Institute of Operating Systems and Computer Networks



#### Secure communication based on noisy input data

Feature extraction and context computing

#### **Stephan Sigg**

May 3, 2011

# Overview and Structure

- Classification methods
- Feature extraction
  - Features from audio
  - Features from RF
- Fuzzy Commitment
- Fuzzy Extractors
- Authentication with noisy data
- Error correcting codes
- Entropy
- Physically unclonable functions



Conclusion

# Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



Bayesian

Non-parametric

Linear discriminant

NN

Sequential Stochastic Conclusion

Pattern recognition and classification

- Classical pattern recognition
  - Mapping of features onto classes by utilisation of prior knowledge
  - What are characteristic features?
  - Which approaches are suitable to obtain these features?



Technische

Universität

Braunschweig

Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

c Conclusion

## Pattern recognition and classification

#### • From features to context

- Measure available data on features
- Context reasoning by appropriate method
  - Syntactical (rule based e.g. RuleML)
  - Bayesian classifier
  - Non-parametric
  - Linear discriminant
  - Neural networks
  - Sequential
  - Stochastic



Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

c Conclusion

# Pattern recognition and classification

- Allocation of sensor value by defined function
  - Correlation of various data sources
  - Several methods possible simple approaches
  - Template matching
  - Minimum distance methods
  - 'Integrated' feature extraction
    - Nearest Neighbour
    - Neural Networks
- Problem
  - Measured raw data might not allow to derive all features required
  - Therefore often combination of sensors



Sufficient feature count



# Pattern recognition and classification

- Methods Syntactical (Rule based)
  - Idea: Description of Situation by formal Symbols and Rules
  - Description of a (agreed on?) world view
  - Example: RuleML
- Comment
  - Pro:
    - Combination of rules and identification of loops and impossible conditions feasible

Contra-

- Very complex with more elaborate situations
- Extension or merge of rule sets typically not possible without contradictions



Bayesian N

Non-parametric

Linear discriminant

t NN Sequential

iential Stochastic Conclusion

# Pattern recognition and classification

- Rule Markup Language: Language for publishing and sharing rules
- Hierarchy of rule-sub-languages (XML, RDF, XSLT, OWL)
- Example:
  - A meeting room was occupied by min 5 people for the last 10 minutes.

ALONE	<	١	t	0	m	>
-------	---	---	---	---	---	---

```
<Rel> occupied </Rel>
<Var> meeting room </Var>
<Ind> min 5 persons </Ind>
<Ind> last 10 minutes </Ind>
</Atom>
```



Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

Conclusion

Pattern recognition and classification

• Also conditions can be modelled

• A Meeting is taking place in a meeting room when it was occupied by min 5 people for the last 10 minutes.





Stephan Sigg  $\mid$  Secure communication based on noisy input data  $\mid$  9

Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

Conclusion

# Pattern recognition and classification

- Logical combination of conditions
  - A Meeting is taking place in a meeting room when it was occupied by min 5 people for the last 10 minutes and the light is on.



# Outline

### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



## Recognition of patterns

- Patterns can be recognised by stating a sufficient number of rules
- Sensor readings are inaccurate
- Very high count of rules to accurately model all variations of sensor readings that belong to one class
- Therefore: Consider machine learning approaches





# Recognition of patterns

- A training set  $x_1 \dots x_N$  of a large number of N samples is utilised
- Classes  $t_1 \dots t_N$  of all samples in this set known in advance
- A machine learning algorithm computes a function y(x) and generates a new target t'





# Polynomial curve fitting

- As an example we assume that a curve shall be approximated by a machine learning approach
- Sample points are created for the function  $sin(2\pi x) + N$  where N is a random noise value





Stephan Sigg | Secure communication based on noisy input data | 14 Institute of O

# Polynomial curve fitting

• We will try to fit the data points into a polynomial function of the form

$$y(x, \vec{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

x



Stephan Sigg | Secure communication based on noisy input data | 15 Ir

0

# Polynomial curve fitting

• We will try to fit the data points into a polynomial function of the form

$$y(x, \overrightarrow{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

• This can be obtained by minimising an error function that measures the misfit between  $y(x, \vec{w})$  and the training data set:

$$E(\overrightarrow{w}) = rac{1}{2} \sum_{i=1}^{N} \left[ y(x_i, \overrightarrow{w}) - t_i 
ight]^2$$

E( w)) is non-negative and zero if and only if all points are covered by the function

Technische Universität Braunschweig

# Polynomial curve fitting

- One problem is the right choice of the dimension *M* of the polynomial
- When M is too small, the accuracy of the approximation might be bad





Stephan Sigg | Secure communication based on noisy input data | 17 Inst

# Polynomial curve fitting

- One problem is the right choice of the dimension *M* of the polynomial
- In particular, when *M* becomes too big, the resulting polynomial will cross all points exactly
- When *M* reaches the count of samples in the training data set, it is always possible to create a polynomial of order *M* that contains all values in the data set exactly.





