**10th VIPS Advanced School on
Computer Vision and Pattern Recognition**

**Dissimilarity-based Representation for
Pattern Recognition**

*Robert P.W. Duin, Delft University of Technology*

*Pattern Recognition Lab
Delft University of Technology
The Netherlands*

*//rduin.nl*

23 September 2013    Representation and Generalization    1

**T**U Delft

---

## Participants

Mohamed Lamine Mekhalfi
Attaullah Buriro
Mohammad Bilal
John Kenneth Thiessen
Andrea Gasparetto
Andres Mendez
Cristina Segalin
Davide Conigliaro

Robert (Bob) Duin
Manuele Bicego
Umberto Castellani

Sami Abduljalil Abdulhak Naji
Pietro Lovato
Ricardo Henrique Gracini Guiraldelli
Anna Pesarin
Filippo Bistaffa
Valeria Garro
Denis Peruzzo

23 September 2013    Representation and Generalization    2

**T**U Delft

---

## Program

1. Representation and Generalization

2. The Dissimilarity Space

3. Pseudo-Euclidean embedding

4. Applications

23 September 2013    Representation and Generalization    3

**T**U Delft

---

## Daily Schedule

10:00 - 12:00 Lectures
        http://www.37steps.com/disrep-course/

13:30 – 16:30 Lab course
        http://www.37steps.com/disrep_exercises/

23 September 2013    Representation and Generalization    4

**T**U Delft

---

**Pattern Recognition Problems**

**T**U Delft

---

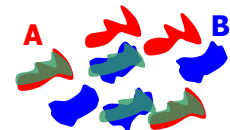## Question



How to represent real world objects,
(with a size and a shape)
given a set of examples
such that we can generalize?

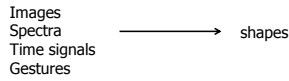23 September 2013    Representation and Generalization    6

**T**U Delft

## Real world objects and events

Images
Spectra
Time signals ───────→ shapes
Gestures

### How to build a representation?
### Features ←→ Structure

**T̃U**Delft

## Blob Recognition



BACK
BREAST
DRUMSTICK
THIGH-AND-BACK
WING

446 binary images, varying size, e.g.: 100 x 130
*Andreu, G., Crespo, A., Valiente, J.M.: Selecting the toroidal self-organizing feature maps (TSOFM) best organized to object recogn. In: ICNN. (1997) 1341–1346.*
Shape classification by weighted-edit distances (Bunke)
*Bunke, H., Buhler, U.: Applications of approximate string matching to 2D shape recognition. Pattern recognition 26 (1993) 1797–1812*
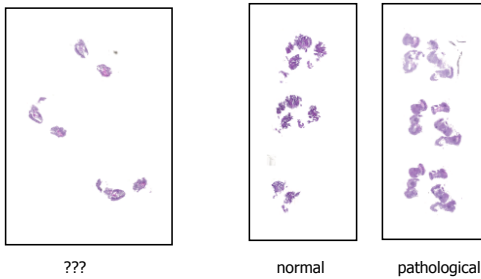
**T̃U**Delft

## Colon Tissue Recognition



???          normal          pathological

**T̃U**Delft

## Volcano / Seismic Signal Classification

Volcano-Tectonic      Long Period



150 000 events (1994 – 2008)
5 volcanos
40 stations
15 classes

J. Makario,
    INGEOMINAS, Manizales, Colombia

M. Orozco-Alzate,
    Nat. Univ. Colombia, Manizales

R. Duin, TUDelft
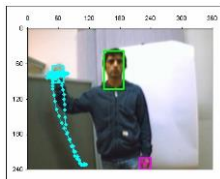
M. Bicego, Univ. of Verona, Italy

Cenatav, Havana, Cuba

**T̃U**Delft

## Gesture Recognition



Is this gesture in the database?

**T̃U**Delft

## Pattern Recognition Problems



To which class belongs an image

To which class (segment) belongs every pixel?

