

Multivariate analysis

Yann Coadou

CPPM Marseille

esipap...

European School of Instrumentation
in Particle & Astroparticle Physics

Archamps, 29 January 2014

Thanks to Harrison Prosper, Balázs Kégl, Jérôme Schwindling, Jan Therhaag



- 1 **Introduction**
- 2 **Optimal discrimination**
 - Bayes limit
 - Multivariate discriminant
- 3 **Machine learning**
 - Supervised and unsupervised learning
 - Example
- 4 **Multivariate discriminants**
 - Random grid search
 - Genetic algorithms
 - Quadratic and linear discriminants
 - Support vector machines
 - Kernel density estimation
 - Neural networks
 - Bayesian neural networks
 - Decision trees
- 5 **Summary**

Typical problems in HEP

- Classification of objects
 - separate real and fake leptons/jets/etc.
- Signal enhancement relative to background
- Regression: best estimation of a parameter
 - lepton energy, \cancel{E}_T value, invariant mass, etc.

Discrimination of signal from background in HEP

- Event level (Higgs searches, ...)
- Cone level (tau-vs-jet reconstruction, ...)
- Lifetime and flavour tagging (*b*-tagging,)
- Track level (particle identification, ...)
- Cell level (energy deposit from hard scatter/pileup/noise, ...)

Input information from various sources

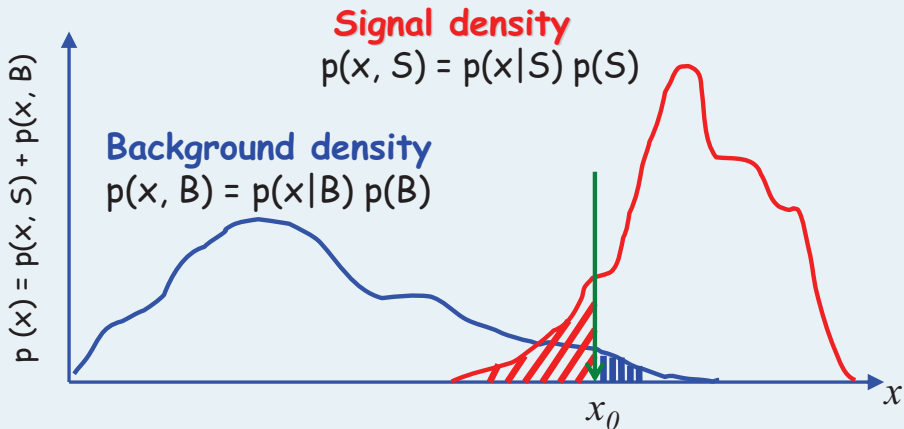
- Kinematic variables (masses, momenta, decay angles, ...)
- Event properties (jet multiplicity, sum of charges, brightness ...)
- Event shape (sphericity, aplanarity, ...)
- Detector response (silicon hits, dE/dx , Cherenkov angle, shower profiles, muon hits, ...)

Most data are (highly) multidimensional

- Use dependencies between $x = \{x_1, \dots, x_n\}$ discriminating variables
- Approximate this n -dimensional space with a function $f(x)$ capturing the essential features
- f is a **multivariate discriminant**
- For most of these lectures, use binary classification:
 - an object belongs to one class (e.g. signal) if $f(x) > q$, where q is some threshold,
 - and to another class (e.g. background) if $f(x) \leq q$

For simplicity: 1-dimension case

- Where to place a cut x_0 on variable x ?



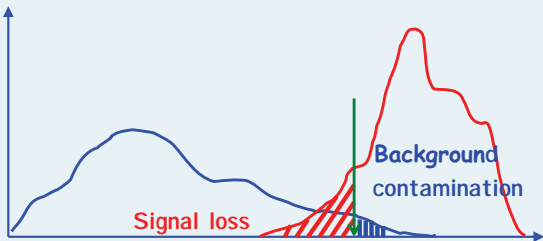
- Optimal choice: minimum misclassification cost at decision boundary
 $x = x_0$

Cost of misclassification

$$C(x_0) = C_S \int H(x_0 - x)p(x, S)dx \quad \text{signal loss}$$
$$+ C_B \int H(x - x_0)p(x, B)dx \quad \text{background contamination}$$

C_S = cost of misclassifying signal as background

C_B = cost of misclassifying background as signal



- $H(x)$: Heaviside step function
- $H(x) = 1$ if $x > 0$,
0 otherwise

- Optimal choice: when cost function C is minimum

Minimising the cost

- Minimise

$$C(x_0) = C_S \int H(x_0 - x)p(x, S)dx + C_B \int H(x - x_0)p(x, B)dx$$

with respect to the boundary x_0 :

$$\begin{aligned} 0 &= C_S \int \delta(x_0 - x)p(x, S)dx - C_B \int \delta(x - x_0)p(x, B)dx \\ &= C_S p(x_0, S) - C_B p(x_0, B) \end{aligned}$$

- This gives the **Bayes discriminant**:

$$BD = \frac{C_B}{C_S} = \frac{p(x_0, S)}{p(x_0, B)} = \frac{p(x_0|S)p(S)}{p(x_0|B)p(B)}$$

Probability relationships

- $p(A, B) = p(A|B)p(B) = p(B|A)p(A)$
- Bayes theorem: $p(A|B)p(B) = p(B|A)p(A)$
- $p(S|x) + p(B|x) = 1$

Generalising to multidimensional problem

- The same holds when x is an n -dimensional variable:

$$BD = B \frac{p(S)}{p(B)} \quad \text{where} \quad B = \frac{p(x|S)}{p(x|B)}$$

- B is the **Bayes factor**, identical to the likelihood ratio when class densities $p(x|S)$ and $p(x|B)$ are independent of unknown parameters

Bayes limit

- $p(S|x) = BD/(1 + BD)$ is what should be achieved to minimise cost, achieving classification with the fewest mistakes
- Fixing relative cost of background contamination and signal loss $q = C_B/(C_S + C_B)$, $q = p(S|x)$ defines decision boundary:
 - signal-rich if $p(S|x) \geq q$
 - background-rich if $p(S|x) < q$
- Any function that approximates conditional class probability $p(S|x)$ with negligible error reaches the **Bayes limit**

How to construct $p(S|x)$?

- $k = p(S)/p(B)$ typically unknown
- Problem: $p(S|x)$ depends on k !
- Solution: it's not a problem...
- Define a **multivariate discriminant**:

$$D(x) = \frac{s(x)}{s(x) + b(x)} = \frac{p(x|S)}{p(x|S) + p(x|B)}$$

- Now:

$$p(S|x) = \frac{D(x)}{D(x) + (1 - D(x))/k}$$

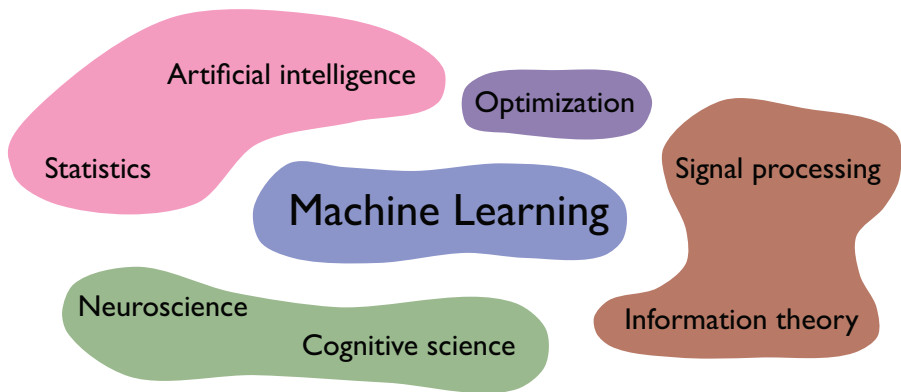
- Cutting on $D(x)$ is equivalent to cutting on $p(S|x)$, implying a corresponding (unknown) cut on $p(S|x)$

Several types of problems

- Classification/decision:
 - signal or background
 - type Ia supernova or not
 - will pay his/her credit back on time or not
- Regression (mostly ignored in these lectures)
- Clustering (cluster analysis):
 - in exploratory data mining, finding features

Our goal

- Teach a machine to learn the discriminant $f(x)$ using examples from a training dataset
- Be careful to not learn too much the properties of the training sample
 - no need to memorise the training sample
 - instead, interested in getting the right answer for new events
⇒ generalisation ability



©Balázs Kégl

Supervised learning

- Training events are labelled: N examples $(x, y)_1, (x, y)_2, \dots, (x, y)_N$ of discriminating variables x (also called **feature variables**) and **class labels** y
- The learner uses example classes to know how good it is doing

Unsupervised learning

- Clustering: find similarities in training sample, without having predefined categories (how Amazon is recommending you books or DVDs...)
- Instead of categories, some sort of reward system. May not even “learn” anything from data, but remembers what triggers reward or punishment
- Not biased by pre-determined classes \Rightarrow may discover unexpected features!

A “giant” neural network

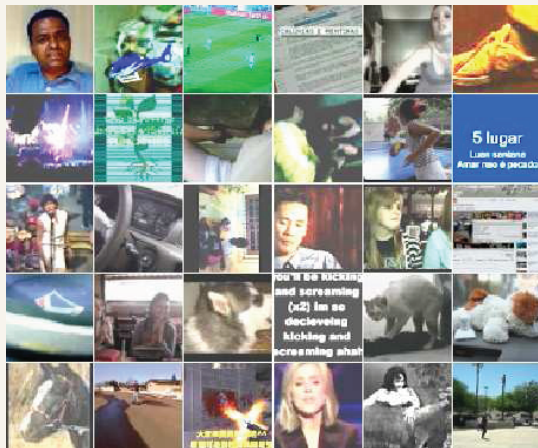
- At Google they trained a 9-layered NN with 1 billion connections
 - trained on 10 million 200×200 pixel images from YouTube videos
 - on 1000 machines (16000 cores) for 3 days, unsupervised learning
- Sounds big? The human brain has 100 billion (10^{11}) neurons and 100 trillion (10^{14}) connections...

