

# Machine Learning and Pattern Recognition in Bioinformatics

Peter Sykacek<sup>1</sup>

Vienna Science Chair of Bioinformatics  
Department of Biotechnology  
BOKU University  
peter.sykacek@boku.ac.at

[Jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.1/44

## Lecture 1 & 2

- Concepts in Machine Learning and Pattern Recognition (ML & PRN)
- Classification of Problems and Methods in ML & PRN
- Basics of Matrices in ML & PRN and Introduction to MatLab

[Jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.2/44

# Machine Learning

and Pattern Recognition (ML & PRN)

Discipline which emerged at the intersection of

- Computer Science and
- Applied Mathematics / Statistics

Advantage: Same objectives but less formal than Statistics; Thus accessible with less mathematics.

Idea: Data analysis with clever algorithms.

– > inevitably requires qualitative understanding of data analysis and implication of algorithms.

## Brief Guideline to Success

ML & PRN provide powerful but dangerous tools.

Avoiding a self inflicted chain saw massacre requires:

- Understand general concepts in data analysis.
- Know the ever growing repertoire of ML & PRN methods.
- Understand implications of all methods you apply.
- Compare competing methods to avoid bias.

Do not apply published approaches in interdisciplinary fields like computational biology, etc. uncritically: Data analysis is not necessarily done by experts!

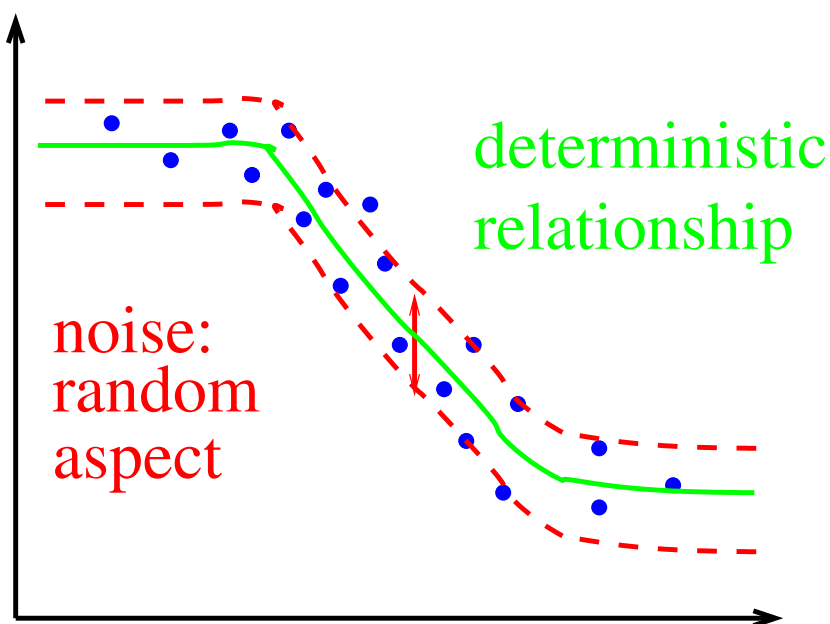
# Problem 1: Noise

ML & PRN tries to “learn” reasonably good **models** from measured data. (in statistics they call this **inference**).

Model - a possibly implicit quantitative abstraction of nature implying:

- imperfection due to unknown or unmodeled factors influencing nature.
- problem of separating the data generating mechanism into:
  - explicable deterministic aspects
  - noise component which contains the rest

## Noise

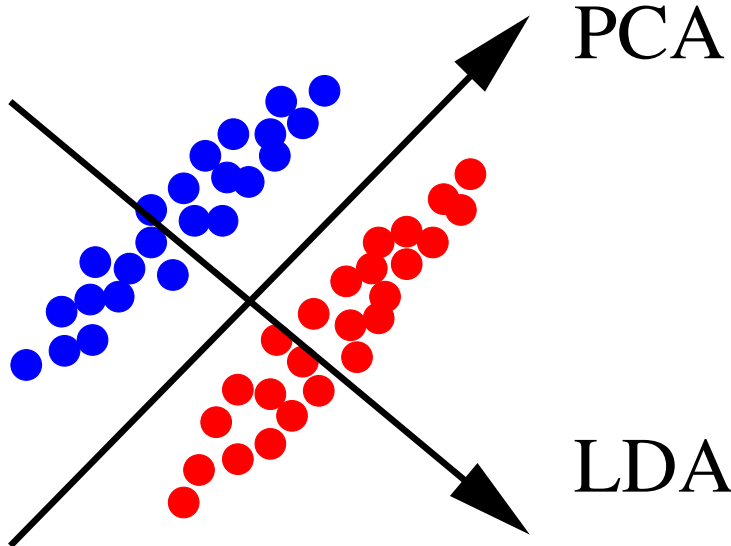


Inexplicable aspects are summarised by a noise component (red error bars).

## Problem 2: Modelling Goal

Result = Data + Model!

Linear discriminant (LDA) and principle component analysis (PCA) give different projections of the same data.



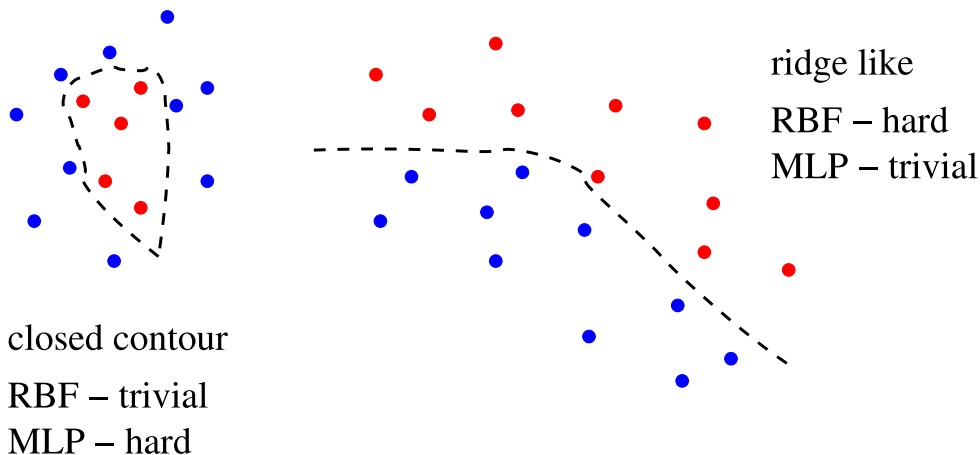
Both use linear projections!

$$t_{\text{PCA}} = \theta_{\text{PCA}}^T x$$

$$t_{\text{LDA}} = \theta_{\text{LDA}}^T x$$

## Problem 3: Methods Bias

Example: dichotomous classification (2 labels)



Methods bias leads to results being influenced by the analysis method, motivating the **no free lunch theorem**. Lack of overall winning algorithms requires comparing competing methods to ensure acceptable performance.

# The Nature of Data

- Discrete valued observations (e.g. class labels).
- Continuous valued observations (e.g. measurements).

Measurement processes involve **errors** which arise from **noise** (fluctuations) that are or can not be captured:

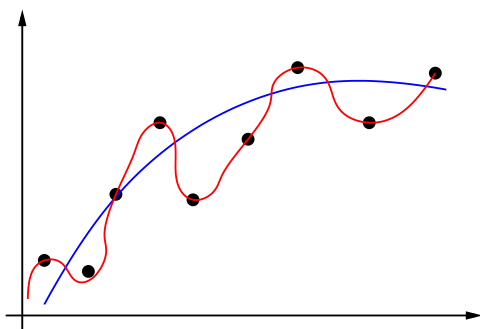
- Measurement noise.
- Wrong classifications (e.g. disease state).
- Simplified Models.

Individual data points do thus not reflect ground truth. Data analysis uses **replicates to estimate noise** and model the remaining aspects as good as possible.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.9/44

# Fundamental Principle in Data Analysis

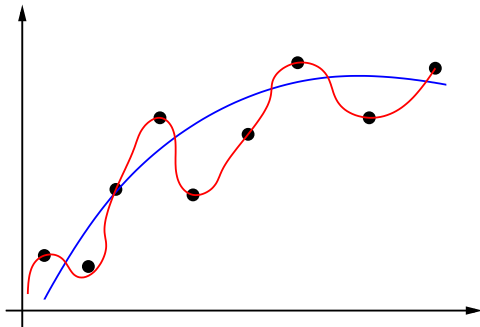


Which is the better model? Why is that the case?

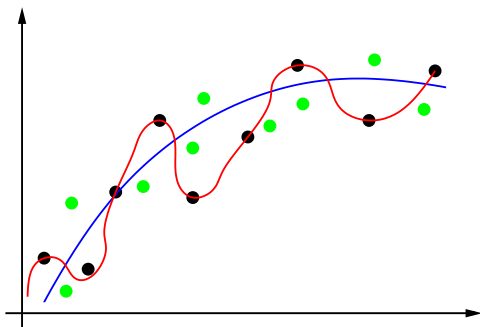
jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.10/44

# Fundamental Principle in Data Analysis



Which is the better model? Why is that the case?



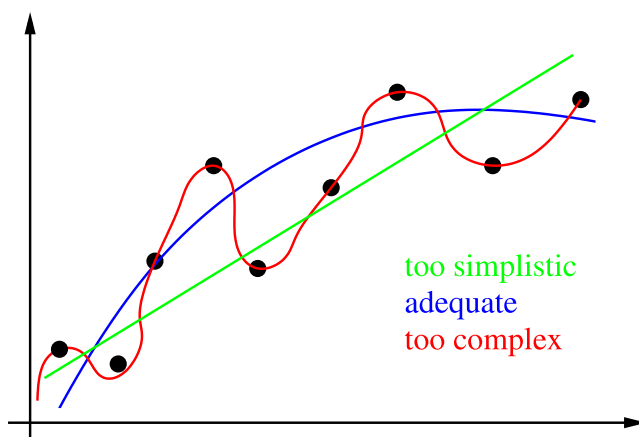
Find (or abstract from) the **underlying model** that generated the data.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.10/44

## Adequate models

Capture underlying structure and avoid **overfitting**. “Fiddle parameters” affecting model complexity can have adverse effects.



Idea: overfitting is a result of tuning the model towards the training data. Over or under-complex models that do not capture the underlying data generating mechanism will perform worse on novel data obtained from the generating model than an appropriate model.

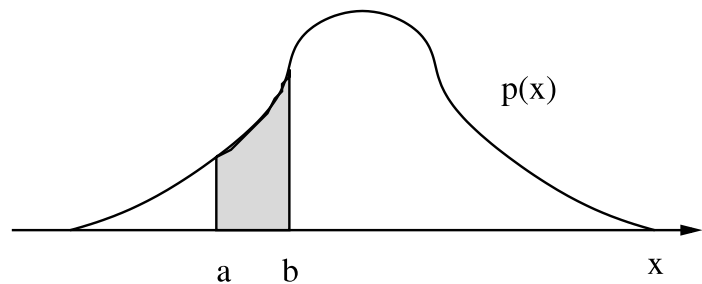
jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.11/44

# Definition: PDF

Data analysis is inherently connected with the concept of **random variables**. A random variable is a non deterministic quantity where repeated observations differ and are generated according to some overall property. Properties of random variables are for example captured by an associated **probability density function (pdf),  $p(x)$** .

The pdf allows deducing the probability that a new realisation falls into a particular set,  $P(x \in [a, b]) = \int_{x=a}^b p(x) dx$ .



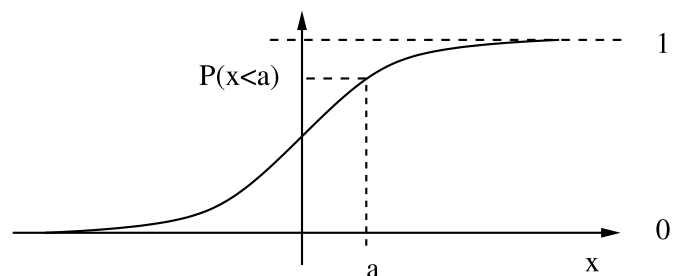
[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.12/44

# Definition: CDF

An equivalent characterisation for univariate random variables is provided by the so called **cumulative distribution function (cdf),  $F(x)$** .

The cdf  $F(x)$  denotes the probability that a realisation of the random variable is smaller than  $x$ .  $F(a)$  is thus the probability  $P(x < a)$ .



– > the pdf  $p(x) = \left. \frac{dF(\xi)}{d\xi} \right|_{\xi=x}$  is the derivative of the cdf at  $x$ .