Where is an object of interest (detection); What is it (classification)?
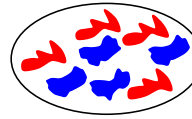
**T̃U**Delft

## Pattern Recognition: Shape Recognition

Pattern Recognition is very often Shape Recognition:
- Images: B/W, grey value, color, 2D, 3D, 4D
- Time Signals
- Spectra

**T**U Delft

## Pattern Recognition: Shapes



Examples of objects for different classes

Object of unknown class to be classified

**A  ?  B**

**T**U Delft

## Vector Representation

**T**U Delft

## Pattern Recognition System



area

perimeter

Feature Representation

**T**U Delft

## Pattern Recognition System



Pixel Representation

**T**U Delft

## Pattern Recognition System



Dissimilarity Representation

**T**U Delft

3

## Pattern Recognition System

Sensor → Representation → Generalization

Combining Classifiers

Classifier_2 / Classifier_1 with classes A and B

*T*UDelft

## Good Representations

- Class specific
  Different classes should be represented in different positions in the representation space.

- Compact
  Every class should be represented in a small set of finite domains.

*T*UDelft

## Compactness

Representations of real world similar objects are close.
There is no ground for any generalization (induction) on representations that do not obey this demand.

*(A.G. Arkedev and E.M. Braverman, Computers and Pattern Recognition, 1966.)*

$x_2$ (area) / (perimeter) $x_1$

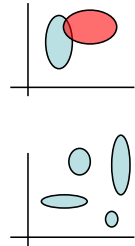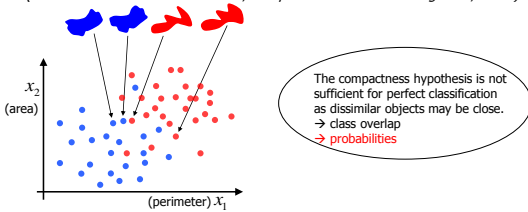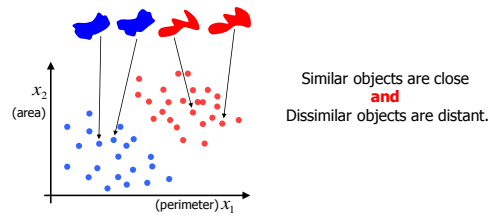The compactness hypothesis is not sufficient for perfect classification as dissimilar objects may be close.
→ class overlap
→ probabilities

*T*UDelft

## True Representations

$x_2$ (area) / (perimeter) $x_1$

Similar objects are close
**and**
Dissimilar objects are distant.

→ no probabilities needed, domains are sufficient!

*T*UDelft

## Distances and Densities

- **?** to be classified as
- **B** – because it is most close to an object B
- **A** – because the local density of A is larger.

$x_2$ (area) / (perimeter) $x_1$ with classes A and B

*T*UDelft

## Features Reduce

objects

$x_2$ / (perimeter) $x_1$ with classes A and B

Due to reduction essentially different objects are represented identically.
→ The feature space representation needs a statistical, probabilistic generalization

*T*UDelft

## Classification

## Classification error

Non-optimal classifier, e.g. based on wrong density estimates



$p(x\,|\,A)p(A)$    ε    $p(x\,|\,B)p(B)$

$S(x) > 0 \rightarrow$ Class A    $S(x) = 0$    $S(x) \le 0 \rightarrow$ Class B

$$\varepsilon = p(S(x) > 0, B) + p(S(x) < 0, A)$$
$$\varepsilon = p(S(x) > 0\,|\,B)p(B) + p(S(x) < 0\,|\,A)p(A)$$
$$\varepsilon = \int_{S(x)>0} p(x\,|\,B)p(B)dx + \int_{S(x)\le0} p(x\,|\,A)p(A)dx$$
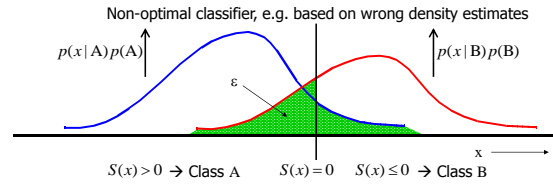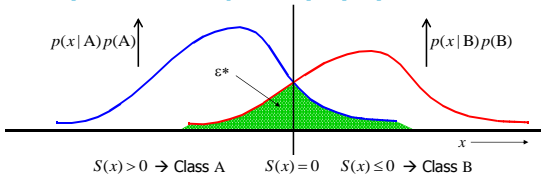
## Bayes error for optimal (Bayes) classifier