Stephan Sigg | Secure communication based on noisy input data | 18

# Polynomial curve fitting

- This event is called overfitting
- The polynomial now trained too well to the training data
- It will therefore perform badly on another sample of test data for the same phenomenon
- We can visualise this by the Root of the Mean Square (RMS) of  $E(\overrightarrow{w})$





Universität Braunschweig Stephan Sigg | Secure communication based on noisy input data | 19

### Polynomial curve fitting

• With increasing number of data points, the problem of overfitting becomes less severe for a given value of M





Stephan Sigg | Secure communication based on noisy input data | 20

# Polynomial curve fitting

- One solution to cope with overfitting is regularisation
- A penalty term is added to the error function
- This term discourages the coefficients of  $\vec{w}$  from reaching large values

$$\overline{E}(\overrightarrow{w}) = \frac{1}{2} \sum_{i=1}^{N} \left[ y(x_i, \overrightarrow{w}) - t_i \right]^2 + \frac{\lambda}{2} ||\overrightarrow{w}||^2$$

with

$$||\overrightarrow{w}||^2 = \overrightarrow{w}^T \overrightarrow{w} = w_0^2 + w_1^2 + \dots + w_M^2$$



# Polynomial curve fitting

• Depending on the value of  $\lambda$ , overfitting is controlled



$$\overline{E}(\overrightarrow{w}) = rac{1}{2}\sum_{i=1}^{N} \left[y(x_i, \overrightarrow{w}) - t_i
ight]^2 + rac{\lambda}{2}||\overrightarrow{w}||^2$$



Stephan Sigg | Secure communication based on noisy input data | 22

# Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



### Bayesian decision theory

- With probability theory, the probability of events can be estimated by repeatedly generating events and counting their occurrences
- When, however, an event only very seldom occurs or is hard to generate, other methods are required
- Example: Probability that the Arctic ice cap will have disappeared by the end of this century
- In such cases, we have to find a way to model uncertainty.
- In fact, it is possible to represent uncertainty by probability



# Conditional probability

### Conditional probability

The conditional probability of two events  $\chi_1$  and  $\chi_2$  with  $P(\chi_2) > 0$  is denoted by  $P(\chi_1|\chi_2)$  and is calculated by

$$P(\chi_1|\chi_2) = \frac{P(\chi_1 \cap \chi_2)}{P(\chi_2)}$$

 $P(\chi_1|\chi_2)$  describes the probability that event  $\chi_2$  occurs in the presence of event  $\chi_2$ .



Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion Bayesian decision theory

• With the notion of conditional probability we can express the effect of observed data  $\overrightarrow{t} = t_1, \ldots, t_N$  on a probability distribution of  $\overrightarrow{w}$ :  $P(\overrightarrow{w})$ :

$$P(\overrightarrow{w}|\overrightarrow{t}) = \frac{P(\overrightarrow{t}|\overrightarrow{w})P(\overrightarrow{w})}{P(\overrightarrow{t})}$$

- Thomas Bayes described a way to evaluate the uncertainty of  $\overrightarrow{w}$  after observing  $\overrightarrow{t}$
- $P(\overrightarrow{t}|\overrightarrow{w})$  expresses how probable a value for  $\overrightarrow{t}$  is given a fixed choice of  $\overrightarrow{w}$



### Bayesian decision theory

- A principle difference between the Bayesian viewpoint and a frequentist viewpoint is that prior assumptions are provided for a model
- Example:
  - Consider a fair-looking coin that scores heads in three consecutive coin tosses.
  - A classical maximum likelihood estimate will result in predict head for future coin tosses with probability 1
  - A Bayesian approach will include prior assumptions on the probability of events and would result in a less extreme conclusion



and Computer Networks

# Bayesian curve fitting

Technische

Braunschweig

- In the curve fitting problem, we are given  $\overrightarrow{x}$  and  $\overrightarrow{t}$  together with a new sample  $x_{M+1}$
- The task is to find a good estimation of the value  $t_{M+1}$
- This means that we want to evaluate the predictive distribution

$$p(t_{M+1}|x_{M+1}, \overrightarrow{x}, \overrightarrow{t})$$

 To account for measurement inaccuracies, typically also a probability distribution (e.g. Gauss) is underlying the sample vector <del>x</del>

# Bayesian curve fitting

• This means that we want to evaluate the predictive distribution

$$p(t_{M+1}|x_{M+1}, \overrightarrow{x}, \overrightarrow{t})$$

• After consistent application of the sum and product rules of probability we can rewrite this as

$$p(t_{M+1}|x_{M+1},\overrightarrow{x},\overrightarrow{t}) = \int p(t_{M+1}|x_{M+1},\overrightarrow{w})p(\overrightarrow{w}|\overrightarrow{x},\overrightarrow{t})d\overrightarrow{w}$$



## Bayesian curve fitting





and Computer Networks

### Example





### Example





## Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



# Histogram methods

- An alternative approach to function estimation are histogram methods
- In general, the probability density of an event is estimated by dividing the range of N values into bins of size Δ<sub>i</sub>
- Then, count the number of observations that fall inside bin  $\Delta_i$
- This is expressed as a normalised probability density





Stephan Sigg | Secure communication based on noisy input data | 34

$$p_i = \frac{n_i}{N\Delta_i}$$

# Histogram methods

- The accuracy of the estimation is dependent on the width of the bins
- The approach is well suited also for huge amounts of data since the data items can be discarded once the histogram is created





Stephan Sigg | Secure communication based on noisy input data | 35

## Histogram methods

- Issues:
  - Due to the edges of the bins, the modelled distribution is characterised by discontinuities not present in the underlying distribution observed
  - The method does not scale well with increasing dimension (Curse of dimensionality)





Stephan Sigg | Secure communication based on noisy input data | 36
#### Parzen estimator methods

