h: \text{VC Dimension of learning method (complexity)}$$

→ Create nested subsets of function spaces with rising complexity



Support Vector Machines



Separating hyperplane with maximum distance to each datapoint: Maximum margin classifier

Found by setting up condition for correct classification $y_i(\vec{w}\vec{x}_i + b) \geq 1$ and minimizing $\|\vec{w}\|^2$ which leads to the Lagrangian

$$L = \frac{1}{2}\|\vec{w}\|^2 - \sum \alpha_i (y_i(\vec{w}\vec{x}_i + b) - 1)$$

Necessary condition for a minimum is $\vec{w} = \sum \alpha_i y_i \vec{x}_i$

So the output becomes $out = \text{sgn}\left(\sum \alpha_i y_i \vec{x}_i \vec{x} + b\right)$ KKT: only SV have $\alpha_i \neq 0$

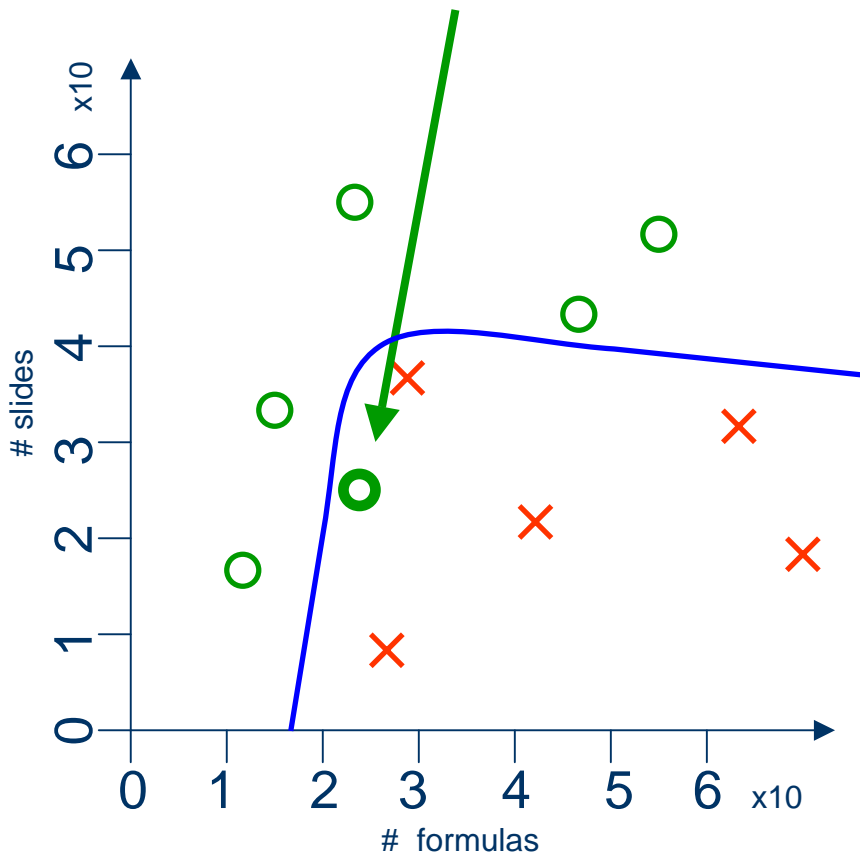
Only linear separation? No! Replace dot products: $\vec{x}\vec{y} \rightarrow \Phi(\vec{x})\Phi(\vec{y})$

The mapping to feature space $\Phi: R^d \rightarrow F$
is hidden in a kernel $K(\vec{x}, \vec{y}) = \Phi(\vec{x})\Phi(\vec{y})$

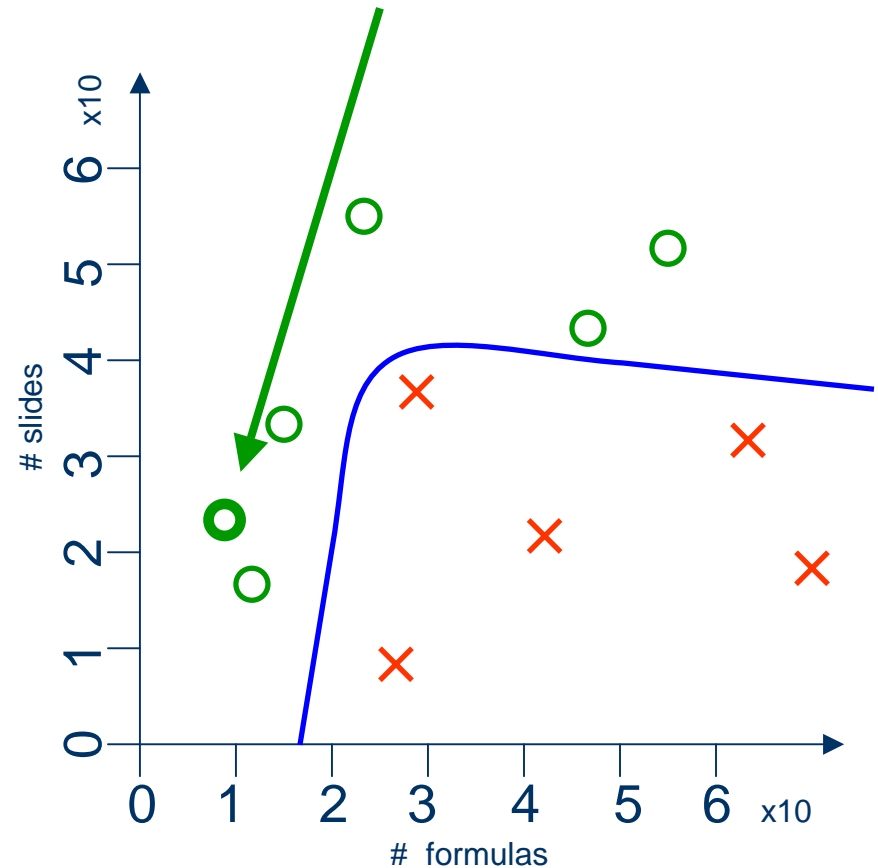
Non-separable case:
 $\frac{1}{2}\|\vec{w}\|^2 \rightarrow \frac{1}{2}\|\vec{w}\|^2 + C \sum \xi_i$