$p(x\,|\,A)p(A)$   ε*   $p(x\,|\,B)p(B)$

$S(x) > 0 \rightarrow$ Class A    $S(x) = 0$    $S(x) \le 0 \rightarrow$ Class B

Classification error is minimal, ε*, if the decision function is optimal:

$$S(x)^* = p(x\,|\,A)p(A) - p(x\,|\,B)p(B)$$
$$\varepsilon^* = \int \min\{p(x\,|\,A)p(A), p(x\,|\,B)p(B)\}dx$$
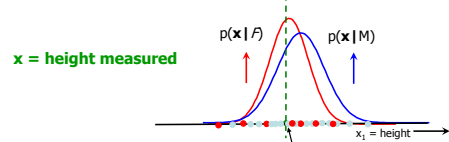
Only possible if true distributions are known

## Probabilistic Generalization

**x = height measured**



$p(\mathbf{x}\,|\,F)$   $p(\mathbf{x}\,|\,M)$   $x_1$ = height

**What is the gender of a person with this height?**

Best guess is to choose the most 'probable' class ($\rightarrow$ small error).

$\Rightarrow$ Good for overlapping classes.

$\Rightarrow$ Assumes the existence of a probabilistic class distribution and a representative set of examples.

## Bayes decision rule, formal

p(A|x)   >   p(B|x)   $\rightarrow$   A else B

Bayes:   $\dfrac{p(x|A)\,p(A)}{p(x)}$   >   $\dfrac{p(x|B)\,p(B)}{p(x)}$   $\rightarrow$   A else B

p(x|A) p(A)   >   p(x|B) p(B)   $\rightarrow$   A else B

2-class problems: S(x) = p(x|A) p(A) - p(x|B) p(B) > 0 $\rightarrow$ A else B

n-class problems: Class(x) = argmax$_\omega$(p(x|ω) p(ω))

## Density estimation

- The density is defined on the whole feature space.
- Around object $x$, the density is defined as:

$$p(x) = \frac{dP(x)}{dx} = \left(\frac{\text{fraction of objects}}{\text{volume}}\right)$$

- Given $n$ measured objects, e.g. person's height (m) how can we estimate $p(x)$?
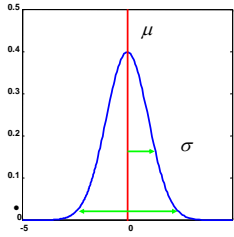
5

## The Gaussian distribution (3)



- Normal distribution = Gaussian distribution

- Standard normal distribution:
  $\mu = 0,\ \sigma^2 = 1$

- 95% of data between $[\,\mu - 2\sigma,\ \mu + 2\sigma\,]$ (in 1D!)

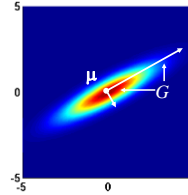$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}\right)$$

23 September 2013　　Representation and Generalization　　31

**T**UDelft

## Multivariate Gaussians



$$G = \begin{bmatrix} 3 & 1\frac{1}{2} \\ 1\frac{1}{2} & 2 \end{bmatrix}$$

- $k$ - dimensional density:

$$p(x) = \frac{1}{\sqrt{2\pi^k \det(G)}} \exp\left(-\frac{1}{2}(x-\mu)^{\mathrm{T}} G^{-1} (x-\mu)\right)$$
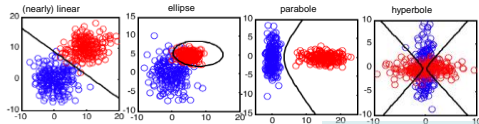
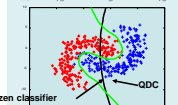23 September 2013　　Representation and Generalization　　32

**T**UDelft

## Quadratic discriminant functions

$$R(x) = -\frac{1}{2}(x - \hat{\mu}_A)^{\mathrm{T}} \hat{\Sigma}_A^{-1}(x - \hat{\mu}_A) + \frac{1}{2}(x - \hat{\mu}_B)^{\mathrm{T}} \hat{\Sigma}_B^{-1}(x - \hat{\mu}_B) + \mathrm{const}$$



QDC assumes that classes are normally distributed. Wrong decision boundaries are estimated if this does not hold.