- Assume an unknown probability density  $p(\cdot)$
- We want to estimate the probability density  $p(\vec{x})$  of  $\vec{x}$  in a  $\mathcal{D}$ -dimensional Euclidean space
- We consider a small region  $\mathcal{R}$  around  $\overrightarrow{x}$ :

$$P = \int_{\mathcal{R}} p(\overrightarrow{x}) d\overrightarrow{x}$$



#### Parzen estimator methods

- We utilise a data set of N observations
- Each observation has a probability of P to fall inside  $\mathcal R$
- With the binomial distribution we can calculate the count *K* of points falling into *R*:

$$\mathsf{Bin}(K|N,P) = \frac{N!}{K!(N-K)!} P^{K} (1-P)^{N-K}$$



#### Parzen estimator methods

- We utilise a data set of N observations
- Each observation has a probability of P to fall inside  $\mathcal R$
- With the binomial distribution we can calculate the count *K* of points falling into *R*:

$$\mathsf{Bin}(K|N,P) = \frac{N!}{K!(N-K)!}P^{K}(1-P)^{N-K}$$

• For large N we can show

 $K \approx NP$ 

• With sufficiently small  ${\mathcal R}$  we can also show for the volume V of  ${\mathcal R}$ 

$$P \approx p(\overrightarrow{x})V$$

Therefore, we can estimate the density as

$$p(\overrightarrow{x}) = \frac{K}{NV}$$



## Parzen estimator methods

- ullet We assume that  ${\mathcal R}$  is a small hypercube
- In order to count the number K of points that fall inside  $\mathcal{R}$  we define

$$k(\overrightarrow{u}) = \left\{ egin{array}{cc} 1, & |u_i| \leq rac{1}{2}, & i=1,\ldots,D, \ 0, & ext{otherwise} \end{array} 
ight.$$

- This represents a unit cube centred around the origin
- This function is an example of a kernel-function or Parzen window



Introduction Recognition

Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic Conclusion

#### Parzen estimator methods

$$k(\overrightarrow{u}) = \left\{ egin{array}{cc} 1, & |u_i| \leq rac{1}{2}, & i=1,\ldots,D, \\ 0, & ext{otherwise} \end{array} 
ight.$$

• When the measured data point  $\overrightarrow{x_n}$  lies inside a cube of side *h* centred around  $\overrightarrow{x}$ , we have

$$k\left(\frac{\overrightarrow{x}-\overrightarrow{x_n}}{h}\right)=1$$

• The total count K of points that fall inside this cube is

$$K = \sum_{n=1}^{N} k\left(\frac{\overrightarrow{x} - \overrightarrow{x_n}}{h}\right)$$

Stephan Sigg | Secure communication based on noisy input data | 41 Instii

Institute of Operating Systems and Computer Networks



#### Parzen estimator methods

• The total count K of points that fall inside this cube is

$$K = \sum_{n=1}^{N} k \left( \frac{\overrightarrow{x} - \overrightarrow{x_n}}{h} \right)$$

• When we substitute this in the density estimate derived above

$$p(\overrightarrow{x}) = \frac{K}{NV}$$

• with volume  $V = h^D$  we obtain the overall density estimate as

$$p(\overrightarrow{x}) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{h^{D}} \left( \frac{\overrightarrow{x} - \overrightarrow{x_{n}}}{h} \right)$$

Stephan Sigg | Secure communication based on noisy input data | 42 Institute of Operating Systems and Computer Networks



Introduction Recognition

Bayesian N

Non-parametric

Linear discriminant NN

Sequential Stochastic Conclusion

# Parzen estimator methods

$$p(\overrightarrow{x}) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{h^{D}} \left( \frac{\overrightarrow{x} - \overrightarrow{x_{n}}}{h} \right)$$

- Again, this density estimator suffers from artificial discontinuities
  Due to the fixed boundaries of the cubes
- This problem can be overcome by choosing a smoother kernel function
  - ${\, {\bullet}\,}$  A common choice is a Gaussian kernel with a standard deviation  $\sigma$

$$p(\overrightarrow{x}) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{(2\pi\sigma^2)^{\frac{D}{2}}} e^{-\frac{||\overrightarrow{x}-\overrightarrow{x_n}||^2}{2\sigma^2}}$$



#### Parzen estimator methods

 $\bullet\,$  Density estimation for various values of  $\sigma\,$ 





Stephan Sigg | Secure communication based on noisy input data | 44

Institute of Operating Systems and Computer Networks

#### Nearest neighbour methods

- A problem with Parzen estimator methods is that the parameter governing the kernel width (h or σ) is fixed for all values x
  - In regions with high data density, a high kernel width might lead to over-smoothing
  - In regions with low data density, the same width may lead to noisy estimates





Stephan Sigg | Secure communication based on noisy input data | 45 In

 Institute of Operating Systems and Computer Networks

# Nearest neighbour methods

- Nearest neighbour methods address this problem by adapting the width to the data density
  - Parzen estimator methods fix V and determine K from the data
  - Nearest neighbour methods fix K and choose V accordingly
- Again, we consider a point  $\vec{x}$  and estimate the density  $p(\vec{x})$
- The radius of the sphere is increased until K data points (the nearest neighbours) are covered



#### Nearest neighbour methods

- The value K then controls the amount of smoothing
- Again, an optimum value for K exists





# Nearest neighbour methods

- For classification we apply the KNN-density estimation for each class
- Then, we utilise the Bayes theorem
- Assume a data set of N points with  $N_k$  points in class  $C_k$
- In order to classify a sample  $\vec{X}$ , we draw a sphere containing K points around  $\overrightarrow{x}$
- The sphere can contain other points regardless of their class
- Assume that the sphere has volume V and contains  $K_k$  points from class  $C_k$



and Computer Networks

Introduction Recognition

Bayesian

Non-parametric

Linear discriminant NN

Sequential Stochastic Conclusion

# Nearest neighbour methods

- Assume: Sphere of volume V contains  $K_k$  points from class  $C_k$
- We estimate the density of class  $C_k$  as

$$p(\overrightarrow{x}|C_k) = \frac{K_k}{N_k V}$$

• The unconditional density is given as

$$p(\overrightarrow{x}) = \frac{K}{NV}$$

• The probability to experience a class  $C_k$  is given as

$$p(C_k) = \frac{N_k}{N}$$

• With the Bayes theorem we can combine this to achieve

$$p(C_k | \vec{x}) = \frac{p(\vec{x} | C_k) p(C_k)}{p(\vec{x})} = \frac{K_k}{K}$$