Finally

Include Statistical Learning Theory:
~25 formulas on 27 slides



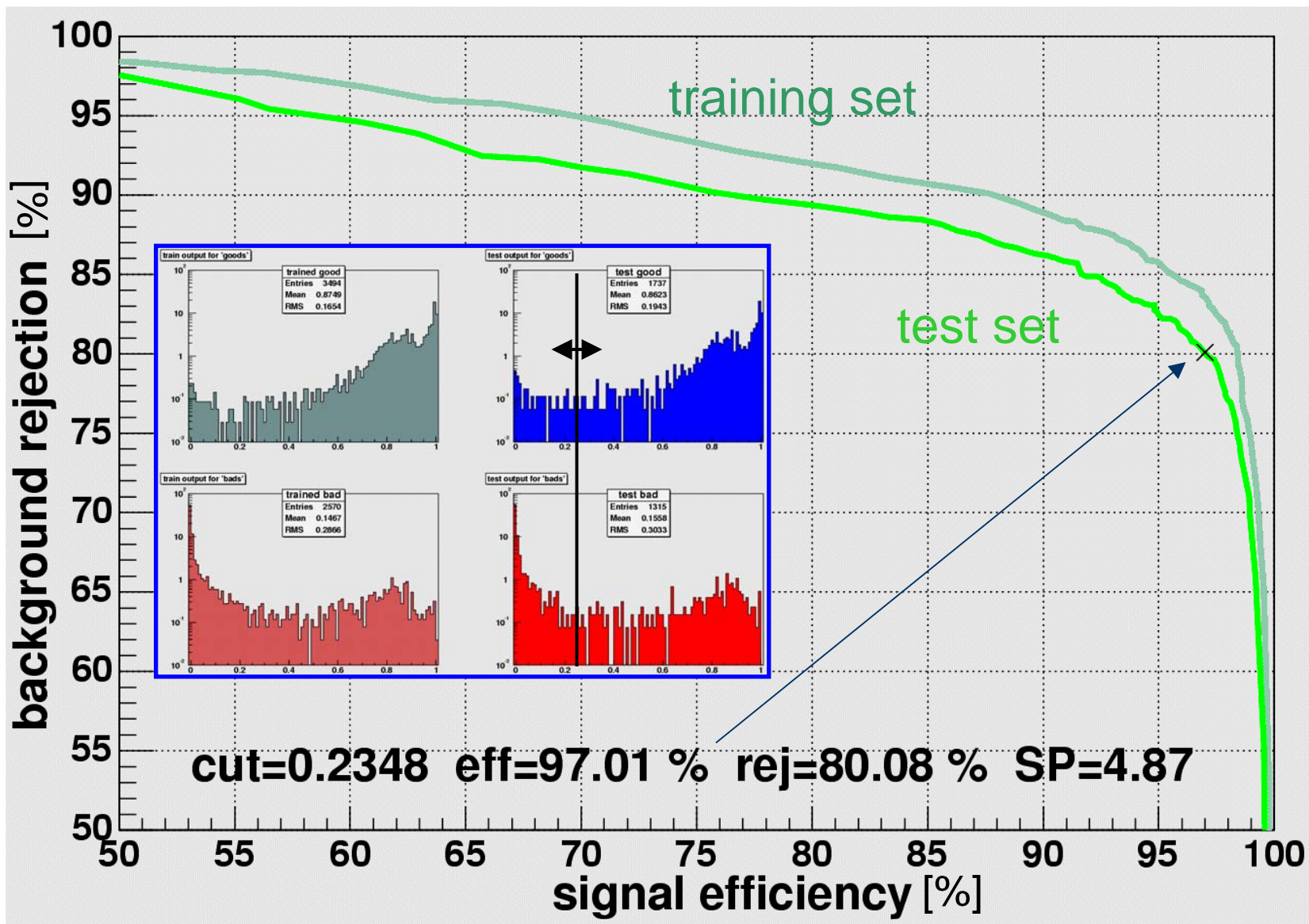
Skip Theory:
~7 formulas on 24 slides



Backups

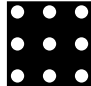


Deeply Virtual Compton Scattering (S41) - Efficiency/Rejection



Supervised Learning needs Training Data

Consequences

- ◆ Correct output t_i is known for every learning sample because
 - Data was generated by simulation
 - Data was taken from experiment under specific conditions 
- ◆ Correct output t_i is not known but very well approximated by some complicated analysis