What it did

- It learned to recognise faces, one of the original goals
- ... but also cat faces (among the most popular things in YouTube videos) and body shapes



Features extracted from such images



- Results shown to be robust to
 - colour
 - translation
 - scaling
 - out-of-plane rotation

Finding the multivariate discriminant $y = f(x)$

- Given our N examples $(x, y)_1, \dots, (x, y)_N$ we need
 - a function class $\mathbb{F} = \{f(x, w)\}$ (w : parameters to be found)
 - a constraint $Q(w)$ on \mathbb{F}
 - a loss or error function $L(y, f)$, encoding what is lost if f is poorly chosen in \mathbb{F} (i.e., $f(x, w)$ far from the desired $y = f(x)$)
- Cannot minimise L directly (would depend on the dataset used), but rather its average over a training sample, the **empirical risk**:

$$R(w) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i, w))$$

subject to constraint $Q(w)$, so we minimise the **cost function**:

$$C(w) = R(w) + \lambda Q(w)$$

- At the minimum of $C(w)$ we select $f(x, w_*)$, our estimate of $y = f(x)$

Loss function in regression

- Goal: set $f(x, w)$ as close as possible to y
- Therefore, loss increases with difference between $f(x, w)$ and y
- Most widely used loss function is **quadratic loss**:

$$L(y, f) = (f(x, w) - y)^2$$

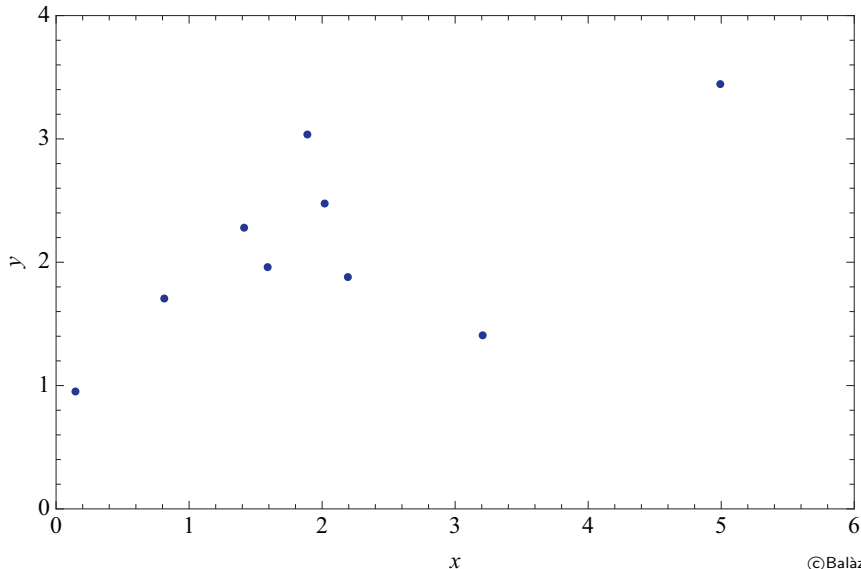
Loss function in classification

- There is no “distance” between classes
- Goal: $f(x, w)$ predicts properly class y
- Usual loss function is **one-loss** or **zero-one loss**:

$$L(y, f) = \mathbb{I}(f(x, w) \neq y)$$

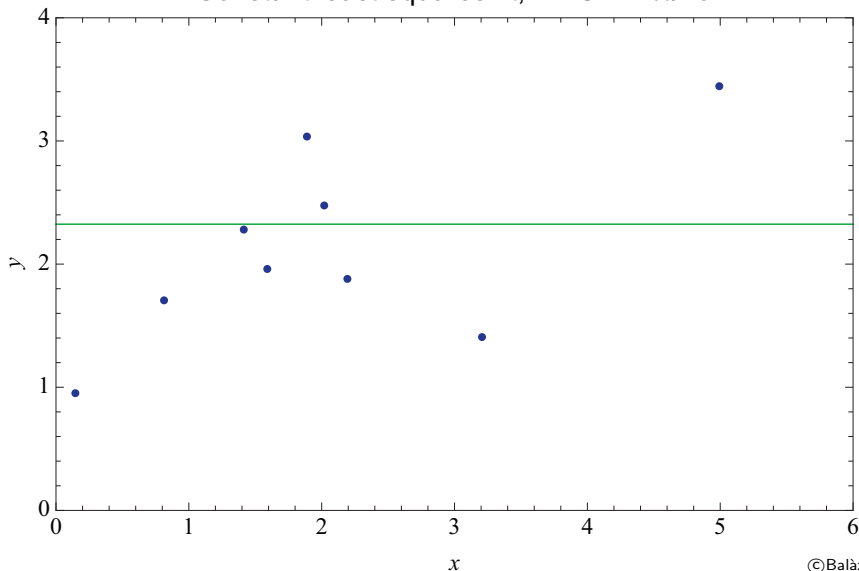
where indicator function $\mathbb{I}(X) = 1$ if X is true, 0 otherwise

Data generated from an unknown function with unknown noise



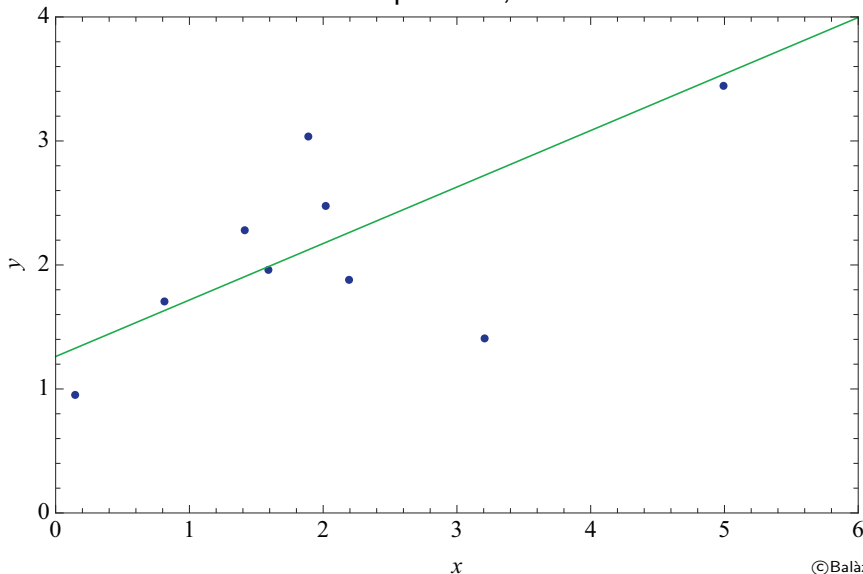
©Balázs Kégl

Constant least squares fit, RMSE = 0.915



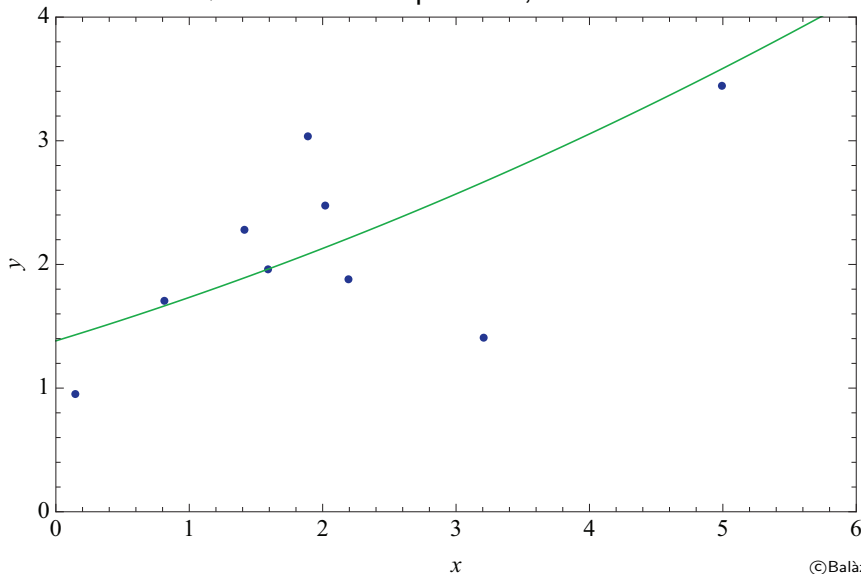
©Balázs Kégl

Linear least squares fit, RMSE = 0.581



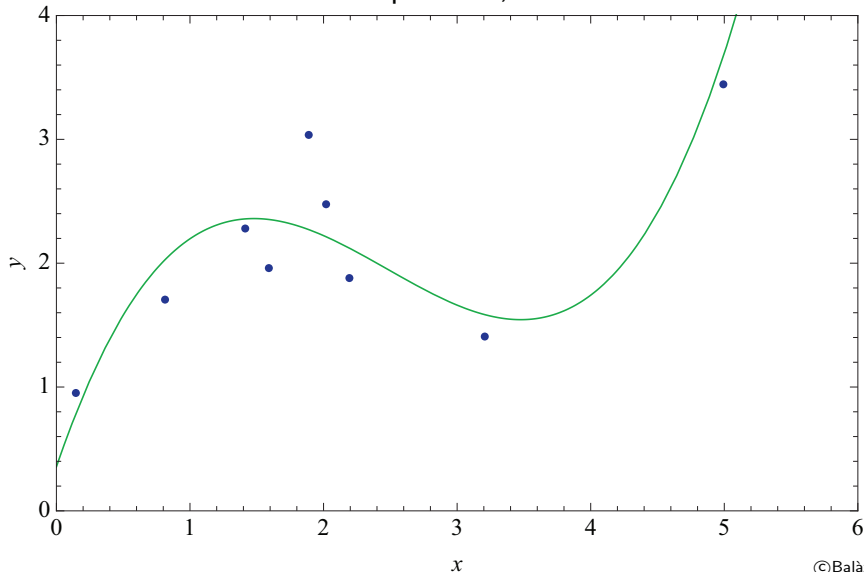
©Balázs Kégl

Quadratic least squares fit, RMSE = 0.579



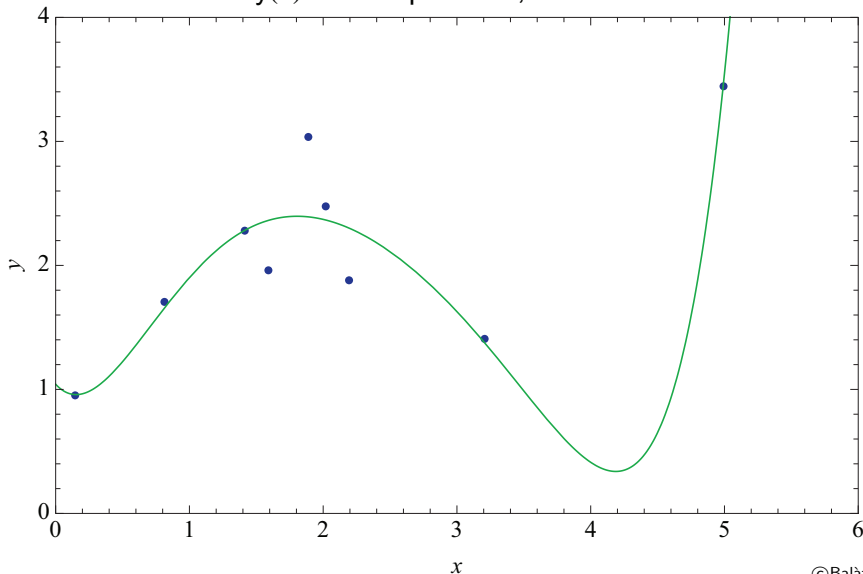
©Balázs Kégl

Cubic least squares fit, RMSE = 0.339



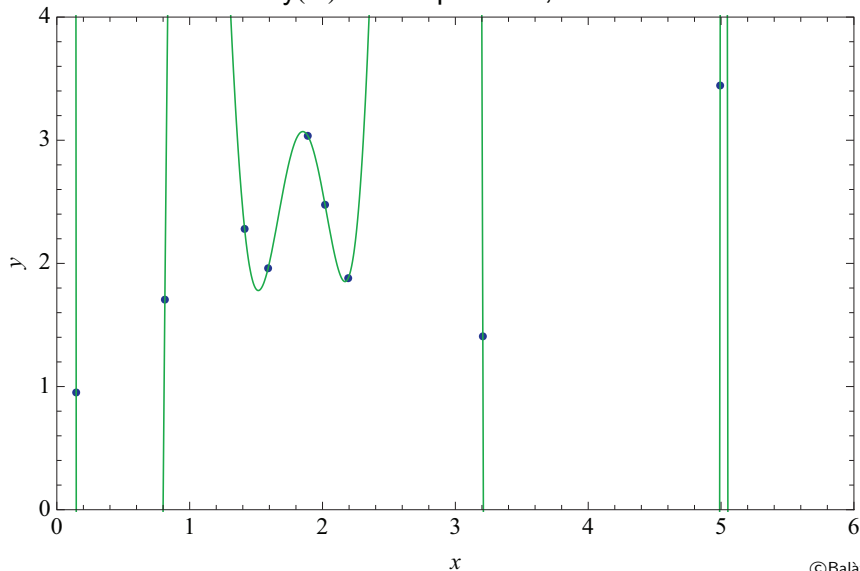
©Balázs Kégl

Poly(6) least squares fit, RMSE = 0.278



©Balázs Kégl

Poly(9) least squares fit, RMSE = 0



©Balázs Kégl

Quality of fit

- Increasing degree of polynomial increases flexibility of function
- Higher degree \Rightarrow can match to more features
- If degree = $\#$ points, polynomial passes through each point: perfect match!