Stephan Sigg | Secure communication based on noisy input data | 49 Institute of Opera

Institute of Operating Systems and Computer Networks Introduction Recognition

Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic Conclusion

# Nearest neighbour methods

$$p(C_k | \overrightarrow{x}) = \frac{p(\overrightarrow{x} | C_k) p(C_k)}{p(\overrightarrow{x})} = \frac{K_k}{K}$$

- To minimise the probability of misclassification, assign  $\overrightarrow{x}$  to class with the largest probability
- This corresponds to the largest value of

$$\frac{K_k}{K}$$



#### Nearest neighbour methods

- To classify a point, we identify the K nearest points
- And assign the point to the class having most representatives in this set
- The choice K = 1 is called the nearest neighbour rule
  - For this choice, the error rate is never more than twice the minimum achievable error rate of an optimum classifier<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>T. Cover and P. Hart: Nearest neighbour pattern classification. IEEE Transactions on Information Theory, IT-11,



Technische

Braunschweig

## Nearest neighbour methods

• Classification of points by the K-nearest neighbour classifier





Stephan Sigg | Secure communication based on noisy input data | 52 Inst

Institute of Operating Systems and Computer Networks

#### Nearest neighbour methods

• Classification of points by the K-nearest neighbour classifier





Stephan Sigg | Secure communication based on noisy input data | 53 In

Institute of Operating Systems and Computer Networks

#### Nearest neighbour methods

• The KNN-method and the Parzen-method are not well suited for large data sets since they require the entire data set to be stored



# Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



Recognition Bavesian Non-parametric Linear discriminant Sequential Stochastic Conclusion

# Support vector machines

- In classification problems we want to take an input  $\overrightarrow{x}$  and assign it to one of K discrete classes  $C_k$
- The input is divided by decision boundaries
- Here we assume that the decision boundaries are linear functions of  $\overrightarrow{\mathbf{x}}$
- Data sets that can be separated exactly by linear decision surfaces are linear separable
- With sufficiently high dimension, a data set of two classes is always linearly separable



and Computer Networks





- A support vector machine pre-processes data to represent patterns in a high dimension
- Dimension often much higher than original feature space
- Then, insert hyperplane in order to separate the data



# Support vector machines

- A pattern  $\overrightarrow{x_k}$  is transformed to  $\overrightarrow{y_k} = \varphi(\overrightarrow{x_k})$
- Also, each  $\overrightarrow{x_k}$  is associated with  $z_k \in \{-1,1\}$
- A linear discriminant in an augmented  $\overrightarrow{y}$  space is  $g(\overrightarrow{y}) = \overrightarrow{a}^t \overrightarrow{y}$
- A separating hyperplane ensures for  $y_0 = 1, a_0 \ge 1$

 $z_k g(y_k) \geq 1$ 



# Support vector machines

• The goal for support vector machines is to find a separating hyperplane with the largest margin *b* to the outer points in all sets

$$\frac{z_k g(y_k)}{||\overrightarrow{a}||} \ge b, \ k = 1, \dots, n$$

- If no such hyperplane exists, map all points into a higher dimensional space until such a plane exists
- Support vectors satisfy ' $\cdot = b'$





Stephan Sigg | Secure communication based on noisy input data | 60

Institute of Operating Systems and Computer Networks

- Simple application to several classes by iterative approach:
  - belongs to class 1 or not?
  - belongs to class 2 or not?
  - ...
- Search for optimum mapping between input space and feature space complicated (no optimum approach known)



- Simple learning approach to find the correct hyperplane:
  - Starting from an initial separating hyperplane
  - Find worst classified pattern (on the wrong side of the hyperplane)
  - Design a new hyperplane with this pattern as one of the support vectors
  - Iterate until all patterns are correctly classified









# Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



# Neural networks

#### Neural networks

- Learn the correct mapping from an input vector to an output vector
- Representation of the mapping function by an edge-weighted graph
- Distinction between
  - Input neurons
  - Output neurons
  - Hidden nodes



#### Neural networks



Bayesian Non-parametric

Linear discriminant

NN

Sequential Stochastic

#### c Conclusion

# Neural networks

- Neural networks
  - Input neurons are only equipped with outgoing edges
  - Hidden nodes 'fire' (assign the value 1 to their outgoing edges) with respect to the configuration of the weighted ingoing edges and a threshold function  $\Theta$



$$y_i = \begin{cases} 1, & \text{if } \sum_{i=1}^n x_i w_i \ge \Theta \\ 0, & \text{else.} \end{cases}$$

# Neural networks

Learning with back-propagation (schematic):

- Iterate until the error is sufficiently small
  - Choose a training-pair and copy it to the input layer
  - 2 Propagate it through the network
  - Calculate error between computed and expected output
  - Propagate the sum product of the weights back into the network in order to calculate the error in internal layers
  - 6 Adapt weights to the error



# Neural networks

Discussion

- With only one layer of hidden nodes it is possible to represent arbitrary multi-dimensional functions
- Is well suited to work on noisy input data
- Implicit clustering of input data possible ۰
- Issues:
  - Complicated to extend existing network (e.g. when new features are added)





- A Neural network consists of a series of functional transformations
- For the input layer, we construct *M* linear combinations of the input variables *x*<sub>1</sub>,..., *x*<sub>D</sub> and weights *w*<sub>1</sub>,..., *w*<sub>D</sub>

$$a_j = \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

• Each *a<sub>j</sub>* is transformed using a differentiable, non-linear activation function

Technische

Braunschweig

$$z_j = h(a_j)$$



#### Neural networks

• Input layer: *M* linear combinations of  $x_1, \ldots, x_D$  and  $w_1, \ldots, w_D$ 