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.13/44

# Classification of Methods & Problems

Data analysis can be grouped into two categories:

1. **Supervised Learning** methods are used for regression problems.
2. **Unsupervised Learning** methods are used for exploratory data analysis.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.14/44

## Regression Problems

Regression is concerned with data sets of size  $N$  of the type  $\mathcal{Z} = \{(y_1, x_1), (y_2, x_2), \dots, (y_N, x_N)\}$ , with the tuples  $(y_n, x_n)$  drawn from an unknown joint pdf  $p(y, x)$ . The learning task is modelling the dependent variable,  $y$ , as a function of the independent variable  $x$ .

Typical regression models:

**Linear regression:**  $p(y|x) = \mathcal{N}(kx + d, \lambda)$  and  $y \in \mathbb{R}$ .

**Logistic regression:**  $P(y|x) = \text{cdf}_{\text{logistic}}(kx + d)$  and  $y \in \{0, 1\}$ .

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.15/44



# Exploratory Data Analysis

Exploratory data analysis searches for unknown structure in a data set of size  $N$ ,  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ , with  $x_n \in \mathbb{R}^d$  drawn from an unknown pdf  $p(x)$ . The learning task is modelling  $x$  as a function of an **unobserved (latent)** variable  $t$ , providing an easy to grasp summary of the data.

Typical models for exploratory data analysis:

**Mixture density models:**

$$p(x_n) = \sum_k P(t_n = k)p(x_n|t_n = k), \text{ and } t_n \in \{1, \dots, K\}.$$

**Continuous latent variable models:**

$$p(x_n) = \int p(x_n|t) p(t) dt, \quad t \in \mathbb{R}^k \text{ and } d > k$$

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.16/44

## Summary Part 1

- ML & PRN is concerned with **algorithms** which are used for **data analysis**.
- Data analysis **learns** models from **noisy data**.
- We know two types of problems:
  - Data which should be analysed for **simplifying explanation**. – > **unsupervised learning**.
  - Data which requires predicting **responses from independent variables**. – > **supervised learning**.
- Fundamental principle: Search for **adequate models of appropriate complexity**. Also requires comparing different models!

# Matrices for Data Analysis

Important to simplify notation!

Definition  $n$ -dimensional Euclidean Space  $\mathbb{R}^n$ :

Definition of a matrix ( $n$  rows,  $m$  columns):

$$M = \begin{pmatrix} m_{1,1} & \cdots & m_{1,m} \\ \vdots & \ddots & \vdots \\ m_{n,1} & \cdots & m_{n,m} \end{pmatrix} = (\mathbf{m}_1, \cdots, \mathbf{m}_m) \text{ and } \mathbf{m}_i \in \mathbb{R}^n$$

MatLab: `>> M = [[a, b, c]; [d, e, f]; ...];` What are the  $\mathbf{m}_i$ ?

Important: `>> help < commandname >`, where one may type any MatLab command name, **provides help** for the corresponding command.

# Matrix Operations

Transposition:  $B = A^T, \forall n, m : b_{m,n} = a_{n,m}$

MatLab: `>> B = A'`;

Addition ( $A, B$  equal size):

$$C = A + B, \forall n, m : c_{n,m} = a_{n,m} + b_{n,m}$$

MatLab: `>> C = A + B`;

Associative and commutative?

Matrix times constant:

$$B = \lambda A \forall n, m : b_{n,m} = \lambda a_{n,m}$$

MatLab: `>> B = lambda * A`;

Associative and commutative?

# Matrix Multiplication

Matrix Inner Product ( $A$ 's column no. equals  $B$ 's row no.):  $C = AB \forall n, m : c_{n,m} = \sum_i a_{n,i} b_{i,m}$

MatLab: `>> C = A * B;`

Associative and commutative?

Note:  $(AB)^T = B^T A^T$

**However:**  $(A + B)^2 = A^2 + AB + BA + B^2$

Hadamard Product ( $A, B$  equal size):

$C = A \cdot B, \forall n, m : c_{n,m} = a_{n,m} b_{n,m}$

MatLab: `>> C = A .* B;`

Associative and commutative?

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.20/44

# Matrices and Data Analysis

Just to make sure that you see the connection between matrices and data analysis, here an example:

Assume  $N$  samples of  $k$  “input” measurements collected in  $\mathbf{x}_n$  and one dependent variable  $y_n$ , which we intend modelling as a function  $f(\mathbf{x}_n, \Theta)$  parameterised by  $\Theta$ . This type of modelling is called *regression*.

We can only move on deciding on a particular  $f(\mathbf{x}_n, \Theta)$ .

For simplicity we assume that the best guess of  $y_n$  is obtained as linear combination of  $\mathbf{x}_n$ . This allows writing:

$$y_n = \sum_k \mathbf{x}_n[k] \Theta[k], \text{ or } y_n = \mathbf{x}_n^T \Theta$$

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.21/44

# Model fitting I

is often done by **minimising least squares differences**

$$\hat{\Theta} = \operatorname{argmin}_{\Theta} \left( \sum_n (y_n - \mathbf{x}_n^T \Theta)^2 \right)$$

Importance of thinking in terms of matrices:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} \quad \text{and} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}$$

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.23/44

# Model fitting II

We can then immediately write

$$\text{lsd} = \sum_n (y_n - \mathbf{x}_n^T \Theta)^2 = (\mathbf{X}\Theta - \mathbf{y})^T (\mathbf{X}\Theta - \mathbf{y})$$

which contains no sums any more and is an extremely convenient method for deriving model fitting procedures and code for numerical tools like MatLab.

```
MatLab:>> y_d = X * theta - y;  
>> LSD = y_d' * y_d;
```

two lines to be more efficient!

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.23/44

# The Least Squares Optimum

is, as previously discussed, obtained as  $\hat{\Theta} = \operatorname{argmin}_{\Theta}(\text{lsd})$ . Similar to unconditional optimisation of functions:

- We take the derivative with respect to the quantity we want to optimise for.
- and set the derivative equal to zero.

$$\text{lsd} = \Theta^T \mathbf{X}^T \mathbf{X} \Theta - 2\Theta^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$$

Special: we take derivatives w.r.t the vector  $\Theta$ !

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.24/44

# Gradient Vectors of LSD Expression

$$\nabla_{\Theta}(\Theta^T \mathbf{X}^T \mathbf{X} \Theta) = 2\mathbf{X}^T \mathbf{X} \Theta$$

$$\nabla_{\Theta}(-2\Theta^T \mathbf{X}^T \mathbf{y}) = -2\mathbf{X}^T \mathbf{y}$$

and

$$\nabla_{\Theta}(\mathbf{y}^T \mathbf{y}) = 0$$

The gradient vector w.r.t.  $\Theta$  is thus given by:

$$\nabla_{\Theta}(\text{lsd}) = 2(\mathbf{X}^T \mathbf{X} \Theta - \mathbf{X}^T \mathbf{y}).$$

Solution:  $2(\mathbf{X}^T \mathbf{X} \hat{\Theta} - \mathbf{X}^T \mathbf{y}) = 0$  or  $\mathbf{X}^T \mathbf{X} \hat{\Theta} = \mathbf{X}^T \mathbf{y}$

– > we have to remove  $\mathbf{X}^T \mathbf{X}$  from the left side of the equality.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.25/44

# Square Matrices

Rank of a matrix  $r(\mathbf{A})$ : number of linearly independent columns (or rows, that's the same) of  $\mathbf{A}$ .

Square matrix  $\mathbf{A}$  is a square matrix if no. rows equals no. cols, that is:  $n = m$ .

Square matrix  $\mathbf{A}$  is non-singular if rank  $r(\mathbf{A}) = n$ .  
The Determinant and the inverse are defined for square matrices.

**Note:**  $\mathbf{X}^T \mathbf{X}$  from the previous slide is a square matrix!

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.26/44

# The Determinant

The determinant  $|\mathbf{A}|$  (MatLab: `>> det(A)`) converts a square matrix to a real number.

$|\mathbf{A}| = 0$  implies that  $\mathbf{A}$  is singular

$|\mathbf{A}| = \prod_n \lambda_n$ , where  $\lambda_n$  are the eigenvalues of  $\mathbf{A}$

For  $\mathbf{A}$  upper/lower diag. matrix:  $|\mathbf{A}| = \prod_{i=1}^n a_{i,i}$ .

Recursive definition w.r.t  $j$ -th row (high school!):

$$|\mathbf{A}| = \sum_i (-1)^{i+j} a_{i,j} |\mathbf{A}_{i,j}|, \text{ where } \mathbf{A}_{i,j} = \begin{pmatrix} \cdots & | & \cdots \\ \cdots & a_{i,j} & \cdots \\ \cdots & | & \cdots \end{pmatrix}$$

$|\mathbf{A}_{i,j}|$  is the determinant of the submatrix when removing the  $j$ -th row and  $i$ -th column.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.27/44

# Diagonal and Identity Matrices

Diagonal matrix:  $A = \text{diag}(a_{1,1}, \dots, a_{n,n})$ , defines a matrix with the only non zero elements located on the main diagonal

MatLab: `>> A = diag(a); % places a into main diagonal of A`  
`>> a = diag(A); % places main diagonal of A into a`

Identity matrix:  $I = \text{diag}(1, \dots, 1)$

neutral element of matrix multiplication:

$$IA = AI = A$$

MatLab: `>> I = eye(n); % generates an  $[n \times n]$  identity matrix.`

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.28/44

# Inverse Matrix

If matrix  $A$  is non-singular, we get a non-singular matrix  $B = A^{-1}$ , such that,  $BA = AB = I$ .

Matrix  $B$  is the inverse of  $A$

MatLab: `>> B = A^(-1)`

Remarks:  $(A^{-1})^T = (A^T)^{-1}$  and  
 $(AB)^{-1} = B^{-1}A^{-1}$

Matrix  $A$  is orthonormal if  $A^T A = I$ , hence  
 $A^T = A^{-1}$

Examples: projection to principal axis (PCA)

Permutation matrix  $P$  in every row and column one 1 entry, otherwise 0

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.29/44

# Using Inverse Matrices

Consider resolving:

$$\mathbf{Ax} = \mathbf{b}$$

w.r.t  $x$  for  $A$  square  $[n \times n]$ , then:

$$\mathbf{A}^{-1} \times | \mathbf{Ax} = \mathbf{b} \text{ leads to } \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

Such operations are often found in data analysis, e.g. in finding least squares solutions.

**Most important:** Unlike in one dimensional operations, you must here multiply from the left! About 80% get this wrong in the exam!

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.30/44

# Moore Penrose Pseudo Inverse

Inverting ill conditioned (close to singular) matrices involves quantities close to machine precision. Results derived from such inverse matrices result in large numerical errors.