Quality of fit

- Increasing degree of polynomial increases flexibility of function
- Higher degree \Rightarrow can match to more features
- If degree = $\#$ points, polynomial passes through each point: perfect match!

Is it meaningful?

- It could be:
 - if there is no noise or uncertainty in the measurement
 - if the true distribution is indeed perfectly described by such a polynomial
- . . . not impossible, but not very common . . .

Quality of fit

- Increasing degree of polynomial increases flexibility of function
- Higher degree \Rightarrow can match to more features
- If degree = $\#$ points, polynomial passes through each point: perfect match!

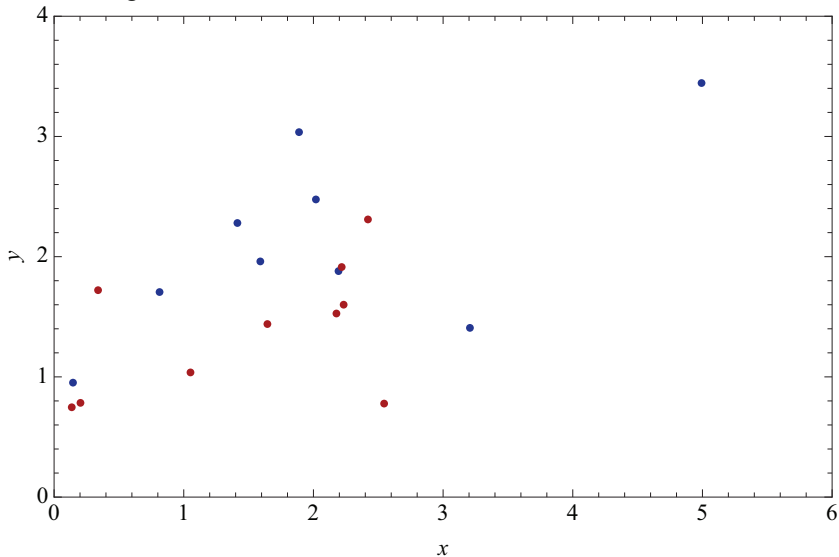
Is it meaningful?

- It could be:
 - if there is no noise or uncertainty in the measurement
 - if the true distribution is indeed perfectly described by such a polynomial
- ... not impossible, but not very common...

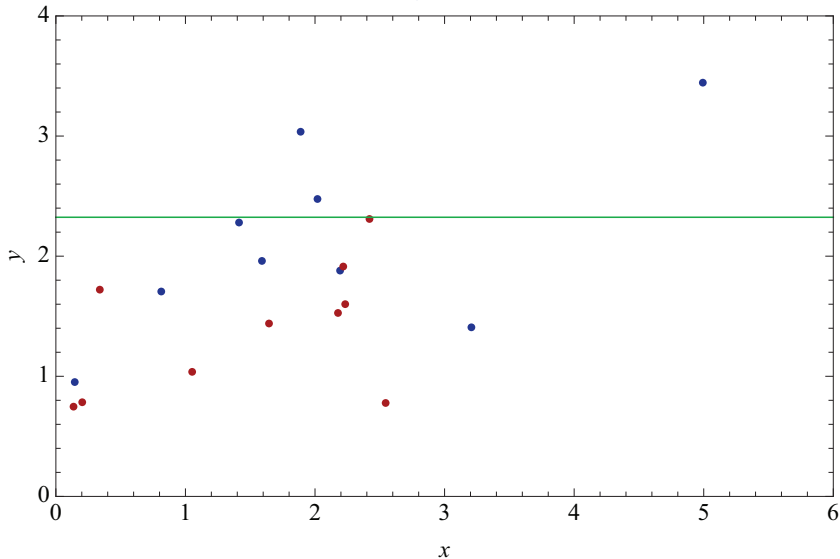
Solution: testing sample

- Use independent sample to validate the result
- Expected: performance will also increase, go through a maximum and decrease again, while it keeps increasing on the training sample

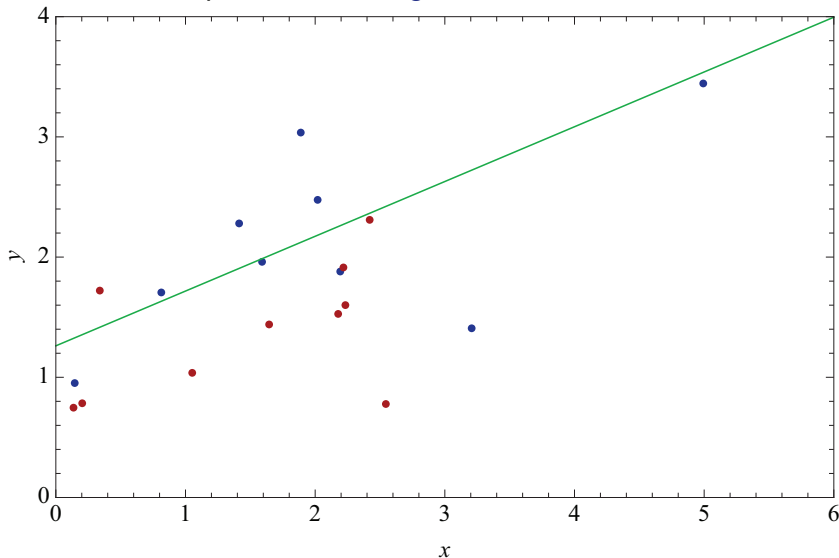
Data generated from an unknown function with unknown noise



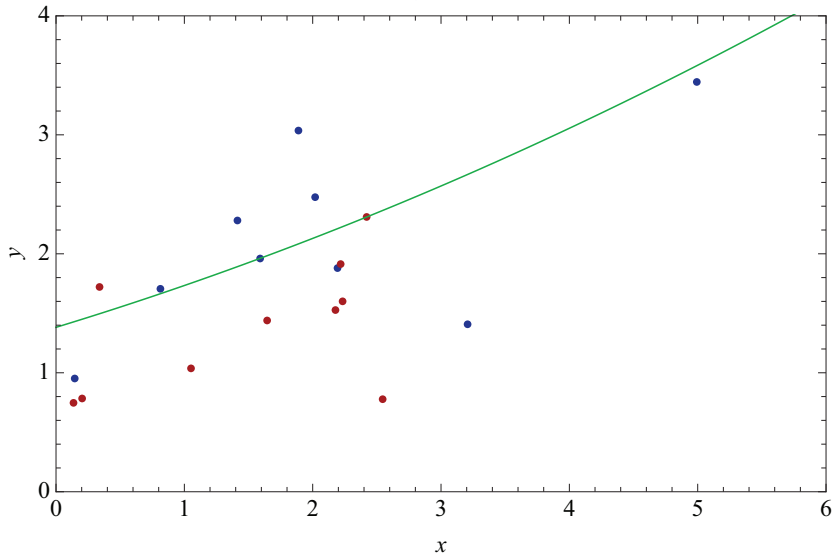
Const. least squares fit, **training** RMSE = 0.915, **test** RMSE = 1.067