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

• Differentiable, non-linear activation function:

$$z_j=h(a_j)$$

• The  $h(\cdot)$  function is usually a sigmoidal function or tanh





Stephan Sigg | Secure communication based on noisy input data | 71 Institute of Operating Systems and Computer Networks

#### Neural networks

• The values  $z_j$  are then again linearly combined in the hidden layers:

$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)}$$

- with  $k = 1, \ldots, K$  describing the total number of outputs
- Again, these values are transformed using a sufficient transformation function σ to obtain the network outputs y<sub>k</sub>

$$y_k = \sigma(a_k)$$

• For multi-class problems, we use a function such as

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

Stephan Sigg | Secure communication based on noisy input data | 72 Institute of Operating Systems and Computer Networks




• We can combine these various stages to give the overall network function in the form

$$y_k(\overrightarrow{x},\overrightarrow{w}) = \sigma \left( \sum_{j=1}^M w_{kj}^{(2)} h\left( \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

• When the network contains multiple hidden layers, these are added analogously



## Neural networks

- When the activation functions of all hidden units in a network are linear, we can always find an equivalent network without hidden units
  - Since the composition of successive linear transformations is itself a linear transformation





Stephan Sigg | Secure communication based on noisy input data | 74

Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion

#### Neural networks

- When the number of hidden units is smaller than the number of input or output units, not all linear functions can be generated
  - Since Information is lost in the dimensionality reduction at the hidden units





Stephan Sigg | Secure communication based on noisy input data | 75

Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion

#### Neural networks

- Sometimes, skip-layer connections are utilised
  - In principle, a network with sigmoidal hidden units can always mimic skip layer connections
  - By using sufficiently small first-layer weight
  - and then compensating with a large weight from the hidden unit to the output



#### Neural networks

- Neural networks are Universal approximators<sup>2 3 4 5 6 7 8 9</sup>
  - A two-layer network with only linear outputs can uniformly approximate any continuous function on a compact input domain to arbitrary accuracy

<sup>4</sup>K. Hornik, M. Sinchcombe, H. White: Multilayer feed-forward networks are universal approximators. Neural Networks, 2(5), 359-366, 1989

<sup>5</sup>N.E. Cotter: The stone-Weierstrass theorem and its application to neural networks. IEEE Transactions on Neural Networks 1(4), 290-295, 1990

 $^{6}$ Y. Ito: Representation of functions by superpositions of a step or sigmoid function and their applications to neural network theory. Neural Networks 4(3), 385-394, 1991

<sup>7</sup>K. Hornik: Approximation capabilities of multilayer feed forward networks: Neural Networks, 4(2), 251-257, 1991

<sup>8</sup>Y.V. Kreinovich: Arbitrary non-linearity is sufficient to represent all functions by neural networks: a theorem. Neural Networks 4(3), 381-383, 1991

Pattern Recognition and Neural Networks. Cambridge University Press, 1996



Technische Universität Braunschweig

Stephan Sigg | Secure communication based on noisy input data | 77 Institute of Operating Systems and Computer Networks

 $<sup>^2{\</sup>rm K}.$  Funahashi: On the approximate realisation of continuous mappings by neural networks, Neural Networks, 2(3), 183-192, 1989

 $<sup>^3</sup>$  G. Cybenko: Approximation by superpositions of a sigmoidal function. Mathematics of control, signals and systems, 2, 304-314, 1989



- Remaining issue in neural networks
  - How to find suitable parameter values given a set of training data
  - Several learning approaches have been proposed





Stephan Sigg | Secure communication based on noisy input data | 78



- A simple approach to the problem of determining the network parameters is to minimise a sum-of-squared error function
  - Given a training set  $\overrightarrow{x_n}$  with  $n \in \{1, \ldots, N\}$
  - And a corresponding set of target vectors  $\vec{t_n}$
  - We minimise the error function

$$E(\overrightarrow{w}) = \frac{1}{2} \sum_{n=1}^{N} (y(\overrightarrow{x_n}, \overrightarrow{w}) - \overrightarrow{t_n})^2$$



# Neural networks

Classification with neural networks

- Classification of two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$ 
  - We consider a network with a single output

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

- The output is interpreted as the conditional probability  $p(\mathcal{C}_1 | \overrightarrow{x})$
- Analogously, we have  $p(\mathcal{C}_2|\overrightarrow{x}) = 1 p(\mathcal{C}_1|\overrightarrow{x})$



Introduction Recognition Bayesian Non-parametric Linear discriminant **NN** Sequential Stochastic Conclusion

#### Neural networks

Classification with neural networks

- Classification of K classes  $C_1, \cdots, C_K$ 
  - Binary target variables  $t_k \in \{0, 1\}$
  - Network outputs are interpreted as  $y_k(\overrightarrow{x}, \overrightarrow{w}) = p(t_k = 1 | \overrightarrow{x})$



NN

# Introduction to self organising maps

- Self organising map (SOM) algorithm proposed by Teuvo Kohonen<sup>10</sup>
- Presented it as a model of the self-organisation of neural connections.
- Map high dimensional input data to low dimensional representation
- Based on neural network learning of the underlying mapping



Technische Universität Braunschweig en, Self-Organizing Maps, Springer, 2001.