23 September 2013　　Representation and Generalization　　33

**T**UDelft

## Linear discriminant function (summary) [G]



Normal distributions with equal covariance matrices Σ are optimally separated by a linear classifier

$$R(x) = (\mu_A - \mu_B)^{\mathrm{T}} \Sigma^{-1} x + \mathrm{const}$$

Optimal classifier for normal distributions with unequal covariance matrices $\Sigma_A$ and $\Sigma_B$ can be approximated by:

$$R(x) = (\mu_A - \mu_B)^{\mathrm{T}} (p(A)\Sigma_A + p(B)\Sigma_B)^{-1} x + \mathrm{const}$$

23 September 2013　　Representation and Generalization　　34

**T**UDelft

## Parzen density estimation (1)

Fix volume of bin, vary positions of bins, add contribution of each bin
Define 'bin'-shape (kernel):

$$K(\mathbf{r}) > 0$$

$$\int K(\mathbf{r})\, d\mathbf{r} = 1$$

For test object $z$ sum all bins

$$p(z) = \frac{1}{hn} \sum_i K\left(\frac{z - x_i}{h}\right)$$



23 September 2013　　Representation and Generalization　　35

**T**UDelft

## Parzen density estimation (2)

- With Gaussian kernel: $K(x) = \frac{1}{h\sqrt{2\pi}} \exp\left(-\frac{x^2}{2h^2}\right)$



Parzen:

$\hat{p}(x)$

$x \longrightarrow$

23 September 2013　　Representation and Generalization　　36

**T**UDelft

## Parzen: density estimates vs the smoothing parameter

1D

Small h    Optimal h    Large h

Increasing smoothing parameter h

2D

Small h    Optimal h    Large h

**T**UDelft

---

## Parzen classifier performance

`knnc`

Classification error $\varepsilon$

$\varepsilon_{1NN}$

$\varepsilon_{NMC}$

`nmc`

May be approximated by leave-one-out optimization of the error on of the training set

Smoothing parameter

`parzenc`

Small smoothing parameters: 1-NN performance

Large smoothing parameters: Nearest mean performance

**T**UDelft

---

## Nearest neighbor rule (1-NN rule)

**Assign a new object to the class of the nearest neighbor in the training set.**

1-NN rule:

• Often relies on the Euclidean distance. Other distance measures can be used.
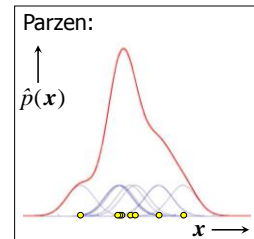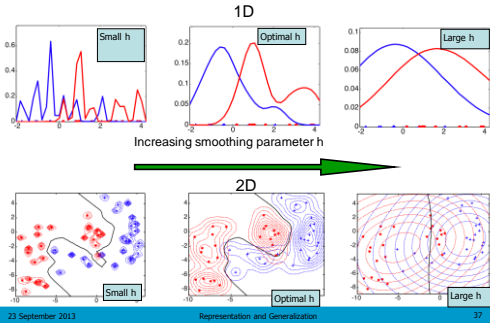
• Insensitive to prior probabilities!

• Scaling dependent. Features should be scaled properly.

There are no errors on the training set. The classifier is overtrained.

**T**UDelft

---

## Nearest neighbor examples

Simple Problem    Banana Set

Good for almost separable classes.
Useful to shape non-linear decision functions.
No training time. Long execution time.
All data should be stored.

**T**UDelft

---

## Nearest neighbor error

Asymptotically (very large training sets):

$$\varepsilon \leq 2\varepsilon^*(1-\varepsilon^*)$$

$$\varepsilon \leq 2\varepsilon^*$$

The nearest neighbor classifier will not perform worse than twice the best possible classifier

**T**UDelft

---

## K-nearest neighbor classifier

`knnc`

Assign new objects to the class of the majority of the k nearest neighbors in the training set.

More smooth.
Less local.

**T**UDelft

7

## K-nearest-neighbor performance

$e = 1 - \min_i\{p(w_i)\}$

Classification error $e$

$e_{1NN}$

May be approximated by leave-one-out optimization of the error on the training set

$k \longrightarrow$

**T**UDelft

## Prototype selection

`edicon`

$x_2$

$x_1$

**Editting**:
Removing some objects may be more accurate.

**T**UDelft

## Prototype selection (2)

`edicon`

$x_2$

$x_1$

**Condensing**:
Removing more objects may be faster.

**T**UDelft

## Support vector machine (SVM)

`svc`

1995-2005

$x_2$

Find for linear separable classes the few objects that determine the classifier: the support objects.

They have the same distance to the classifier: the margin.