Linear least squares fit, **training** RMSE = 0.581, **test** RMSE = 0.734

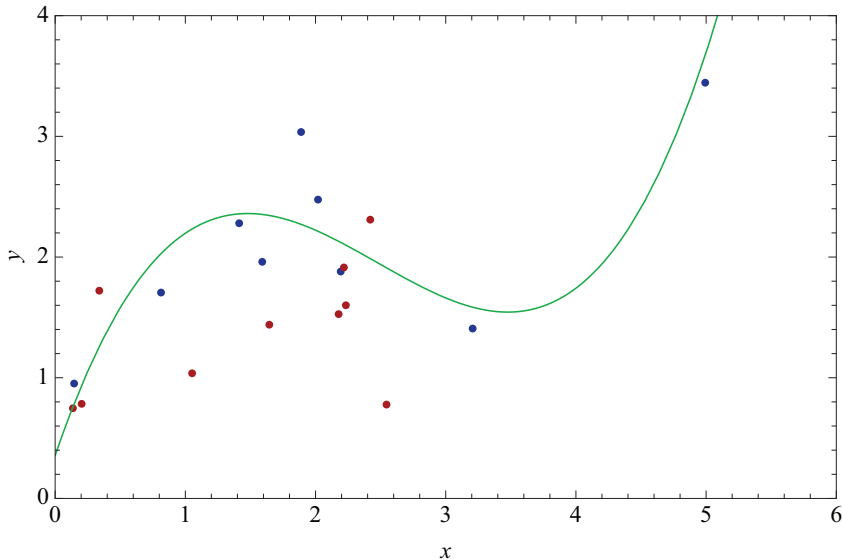


Quadr. least squares fit, **training** RMSE = 0.579, **test** RMSE = 0.723

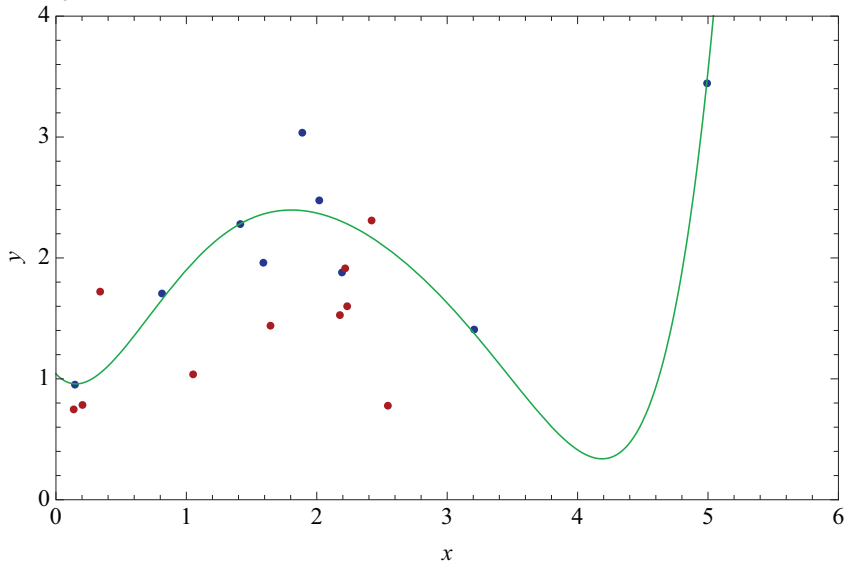


©Balázs Kégl

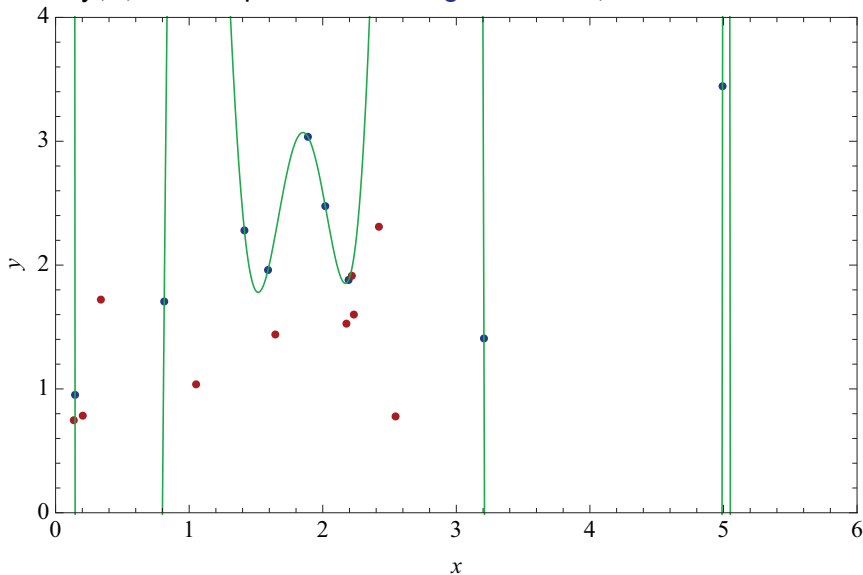
Cubic least squares fit, **training** RMSE = 0.339, **test** RMSE = 0.672



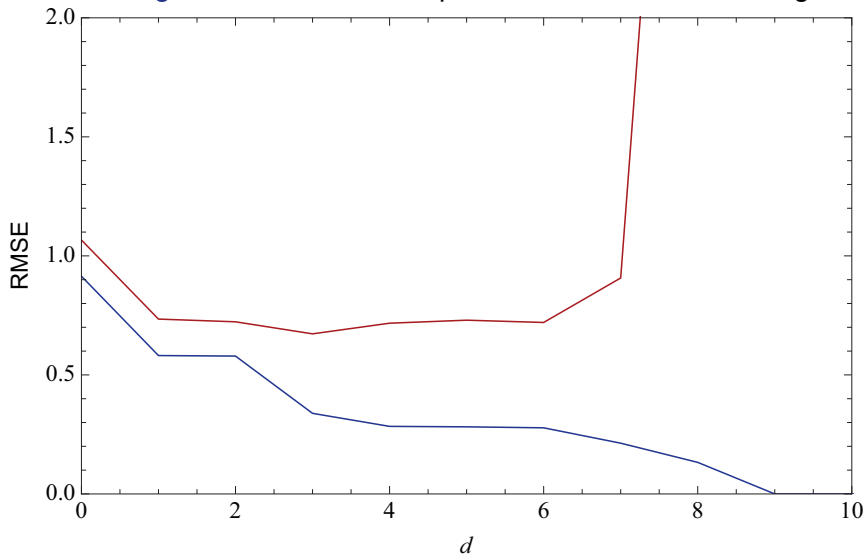
Poly(6) least squares fit, training RMSE = 0.278, test RMSE = 0.72



Poly(9) least squares fit, training RMSE = 0, test RMSE = 46.424



Training and test RMSE's for polinomial fits of different degrees



Non-parametric fit

- Minimising the training cost (here, RMSE) does not work if the function class is not fixed in advance (e.g. fix the polynomial degree): complete loss of generalisation capability!
- But if you do not know the correct function class, you should not fix it! Dilemma...

Capacity control and regularisation

- Trade-off between approximation error and estimation error
- Take into account sample size
- Measure (and penalise) complexity
- Use independent test sample
- In practice, no need to correctly guess the function class, but need enough flexibility in your model, balanced with complexity cost

1 Introduction

2 Optimal discrimination

- Bayes limit
- Multivariate discriminant

3 Machine learning

- Supervised and unsupervised learning
- Example

4 Multivariate discriminants

- Random grid search
- Genetic algorithms
- Quadratic and linear discriminants
- Support vector machines
- Kernel density estimation
- Neural networks
- Bayesian neural networks
- Decision trees

5 Summary

Reminder

- To solve binary classification problem with the fewest number of mistakes, sufficient to compute the multivariate discriminant:

$$D(x) = \frac{s(x)}{s(x) + b(x)}$$

where:

- $s(x) = p(x|S)$ signal density
- $b(x) = p(x|B)$ background density
- Cutting on $D(x)$ is equivalent to cutting on probability $p(S|x)$ that event with x values is of class S

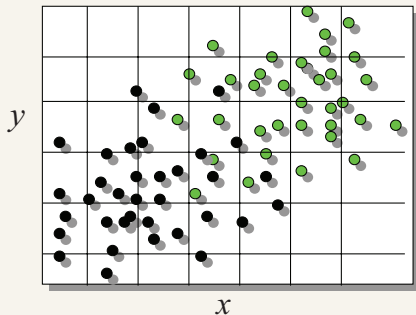
Which approximation to choose?

- Best possible choice: cannot beat Bayes limit (but usually impossible to define)
- No single method can be proven to surpass all others in particular case
- Advisable to try several and use the best one

Cut-based analysis

- Simple approach: cut on each discriminating variable
- Difficulty: how to optimise the cuts?

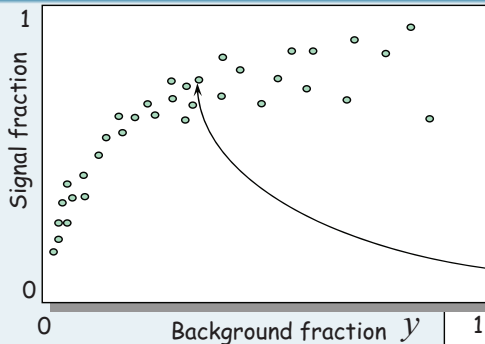
Grid search



©Harrison Prosper

- Split each variable in K values
- Apply cuts at each grid point:
 $x > x_i, y > y_i$
- Number of points scales with K^n : **curse of dimensionality**

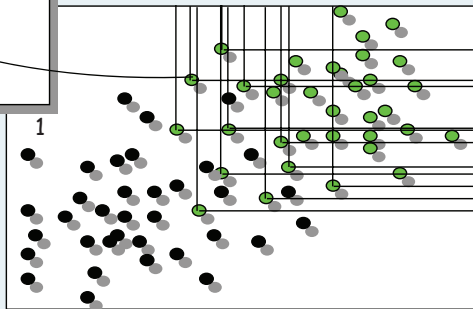
RGS



- Use each point in signal sample as grid point:

$$x > x_i$$

$$y > y_i$$

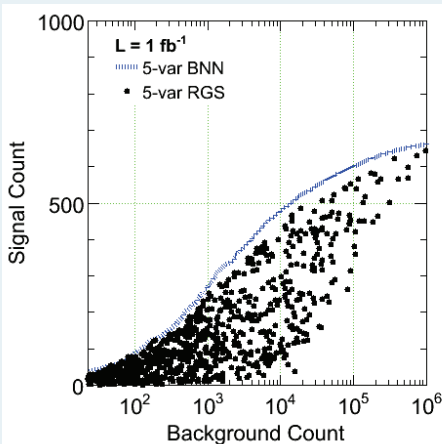


- Number of cut points independent of dimensionality
- Sampled points density follows signal density

x

©Harrison Prosper

Comparison to BNN



©Harrison Prosper

- Blue: 5-dim Bayesian neural network discriminant (see later)
- Points: each cut point from a 5-dim RGS calculation
- Conclusions:
 - RGS can find very good criteria with high discrimination
 - but it usually cannot compete with a full-blown multivariate discriminant
 - and never outsmarts it

Survival of the fittest

- Inspired by biological evolution
- Model: group (population) of abstract representations (genome/discriminating variables) of possible solutions (individuals/list of cuts)
- Typical processes at work in evolutionary processes:
 - inheritance
 - mutation
 - sexual recombination (a.k.a. crossover)
- Fitness function: value representing the individual's goodness, or comparison of two individuals
- For cut optimisation:
 - good background rejection and high signal efficiency
 - compare individuals in each signal efficiency bin and keep those with higher background rejection

Algorithm

- Better solutions more likely to be selected for mating and mutations, carrying their genetic code (cuts) from generation to generation
- Algorithm:
 - 1 Create initial random population (cut ensemble)
 - 2 Select fittest individuals
 - 3 Create offsprings through crossover (mix best cuts)
 - 4 Mutate randomly (change some cuts of some individuals)
 - 5 Repeat from 2 until convergence (or fixed number of generations)
- Good fitness at one generation \Rightarrow average fitness in the next
- Algorithm focuses on region with higher potential improvement

Gaussian problem

- Suppose densities $s(x)$ and $b(x)$ are multivariate Gaussians:

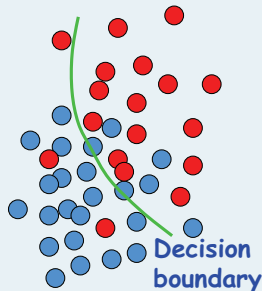
$$\text{Gaussian}(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

with vector of means μ and covariance matrix Σ

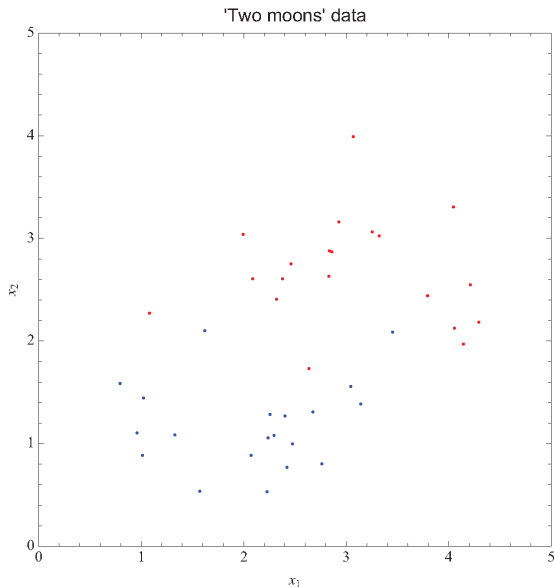
- Then Bayes factor $B(x) = s(x)/b(x)$ (or its logarithm) can be expressed explicitly:

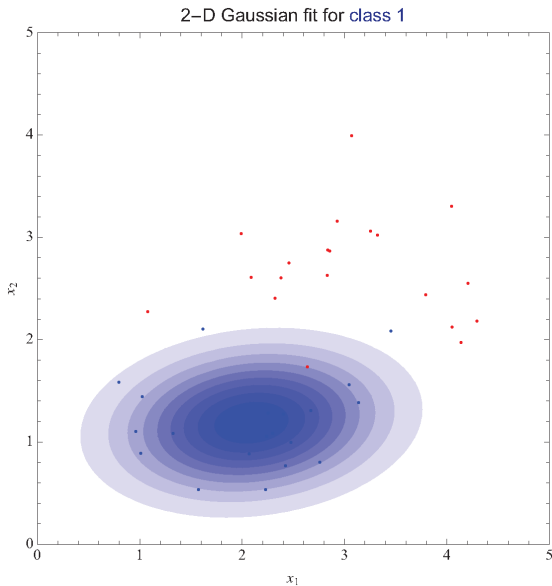
$$\ln B(x) = \lambda(x) \equiv \chi^2(\mu_B, \Sigma_B) - \chi^2(\mu_S, \Sigma_S)$$