NN

#### tic Conclusion

# Introduction to self organising maps

- Idea:
  - Present all points in a source space by points in a target space
  - Given a sequence of points in a sample space,
  - Create a mapping of these points into a target space that respects the neighbourhood relation in the sample space





# Introduction to self organising maps



- SOM is a topology preserving lattice of a predefined number of nodes that represents a topology of elements in the input space.
- Algorithm inherits self-organisation property
  - Able to produce organisation starting from possibly total disorder.
  - SOM algorithm defines and preserves neighbourhood structure between all nodes of the map.
- Learning by two layer neural network



Bayesian I

Non-parametric

Linear discriminant

NN Sequential St

Stochastic Conclusion

#### Introduction to self organising maps





Stephan Sigg | Secure communication based on noisy input data | 85

Linear discriminant

Sequential Stochastic

NN

c Conclusion

# Introduction to self organising maps



- When a pattern  $\overrightarrow{\phi_i}$  is presented, each node (represented by outer neurons) in the target space computes its activation  $\overrightarrow{\phi_i^t} \overrightarrow{w}$ .
- The most activated node  $y^*$  and the weights to its neighbours are updated according to a learning rate  $\rho(t)$

$$w_{ki}(t+1) = w_{ki}(t) + 
ho(t)\Lambda(|y-y^*|)(\overrightarrow{\phi_i} - w_{ki}(t))$$

 A(·) defines a non-increasing neighbourhood function and |y − y\*| describes the distance of nodes in the neighbourhood



Bayesian

Non-parametric

Linear discriminant

t NN Sea

Sequential Stochastic Conclusion

# SOM – Self organisation



# SOM – Definition

 $\bullet$  We recapitulate a condensed definition of the SOM algorithm that can be found  $\mbox{in}^{11}$ 

#### Self organising maps

- Let I = { n₁,..., n₁s₁} be a set of km-dimensional vectors that are associated with nodes in a lattice.
- Neighbourhood structure provided by symmetrical, non-increasing neighbourhood function  $d: I \times I \to \mathbb{R}$  which depends on the distance between two nodes  $\overline{\eta_i}$  and  $\overline{\eta_j} \in I$ .
- The state of the map at time t is given by

$$\eta(t) = \left(\overrightarrow{\eta_1(t)}, \overrightarrow{\eta_2(t)}, \dots, \overrightarrow{\eta_{|S|}(t)}\right),$$

Institute of Operating Systems

and Computer Networks

<sup>11</sup>M. Cottrell J.C. Fort and G. Pages, *Theoretical aspects of the SOM algorithm*, Neurocomputing, pp. 119-138, vol

Stephan Sigg | Secure communication based on noisy input data | 88

Technische Universität Braunschweig Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion

# SOM – Definition

Technische Universität

Braunschweig

#### Self organising map algorithm

The SOM algorithm is recursively defined by

$$i_{c}\left(\overrightarrow{v(t+1)}, \overrightarrow{\eta(t)}\right) = argmin\left\{\left\|\overrightarrow{v(t+1)} - \overrightarrow{\eta_{i}(t)}\right\|, \overrightarrow{\eta_{i}(t)} \in \eta(t)\right\}, \\ \overrightarrow{\eta_{i}(t+1)} = \overrightarrow{\eta_{i}(t)} - \varepsilon_{t}d\left[i_{c}\left(\overrightarrow{v(t+1)}, \overrightarrow{\eta(t)}\right), \overrightarrow{\eta_{i}}\right] \\ \cdot \left(\overrightarrow{\eta_{i}(t)} - \overrightarrow{v(t+1)}\right), \forall \overrightarrow{\eta_{i}} \in I.$$

- In this formula,  $i_c(\overrightarrow{v(t+1)}, \overrightarrow{\eta(t)})$  corresponds to the node in the network that is closest to the input vector.
- Parameter  $\varepsilon_t$  controls the adaptability of the self organising map.

Bayesian

Non-parametric

Linear discriminant

NN Sequential Stochastic

c Conclusion

and Computer Networks

# SOM – Operational principle



• Input values  $v_i(t)$  are to be mapped onto the target space



Bayesian

Non-parametric

Linear discriminant

NN Sequential

ntial Stochastic Conclusion

# SOM – Operational principle



• The node with the lowest distance is associated with the input value:

$$i_c\left(\overrightarrow{v(t+1)},\overrightarrow{\eta(t)}
ight) = \operatorname{argmin}\left\{\left\|\overrightarrow{v(t+1)} - \overrightarrow{\eta_i(t)}\right\|, \overrightarrow{\eta_i(t)} \in \eta(t)
ight\}$$



Bayesian

Non-parametric

Linear discriminant

NN Sequential Stochastic

c Conclusion

# SOM – Operational principle



 Nodes in the neighbourhood of the associated node are moved closer to the input value



Introduction

Recognition Bayesian

Non-parametric

Linear discriminant

NN Sequential Stochastic

c Conclusion

# SOM – Operational principle



 Nodes in the neighbourhood of the associated node are moved to the input value

$$\overrightarrow{\eta_i(t+1)} = \overrightarrow{\eta_i(t)} - \varepsilon_t d \left[ i_c \left( \overrightarrow{v(t+1)}, \overrightarrow{\eta(t)} \right), \overrightarrow{\eta_i} \right] \\ \cdot \left( \overrightarrow{\eta_i(t)} - \overrightarrow{v(t+1)} \right), \forall \overrightarrow{\eta_i} \in I.$$



Introduction

Bayesian

Recognition

Non-parametric

Linear discriminant

NN Sequential S

Stochastic Conclusion

# SOM – Example application: TEA





Stephan Sigg | Secure communication based on noisy input data | 94

Introduction

Recognition Bayesian Non-parametric

Linear discriminant

Sequential Stochastic

NN

Conclusion

# SOM – Example application: TEA





Stephan Sigg | Secure communication based on noisy input data | 95



# SOM – Remarks

- It has been proven that the SOM algorithm always converges<sup>12</sup>
- Normalisation of input vectors might improve numerical accuracy
- Not guaranteed that self-optimisation will always occur (Dependent on choice of parameters)
- The SOM is not optimising any well-defined function, so that it is difficult to set the parameters of the model <sup>13</sup>

<sup>&</sup>lt;sup>12</sup>Y. Cheng, Neural Computation, 9(8), 1997.