Practical rule - never use matrix inversion, always use the Moore Penrose pseudo inverse.

$$\mathbf{A}^+ = \lim_{\delta \rightarrow 0} (\mathbf{A}^T \mathbf{A} + \delta \mathbf{I})^{-1} \mathbf{A}^T$$

MatLab: `>> A_plus = pinv(A);`

If  $A$  square and not ill-conditioned:

$$\mathbf{A}^+ \mathbf{A} = \mathbf{A}^{-1} (\mathbf{A}^T)^{-1} \mathbf{A}^T \mathbf{A} = \mathbf{I}$$

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.31/44



# Solution for Least Squares Optimum

We have to solve:

$$\mathbf{X}^T \mathbf{X} \hat{\Theta} = \mathbf{X}^T \mathbf{y}$$

for  $\hat{\Theta}$ . We thus multiply both sides of the equality from the left with  $(\mathbf{X}^T \mathbf{X})^{-1}$ . Left hand side:

$$(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \hat{\Theta} = \mathbf{I} \hat{\Theta} = \hat{\Theta}$$

The solution is thus:

$$\hat{\Theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

in MatLab: `theta_hat = pinv(X' * X) * X' * y;`

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.32/44

# Linear Regression in MatLab I

For efficiency express predictions simultaneously as  $\mathbf{y} = \mathbf{X}\theta$ . MatLab code:

```
function [y]=linpred(Xtst, W)
% function [y]=linpred(Xtst, W) predicts outputs for linear regression
% Note: we do assume that X is properly augmented allowing for an intercept
% term (i.e. an input independent constant offset)
% Xtst: a [nsmpl x nrinputs] matrix of independent variables (test data).
% W: a [nrinputs x 1] column vector of linear regressor.
% y: predicted responses
% BTW : % is MatLabs comment symbol which makes it ignore the rest of the l
y=Xtst*W;
```

This function may be called in MatLab as:

```
% Assuming that xtst and W are properly defined and matching in dimension,
% y_pred contains the predictions at all samples stored in xtst:
[y_pred]=linpred(xtst, W);
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.33/44

# Significance of Matrix Expressions

$$\begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ y_n \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \\ X_{n1} X_{n2} X_{n3} \\ \bullet \\ \bullet \end{bmatrix} * \begin{bmatrix} W_1 \\ W_2 \\ W_3 \end{bmatrix}$$

$y_n = \sum_c X_{n,c} * W_c$ ,  $n$ : row index in  $y$  and  $X$ ;  $c$ : column index in  $X$  and row index in  $W$ . MatLab calculates:

```
[nrw,mcol]=size(Xtst);
for n = 1:nrw
    y(n) = 0;
    for m = 1:mcol
        y(n)=y(n)+Xtst(n,m)*W(m);
    end
end
% y=Xtst*W is MUCH faster!
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.34/44

# More on MatLab Matrices

## Defining Matrices:

```
xtst=[[0.1, 2.3, -5]; [1, -0.3, 4.5]]; % 2 row by 3 column matrix
xtst=rand(100, 2); % a 100 by 2 matrix of uniform random numbers
% loads tab delimited file (no headers!) and generates matrix data
load -ascii data.txt
```

## Accessing matrices (indices are 1 - based):

```
xtst(3,2)=5; % set the entry in row 3 and column 2 to 5
cv1=xtst(:,1); % writes the first column vector in xtst into cv1
xtst(5,:) =rw5; % stores the elements of rw3 as 5-th row in xtst
```

Calculate  $\sum_n x_n^T A x_n \forall n$  assuming  $x_n^T = X(n, :)$  and avoid a loop:

```
sum_xTAX=sum(sum((X*A).*X, 2)); % .*: element wise matrix multiplication
(X * A) .* X has the  $n$ -th row vector:
```

```
[ $x_{n,1} * A(1, :)$  *  $x_n$ ,  $x_{n,2} * A(2, :)$  *  $x_n$ , ...,  $x_{n,m} * A(m, :)$  *  $x_n$ ]
```

$\text{sum}((X * A) .* X, 2)$  is thus a column vector containing in the  $n$ -th row the expression  $x_n^T A x_n$  and the final sum completes the calculation.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.35/44

# Linear Regression in MatLab II

Calculation of maximum likelihood coefficients as  $\hat{\Theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  and MatLab code:

```
function [W]=lintran(X, t)
% function [W]=lintran(X, t) calculates lsd optimum of linear regression:
% Note: we do assume that X is properly augmented allowing for an intercept
% term (i.e. an input independent constant offset)
% X: a [nsmpl x nrinputs] matrix of independent variables.
% t: a [nsmpl x 1] column vector of corresponding target values.
% W: a [nrinputs x 1] vector containing the optimal regression
% coefficients.
W=pinv(X'*X)*X'*t;
```

This function may be called in MatLab as:

```
% Assuming that xtran and t are properly defined and matching in dimension,
% W contains the optimal regression coefficients:
[W]=lintran(xtran, t);
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.36/44

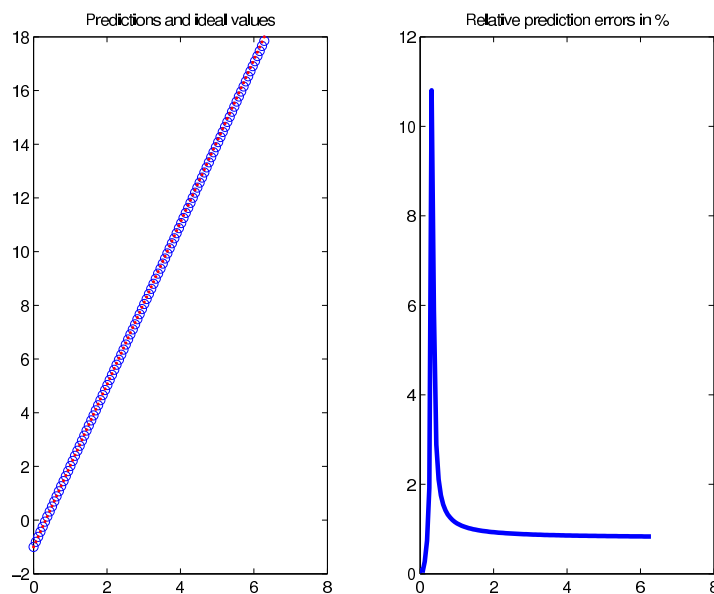
# Test script combining both

```
xtrn=rand(100,1)*2*pi;
theta=[3;-1];
sdns=0.5;
t=xtrn*theta(1)+theta(2)+randn(size(xtrn))*sdns; % noisy targets
% we do now fit the linear regression model
xtrn=[xtrn, ones(size(xtrn))]; % add column of ones for intercept
[W]=lintran(xtrn, t);
% W can now be applied to the test data
xtst=[0:0.01:1]'*2*pi;
xtst=[xtst, ones(size(xtst))]; % add column of ones for intercept
[y_pred]=linpred(xtst, W);
% visualize quality of predictions:
y_ideal=xtst*theta; % using the known correct coefficients
subplot(1,2,1)
plot(xtst(:,1), y_ideal, 'bo', xtst(:,1), y_pred, 'r.')
title('Predictions and ideal values')
subplot(1,2,2)
plot(xtst(:,1), 100*abs(y_pred-y_ideal)./abs(y_ideal), 'b-', 'LineWidth', 3)
title('Relative prediction errors in %')
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.37/44

# Graphical output of script



The current figure (see figure command) can be stored as graphics file:

```
print -depsc lineval.eps % stores the current figure as colour eps file.
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.38/44

# MatLab Functions

MatLab functions like `lintran` and `linpred` are stored as m-files. MatLab's rule is to know a function under the m-file name! For calling

```
[W]=lintran(xtran, t);
```

the corresponding MatLab code **must therefore be stored as file `lintran.m`**.

MatLab functions use local variables.

All parameters are pass by value.

MatLab uses a **search path** for accessing functions. This path is modified using the command **`addpath`**.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.38/44

# MatLab Workspace and Scripts

Statements like:

```
xtrn=rand(100,1)*2*pi;
theta=[3;-1];
sdns=0.5;
t=xtrn*theta(1)+theta(2)+randn(size(xtrn))*sdns; % noisy targets
% we do now fit the linear regression model
xtrn=[xtrn, ones(size(xtrn))]; % add column of ones for intercept
[W]=lintran(xtrn, t);
```

are executed in MatLab

- by either typing each command into MatLabs command line
- or by writing all commands into a script file and executing the script.

Both cases generate all variables (above: xtrn, theta, sdns, t and W) in MatLabs workspace. After execution they are available for further use.

Workspace commands:

who, whos : display all variables at the workspace

load, save: allow loading and saving variables into/from the workspace

clear: remove variables from the workspace.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.40/44

## Scope (“lifetimes”) of Variables I

In general, MatLab variables which are generated within a function have local scope with the following implications:

- They cease to exist as soon as the corresponding function terminates.
- Variables in MatLabs workspace which exist at the time the function gets executed are invisible.

Creating variable “foo” at the workspace:  $foo = [1, 2, 3]$ ; we can thus always create a second variable “foo” inside any function without implications for variable “foo” in MatLabs workspace.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.41/44

## Scope (“lifetimes”) of Variables II

Exception to the previous rule are variables **specifically** defined to have **global scope**:  
global foo; foo = [1, 2, 3]; To force global scope, functions need to specify variables as such.

```
function useglobal % demonstrate global variable
global g_var
if isempty(g_var) % we initialise it
    g_var=1
end
g_var=g_var*2;
```

Calling useglobal  $n$  times results in `g_var` being  $2^n$ . Using global variables at the MatLab prompt or in other functions requires defining them there as global as well.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.42/44

## Handling global variables

```
useglobal; useglobal; useglobal;
global g_var
disp(g_var); % prints 8
clear g_var % only removes the reference!
whos global % prints g_var (which still exists!)
useglobal;
global g_var
disp(g_var); % prints thus 16 !!
clear global g_var % removes the variable
```

`clear global`: removes all global variables

`clear all`: removes all variables including globals

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.43/44

## Summary Part 2

Many **models and algorithms in ML & PRN** are best formulated and **implemented via matrices**.

**MatLab** is very efficient for dealing with **matrices**.

An important **rule for getting fast** implementations is **avoiding loops** and **coding with matrix operations** instead.

Computer science recipes like using **modularised code** and **avoiding side effects** are important.

**Global variables** can help speeding up code but **violate the previous requirement** and should only be used in a separate (final) implementation.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.44/44

# Machine Learning and Pattern Recognition in Bioinformatics

Peter Sykacek<sup>1</sup>

Vienna Science Chair of Bioinformatics

Department of Biotechnology

BOKU University

[peter.sykacek@boku.ac.at](mailto:peter.sykacek@boku.ac.at)

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.1/42

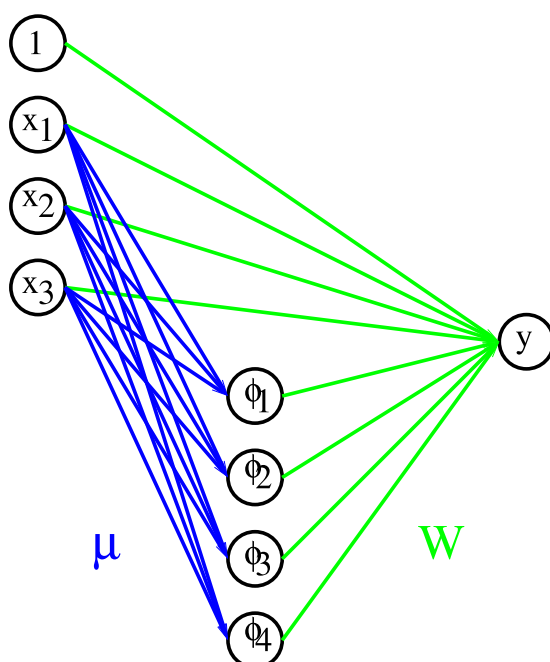
# Modelling Aspects in ML & PRN

ML & PRN are concerned with **algorithms** which are useful in various applications of data analysis. Some properties:

- Like linear regression, ML & PRN **fit models to data**.
- ML & PRN typically use however more **flexible models**.
- Flexible models are advantageous for capturing **nonlinear relationships** often found in biology.
- ML & PRN approaches often suffer however from **identifiability issues and overfitting** (tendency of adding noise to the deterministic part of the model).

**Greater flexibility inevitably renders fitting more difficult!**

## Towards flexible regression



Increase flexibility by introducing a nonlinear mapping  $(\Phi_1, \dots, \Phi_4)$  parameterised by  $\mu = [\mu_1, \dots, \mu_4]$ .  $\Phi_k$  is for example a Gaussian or a *thin plate spline* centred at  $\mu_k$ . Model fitting adjusts kernel parameters  $\mu$  and the regression coefficients  $W$ .

Known as **radial basis function network (RBF)**.



# Thin plate splines

$\Phi_k(\mathbf{x}, \boldsymbol{\mu}_k) = r^2 \log(r)$  with  $r^2 = (\mathbf{x} - \boldsymbol{\mu}_k)^T (\mathbf{x} - \boldsymbol{\mu}_k)$  (square Euclidean distance)

```
function [out]=thinplate(in, K)
% function [out]=thinplate(in, K) calculates size(K,1) output activations
% (columns of out) from each input (in) with each kernel (K) using thin
% plate splines. out has size(K,1) columns and size(in,1) rows.
% (C) P. Sykacek, 2009 <peter@sykacek.net>
nrin=size(in,1);
out=zeros(size(in,1), size(K,1));
for i=1:size(K,1)
    dist=[in-K(i*ones(nrin,1), :)]';
    dist=sum(dist.*dist,1)';
    out(:,i)=dist+eps; % prevent log of zeros
    lgout=log(out(:,i));
    % lgout(find(isinf(lgout)))=0; % this is some debugging output!
    out(:,i)=out(:,i).*out(:,i).*lgout;
end
out=[out, in, ones(nrin, 1)]; % concatenate with original data and ones.
```

# Gaussian Kernel

$\Phi_k(\mathbf{x}, \boldsymbol{\mu}_k) = \exp(-0.5(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu}_k))$ , with  $\boldsymbol{\Lambda}$  ..  
diagonal precision matrix.