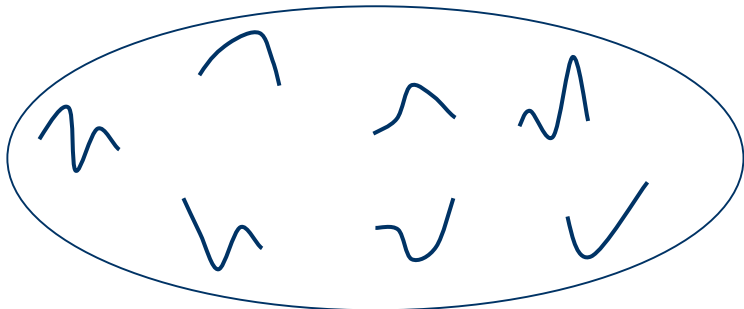
NN may train on unwanted simulation-features

Difficult / bad precision

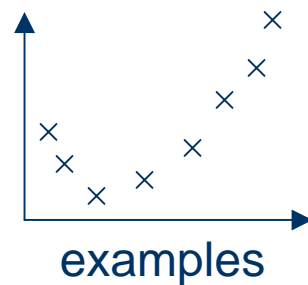
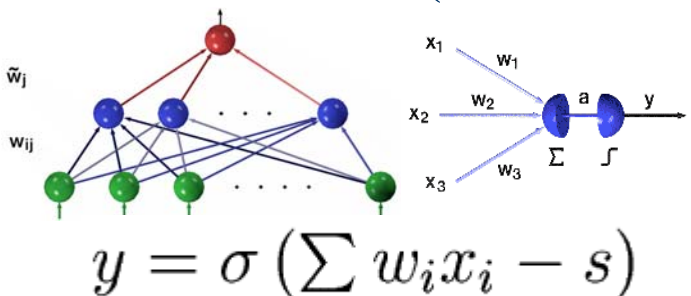
NN will never outperform this analysis

Statistical Learning vs. Parametric Models

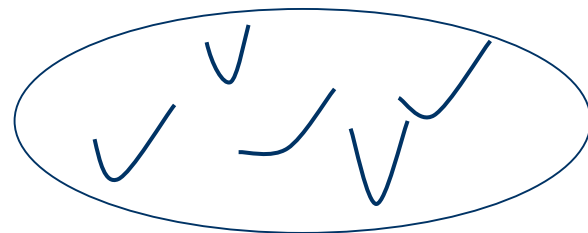
large function space
given by learning method
„Training“



Gradient descent (NN)
QP-Solver (SVM)
Weighted counting (LDE)
Find Max/Min (ML,k-NN)



small function space
given by user (**knowledge!**)
„Fitting“

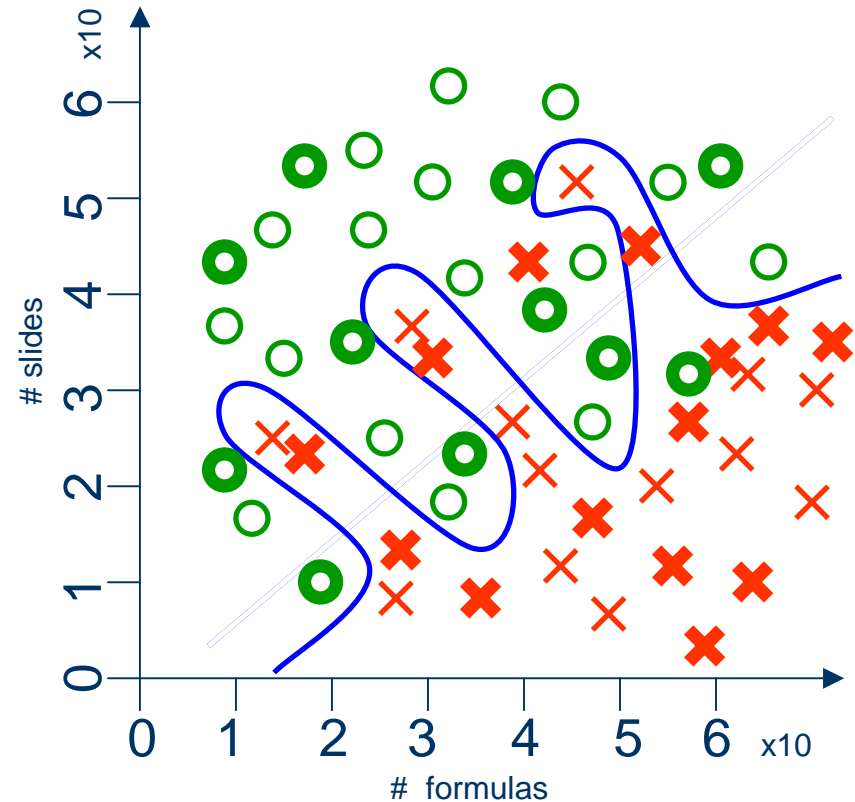
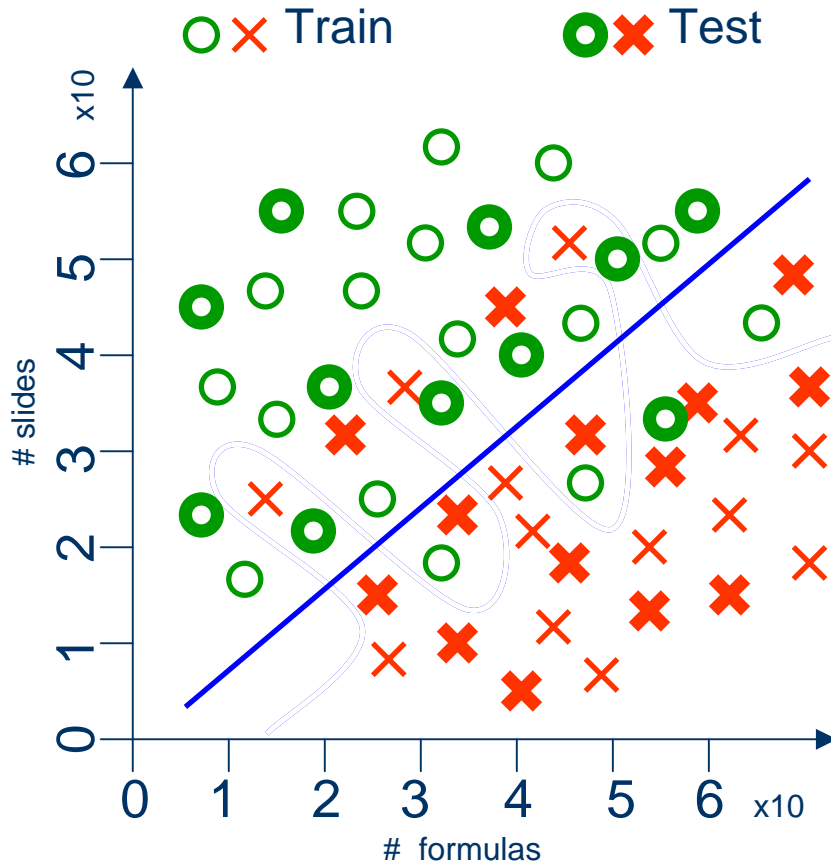