<sup>&</sup>lt;sup>13</sup>E. Erwin, K. Obermayer, K. Schulten: Self-organising maps: Ordering, convergence properties and energy functions. In the Columnatics, 67, 47-55, 1992

Introduction Recognition Linear discriminant Sequential Stochastic Bayesian Non-parametric NN Conclusion

## Problems of SOMs

 If the neighbourhood is chosen to be too small, the map will not be ordered globally.



# Problems of SOMs



- The map created as target space might have several orientations
- It is possible that one part of the map is created following one orientation, while other parts are created following other orientations.



Bayesian No

Non-parametric

Linear discriminant

NN

Sequential Stochastic Conclusion

# Software tools for SOMs – The SOM\_PAK

SOM\_PAK First public-domain software package (1990)<sup>14</sup>

- Released by Laboratory of Computer and Information Science of Helsinki University of Technology
- Source code (ANSI C) and documentation completely available
- Available for UNIX or MS DOS
- Features:
  - Standard incremental-learning SOM
  - (simple) graphics programs included
  - Map size and vector dimension not restricted
  - Several neighbourhood functions available
  - Able to handle largest scale problems

is.hut.fi/research/software.shtml



14

Stephan Sigg | Secure communication based on noisy input data | 99 Institute of Operating Systems

and Computer Networks

Recognition

Bavesian

Non-parametric

Linear discriminant

NN

Sequential Stochastic

Conclusion

#### Software tools for SOMs – The SOM Toolbox

SOM Toolbox Toolbox for MatLab (1996)<sup>15</sup>

- Released by Laboratory of Computer and Information Science of Helsinki University of Technology
- MatLab version 5 or higher required
- Slower than SOM PAK
- Features:
  - Standard incremental-learning SOM and Batch Map SOM
  - Map size and vector dimension not restricted
  - Neighbourhood and training sequences identical to SOM PAK
  - Improved visualisation and analysis capabilities

is.hut.fi/software.shtml

15

Bayesian No

Non-parametric I

Linear discriminant

NN

Sequential Stochastic Conclusion

#### Software tools for SOMs – The Neural Networks Tool

Nenet User friendly ANN Toolbox (1997)<sup>16</sup>

- Released by the Neural network team of Helsinki University of Technology
- 32-bit Windows 95/NT recommended
- Suited for small scale problems only
- Features:
  - Standard incremental SOM
  - Several visualisation options:
    - Component planes with trajectories
    - U-matrix
    - 3D hit histograms
    - Display of active neuron coordinates
  - Easy to use; Good graphics programs

nbnet.fi/~phodju/nenet/Nenet/Gneral.html



16

Stephan Sigg | Secure communication based on noisy input data | 101 Institute of Operating Systems and Computer Networks

NN

## Software tools for SOMs – Viscovery SOMine

Viscovery SOMine Commercial SOM software package <sup>17</sup>

- Released by Eudaptics GmbH in Austria
- Windows NT and Windows NT 4.0
- Features:
  - User-friendly, flexible and powerful package
  - Interfaces for GUI, OLE, SQL and DB2
  - Batch Map algorithm
  - High computing speed
  - Unlimited map size and vector dimension
  - Several visualisation options:
    - Component planes with trajectories
    - U-matrix
    - Cluster windows
    - Iso-contours of hit density

udaptics.co.at



17

Technische Braunschweig

# Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion

#### Markov chains

- Markov processes
  - Intensively studied
  - Major branch in the theory of stochastic processes
- A. A. Markov (1856 1922)
- Extended by A. Kolmogorov by chains of infinitely many states
  - 'Anfangsgründe der Theorie der Markoffschen Ketten mit unendlich vielen möglichen Zuständen' (1936) <sup>18</sup>

18



v,Anfangsgründe der Theorie der Markoffschen Ketten mit unendlich vielen möglichen Zuständen, 1936.

#### Markov chains

#### Markov chains

- Theory of Markov chains applied to a variety of algorithmic problems
- Standard tool in many probabilistic applications
- Intuitive graphical representation
  - Possible to illustrate properties of stochastic processes graphically
- Popular for their simplicity and easy applicability to huge set of problems<sup>19</sup>

19.



Technische Universität Braunschweig

<sup>,</sup> An introduction to probability theory and its applications, Wiley, 1968.

#### Markov chains

#### Independent trials of events

- Set of possible outcomes of a measurement *E<sub>i</sub>* associated with occurrence probability *p<sub>i</sub>*
- Probability to observe sample sequence:
  - $P\{(E_1, E_2, ..., E_i)\} = p_1 p_2 \cdots p_i$
- Dependent trials of events
  - Probability to observe specific sequence  $E_1, E_2, \ldots, E_i$  obtained by conditional probability:

$$P(E_i|E_1,E_2,\ldots,E_{i-1})$$

• In general:

$$P(E_i|E_1, E_2, \ldots, E_{i-1}) \neq P(E_i|E_2, E_1, E_3, E_4, \ldots, E_{i-1})$$

Non-parametric Linear discriminant Recognition Bavesian NN Sequential Stochastic Conclusion

#### Markov chains

- Independent random variables
  - Number of coin tosses until 'head' is observed
  - Radioactive atoms always have the same probability of decaying at the next trial
- Dependent random variables
  - The knowledge that no car has passed for five minutes increases our expectation that it will come soon.
  - Coin tossing:
    - Probability that the cumulative numbers of heads and tails will equalize at the second trial is  $\frac{1}{2}$
    - Given that they did not, the probability that they equalize after two additional trials is only  $\frac{1}{4}$



## Markov chains

#### Markov property

In the theory of stochastic processes the described lack of memory is connected with the Markov property.