Identical to "maximum margin classifier"

$S(x) = \sum_i \alpha_i (x_i^T x)$
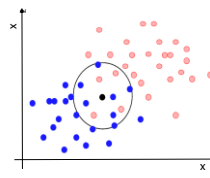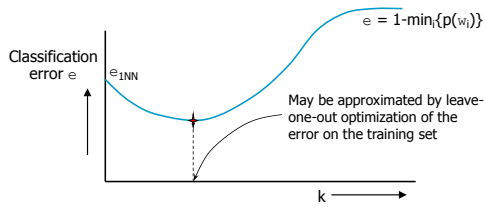
$S(x) = w^T x, \min(w^T w)$

$x_1$

**T**UDelft

## Support vector machine (2)

`svc`

$S(x) = \sum_{x_i \in S} \alpha_i (x_i^T x)$     Depends on support objects $S$ only

$S(x) = w^T x, \min(w^T w)$     Minimum norm → maximum margin

$S(x) = w^T x, \min(w^T w) + \sum_{x_j \in E} \xi(x_j)$     Allow some errors $E$

`svc`: Linear support vector classifier

**T**UDelft

## Non-linear support vector machine: The kernel trick

$S(x) = \sum_{x_i \in S} \alpha_i (x_i^T x) \longrightarrow S(x) = w^T x$    Linear classifier    `svc`

$S(x) = \sum_{x_i \in S} \alpha_i K(x_i^T x) \longrightarrow$    Non-linear classifier

$K(\bullet)$ is a non-linear function of an inner product. A linear classifier in a high-dimensional 'kernel space' is computed, resulting in a non-linear classifier in the feature space.

$K(x_i^T x) = (x_i^T x)^p$    Polynomial classifier    `svc`

$K(x_i^T x) = \Phi(\frac{x - x_i}{s})$    Radial basis SVM (about Parzen)

`rbsvc`    `parzenc`

**T**UDelft

## SVM: Examples



Gaussian Data



Banana Set

`svc`   `ldc`

`svc`   `rbsvc`

`svc` versus `ldc` for classes with very different domains

Linear (`svc`) versus nonlinear support vector classifier (`rbsvc`)

**T**U Delft

---

## Decision trees

`treec`



Implementation of a piece-wise linear classifier

Fast.

Moderate performance

**T**U Delft

---

## Perceptron

`perlc`

$y = \mathbf{w} \bullet \mathbf{x}$



Linear classifier. The weights are corrected for erroneously classified objects only
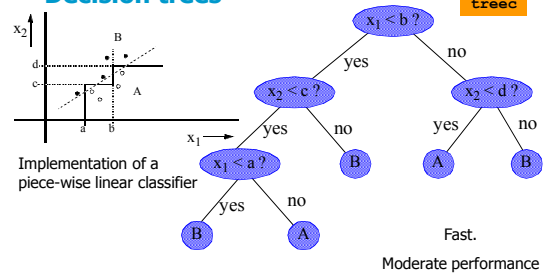
$\mathbf{w}_{n+1} = \mathbf{w}_n + \Delta\mathbf{w}(\mathbf{w}_n, \mathbf{x})$

correction $\Delta\mathbf{w} \bullet \mathbf{x}$

error    correct

Classifier outcome $\mathbf{w} \bullet \mathbf{x}$

$w_0$  $w_1$  $w_n$

1  $x_1$  ........  $x_n$

**T**U Delft

---

## Neural network classifiers

1985-1995

Number of layers with identical neurons (simple linear classifiers) with non-linear transitions in between (sigmoids). Results is a moderately non-linear classifier.

Trained object by object to minimize the MSE on the output compared with targets (labels).

Tricky training procedure.

Slow training and execution, unless special hardware is used.

Good performance, danger of overtraining.