```
function [out]=fastgauss(in, K, l)
% function [out]=fastgauss(in, K, l) calculates size(K,1) output
% activations (columns in out) from each input (in) with each
% kernel (K). out has size(K,1) columns and size(in,1) rows. The
% number of columns in in and K must be identical. l are the
% square roots of the precisions in each dimension.
% (C) P. Sykacek, 2009, <peter@sykacek.net>
nrin=size(in,1);
out=zeros(size(in,1), size(K,1));
l=l.^2;
for i=1:size(K,1)
    dist=[(in-K(i*ones(nrin,1), :))];
    dist=sum(dist.*dist.*l(ones(nrin,1), :), 2);
    out(:,i)=exp(-0.5*dist);
end
out=[out, in, ones(nrin, 1)];
```

# Properties of RBF networks

RBF's are **universal function approximators** (e.g. Hartman *et al.*(1990)).

Meaning: by increasing the number of nonlinear kernels,  $\Phi_k$ , “well behaved” functions can be approximated with arbitrary accuracy.

Unless we do fancy things with the nonlinearities, increasing the number of hidden units will increase the flexibility of the model.

# Fitting RBF networks

Typically, RBF networks are trained by two separate steps:

Step 1: adjust the  $\mu$  parameters of the nonlinear expansion using unsupervised methods (e.g. kernels resulting from mixture density estimation).

Step 2: use the data arising from the nonlinear projection to adjust the regression coefficients  $W$ . This is exactly the procedure we know from linear regression.

A proposal for step 1: don't bother with fancy approaches and use a number of **randomly drawn input vectors as kernel centres**. Similarity to **kernel regression** and the **Nadaraya-Watson kernel smoother**.

# Fast RBF training

```
function [net]=lin2nonlin_tran(Xtrn, ttrn, nrkernels, nrnets, nettyp, krnfk
% returns net, a network structure with a committee of nrnets
% RBF networks using thin plate spline (or gaussian) activation.
% If net is 't' we use thin plate spline if net is 'g' we use Gaussians.
% Use lin2nonlin_pred for prediction.
% Xtrn:      [nsamples x nrin] matrix of independent variables (regressors).
% ttrn:      [nsamples x 1] vector of corresponding target values (responses)
% nrkernels: number of kernels to be used in every network
% nrnets:   number of networks in the committee we average over
% nettyp:    type of RBF function, 't' .. thin plate spline, 'g' .. Gaussian,
%            default .. linear
% krnfkt:    factor for defining the precision of Gaussian kernels
% net:       [nrkernels x 1] vector of data structures describing a committee
%            of trained RBF networks.
% net.typ:   type of nonlinear activation (copy of nettyp)
% net.K:     kernel modes
% net.l:     square roots of precisions, used for Gaussians only.
% net.W:     corresponding regression coefficients.
```

MatLab provides the **dot symbol** for denoting **components of data structures**.

# Fast RBF training - Code

```
function [net]=lin2nonlin_tran(Xtrn, ttrn, nrkernels, nrnets, nettyp, krnfk
nrsamples=size(Xtrn,1);
l=krnfkt/(max(Xtrn) - min(Xtrn));
for i=1:nrnets
    net(i).typ=nettyp;
    if nettyp == 't' || nettyp == 'g'
        net(i).K=uniquerndsub(Xtrn, nrkernels);
        net(i).l=l;
    end
    if nettyp == 't'
        out=thinplate(Xtrn, net(i).K);
    elseif nettyp == 'g'
        out=fastgauss(Xtrn, net(i).K, net(i).l);
    else % default linear, where we only augment with '1' for the intercept
        out=[Xtrn, ones(nrsamples, 1)];
    end
    net(i).W=lintran(out, ttrn);
end
```

# RBF predictions

```
function [y]=lin2nonlin_pred(net, Xtst)
% function [y]=lin2nonlin_pred(net, Xtst) predicts targets for Xtst inputs
% net:   array of trained network structures (see lin2nonlin_tran).
% Xtst:  independent test data used as inputs in the regression model.
% y:     predictions obtained by averaging over all committee members.
nrnets=length(net);
nrsmpl=size(Xtst, 1);
y=zeros(nrsmpl, size(net(1).W,2));
for i=1:nrnets
    if net(i).typ=='t'
        out=thinplate(Xtst, net(i).K);
    elseif net(i).typ=='g'
        out=fastgauss(Xtst, net(i).K, net(i).l);
    else
        out=[Xtst, ones(nrsmpl, 1)];
    end
    y=y+out*net(i).W/nrnets;
end
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.10/42

# Assessing the RBF

A) How do we assess the RBF?

One possibility is monitoring how well  $\hat{y}$  predicts the targets  $y$  by evaluating the sum of squares which we minimise during training.

$$\text{SSD}(\text{RBF}_k) = (\mathbf{y} - \hat{\mathbf{y}}_k)^T (\mathbf{y} - \hat{\mathbf{y}}_k)$$

B) Which RBF parameters do influence performance?

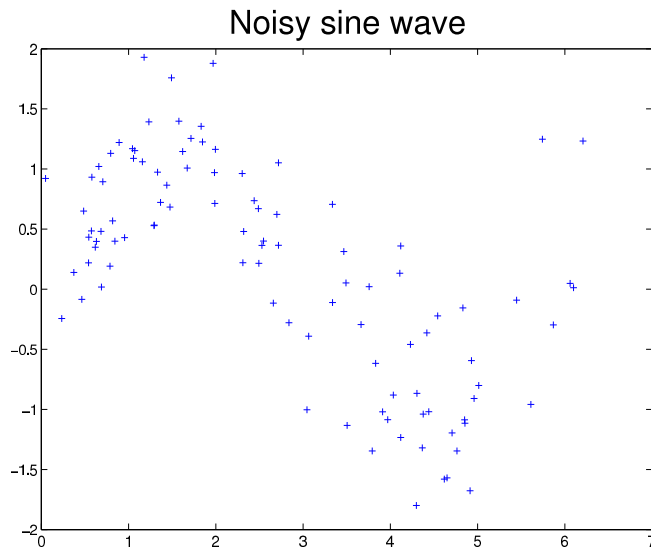
- 1) number of kernels, nrkernels
- 2) network type, nettyp
- 3) Gaussian kernel precision, krnfkt
- 4) number of networks in committee, nrnets

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.11/42

# Assessment on toy data

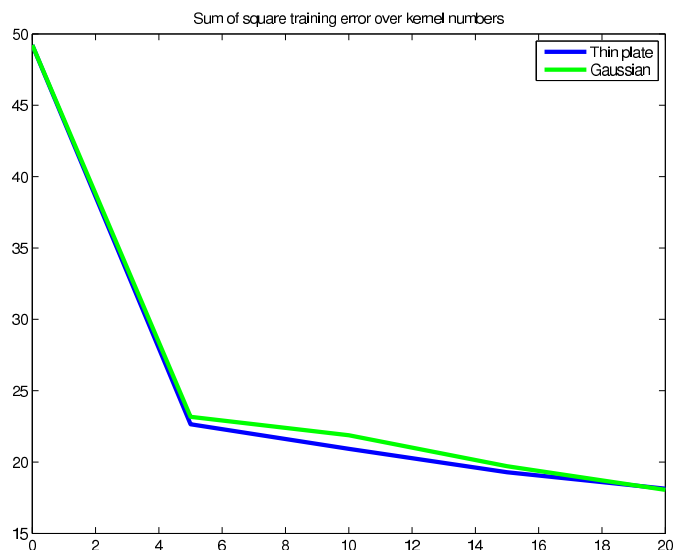
Generate 100 samples of “noisy sine wave”.



```
x=rand(100,1)*2*pi;  
sdn=0.5;  
y=sin(x)+randn(100,1)*sdn;  
plot(x,y,'b+', 'MarkerSize',5);  
title('Noisy sine wave',...  
      'FontSize',24)  
print -depsc noisysine.eps  
save trandata.mat x y
```

jump 2 TOC

D over kernel numbers@

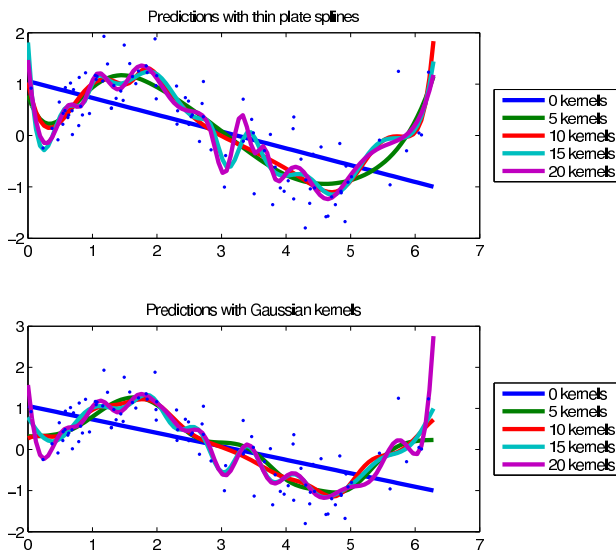


Based squares errors, bc fit the tr when i number the moc

Is this figure good news?

jump 2 TOC

# Visualising the models we fit



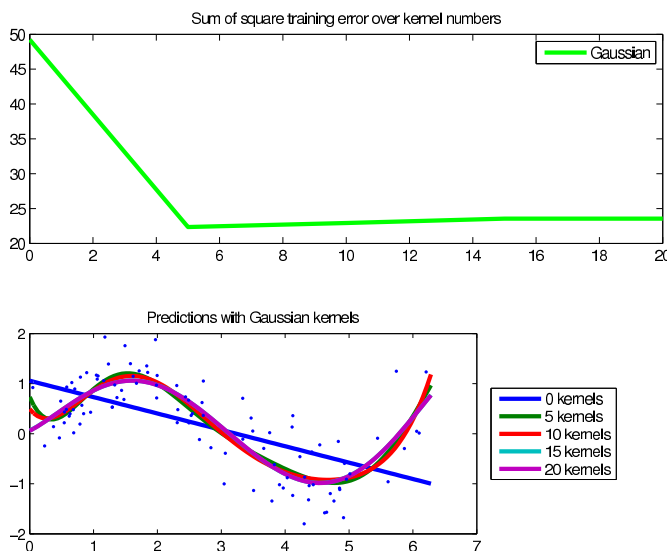
To get these figures we use as inputs data from the range  $0$  to  $2\pi$  and visualise the corresponding predictions.

These visualisations show that **decreasing the sum of squares of prediction errors** below reasonable values is not really a good idea because it **corresponds to fitting noise!**

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.14/42

# Model Complexity



A simple modification of the Gaussian kernel RBF shows that the number of kernels is not representative for model complexity. Exchange (line 2 in `lin2nonlin_tran`):

```
l=krnfkt/((max(Xtrn) - ...
           min(Xtrn)));
```

with:

```
l=krnfkt/((max(Xtrn) - ...
           min(Xtrn))*nkernels);
```

Model complexity depends on the **effective number of parameters**. Increasing the standard deviation of the kernels **regularises**.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.15/42

# Statistical Interpretation of ML & PRN

An experiment provides noisy data

$$\mathcal{Z} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}.$$

Note:  $\mathbf{x}_n$  possibly multivariate i.e. vectors.

Based on  $\mathcal{Z}$ , we have an **inference** (or learning) problem of finding an “optimal” relation between  $x$  and  $y$ :

$$p(y|\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}) + \epsilon(\lambda)$$

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.16/42

# Statistical Interpretation of ML & PRN

An experiment provides noisy data

$$\mathcal{Z} = \{(y_1, \mathbf{x}_1), \dots, (y_N, \mathbf{x}_N)\}.$$

Note:  $\mathbf{x}_n$  possibly multivariate i.e. vectors.

Based on  $\mathcal{Z}$ , we have an **inference** (or learning) problem of finding an “optimal” relation between  $x$  and  $y$ :

$$p(y|\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}) + \epsilon(\lambda)$$

**Noise** requires a **deterministic** and a **random** component.

– > **Inherent uncertainty,  $y$  is a random variable.  
We usually predict its expectation.**

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.16/42

# Two Stages of Inference

Parameter Inference:

Implies knowing  $f(\mathbf{x}; \boldsymbol{\theta})$  and the noise model  $\epsilon(\lambda)$  up to unknown parameters ( $\boldsymbol{\theta}$  and  $\lambda$ ) which we will be **inferring from data**.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.17/42

# Two Stages of Inference

Parameter Inference:

Implies knowing  $f(\mathbf{x}; \boldsymbol{\theta})$  and the noise model  $\epsilon(\lambda)$  up to unknown parameters ( $\boldsymbol{\theta}$  and  $\lambda$ ) which we will be **inferring from data**.