Minimize error:
e.g. least squares fit
→ use e.g. Minuit

$$E_r = \left(a + b \cdot \text{SIZE} + c \cdot \text{IMPACT} + d \cdot \text{LENGTH} + e \cdot \text{WIDTH} + f \cdot \text{SIZE} \times \text{LENGTH} + g \cdot \text{SIZE} \times \text{WIDTH} \right) \cdot \frac{\cos^{-2} \theta + h \cdot \cos^{-4} \theta}{1 + h}$$

$$\frac{1}{N} \sum_{i=1}^N \frac{(E_i - E_r)^2}{E_i}$$

See Overtraining - Want Generalization → Need Regularization

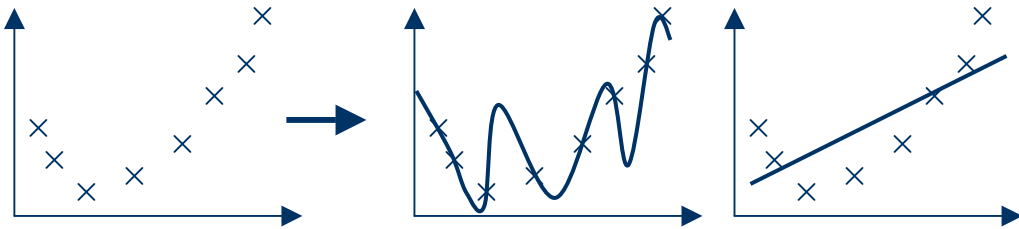


$$E = \frac{1}{N} \sum_{i=1}^N (y_i - out(x_i))^2$$

We want to tune a parameter of the learning algorithm depending on the overtraining we see!

Overtraining

- Too many free parameters
- Too large function space
- Learning „by heart“

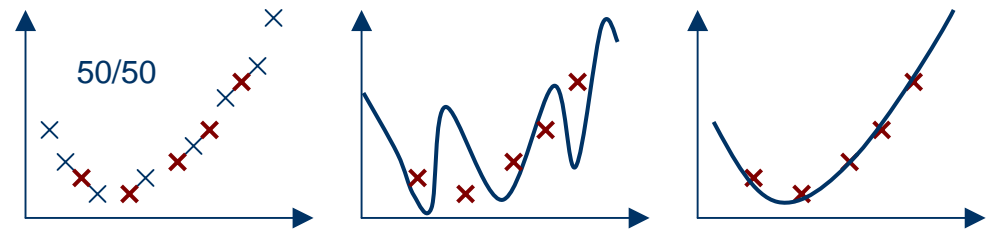


Too high degree of polynomial results in interpolation but too low degree means bad approximation

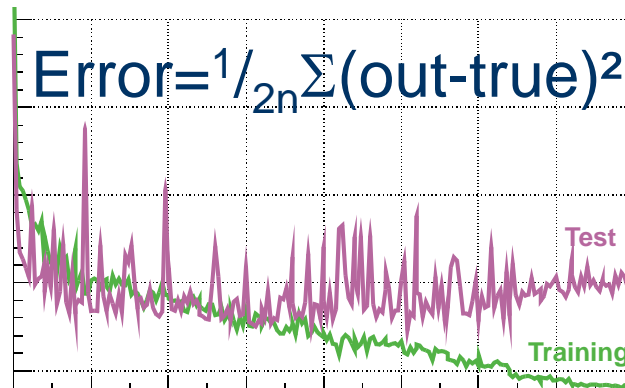
→ MDLP

Control the behaviour:

Split learning samples randomly into training samples \times and test (validation) samples \times



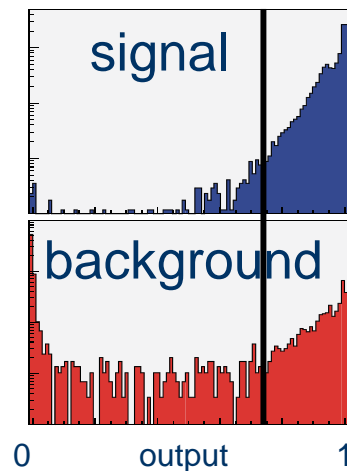
XEUS –
Classification of
photon events with
limited statistics