$f(\mathbf{x}, W)$

Output unit

Weights $w_{2j}$

Hidden units (hidden layer)

Weights $w_{11j}$    Weights $w_{12j}$

Input units

$x_1$  $x_2$

`rnnc`   `neurc`   `bpxnc`   `lmnc`   `rbnc`

**T**U Delft

---

## Neural network overtraining example



Levenberg-Marquardt Optimization, 10 hidden units

2 epochs    10 epochs

20 epochs    50 epochs    100 epochs

**T**U Delft

---

## Classifier outputs

What are the possible outcomes of y = classifier(x)?
- Label, $y_1 \in \{$'apple','banana'$\}$.
- $y_2 \in \{0,1\}$ as crisp numeric labels
- $y_3 \in [0,1]$ for soft labels (confidences)
- $y_4 \in [0,\infty)$ for distances to a class
- $y_5 \in (-\infty,+\infty)$ for distances to a classifier

Conversions are often made, e.g.:

```
y2 = (y1 == 'apple')
y2 = round(y3)
y3 = sigm(y5)
y5 = invsigm(y3)
```

**T**U Delft

9

## Combining classifiers

| parallel | stacked |

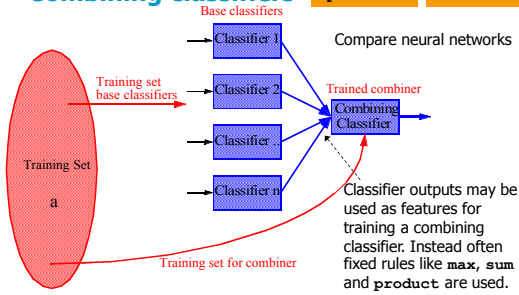Base classifiers

Classifier 1

Classifier 2

Classifier 3

Classifier n

Training set base classifiers

Training Set a

Training set for combiner

Compare neural networks

Trained combiner

Combining Classifier

Classifier outputs may be used as features for training a combining classifier. Instead often fixed rules like `max, sum` and `product` are used.

**T**U Delft

## Evaluation

**T**U Delft

## Classifier evaluation

Sensor → Representation → Generalization

**Feature Space** | **Classification**

**Test object classified as 'A'**

$x_2$

A    B

**How to estimate classifier performance?**

**Learning curves**

**Feature curves**

$x_1$

**T**U Delft

## Learning Curve

| cleval |
| testc |

True classification error ε

**Sub-optimal classifier**

Bayes error ε*

**Bayes consistent classifier**

Size training set

**T**U Delft

## The Apparent Classification Error

The apparent (or resubstitution error) of the training set is positively biased (optimistic).

Classification error

True error ε

bias

Apparent error $\varepsilon_A$ of training set

Size training set

**An independent test set is needed!**

**T**U Delft

## Error Estimation by Test Set

| gendat |
| testc |

Training

Testing → Classifier → Classification Error Estimate $\hat{\varepsilon}$

Design Set

Other training set → other classifier
Other test set → other error estimate $\hat{\varepsilon}$

$$\sigma_{\hat{\varepsilon}}^2 = \text{Var}(\hat{\varepsilon} \mid \text{test set size } N) = \frac{\varepsilon(1-\varepsilon)}{N} \qquad \sigma_{\hat{\varepsilon}} = \sqrt{\frac{\varepsilon(1-\varepsilon)}{N}}$$

| $N$ | $\varepsilon$ 0.01 | 0.03 | 0.1 |
|---|---|---|---|
| 10 | 0.031 | 0.054 | 0.095 |
| 100 | 0.010 | 0.017 | 0.003 |
| 1000 | 0.003 | 0.005 | 0.009 |

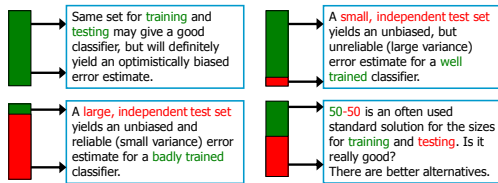**T**U Delft

## Training Set Size ←→ Test Set Size

- Training set should be large for good classifiers.
- Test set should be large for a reliable, unbiased error estimate.

  `gendat`
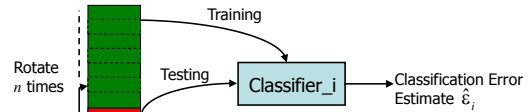- In practice often just a single design set is given



| | |
|---|---|
| Same set for training and testing may give a good classifier, but will definitely yield an optimistically biased error estimate. | A small, independent test set yields an unbiased, but unreliable (large variance) error estimate for a well trained classifier. |
| A large, independent test set yields an unbiased and reliable (small variance) error estimate for a badly trained classifier. | 50-50 is an often used standard solution for the sizes for training and testing. Is it really good? There are better alternatives. |

**T**UDelft

## Crossvalidation

`crossval`



Rotate $n$ times — Training / Testing → Classifier_i → Classification Error Estimate $\hat{\varepsilon}_i$

Size test set $1/n$ of design set.