Model Inference:

A more realistic assumption is that the model class is unknown and we will be **inferring model class (e.g. number of kernels) and parameters**.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.17/42



# Assessing Model Parameters

Idea: subtract the deterministic part from  $y_n$ :

$$\epsilon_n = y_n - f(\mathbf{x}_n; \boldsymbol{\theta})$$

For convenience introduce  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and  $\mathcal{D} = \{y_1, \dots, y_N\}$ . Assuming that  $\epsilon_n$  are i.i.d samples, we get the **likelihood function**:

$$p(\mathcal{D}|\boldsymbol{\theta}, \lambda, \mathcal{X}) = \prod_n p(y_n|\boldsymbol{\theta}, \lambda, \mathbf{x}_n)$$

which is a suitable objective function to be maximised for  $\boldsymbol{\theta}$  and  $\lambda$ .

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.18/42

# Likelihood and Regression

Assuming  $N$  samples and  $\epsilon_n$  being Gaussian:

$$p(y_n|\mathbf{x}_n; \boldsymbol{\theta}, \lambda) = (2\pi)^{-0.5} \lambda^{0.5} \exp(-0.5\lambda(y_n - \mathbf{x}_n^T \boldsymbol{\theta})^2) \text{ and}$$
$$p(\mathcal{D}|\mathcal{X}; \boldsymbol{\theta}, \lambda) = (2\pi)^{-\frac{N}{2}} \lambda^{\frac{N}{2}} \exp(-0.5\lambda \sum_n (y_n - \mathbf{x}_n^T \boldsymbol{\theta})^2)$$

Taking the log, we get the **log likelihood**:

$$\text{llh} = \frac{N}{2}(\log(\lambda) - \log(2\pi)) - 0.5\lambda(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

which, if we consider maximising for  $\boldsymbol{\theta}$  only, is a familiar expression.

– > **minimising least squares assumes Gaussian noise!**

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.19/42

# Central Problem in ML & PRN

True model - linear regression, Gaussian noise:

$$p(y|\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}) + \epsilon(\lambda)$$
$$f(\mathbf{x}; \boldsymbol{\theta}) = [1, \mathbf{x}^T] \boldsymbol{\theta} \quad \text{and} \quad \epsilon(\lambda) = \mathcal{N}(\epsilon; 0, \lambda)$$

Finite sample size and different model classes:

“**Phone book**”: Perfect memorising of all  $y_n$ ,  
modelling error 0,  $\lambda \rightarrow \infty$ ,  $p(\mathcal{D}|\boldsymbol{\theta}, \lambda, \mathcal{X}) \rightarrow \infty$ .

–  $\rightarrow$  likelihood, ssd and information theoretic  
measures are unsuitable for comparing models!

Central Problem in ML & PRN: **controlling complexity**

## Summary Part 3

ML & PRN are like statistics concerned with  
**fitting models to data**.

The **objective functions** used in ML & PRN for  
model fitting are **largely identical** to those known  
in **statistics**.

**Models** used in ML & PRN are typically more  
**flexible** than those commonly used in statistics.

The core problem in ML & PRN is **controlling  
complexity** and **comparing model classes**.

# Complexity control and testing models

ML & PRN models can be **arbitrary complex** (cf. universal function approximator).

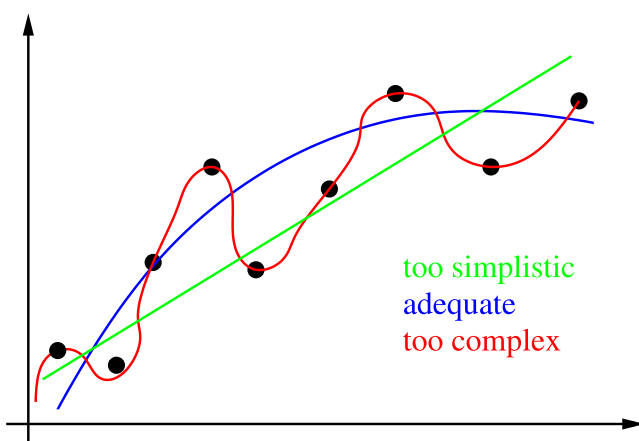
A very important aspect of applied ML & PRN is thus understanding and using methods for

- **choosing models of appropriate dimension** and
- **comparing and assessing models.**

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.22/42

## Assessing model complexity



Two approaches for obtaining appropriate models:

- 1) Testing different model classes (e.g. number of basis functions).
- 2) Allowing for highly flexible models and constraining the model parameters (e.g. regularisation).

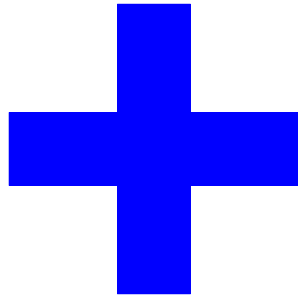
Both approaches can either be dealt with by **augmenting simple objective measures** like SSD or by using **validation data.**

jump 2 TOC

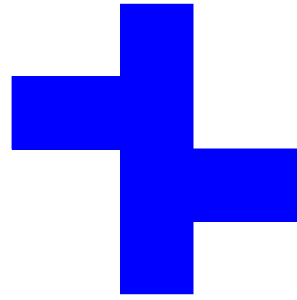
ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.23/42

# Guess the Correct “Model”

How many components?



Object A



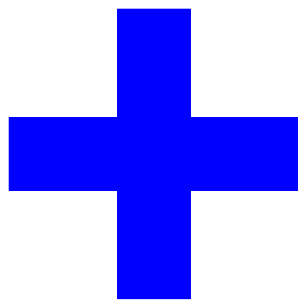
Object B

[jump 2 TOC](#)

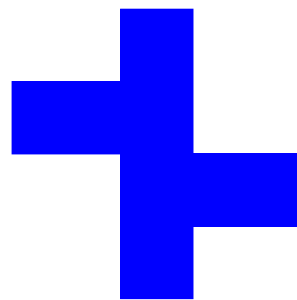
ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.24/42

# Guess the Correct “Model”

How many components?

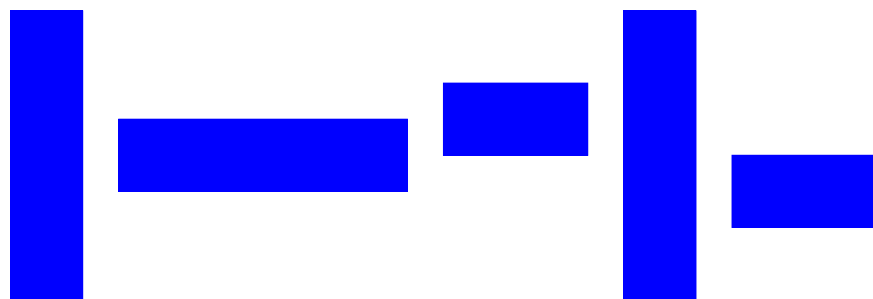


Object A



Object B

Most likely answer:



[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.24/42

# Occam's Razor

We implicitly apply Occam's Razor



William of Occam (or Ockham)  
(1288 - 1348)

*Entia non sunt multiplicanda sine necessitate*: Entities are not to be multiplied without necessity.

Interpretation: One should always opt for an explanation in terms of the fewest possible number of causes, factors, or variables.

Material from [http://en.wikipedia.org/wiki/William\\_of\\_Ockham](http://en.wikipedia.org/wiki/William_of_Ockham).

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.25/42

# Instances of model selection

Problems in context of model selection:

- **Variable selection**: search for input subsets which improve predictive performance or identify important variables (e.g. differentially expressed genes).
- **Changepoint detection**: Separating data into groups which show similar statistical properties.
- **Clustering**: (see above)
- **Determining optimal model orders** (applies to most ML& PRN methods!)
- **Determining suitable noise characteristics**.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.26/42

# Occam's Razor in ML & PRN

Model selection can be done by the following approaches:

- Adding a **complexity penalty** to the objective function. Many choices and active research field: AIC (Akaike information criterion), BIC (Bayesian information criterion), MDL (minimum description length), etc.
- Using **learning** methods like Bayesian inference **with Occam's razor built in**.
- Using **empirical approaches** comparing model classes etc. by **validation testing** (computer simulation using independent data).

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.27/42

## AIC, BIC and MDL

Akaike (1974) suggests using AIC:

$$L_{\text{AIC}}^k = \log(p(\mathcal{D}|\hat{\theta})) - K_k$$

Schwarz (1978) suggested BIC (asymptotic statistics):

$$L_{\text{BIC}}^k = \log(p(\mathcal{D}|\hat{\theta})) - \frac{1}{2}K_k \log(N)$$

Rissanen (1978) suggested MDL (information theory):

$$L_{\text{MDL}}^k = L(\mathcal{D}|k) + L(k)$$

$p(\mathcal{D}|\hat{\theta})$  ... likelihood of size  $N$  dataset  $\mathcal{D}$ ,  $\hat{\theta}$  optimal model coefficients,  $K_k$  .. dimension,  $L(\mathcal{D}|k)$  .. bit length modelling errors,  $L(k)$  .. bit length model  $k$ .  $L_{\text{MDL}}^k$  is equivalent to the Bayesian log posterior model probability plus an additive constant which is independent of the models under consideration.

**AIC, BIC should be maximised, MDL should be minimised.**

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.28/42

# Application to K-means clustering

**K-means clustering** is defined by a simple algorithm which **iteratively minimises the cost function**

$$J = \sum_{k=1} \sum_{x_{n,k} \in \mathcal{S}_k} (x_{n,k} - \mu_k)^2.$$

Upon termination, we get **cluster centres**, **cluster information** and the **optimal  $J$**  as a result.

For using AIC and BIC, we need a **quick and dirty hack** for getting a **likelihood measure** from K-means. Assuming as generative model a Gaussian mixture, hard partitioning and an isotropic precision of  $2 * I$ :

$$\log(p(\mathcal{D}|k)) \approx \sum_k n_k \log(P_k) - J \text{ with } P_k = n_k/N,$$

allows us **calculating  $L_{AIC}^k$  and  $L_{BIC}^k$** .

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.29/42

## MatLab script: calculation

```
% use netlab library. modify path to install dir on your system!
addpath('~/sciwrk/mlbsrc/netlab/')
% load microarray data, etc. from a multi variable .mat file
load mca_tc
% analyse data containing the average of two replicates measured at
% 11 different timepoints stored in the [n-timepoints x n-genes] matrix
% mndat. For clustering similar genes into groups we transpose the data.
clustin=mndat';
% cluster range considered: 15 to 25 clusters.
mink=15; maxk=25;
% 10 repeated fits of k-means and maxit iterations until convergence
nrep=10; maxit=500;
method='A'; % AIC
[alloc_aic, centers_aic, nL_aic]=wrap_kmeans(clustin, mink, maxk, nrep,...
      method, maxit);
method='B'; % BIC
[alloc_bic, centers_bic, nL_bic]=wrap_kmeans(clustin, mink, maxk, nrep,...
      method, maxit);
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.30/42

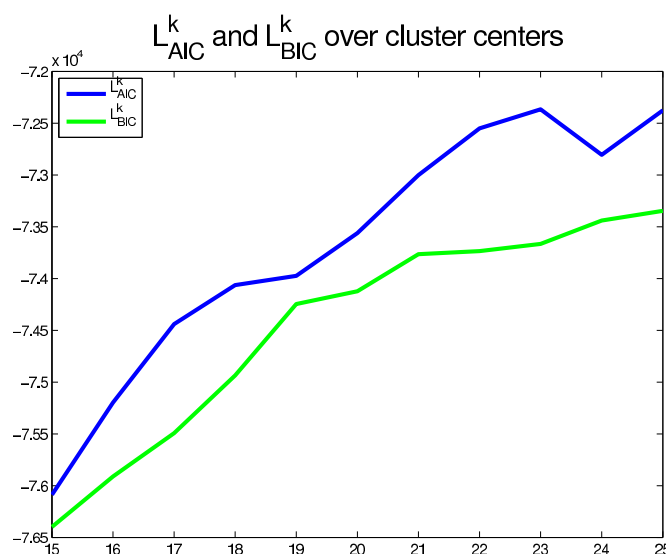
# MatLab script: illustrations