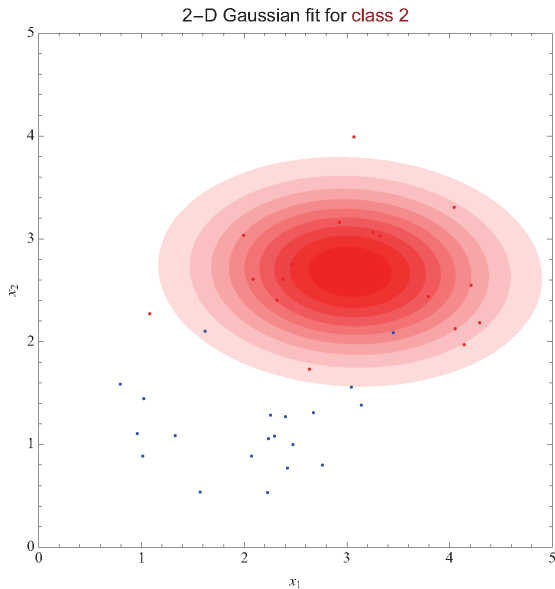
with $\chi^2(\mu, \Sigma) = (x - \mu)^T \Sigma^{-1}(x - \mu)$

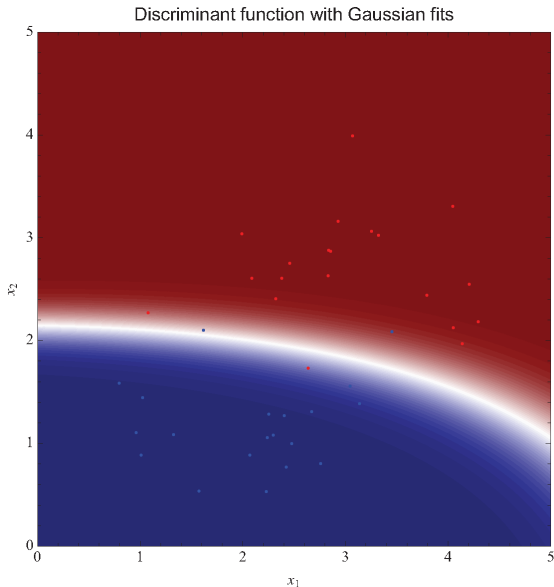


- Fixed value of $\lambda(x)$ defines a quadratic hypersurface partitioning the n -dimensional space into signal-rich and background-rich regions
- Optimal separation if $s(x)$ and $b(x)$ are indeed multivariate Gaussians





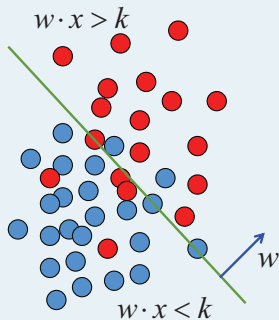




Fisher's discriminant

- If in $\lambda(x)$ the same covariance matrix is used for each class (e.g. $\Sigma = \Sigma_S + \Sigma_B$) one gets **Fisher's discriminant**:

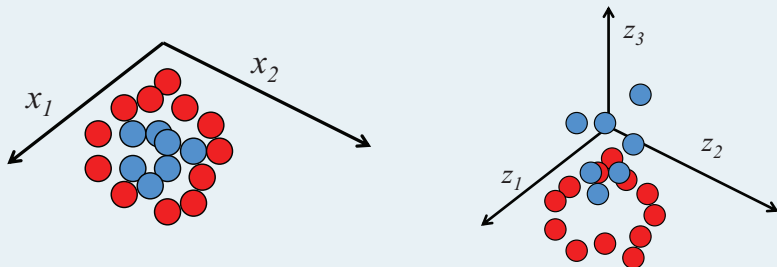
$$\lambda(x) = w \cdot x \quad \text{with} \quad w \propto \Sigma^{-1}(\mu_S - \mu_B)$$



- Optimal linear separation
- Works only if signal and background have different means!
- Optimal classifier (reaches the Bayes limit) for linearly correlated Gaussian-distributed variables

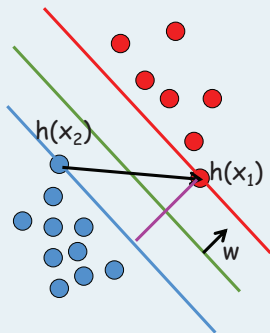
Generalising Fisher discriminant

- Fisher discriminant: may fail completely for highly non-Gaussian densities
- But linearity is good feature \Rightarrow try to keep it
- Idea: data non-separable in n -dim space \mathbb{R}^n , but better separated if mapped to higher dimension space \mathbb{R}^H : $h : x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^H$
- Use hyper-planes to partition higher dim space: $f(x) = w \cdot h(x) + b$
- Example: $h : (x_1, x_2) \rightarrow (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$



Starting simple: separable data

- Consider separable data in \mathbb{R}^H , and three parallel hyper-planes:
 - $w \cdot h(x) + b = 0$ (separating hyper-plane between red and blue)
 - $w \cdot h(x_1) + b = +1$ (contains $h(x_1)$)
 - $w \cdot h(x_2) + b = -1$ (contains $h(x_2)$)



- Subtract blue from red:
 - $w \cdot (h(x_1) - h(x_2)) = 2$
- With unit vector $\hat{w} = w/\|w\|$:
 - $\hat{w} \cdot (h(x_1) - h(x_2)) = 2/\|w\| = m$
- Margin m is distance between red and blue planes
- Best separation: maximise margin
- \Rightarrow empirical risk margin to minimise:
 - $R(w) \propto \|w\|^2$

Constraints

- When minimising $R(w)$, need to keep signal and background separated
- Label red dots $y = +1$ (“above” red plane) and blue dots $y = -1$ (“below” blue plane)

- Since:

$$w \cdot h(x) + b > 1 \text{ for red dots}$$

$$w \cdot h(x) + b < -1 \text{ for blue dots}$$

all correctly classified points will satisfy constraints:

$$y_i (w \cdot h(x_i) + b) \geq 1, \quad \forall i = 1, \dots, N$$

- Using Lagrange multipliers $\alpha_i > 0$, cost function can be written:

$$C(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [y_i (w \cdot h(x_i) + b) - 1]$$

Minimisation

- Minimise cost function $C(w, b, \alpha)$ with respect to w and b :

$$C(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (h(x_i) \cdot h(x_j))$$

- At minimum of $C(\alpha)$, only non-zero α_i correspond to points on red and blue planes: **support vectors**

Kernel functions

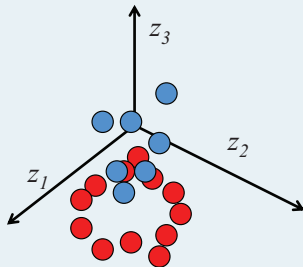
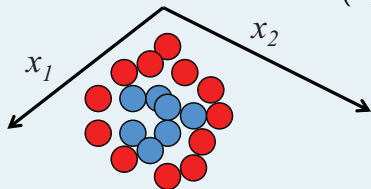
- Issues:
 - need to find h mappings (potentially of infinite dimension)
 - need to compute scalar products $h(x_i) \cdot h(x_j)$
- Fortunately $h(x_i) \cdot h(x_j)$ are equivalent to some kernel function $K(x_i, x_j)$ that does the mapping and the scalar product:

$$C(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j)$$

Example

- $h : (x_1, y_2) \rightarrow (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$

$$\begin{aligned}h(x) \cdot h(y) &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2) \\ &= (x \cdot y)^2 \\ &= K(x, y)\end{aligned}$$



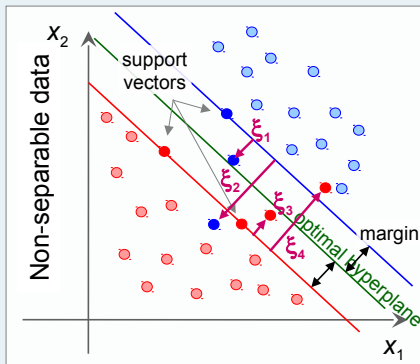
- In reality: do not know a priori the right kernel
- \Rightarrow have to test different standard kernels and use the best one

Real life: non-separable data

- Even in infinite dimension space, data are often non-separable
- Need to relax constraints:

$$y_i(w \cdot h(x_i) + b) \geq 1 - \xi_i$$

with **slack variables** $\xi_i > 0$

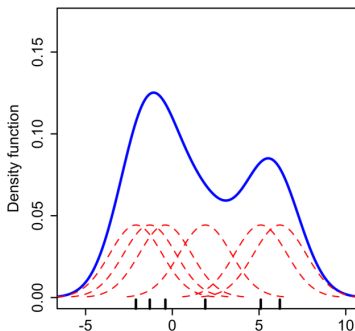
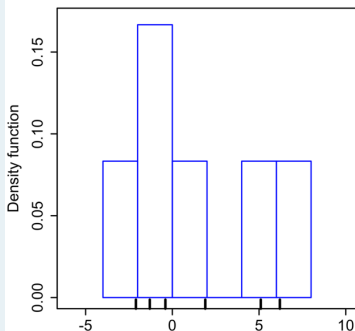


- $C(w, b, \alpha, \xi)$ depends on ξ , modified $C(\alpha, \xi)$ as well
- Values determined during minimisation

Basic principle

- Introduced by E. Parzen in the 1960s
- Place a kernel $K(x, \mu)$ at each training point μ
- Density $p(x)$ at point x approximated by:

$$p(x) \approx \hat{p}(x) = \frac{1}{N} \sum_{j=1}^N K(x, \mu_j)$$



Choice of kernel

- Any kernel can be used
- In practice, often product of Gaussians:

$$K(x, \mu) = \prod_{i=1}^n \text{Gaussian}(x_i | \mu, h_i)$$

each with **bandwidth** (width) h_i

Optimal bandwidth

- Too narrow: noisy approximation
- Too wide: loose fine structure
- In principle found by minimising risk function

$$R(\hat{p}, p) = \int (\hat{p}(x) - p(x))^2 dx$$

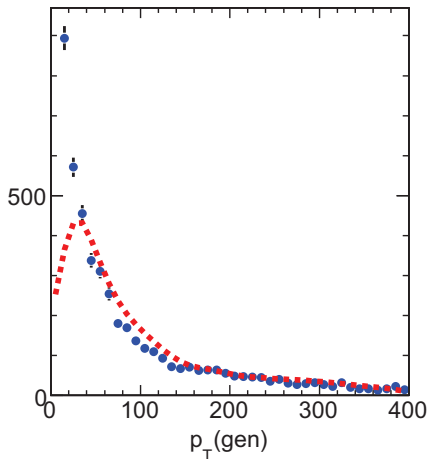
- For Gaussian densities:

$$h = \sigma \left(\frac{4}{(n+2)N} \right)^{1/(n+4)}$$

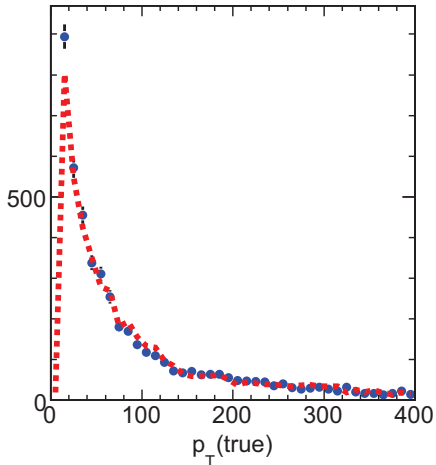
- Far from optimal for non-Gaussian densities

Example

with Gaussian optimal bandwidth



with optimised bandwidth



Why does it work?