Overtraining in different
learning methods:
NN: # hidden neurons
LDE: width of local region
SVM: # support vectors

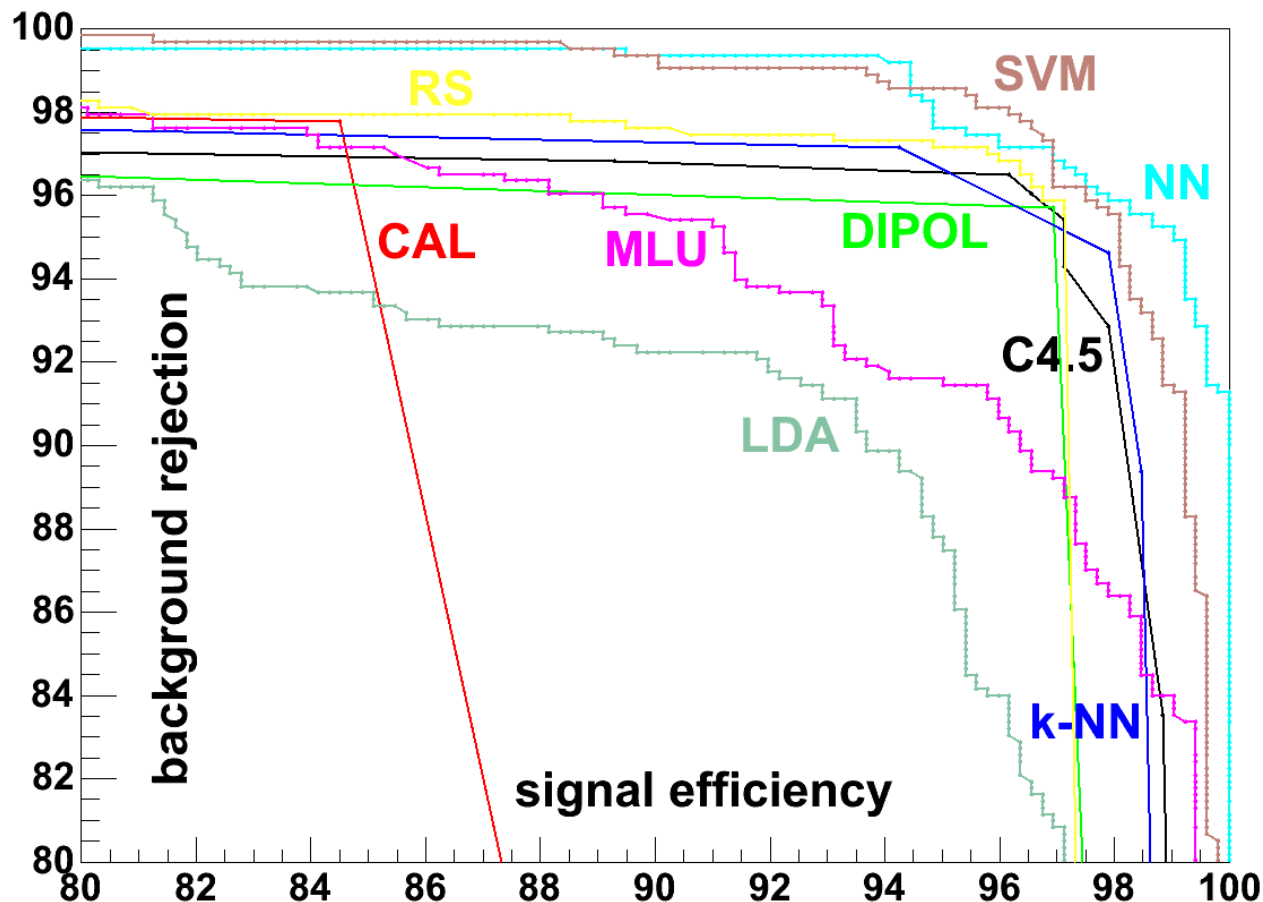
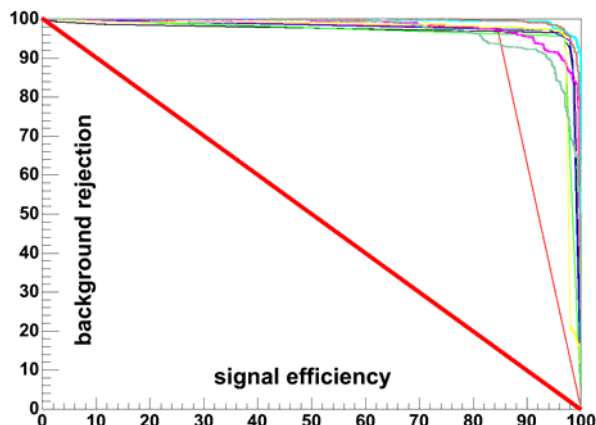
→ VC-Dim