```
% note negation since above calculation provides penalized neg. llh
plot([mink:maxk], -nL_aic, 'b-', [mink:maxk], -nL_bic,...
     'g-', 'LineWidth', 3);
title('L_{AIC}^k and L_{BIC}^k over cluster centers', 'FontSize', 20);
legend('L_{AIC}^k', 'L_{BIC}^k', 'Location', 'NorthWest');
print -depsc kmeans_aicbic.eps
figure(2)
subplot(2,1,1)
plot(centers_aic', 'LineWidth', 3);
title('Cluster centers obtained using AIC penalty', 'FontSize', 18);
subplot(2,1,2)
plot(centers_bic', 'LineWidth', 3);
title('Cluster centers obtained using BIC penalty', 'FontSize', 18);
print -depsc kmeans_cluster.eps
Note: see also MatLab code of functions!
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.31/42

## K-means, AIC and BIC over kernels



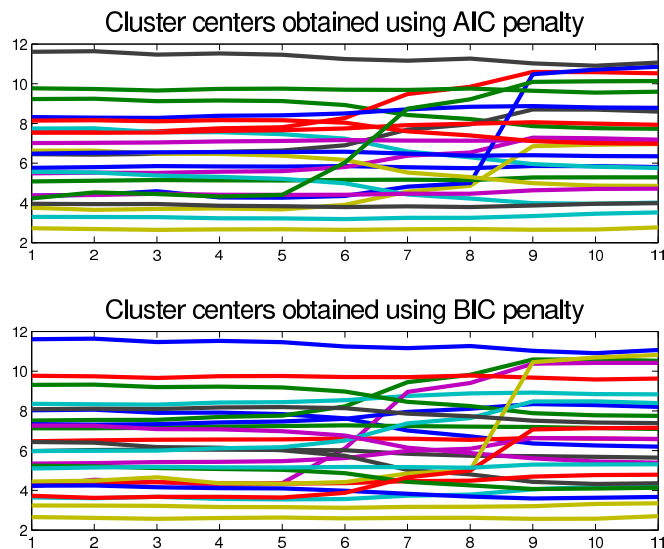
Search for optimal kernel numbers peaks at values larger 25. Apparent difficulty of finding good optima at larger kernel numbers.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.32/42



# Cluster centres



Input: all ( $\approx 12500$ ) genes measured. Several clusters represent house keeping stuff which does not change expression. As a negative side effect of having used unsupervised data analysis, we have thus to analyse the result further for interesting components.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.33/42

# Model selection using validation data

Applying a nonlinear RBF, we saw that increasing complexity will always improve the objective function.

This problem can be overcome by using independent data to assess models with different complexity. Difficulties:

- Samples are usually very precious and thus rare.
- To get good results from model building, we can not afford wasting many samples for validation purposes.
- To get reliable assessments of model complexity, we also need sufficiently many samples.

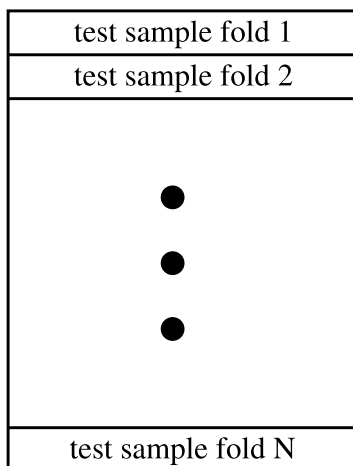
Seems like a deadlock, however: – > iterative schemes which reuse data can help.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.34/42

# N-Fold Cross Testing

## Sketch and MatLab like pseudo code



```
allres=[]; allpred=[];
for n=1:n_folds
    % split into training and test data
    [train, test]=foldsplit(orig_data, n_folds, n);
    % model inference
    [model]=trainfunc(train, fiddleparams);
    % store this folds true targets and predictions
    [res]=truetarg(test);
    [pred]=predtarg(test, model);
    allres=[allres; res];
    allpred=[allpred; pred];
end
```

**Leave one out** has as many folds as samples. An alternative by resampling with replacement is called **Bootstrapping**.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.35/42

# Means for complexity control

In addition to changing model order, a technique called **regularisation** allows to control the complexity during learning. Regularisation penalises non-smooth functions by **adding a suitable penalty term to the objective function**:

- **Multi layer perceptrons (MLPs)** have sigmoidal shaped nonlinearities which, for input values close to zero, are approximately linear. By forcing small coefficients, we can thus avoid excessive nonlinear behaviour:  $\gamma \theta^T \theta$  is thus added as **penalty term** to the objective function. Note that considerations of scale consistency require using one  $\gamma$  per layer.
- If the behaviour of the resulting predictor is not directly connected to the model parameters, the **penalty is formulated via the norm (size) of the second derivatives of the resulting function**. (Bishop 1995), suggests for example using

$$\gamma \sum_{n,i} \left( \frac{\partial^2 f(x_n, \theta)}{\partial x_i^2} \right)^2$$

as penalty term for RBF networks.

- A rather ad-hoc approach to regularisation can be obtained by adjusting the width of the kernel function.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.36/42

# Adjusting kernel widths

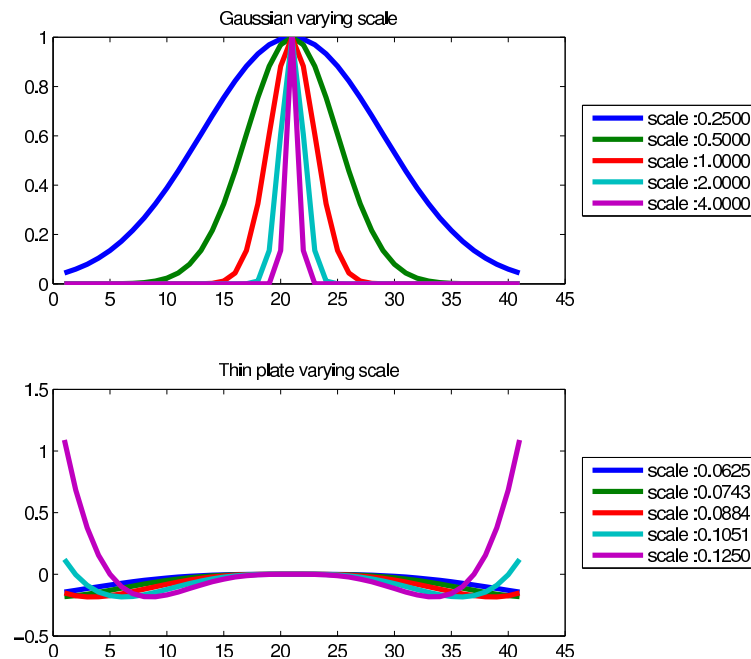


Illustration of the effect of adjusting the scales (variance) of the kernel functions.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.37/42

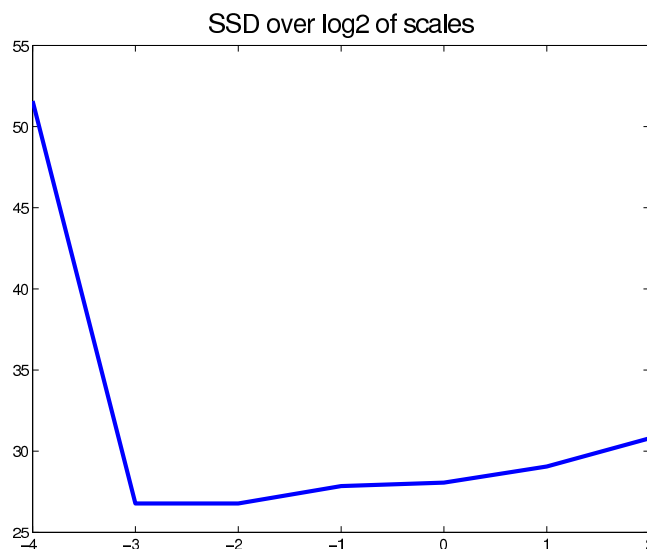
# Validating regularisation penalties

```
for n=1:length(myscaleg)
    csg=myscaleg(n); % get n-th length scale
    % init store for true targets and predictions.
    alltrtrg=[]; allpred=[]; allscales(n)=csg;
    nfold=10;
    for fldno=1:nfold
        % split training data into inputs and targets
        [tr, tst]=nfldspl(data, nfold, fldno);
        xtr=tr(:,1); ttr=tr(:,2);
        alltrtrg=[alltrtrg;tst(:,2)]; xtst=tst(:,1);
        [net, tranpred, ssd]=rbf_optim(xtr, ttr, nkrnl, nrnets, nettyp,...
            nrep, csg);
        % prediction on xtst: unbiased performance measures
        [ytst]=lin2nonlin_pred(net, xtst); allpred=[allpred; ytst];
    end
    % calculate ssd on the independent test data obtained by cv!
    regssd(n)=sum((alltrtrg-allpred).^2);
end
```

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.38/42

# Sum of squares over length scale

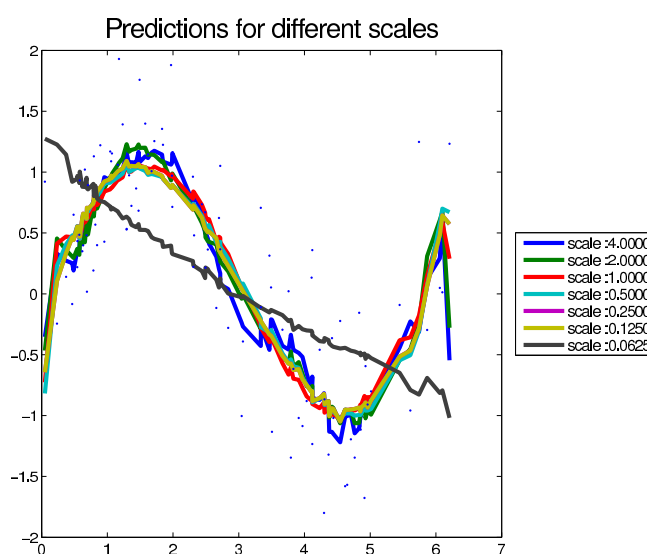


The **most important** aspect is using **independent** validation data. **Predictions** are thus **unbiased** and we obtain **meaningful complexity assessments**.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.39/42

# Corresponding Predictions



Although regularisation helps, further improvements are possible. Further options: **constraining the norm** of model coefficients ( $\theta$ ) or **model averaging** (e.g. bagging).

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.40/42

# Summary Part 4

ML & PRN provide different approaches (**complexity penalties, “direct” regularisation and model averaging**) which help

- getting **model complexities** right and
- obtaining solutions which **generalise well** on independent test data.

Technically these methods allow

- selecting most promising model classes and
- tuning of “fiddle parameters” like the amount of “force” ( $\gamma$ ) to obtain “small” model coefficients.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.41/42

# Bibliography

- Akaike H., A new look at the statistical identification model, in *IEEE Trans. Auto. Control* 19, 716-723, 1974.
- Bishop C. M., *Neural Networks for Pattern Recognition*, Clarendon Press Oxford, UK, 1995.
- Hartman E. J., Keeler J. D. and Kowalski J. M., Layered neural networks with Gaussian hidden units as universal approximations, in *Neural Computation* 2, 210-215, 1990.
- Rissanen J., Modeling by the shortest data description, in *Automatica* 14, 465-471, 1978.
- Schwarz G., Estimating the dimension of a model, in *The annals of Statistics* 6, 461-464, 1978.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.42/42

# Machine Learning and Pattern Recognition in Bioinformatics

Peter Sykacek<sup>1</sup>

Vienna Science Chair of Bioinformatics  
Department of Biotechnology  
BOKU University  
peter.sykacek@boku.ac.at

Jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.1/24

## Likelihood and Classification

“Classification” often used synonymously for regression with discrete outcomes. Likelihood of regression model:

$$P(\mathcal{D}|\mathcal{X}; \boldsymbol{\theta}) = \prod_n P(y_n|\mathbf{x}_n, \boldsymbol{\theta})$$

To enforce  $\sum_{y_n} P(y_n|\mathbf{x}_n, \boldsymbol{\theta})$  is 1, we apply a suitable output transformation, e.g. the cdf of the logistic distribution:

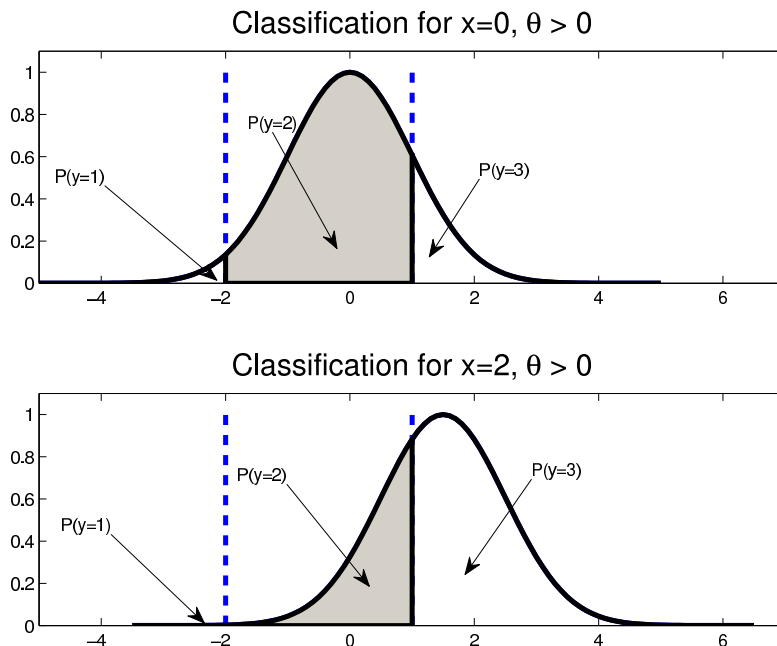
$$P(y_n|\mathbf{x}_n^T \boldsymbol{\theta}) = \frac{1}{1 + \exp((2y_n - 1)\mathbf{x}_n^T \boldsymbol{\theta})}$$

Probabilities are certainty measures about which label to predict. The **linearly parameterised logistic distribution** is optimal for a two class problem where both classes are **Gaussian distributed and have identical covariances**.

Jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykacek – p.2/24

# Ordered multi class problems



Two main types of 1-of- $c$  classification problems ( $c$ : number of classes). For ordered classes like quality or severity indices:

$$\log(\gamma_k(\mathbf{x})/(1-\gamma_k(\mathbf{x}))) = \beta_k - \boldsymbol{\theta}^T \mathbf{x} \text{ with } \gamma_k(\mathbf{x}) = P(y \leq k).$$

All classes share regression coefficients but differ in threshold  $\beta_k$ . To **identify** the model, one  $\beta_k$  is fixed.

The idea is best understood as a logistic distribution with mode  $\boldsymbol{\theta}^T \mathbf{x}$  and probabilities for class corresponding to the area under the pdf between neighbouring thresholds.

# Unordered multi class problems

Are usually dealt with by the output mapping:

$$P(y = k|\mathbf{x}) = \frac{\exp((\Phi(\mathbf{x}))^T \boldsymbol{\theta}_k)}{\sum_k \exp((\Phi(\mathbf{x}))^T \boldsymbol{\theta}_k)}$$

In ML & PRN this is called **softmax transformation**. The  $\Phi$  should remind you that these output transformations can be used **in combination with nonlinear mappings** of the original inputs (RBF, MLP etc.). The softmax transformation can be regarded as **multivariate generalisation of the logistic cdf**. Expressing  $P(y = k|\mathbf{x})$  vs.  $P(y = \neg k|\mathbf{x})$  results consequently in a logistic cdf.

A note on evaluating softmax: for **numerical reasons** one should rather **express it as**  $P(y = k|\mathbf{x}) = 1/(1 + \sum_{i \neq k} \exp(\Phi(\mathbf{x})^T (\boldsymbol{\theta}_i - \boldsymbol{\theta}_k)))$ .

# Error functions for classification

Taking the log of the above likelihood results in:

$$L_{\text{class}} = \sum_n \log P(y_n | \mathbf{x}_n, \boldsymbol{\theta})$$

The negative log likelihood, also called **cross entropy error**, is in classification the **logical counterpart of sum of square differences** in regression:

$$E_{\text{ce}} = - \sum_n \log P(y_n | \mathbf{x}_n, \boldsymbol{\theta})$$

Note however that **SSD** and a zero one target coding is still guaranteed providing **asymptotically posterior probabilities for class**.

# Fitting Classifiers

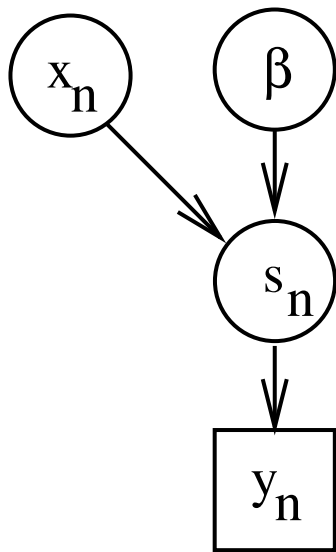
**Nonlinear** output mapping, makes **fitting classifiers more difficult** than fitting linear in the parameter models.

Requires generic nonlinear optimisers (see Press *et al.* (2007)):

- gradient descent:  $\delta \boldsymbol{\theta} \propto -\nabla_{\boldsymbol{\theta}} E_{\text{ce}}(\boldsymbol{\theta})$ .
- quasi Newton:  $\delta \boldsymbol{\theta}_n = -\alpha_n \boldsymbol{\Lambda}(\boldsymbol{\theta}_n) \nabla_{\boldsymbol{\theta}} E_{\text{ce}}(\boldsymbol{\theta}_n)$ , with  $\boldsymbol{\Lambda}$  approx. Hessian  $\nabla \nabla_{\boldsymbol{\theta}}^T E_{\text{ce}}(\boldsymbol{\theta})$  at minimum. Close to optimum,  $E_{\text{ce}}(\boldsymbol{\theta})$  approx. quadratic – > convergence speed close to Newton method.
- conjugate gradient similar to quasi Newton, however only memory requirement for  $\dim(\boldsymbol{\theta})$  vector instead of  $[\dim(\boldsymbol{\theta}) \times \dim(\boldsymbol{\theta})]$  matrix.



# EM and Classification



$$p(s_n | \mathbf{x}_n, \boldsymbol{\theta}) = (2\pi)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(s_n - \mathbf{x}_n^T \boldsymbol{\theta})^2\right)$$

$$P(y_n = 1 | s_n) = \begin{cases} 1, & \text{if } s_n > 0 \\ 0, & \text{if } s_n \leq 0 \end{cases},$$

$$P(y_n = 0 | s_n) = 1 - P(y_n = 1 | s_n)$$

$$\log(P(y_n | \mathbf{x}_n, \boldsymbol{\theta})) \geq \int_{s_n} Q(s_n) \left( \log(P(y_n | s_n)) + \log(p(s_n | \mathbf{x}_n, \boldsymbol{\theta})) - \log(Q(s_n)) \right) ds_n$$

$$E_{ce}(\boldsymbol{\theta}) \leq - \sum_n \left[ \int_{s_n} Q(s_n) \left( \log(P(y_n | s_n)) + \log(p(s_n | \mathbf{x}_n, \boldsymbol{\theta})) - \log(Q(s_n)) \right) ds_n \right]$$

# Probit EM

Expectation maximisation (EM) scheme for  $E_{ce}(\boldsymbol{\theta})$ .

Minimise iteratively w.r.t all  $Q(s_n)$  and the parameters  $\boldsymbol{\theta}$ .

$Q(s_n)$ : truncated Gaussian with mean  $\mathbf{x}_n^T \boldsymbol{\theta}$ , std. 1 and  $y_n$  dependant truncation:

$$Q(s_n) \propto \mathcal{N}(s_n; \mathbf{x}_n^T \boldsymbol{\theta}, 1) \Pi[\gamma_n, \delta_n]$$

with  $\Pi[\gamma_n, \delta_n]$  denoting an improper uniform density.

$\Pi[\gamma_n, \delta_n] = [-\infty, 0]$  if  $y_n = 0$  and  $\Pi[\gamma_n, \delta_n] = [0, \infty]$  if  $y_n = 1$ .

Optimal  $\hat{\boldsymbol{\theta}}$ :

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{s}$$

$\mathbf{s}$  denotes all expectations  $\langle s_n \rangle_{Q(s_n)}$  as column vector.

# Probit EM on toy data

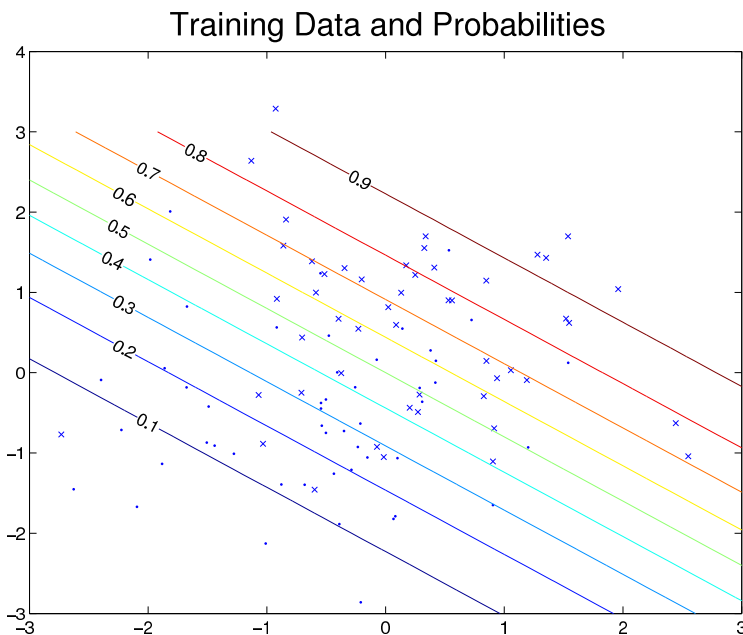


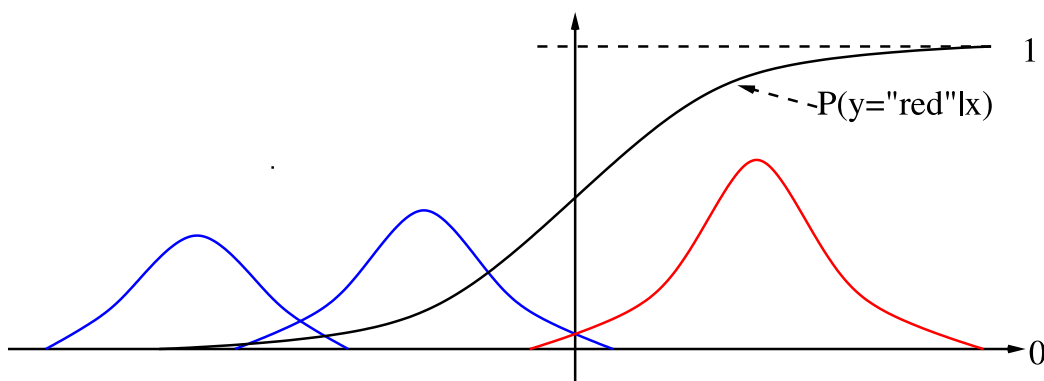
Illustration of data drawn from two Gaussians and the probabilities for class one ('x') obtained by a linear probit model. **Nonlinear decision boundaries** can be obtained without

changing the EM algorithm. Similar to RBF regression, we only have to **project inputs onto a nonlinear basis**.

# Classification and Sampling Paradigm

$$P(y_n | \mathbf{x}_n) = \frac{P(y_n)p(\mathbf{x}_n | y_n)}{p(\mathbf{x}_n)}$$

– > **Bayes theorem** suggests that we can also model **class priors**  $P(y_n)$  and **class conditional densities**  $p(\mathbf{x}_n | y_n)$ .



Advantage: a useful density model, disadvantage: more complicated

# Model Diagnostic Measures

The mean square generalisation error ( $\text{MSE}_{test}$ ) allows assessing regression models.

$$\text{MSE}_{test} = \frac{1}{N} (\mathbf{y}_{test} - \mathbf{f}(\mathbf{X}_{test}; \boldsymbol{\theta}_{train}))^T (\mathbf{y}_{test} - \mathbf{f}(\mathbf{X}_{test}; \boldsymbol{\theta}_{train}))$$

Classification aims at labelling novel samples correctly.

This is tested by the generalisation accuracy  $\text{acc}_{test}$ .

$$\forall n \hat{\mathbf{y}}_{test}[n] = \text{argmax}_k (P(y = k | \mathbf{X}_{test}[n, :] \boldsymbol{\theta}_{train}))$$
$$\text{acc}_{test} = \frac{1}{N} \sum_n \delta(\mathbf{y}_{test}[n], \hat{\mathbf{y}}_{test}[n])$$

We classify such that the most probable class wins and estimate the fraction of correctly classified test cases.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.11/24

# Testing vs. Validation

Validation: **repeated assessments** of ML & PRN methods in order to **calibrate fiddle parameters** of algorithms.