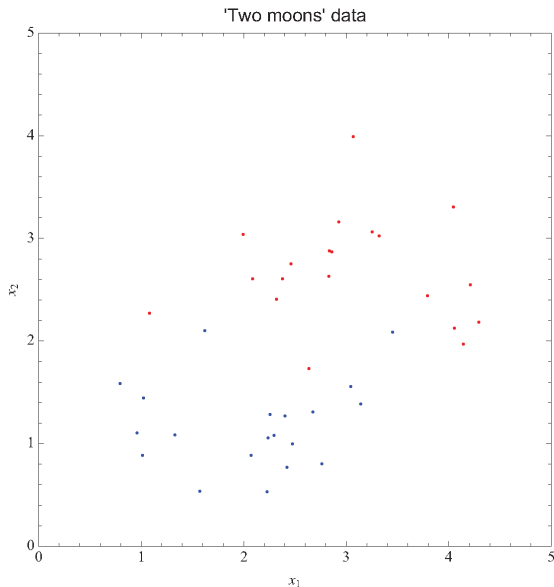
- When $N \rightarrow \infty$:

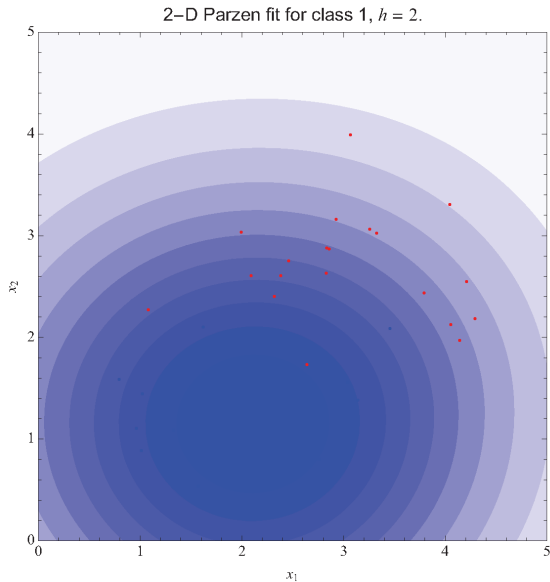
$$\hat{p}(x) = \int K(x, \mu) p(\mu) d\mu$$

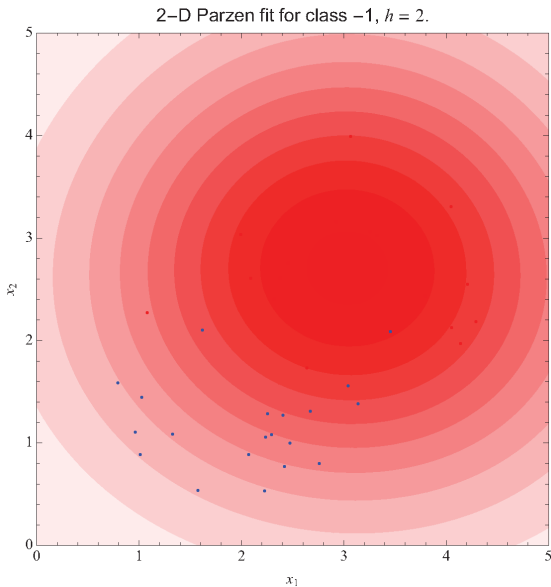
- $p(\mu)$: true density of x
- Kernel bandwidth getting smaller with N , so when $N \rightarrow \infty$, $K(x, \mu) \rightarrow \delta^n(x - \mu)$ and $\hat{p}(x) = p(x)$
- KDE gives consistent estimate of probability density $p(x)$

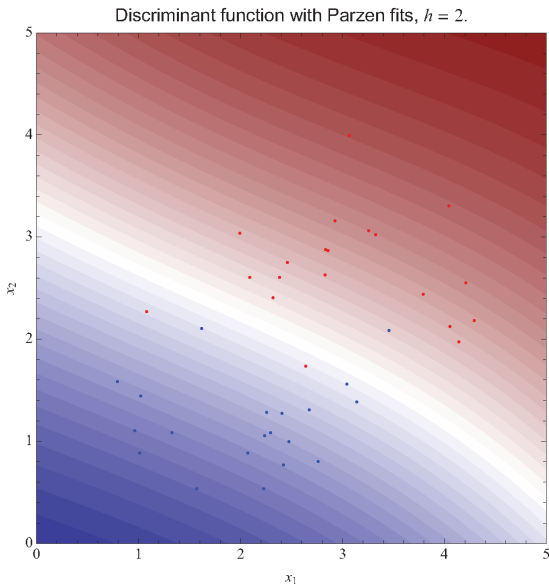
Limitations

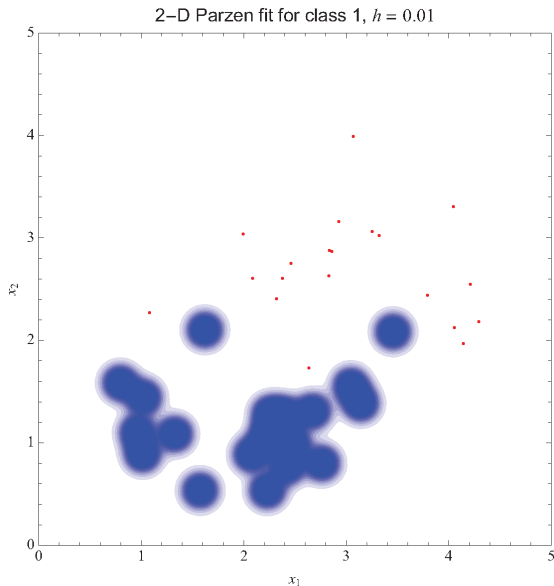
- Choice of bandwidth non-trivial
- Difficult to model sharp structures (e.g. boundaries)
- Kernels too far apart in regions of low point density
- (both can be mitigated with adaptive bandwidth choice)
- Requires evaluation of N n -dimensional kernels