Size training set is $(n - 1)/n$ of design set.

Train and test n test times. Average errors. (Good choice: $n = 10$)

All objects are tested ones → most reliable test result that is possible.

Final classifier: trained by all objects → best possible classifier.
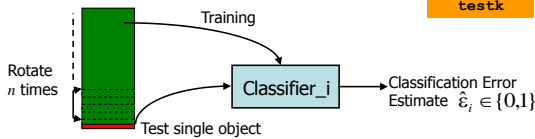
Error estimate is slightly pessimistically biased.

**T**UDelft

## Leave-one-out Procedure

`crossval`

`testk`



Rotate $n$ times — Training / Test single object → Classifier_i → Classification Error Estimate $\hat{\varepsilon}_i \in \{0,1\}$
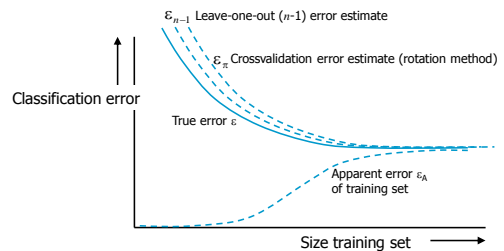
Crossvalidation in which $n$ is total number of objects.
One object tested at a time.
$n$ classifiers to be computed.
In general unfeasible for large $n$.
Doable for k-NN classifier (needs no training).

**T**UDelft

## Expected Learning Curves by Estimated Errors

`cleval`



$\varepsilon_{n-1}$ Leave-one-out ($n$-1) error estimate

$\varepsilon_\pi$ Crossvalidation error estimate (rotation method)

Classification error

True error $\varepsilon$

Apparent error $\varepsilon_A$ of training set

Size training set

**T**UDelft

## Averaged Learning Curve

`cleval`



For obtaining 'theoretically expected' curves many repetitions are needed.

```
a = gendath([200 200]);
e = cleval(a,ldc,[2,3,5,7,10,15,20],500);
plote(e);
```
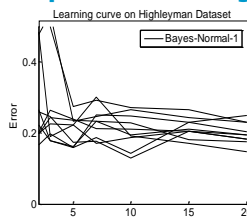
**T**UDelft

## Repeated Learning Curves
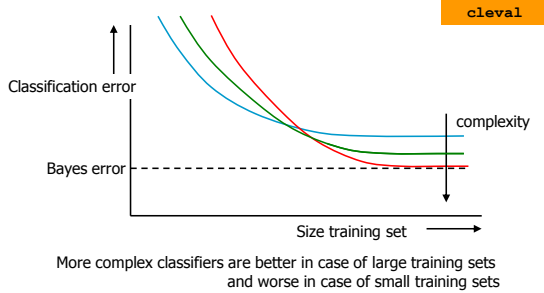
`cleval`

`plote`



Small sample sizes have a very large variability.

```
a = gendath([200 200]);
for j=1:10
  e = cleval(a,ldc,[2,3,5,7,10,15,20],1);
  hold on; plote(e);
end
```

**T**UDelft

11

## Learning Curves for Different Classifier Complexity

`cleval`



More complex classifiers are better in case of large training sets
and worse in case of small training sets

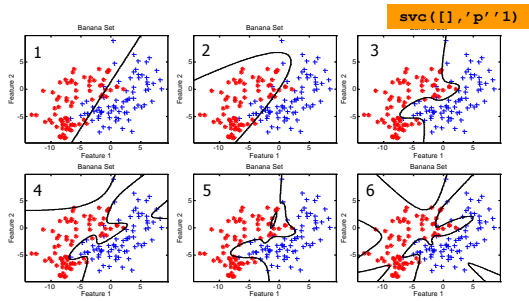23 September 2013   Representation and Generalization   67