Performance Measurement - Classification



Eff@Rej = xx%
Rej@Eff = xx%

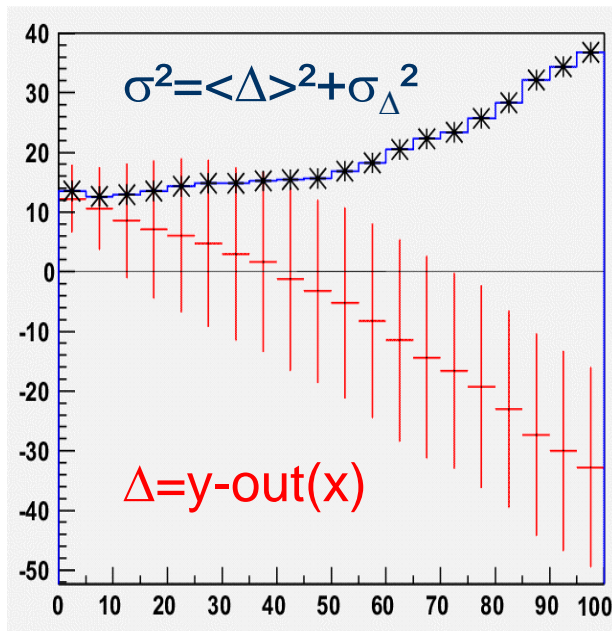
Misclassification =
 200%-Eff-Rej



Performance Measurement - Regression

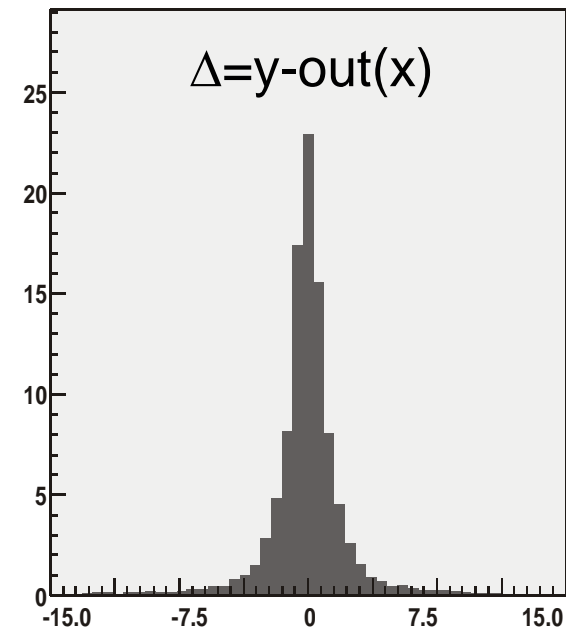
$$E = \frac{1}{N} \sum_{i=1}^N (y_i - out(x_i))^2$$

$$E = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - out(x_i)}{y_i} \right)^2$$