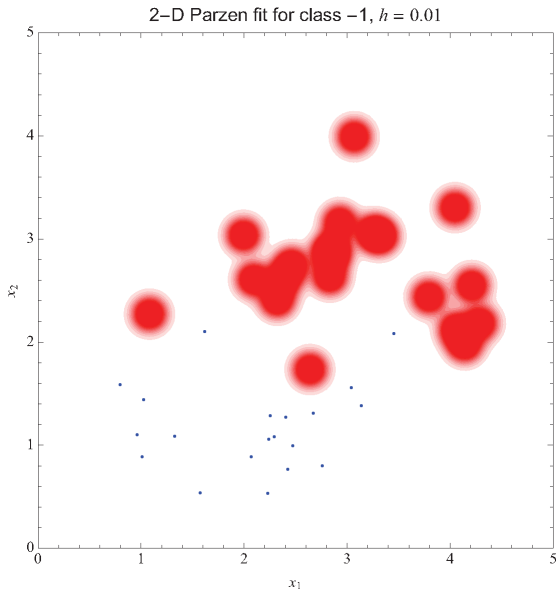


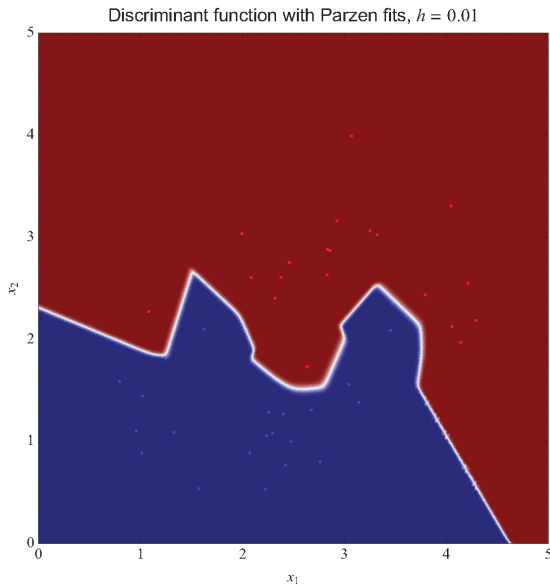


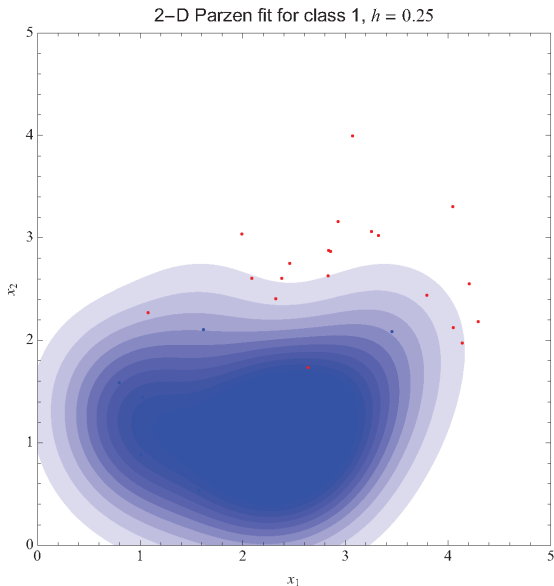


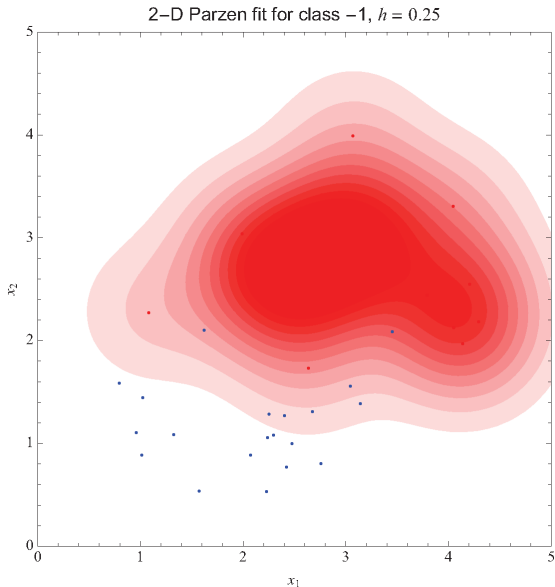


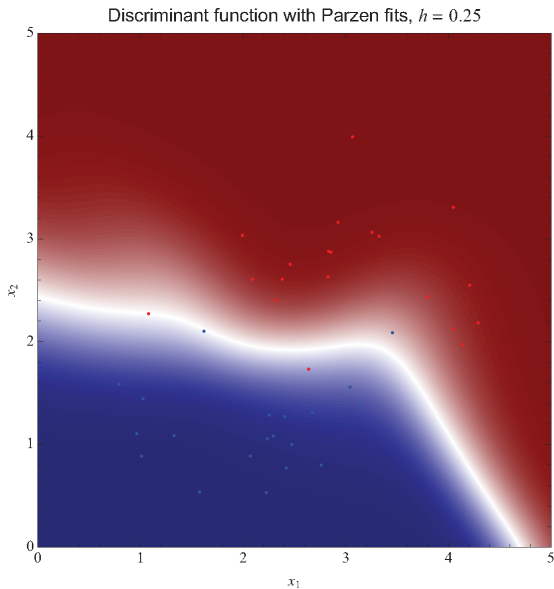




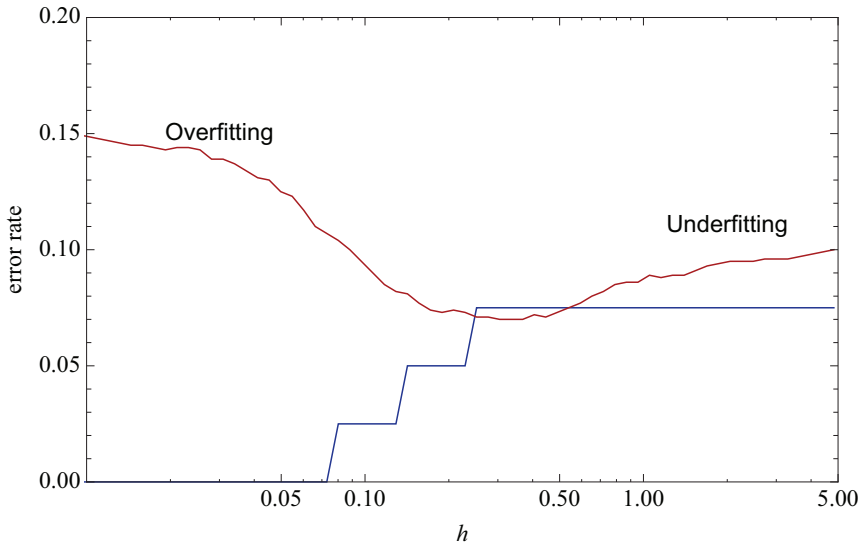


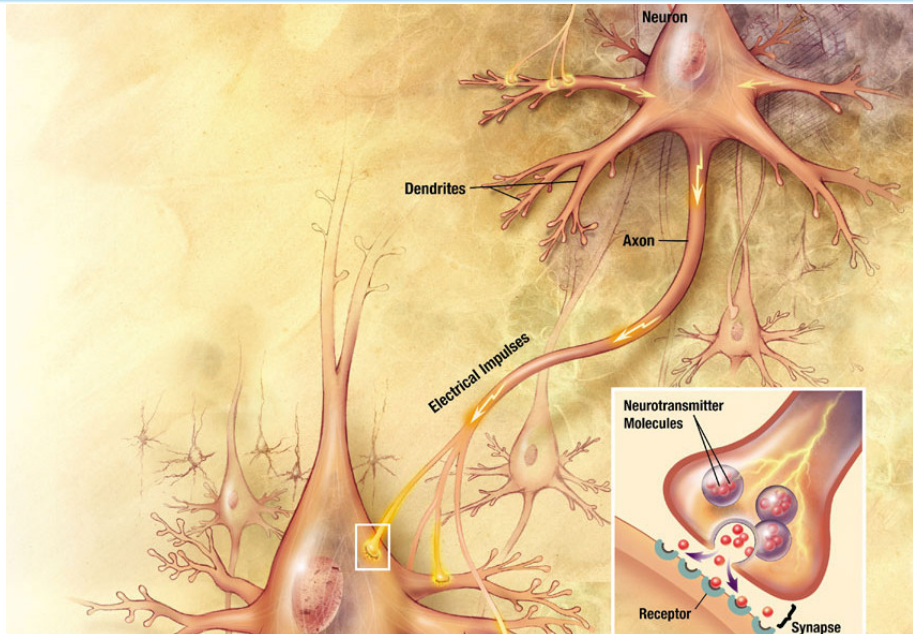






Training and test error rates



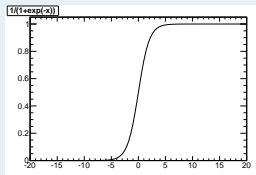


Brief history of artificial neural networks

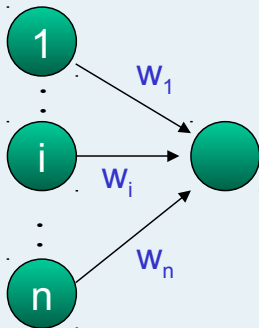
- 1943: W. McCulloch and W. Pitts explore capabilities of networks of simple neurons
- 1958: F. Rosenblatt introduces perceptron (single neuron with adjustable weights and threshold activation function)
- 1969: M. Minsky and S. Papert prove limitations of perceptron (linear separation only) and (wrongly) conjecture that multi-layered perceptrons have same limitations
⇒ ANN research almost abandoned in 1970s!!!
- 1986: Rumelhart, Hinton and Williams introduce “backward propagation of errors”: solves multi-layered learning
- Today: will only talk about multilayer perceptron (MLP), but there are many recent advances in ANN

Single neuron

- Remember linear separation:
 $\lambda(x) = w \cdot x = \sum_{i=1}^n w_i x_i + w_0$
- Boundary at $\lambda(x) = 0$
- Replace threshold boundary by sigmoid:



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

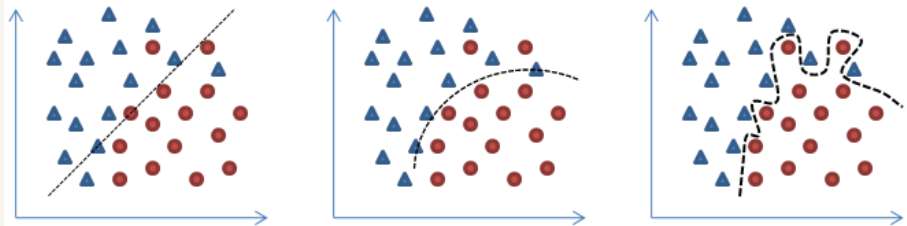


- $\sigma(\lambda)$ is neuron activity, λ is activation
- Neuron behaviour completely controlled by weights $w = \{w_0, \dots, w_n\}$
- Training: minimisation of error/loss function (quadratic deviations, entropy [maximum likelihood]), via gradient descent or stochastic approximation

Training

- Minimise error function $E(w)$
- Gradient descent: $w^{(k+1)} = w^{(k)} - \eta \frac{dE^{(k)}}{dw}$
- $\frac{\partial E}{\partial w_j} = \sum_{n=1}^N -(t^{(n)} - y^{(n)})x_j^{(n)}$ with target $t^{(n)}$ (0 or 1), so $t^{(n)} - y^{(n)}$ is the error on event n
- All events at once (batch learning):
 - weights updated all at once after processing the entire training sample
 - finds the actual steepest decent
 - takes more time
- or one-by-one (online learning):
 - speeds up learning
 - useful in HEP because of redundant datasets (large Monte Carlo samples with many similar events)
 - may avoid local minima with stochastic component in minimisation
 - depends on the order of training events
- One epoch: going through the training data once

Overtraining



- Diverging weights can cause overfitting
- Mitigate by:
 - early stopping (after a fixed number of epochs)
 - monitoring error on test sample
 - regularisation, introducing a “weight decay” term:

$$\tilde{E}(w) = E(w) + \frac{\alpha}{2} \sum_i w_i^2$$

Theorem

Let $\sigma(\cdot)$ be a non-constant, bounded, and monotone-increasing continuous function. Let $\mathcal{C}(I_n)$ denote the space of continuous functions on the n -dimensional hypercube. Then, for any given function $f \in \mathcal{C}(I_n)$ and $\varepsilon > 0$ there exists an integer M and sets of real constants w_j, w_{ij} where $i = 1, \dots, n$ and $j = 1, \dots, M$ such that

$$y(x, w) = \sum_{j=1}^M w_j \sigma \left(\sum_{i=1}^n w_{ij} x_i + w_{0j} \right)$$

is an approximation of $f(\cdot)$, that is $|y(x) - f(x)| < \varepsilon$

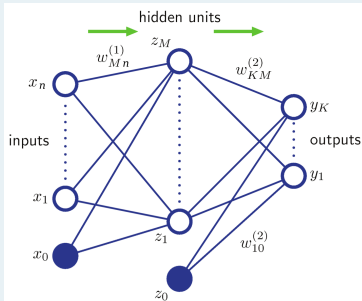
Interpretation

- You can approximate any continuous function to arbitrary precision with a linear combination of sigmoids
- Corollary 1: can approximate any continuous function with neurons!
- Corollary 2: a single hidden layer is enough
- Corollary 3: a linear output neuron is enough

Multilayer perceptron: feedforward network

- Neurons organised in layers
- Output of one layer becomes input to next layer

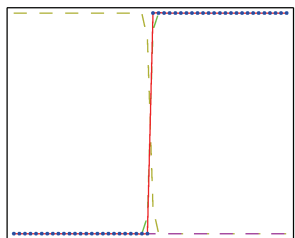
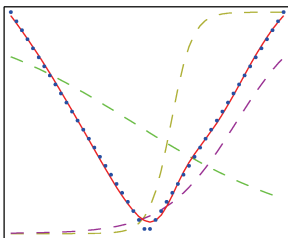
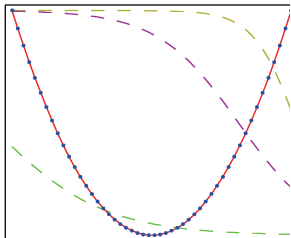
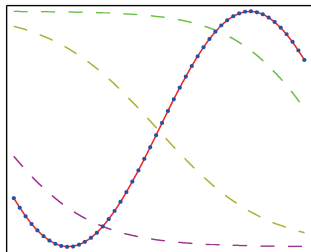
$$y_k(x, w) = \sum_{j=0}^M w_{kj}^{(2)} \underbrace{\sigma \left(\sum_{i=0}^n w_{ji}^{(1)} x_i \right)}_{z_j}$$



Can fit any function: examples

- 1 input (training data), 1 output
- 3 hidden neurons on one hidden layer

©Jan Therhaag



Backpropagation

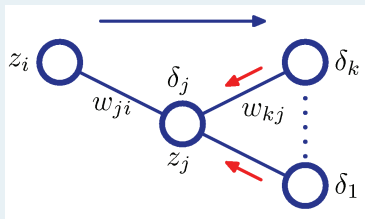
- Training means minimising error function $E(w)$
- For single neuron: $\frac{dE}{dw_k} = (y - t)x_k$
- One can show that for a network:

$$\frac{dE}{dw_{ji}} = \delta_j z_i, \text{ where}$$

$$\delta_k = (y_k - t_k) \text{ for output neurons}$$

$$\delta_j \propto \sum_k w_{kj} \delta_k \text{ otherwise}$$

- Hence errors are propagated backwards

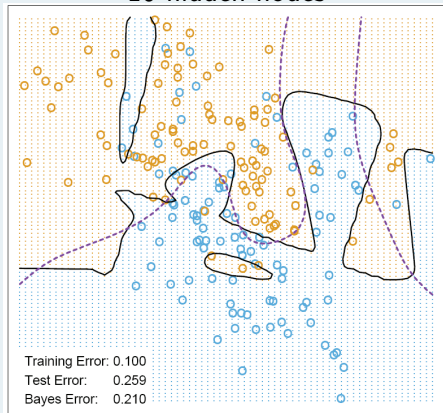


- As before, weights can be regularised:

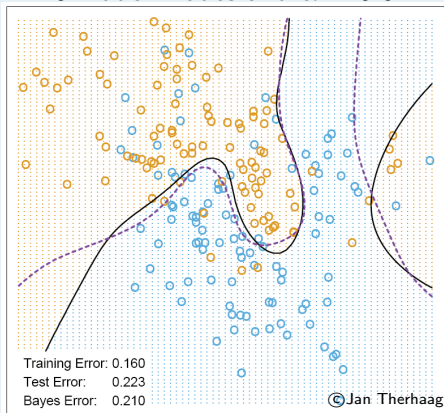
$$\tilde{E}(w) = E(w) + \frac{\alpha}{2} \sum_i w_i^2$$

Regularisation

10 hidden nodes



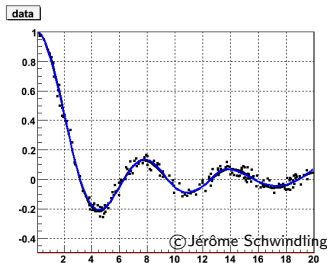
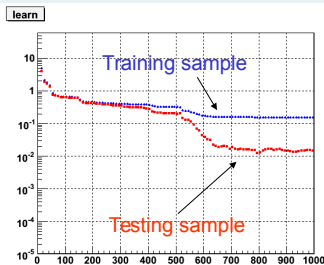
10 hidden nodes and $\alpha = 0.04$



- Much less overfitting, better generalisation properties

Getting confused: testing better than training?

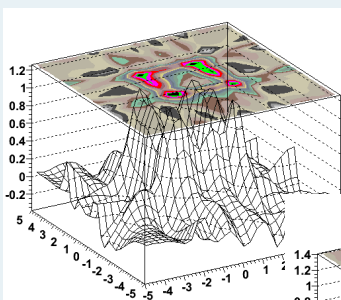
- Train on noisy data centred on true value
- Test on no-noise data
- Testing error becomes better: during training, the NN learned the true distribution (average of noisy inputs)
- \Rightarrow testing converges
- Example:
 $\sin(x)/x + \text{rand}(-0.05, +0.05)$
- Of course doesn't work as well if noise is not symmetric



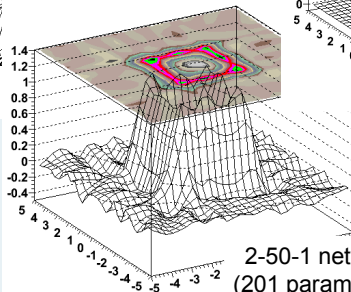
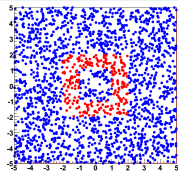
Tricks of the trade

- Preprocess data:
 - if relevant, provide e.g. x/y instead of x and y
 - subtract the mean because the sigmoid derivative becomes negligible very fast
 - normalise variances (close to 0)
 - shuffle training sample (order matters in online training)
- Initial random weights should be small to avoid saturation
- Batch/online training: depends on the problem
- Regularise weights to minimise overtraining. May also help select good variables via Automatic Relevance Determination (ARD)
- Make sure the training sample covers the full parameter space
- No rule (not even valid guestimates) about the number of hidden nodes
- A single hidden layer is enough for all purposes, but multiple hidden layers may allow for a solution with fewer parameters

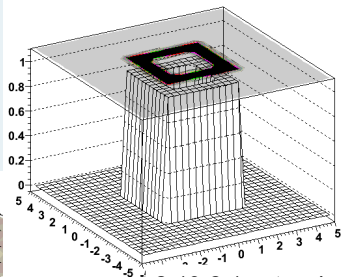
Adding a hidden layer



2-20-1 network
(81 parameters)



2-50-1 network
(201 parameters)



2-10-2-1 network
(55 parameters)

Introduction

- As name says: Bayesian approach, try to *infer* functions $f(x)$
- Training sample T of N examples $(x, y)_1, (x, y)_2, \dots, (x, y)_N$ of discriminating variables x and class labels y
- Each point w corresponds to a function $f(x, w)$
- Assign probability density $p(w|T)$ to it
- If $p(w_1|T) > p(w_2|T)$, then associated function $f(x, w_1)$ more compatible with training data T than function $f(x, w_2)$
- Posterior density $p(w|T)$ is final result of Bayesian inference
- BNN is the predictive distribution

$$p(y|x, T) = \int p(y|x, w)p(w|T)dw$$

where the function class is class of feedforward neural networks with a fixed structure (inputs, layers, hidden nodes, outputs)

In practice

- Take the mean of the predictive distribution:

$$\begin{aligned}y(x) &= \int zp(z|x, T)dz \\ &= \int f(x, w)p(w|T)dw\end{aligned}$$

- Why? For classification $p(y|x, w) = f(x, w)^y(1 - f(x, w))^{1-y}$
 - for $y = 1$: $p(y|x, w) = f(x, w)$
 - for $y = 0$: $p(y|x, w) = 1 - f(x, w)$
 - so only $f(x, w)$ contributes to the mean

- Example usage:

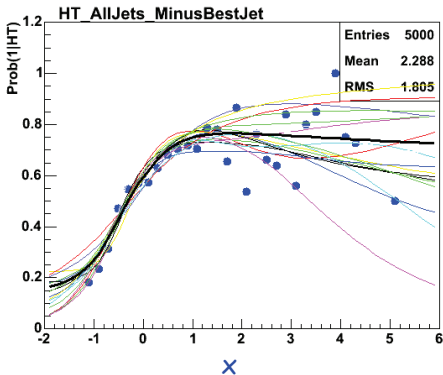
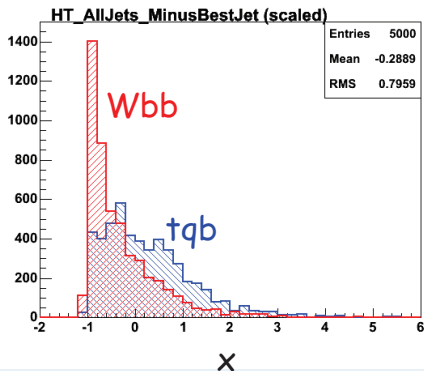
$$\begin{aligned}f(x, w) &= \frac{1}{1 + e^{-g(x, w)}} \\ g(x, w) &= b + \sum_{j=1}^H v_j \tanh(a_j + \sum_{i=1}^n u_{ij}x_i)\end{aligned}$$

with H hidden nodes

Implementation

- Scanning NN parameter space can be daunting
- Can approximate integral in $y(x)$ using Markov chain Monte Carlo method (MCMC)
- Will generate M sample weights w_1, \dots, w_M from posterior density $p(w|T)$
- $y(x) \approx \frac{1}{M} \sum_{m=1}^M f(x, w_m)$
- Use sparse subset of MCMC points to avoid correlations
- Start with “reasonable” guesses for parameters (e.g. zero-centred Gaussians)

Example



- points: bin by bin histogram ratio
- thin curves: each $f(x, w_k)$
- thick curve: average, which approximates $D(x)$

Details tomorrow

Summary of MVA techniques

Criteria		Classifiers								
		Cuts	Likelihood	PDERS / k-NN	H-Matrix x	Fisher	MLP	BDT	RuleFit	SVM
Performance	no / linear correlations	☹️	😊	😊	☹️	😊	😊	☹️	😊	😊
	nonlinear correlations	☹️	☹️	😊	☹️	☹️	😊	😊	☹️	😊
Speed	Training	☹️	😊	😊	😊	😊	☹️	☹️	☹️	☹️
	Response	😊	😊	☹️/☹️	😊	😊	😊	☹️	☹️	☹️
Robustness	Overtraining	😊	☹️	☹️	😊	😊	☹️	☹️	☹️	☹️
	Weak input variables	😊	😊	☹️	😊	😊	☹️	☹️	☹️	☹️
Curse of dimensionality		☹️	😊	☹️	😊	😊	☹️	😊	☹️	☹️
Transparency		😊	😊	☹️	😊	😊	☹️	☹️	☹️	☹️









(according to TMVA authors)







©Andreas Hoecker

- When trying to achieve optimal discrimination one can try to approximate

$$D(x) = \frac{s(x)}{s(x) + b(x)}$$

- Many techniques and tools exist to achieve this
- (Un)fortunately, no one method can be shown to outperform the others in all cases.
- One should try several and pick the best one for any given problem
- Multivariate techniques are at work in your everyday life without your knowing and can easily outsmart you for many tasks
- Try this to convince yourself http://www.phi-t.de/mousegame/index_eng.html

-  V. Vapnik, *The Nature of Statistical Learning Theory*, Springer, New York, 2nd Edition, 2000
-  T. Hastie, R. Tibshirani and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*, Springer-Verlag, New York, 2nd Edition, 2009
-  R.M. Neal, *Bayesian Learning of Neural Networks*, Springer-Verlag, New York, 1996
-  C.M. Bishop, *Pattern Recognition and Machine Learning*, Springer, New York, 2007
-  M. Minsky and S. Papert, "Perceptrons", M.I.T. Press, Cambridge, Mass., 1969
-  H.B. Prosper, "The Random Grid Search: A Simple Way to Find Optimal Cuts", Computing in High Energy Physics (CHEP 95) conference, Rio de Janeiro, Brazil, 1995
-  W.S. McCulloch & W. Pitts, "A logical calculus of the ideas immanent in nervous activity", *Bulletin of Mathematical Biophysics*, 5, 115-137, 1943
-  F. Rosenblatt, "The Perceptron: A Probabilistic Model for Information Storage & Organization in the Brain", *Psychological Review*, 65, pp. 386-408, 1958

-  D.E.Rumelhart et al., “Learning representations by back-propagating errors”, *Nature* vol. 323, p. 533, 1986
-  K.Hornik et al., “Multilayer Feedforward Networks are Universal Approximators”, *Neural Networks*, Vol. 2, pp 359-366, 1989
-  Y. LeCun, L. Bottou, G. Orr and K. Muller, “Efficient BackProp”, in *Neural Networks: Tricks of the trade*, Orr, G. and Muller K. (Eds), Springer, 1998
-  P.C. Bhat and H.B. Prosper, “Bayesian neural networks”, in *Statistical Problems in Particles, Astrophysics and Cosmology*, Imperial College Press, Editors L. Lyons and M. Ünel, 2005
-  Q.V. Le et al., “Building High-level Features Using Large Scale Unsupervised Learning”, in *Proceedings of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, UK, 2012
-  A. Hoecker, P. Speckmayer, J. Stelzer, J. Therhaag, E. von Toerne, and H. Voss, “TMVA: Toolkit for Multivariate Data Analysis”, PoS A CAT 040 (2007) [physics/0703039], <http://tmva.sourceforge.net>