## Peaking Phenomenon, Overtraining
## Curse of Dimensionality, Rao's Paradox

`clevalf`
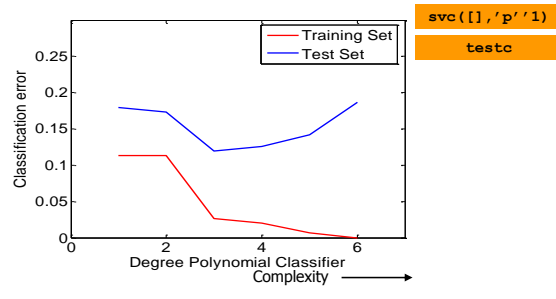


23 September 2013   Representation and Generalization   68
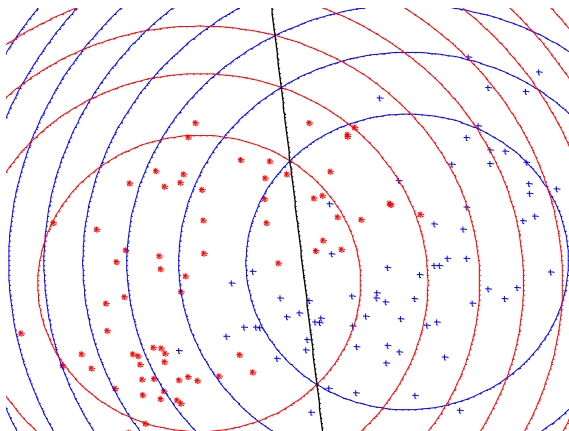
## Example Overtraining, Polynomial Classifier

`svc([],'p''1)`



23 September 2013   Representation and Generalization   69

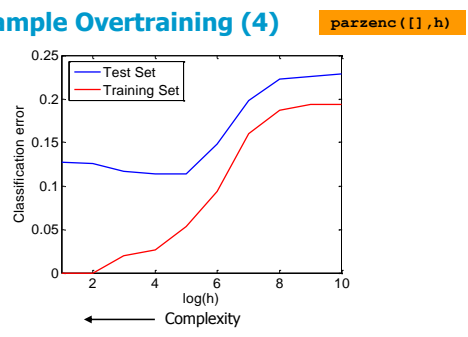## Example Overtraining (2)

`gendat`
`svc([],'p''1)`
`testc`



23 September 2013   Representation and Generalization   70



## Example Overtraining (4)

`parzenc([],h)`



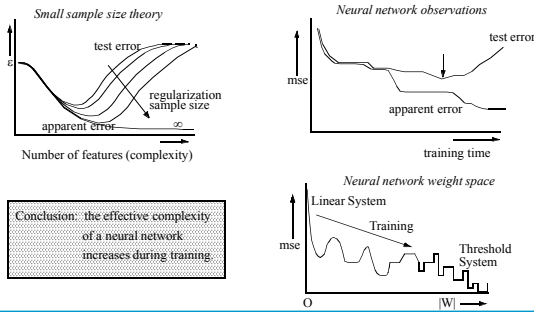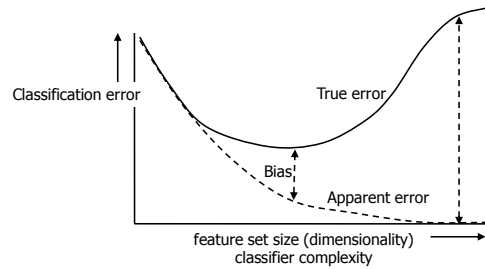23 September 2013   Representation and Generalization   72

## Neural Network Understanding

*Small sample size theory*



test error

regularization
sample size

apparent error

∞

Number of features (complexity)

*Neural network observations*

mse

test error

apparent error

training time

Conclusion: the effective complexity
of a neural network
increases during training

*Neural network weight space*

Linear System

Training

mse

Threshold
System

O                    |W|

**T**U Delft

## Overtraining ←→ Increasing Bias



Classification error

True error

Bias

Apparent error

feature set size (dimensionality)
classifier complexity

**T**U Delft

## Example Curse of Dimensionality



Feature curve for Sonar

Averaged error (50 experiments)

— original feature ranking
— feature ranking 1
— feature ranking 2

Feature size

Fisher classifier for
Various feature rankings

**T**U Delft

## Conclusions on Evaluation

- Larger training sets yield better classifiers.
- Independent test sets are needed for obtaining unbiased error estimates.
- Larger test sets yield more accurate error estimates.
- Leave-one-out crossvalidation seems to be an optimal compromise, but might be computationally infeasible.
- 10-fold cross-validation is a good practice.
- More complex classifiers need larger training sets to avoid overtraining.
- This holds in particular for larger feature sizes, due to the curse of dimensionality.
- For too small training sets, more imple classifiers or smaller feature sets are needed.

**T**U Delft