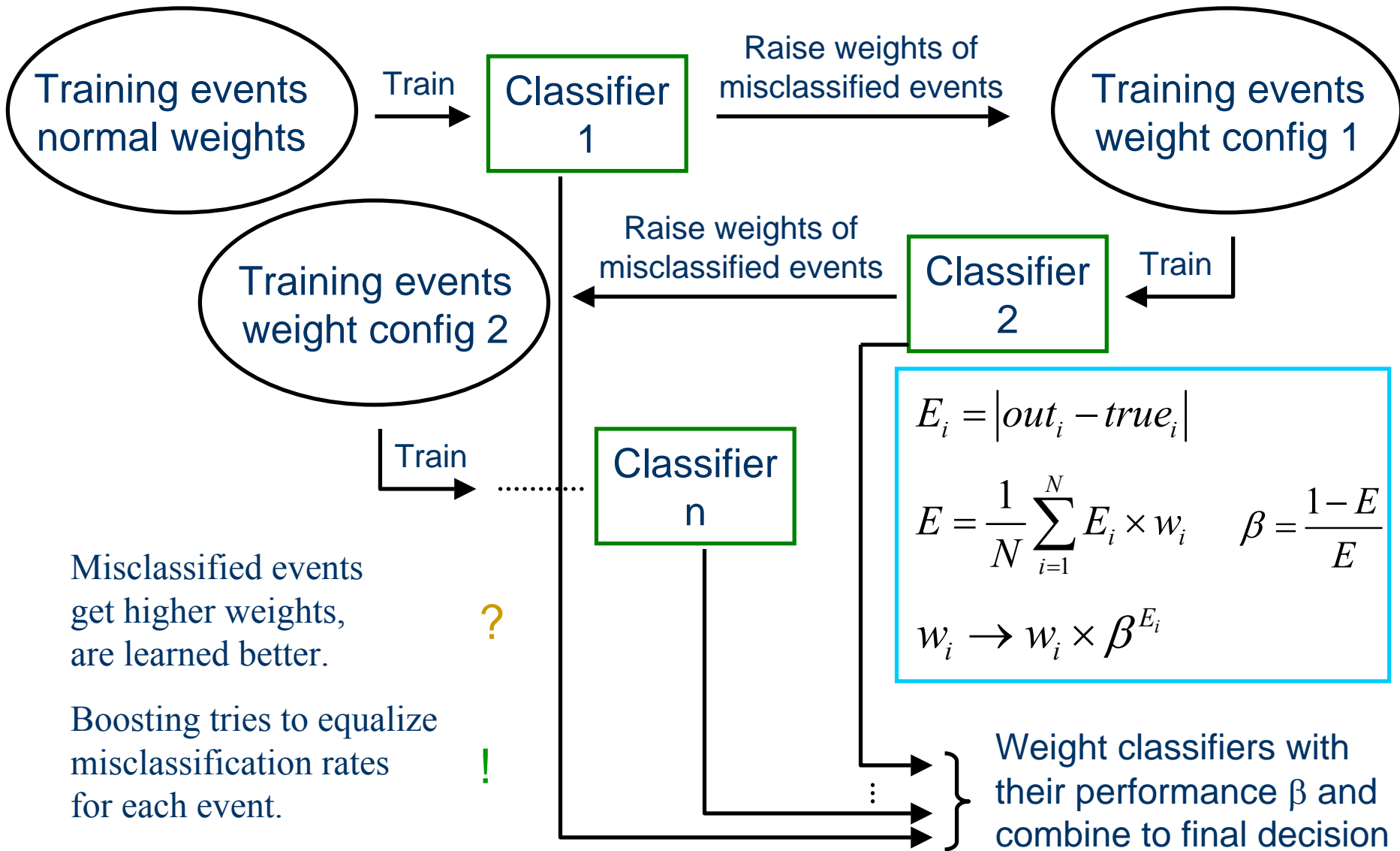


$$m = \langle y_i - out(x_i) \rangle$$

$$s^2 = \langle (y_i - out(x_i))^2 \rangle - m^2$$

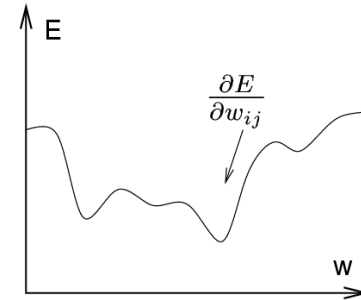


Boosting – Procedure

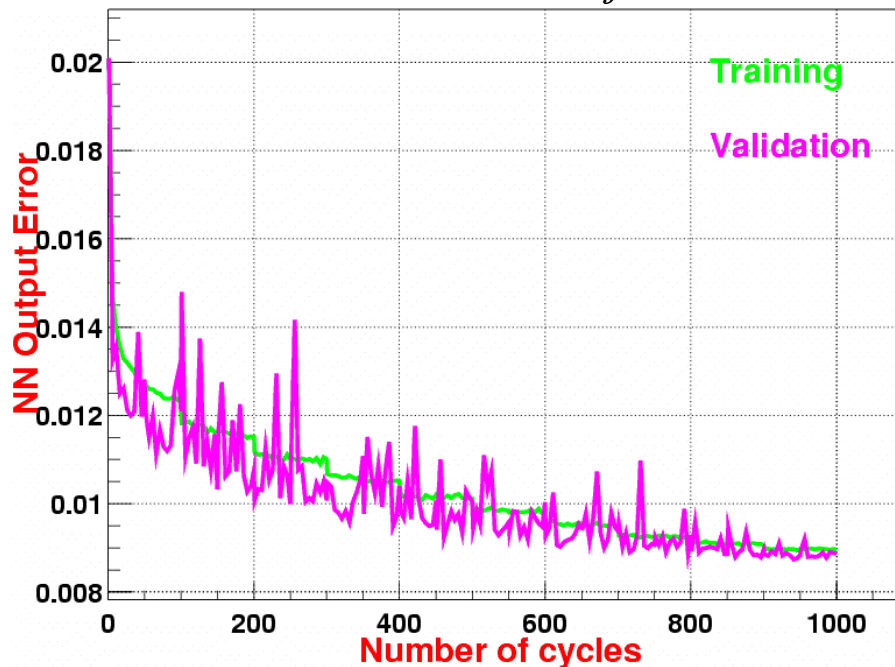


Training Neural Networks – Theory

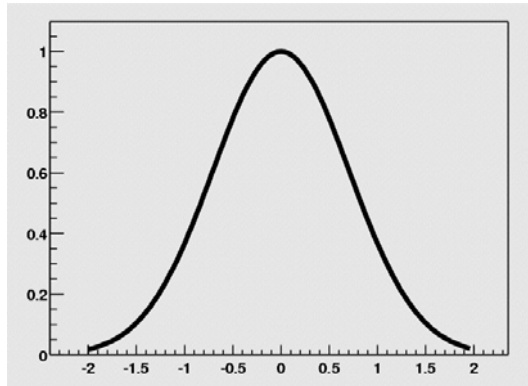
$$E = \frac{1}{N} \sum_{i=1}^N (y_i - \text{out}(x_i))^2$$



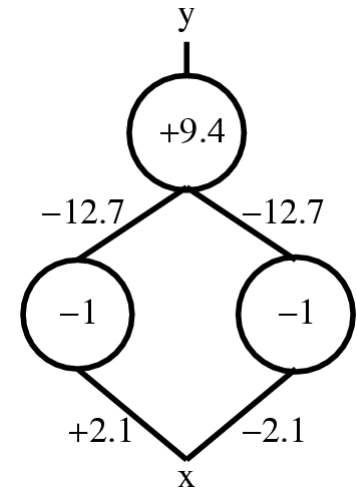
$$\Delta w_{ij}(t+1) = -\eta \frac{\partial E}{\partial w_{ij}} + \alpha \Delta w_{ij}(t)$$