Implication: **tuning** of some algorithmic aspects **towards matching the validation data well**.

The final **results** are thus **no longer unbiased**. We have to be careful here: **too much validation has similar effects as overfitting!**

Validation does not provide valid generalisation accuracies. The method will inevitably look too good.

Testing: assessment of **one and only one** model on **independent test data**.

The idea is obtaining unbiased generalisation results to

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok - p.12/24

# More on Assessing Classifiers

What is the implication of predicting the label which got largest posterior probability?

Consider an individual prediction and that we know the correct posterior  $P(y|\mathbf{x})$ : deciding for label  $y = 1$  implies being correct with probability  $P(y = 1|\mathbf{x})$  and being wrong with probability  $1 - P(y = 1|\mathbf{x})$ .

If we decide for the label  $k = \operatorname{argmax}_k(P(y = k|\mathbf{x}))$ , we will thus minimise the overall number of missclassifications.

What if the missclassifications cost is class dependant?

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.13/24

# Missclassification Cost

Classifying correctly induces zero cost. A false negative for class  $k$  induces cost  $\alpha_k$ .

Predicting  $y = t$  results in an **expected cost** (deciding that the true  $y$  is probably  $t$ , this is the cost we expect to suffer):

$$C = \sum_{k \neq t} P(y = k|\mathbf{x})\alpha_k$$

Missclassification cost  $C$  is thus minimised by classifying  $y = \operatorname{argmax}_k(P(y = k|\mathbf{x})\alpha_k)$ .

This structure is commonly found in life science applications. In cancer screening a false negative is obviously much worse than a false positive.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.14/24

# ROC Curve I

For unknown missclassification cost, the Receiver Operating Characteristic (ROC) curve provides means for assessing binary classifiers.

$$\text{Sensitivity: } s(\gamma) = \frac{\sum_{n|y_n=1} \delta(P(y_n = 1|\mathbf{x}_n) > \gamma)}{N_+}$$

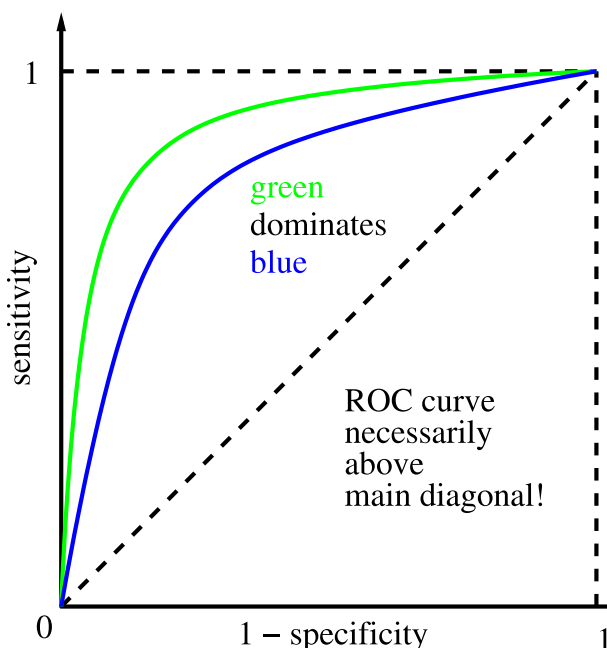
$$\text{Specificity: } p(\gamma) = \frac{\sum_{n|y_n=0} \delta(P(y_n = 1|\mathbf{x}_n) < \gamma)}{N_-}$$

depend on a detection threshold  $\gamma$ ;  $\delta(\cdot)$  maps “true” to 1 and “false” to 0;  $N_+$  is the number of positive and  $N_-$  the number of negative samples.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.15/24

# ROC Curve II



ROC curve:  $s(\gamma)$  over  $1 - p(\gamma)$

Classifier  $a$  dominates  $b$  if

$$\forall s_a, s_b = s : p_b(s_b^{-1}) < p_a(s_a^{-1})$$

Given domination, the area under the ROC curve (AUC) provides a quality measure of the classifier.

AUC and survival probability:  $AUC = P(P(y = 1|x \sim p(x|1)) > P(y = 1|x \sim p(x|0)))$ .

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeek – p.16/24

# Comparing Classifiers

Verification particular choices like the chosen model and fiddle parameters requires performance comparisons.

## Default Accuracy:

The most trivial competitor predicts for every sample it's prior probability and thus always the majority class.

## Competing Classifiers:

Since it is not at all clear, whether a particular model works well in the case at hand, one should use a set of different approaches. A simple, yet still powerful approach is the so called **k nearest neighbour** classifier.

Given several performance measures, we should **investigate, whether differences are significant.**

# Mc Nemars Test

The idea of Mc Nemars test is that differences in classifiers manifest themselves in differently classified samples. This allows producing a 2 by 2 contingency table:

	$C_{2w}$	$C_{2c}$
$C_{1c}$	$n_a$	$n_{bthc}$
$C_{1w}$	$n_{bthw}$	$n_b$

$C_{xc}$  and  $C_{xw}$  refer to samples correctly / wrongly labelled by classifier  $x$ . Differences manifest in the number of samples,  $n_a$ , which  $C_1$  labels correctly and  $C_2$  labels wrongly and the number of samples,  $n_b$ , where the situation is vice versa.

If both classifiers are equal,  $(n_a, n_b)$  is a sample of drawing  $n_a + n_b$  times from a Binomial distribution with probability 0.5. The null hypothesis of Mc Nemars test is the Binomial  $\mathcal{Bn}(n_a + n_b, 0.5)$  and we obtain the p-value by calculating the tail probability from  $(n_a, n_b)$  onwards.

# Continuous latent variable models

Common aspect:  $x = [m+]Wt[+\epsilon]$ ,  $W : [d \times k]$  dimensional coefficients matrix,  $m, \epsilon$  : optional mean and noise term.

**PCA (principle component analysis):**  $t \sim \mathcal{N}(t; 0, \Lambda)$ ,

$\Lambda : [d \times d]$  diagonal cov. matrix,  $\epsilon : 0$ ,  $m$ : data mean.

**Factor analysis:**  $t \sim \mathcal{N}(t; 0, \Lambda)$ ,  $\Lambda : [d \times d]$  general cov. matrix,  $\epsilon_k \sim N(\epsilon_k; 0, \lambda_k)$  and  $m$ : data mean.

**ICA (independent component analysis):**  $t \sim \prod_d p(t_d|\theta_d)$  and  $p(t_d|\theta_d)$  : univariate density functions, at maximum one Gaussian,  $m, \epsilon$  : implementation dependent.

These models provide as summary a lower dimensional representation of the data – > **dimensionality reduction**.

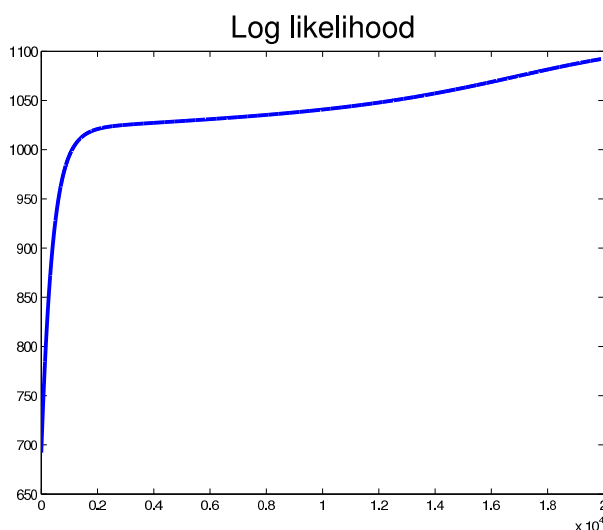
jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.19/24

## ICA Optimisation

ICA and one observation:

$$p(\mathbf{x}_n|W) = \int_{\mathbf{t}_n} \left( p(\mathbf{x}_n|W, \mathbf{t}_n) \prod_k p(t_{n,k}) d\mathbf{t}_n \right).$$



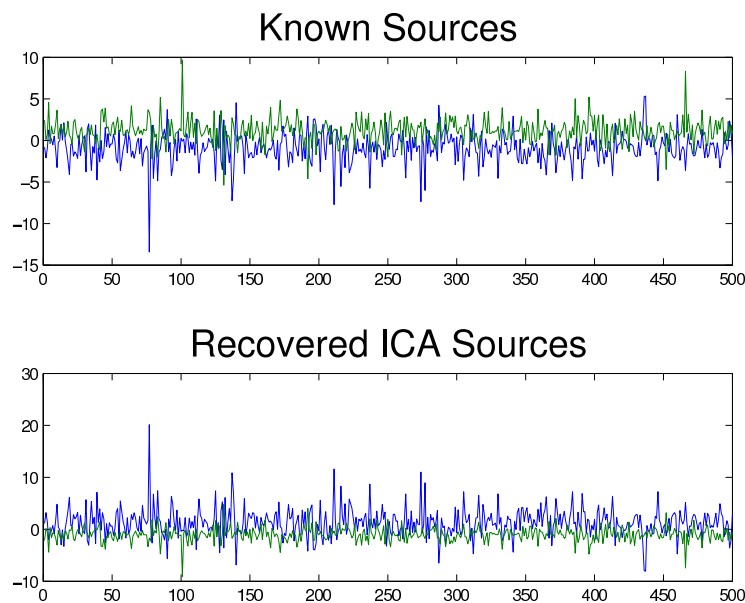
For noiseless observations,  $p(\mathbf{x}_n|W, \mathbf{t}_n) = \prod_j \delta(x_{n,j} - W[j, :]t_n)$ .

As can be seen, convergence of ICA is rather slow.

jump 2 TOC

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.20/24

# Original and Recovered Sources



Note the identifiability problem which loses

- 1) ordering (sources appear in different order)
- 2) scale including sign.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.21/24

# Mixture Density Models

**Gaussian mixture model:**  $p(x|t = k) = \mathcal{N}(x; \mu_k, \lambda_k)$ , i.e. a (possibly multivariate) Gaussian density function.

**K-means clustering:** can be regarded as a mixture density model with  $P(t = k) = 1/K$  and  $p(x|t = k)$  being  $K$  uniform densities with domains emerging from the Voronoi tessellation defined by the  $K$  cluster centres.

**Hidden Markov Model:** assumes a one dimensional ordering (e.g. time) among the latent variables  $t_n$ . We have thus a more complicated prior:  $P(t_n|t_{n-1})$ .

These models infer as summary information which mixture component generated the data point  $x_n \rightarrow$  **Clustering**.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.22/24



# Fitting a Gaussian Mixture - EM

$$\log(p(\mathcal{D}|\{P_k, \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k, \forall k\})) = \sum_n \log\left(\sum_k P_k p(x_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k)\right)$$

Problem: log of sum:

$$\log\left(\sum_k \log(P_k p(x_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k))\right) \geq \sum_k Q(K_n = k) \\ (\log(P_k) + \log(p(x_n|\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k))) - \log(Q(K_n = k))$$

Introducing additional  $K_n$  latent variables into the model allows using this trick – > **expectation maximisation (EM)** algorithm. We have thus to maximise for every sample the left hand side w.r.t.  $Q(K_n)$ , evaluate the sums (i.e. taking expectations) – > **E-step** and then **maximise** the resulting expression w.r.t all parameters  $\{P_k, \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k, \forall k\}$  – > **M-step**. EM-algorithm iterates E-step and M-step until convergence.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.23/24

## Bibliography

- Akaike H., A new look at the statistical identification model, in *IEEE Trans. Auto. Control* 19, 716-723, 1974.
- Bishop C. M., *Neural Networks for Pattern Recognition*, Clarendon Press Oxford, UK, 1995.
- Hartman E. J., Keeler J. D. and Kowalski J. M., Layered neural networks with Gaussian hidden units as universal approximations, in *Neural Computation* 2, 210-215, 1990.
- Press W. H., Teukolsky S. A., Vetterling W. T. and Flannery B. P., *Numerical Recipes*, Cambridge University Press, 2007.
- Rissanen J., Modeling by the shortest data description, in *Automatica* 14, 465-471, 1978.
- Schwarz G., Estimating the dimension of a model, in *The annals of Statistics* 6, 461-464, 1978.

[jump 2 TOC](#)

ML & PRN in Bioinformatics (793.404), Peter Sykaeok – p.24/24