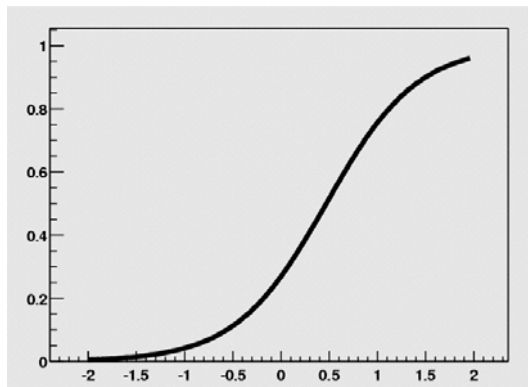
Regression with Neural Networks



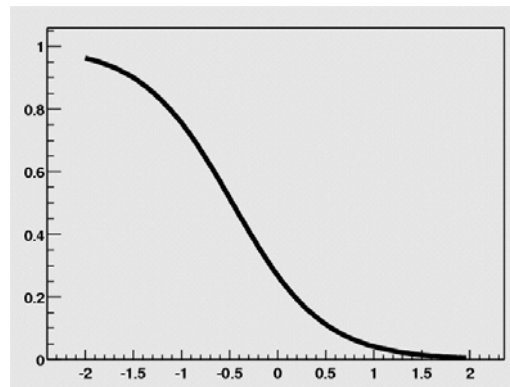
- ◆ Example: Train Gauß with two hidden neurons
- ◆ Minimization of error during training is **scaling and shifting the two sigmoids**
scaling and shifting their sum to fit the function



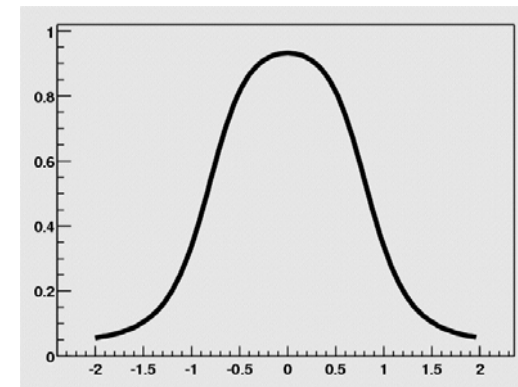
$$\sigma(+2.1x - 1)$$



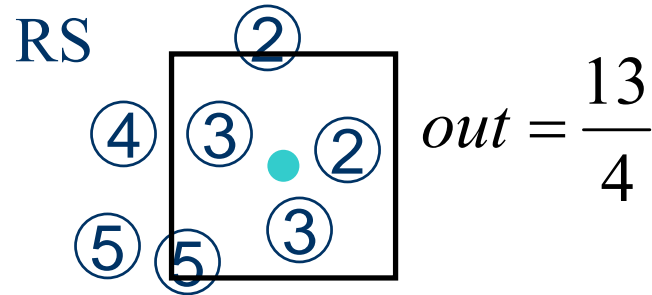
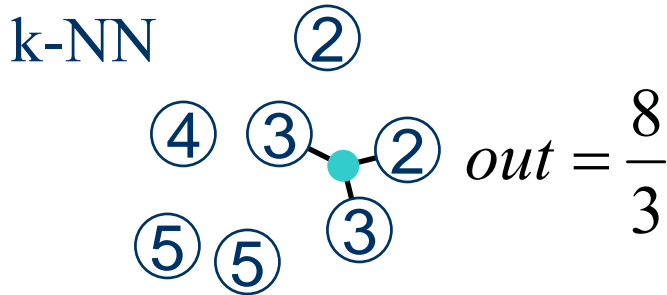
$$\sigma(-2.1x - 1)$$



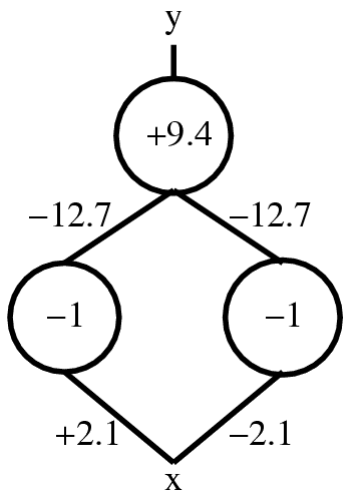
output



From Classification to Regression



NN

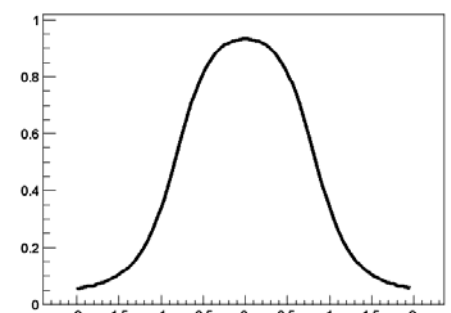
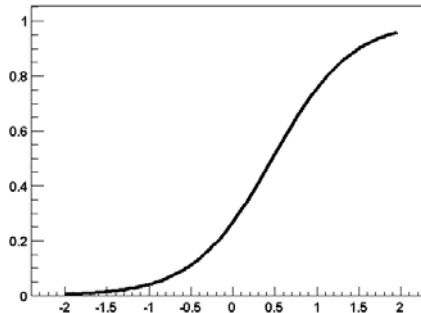
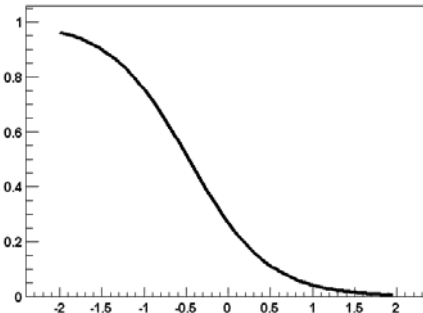


$$E = \frac{1}{N} \sum_{i=1}^N (y_i - out(x_i))^2$$

$$a = \sigma(-2.1x - 1)$$

$$b = \sigma(+2.1x - 1)$$

$$out = \sigma(-12.7a - 12.7b + 9.4)$$



Fit Gauss