- Theory of Markov chains:
  - Outcome of any trial depends exclusively on the outcome of the directly preceding trial
  - Outcome of  $E_k$  is no longer associated with fixed probability  $p_k$ 
    - Instead: With every pair  $(E_i, E_i)$  a conditional probability  $p_{ii}$
    - Probability that  $E_i$  is observed after  $E_i$
    - Additionally: Probability  $a_i$  of the event  $E_i$


### Markov chain

A sequence of observations  $E_1, E_2, \ldots$  is called a Markov chain if the probabilities of sample sequences are defined by

$$P(E_1, E_2, \ldots, E_i) = a_1 \cdot p_{12} \cdot p_{23} \cdot \cdots \cdot p_{(i-1)i}.$$

and fixed conditional probabilities  $p_{ij}$  that the event  $E_i$  is observed directly in advance of  $E_j$ .



### Markov chains

- Theory of Markov chains:
  - $P\{(E_i, E_j)\} = a_i p_{ij}$ •  $P\{(E_i, E_j, E_k)\} = a_i p_{ij} p_{jk}$
  - $P\{(E_i, E_j, E_k, E_l)\} = a_i p_{ij} p_{jk} p_{kl}$
  - $P\{(E_i, E_j, ..., E_m, E_n)\} = a_i p_{ij} ... p_{mn}$



### Markov chains

- Markov chain sometimes modelled as directed graph G = (V, E)
- Labelled edges in E
- states (or vertices) in V.
- Transition probabilities  $p_{ij}$  between  $E_i, E_j \in V$





and Computer Networks



• Markov chain described by probability *a* for initial distribution and matrix *P* of transition probabilities.

$$P = \begin{bmatrix} p_{11} & p_{12} & p_{13} & \cdots \\ p_{21} & p_{22} & p_{23} & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

- *P* is a square matrix with non-negative entries that sum to 1 in each row.
- *P* is called a stochastic matrix



- $p_{ij}^k$  denotes probability that  $E_j$  is observed exactly k observations after  $E_i$  was observed.
- Calculated as the sum of the probabilities for all possible paths  $E_i E_{i_1} \cdots E_{i_{k-1}} E_j$  of length k
- We already know

$$p_{ij}^1 = p_{ij}$$

Consequently:

$$P_{ij}^{2} = \sum_{\nu} p_{i\nu} \cdot p_{\nu j}$$
$$P_{ij}^{3} = \sum_{\nu} p_{i\nu} \cdot p_{\nu j}^{2}$$



Stephan Sigg | Secure communication based on noisy input data | 113 Institute of Operating Systems and Computer Networks

### Markov chains

#### • By mathematical induction:

$$p_{ij}^{n+1} = \sum_{
u} p_{i
u} \cdot p_{
u j}^n$$

$$p_{ij}^{n+m} = \sum_{\nu} p_{i\nu}^m \cdot p_{\nu j}^n = \sum_{\nu} p_{i\nu}^n \cdot p_{\nu j}^m$$



- Similar to matrix P we can create a matrix  $P^n$  that contains all  $p_{ii}^n$
- We obtain  $P_{ij}^{n+1}$  from  $P^{n+1}$  by multiplying all elements of the *i*-th row of P with the corresponding elements of the *j*-th column of  $P^n$  and add all products.

• Symbolically: 
$$P^{n+m} = P^n P^m$$



### Closed set of states

- A set *C* of states is closed if no state outside *C* can be reached from any state *E<sub>i</sub>* in *C*.
- For an arbitrary set C of states the smallest closest set containing C is called the closure of C
- A single state  $E_k$  forming a closed set is called absorbing



### Markov chains

#### Closed sets in stochastic matrices

If in a matrix  $P^n$  all rows and all columns corresponding to states outside a closed set C are deleted, the remaining matrices are again stochastic matrices.



#### Irreducible Markov chain

A Markov chain is irreducible if there exists no closed set other than the set of all states.

### Criterion for irreducible chains

A chain is irreducible if, and only if, every state can be reached from every other state.



### Periodicity of states

- The state  $E_j$  has period t > 1 if  $p_{jj}^n = 0$  unless n = vt is a multiple of t and t is the largest integer with this property.
- The state  $E_j$  is aperiodic if no such t > 1 exists
- A state  $E_j$  to which no return is possible is considered aperiodic





- To deal with a periodic *E<sub>j</sub>* it suffices to consider the chain at the trials *t*, 2*t*, 3*t*
- In this way we obtain a new Markov chain with transition probabilities p<sup>t</sup><sub>ik</sub>
- In this new chain  $E_j$  is aperiodic
- Results concerning aperiodic states can thus be transferred to periodic states



### Persistent and transient states

- The state  $E_j$  is persistent if  $\sum_{n=1}^\infty p_{jj}^n=1$  and transient if  $\sum_{n=1}^\infty p_{jj}^n<1$
- A persistent state  ${\it E}_j$  is called null state if its mean recurrence time  $\mu_j = \infty$

### Ergodic states

• An aperiodic persistent state  $E_j$  with  $\mu_j < \infty$  is called ergodic



# Hidden Markov Models

- Make a sequence of decisions for a process that is not directly observable<sup>20</sup>
- Current states of the process might be impacted by prior states
- HMM often utilised in speech recognition or gesture recognition



and Computer Networks



# Hidden Markov Models

### Computational biology

- Align biological sequences
- Find sequences homologous to a known evolutionary family
- Analyse RNA secondary structure <sup>21</sup>
- Computational linguistics<sup>22</sup>
  - Topic segmentation of text
  - Information extraction

Technische

Universität Braunschweig

<sup>&</sup>lt;sup>21</sup> R. Durbin, S. Eddy, A. Krogh and G. Mitchison, *Biological sequence analysis: Probabilistic models of proteins and nucleic acids*, Cambridge University Press, 1998.

and H. Schütze, Foundations of statistical natural language processing, MIT Press, 1999.

Introduction Recognition

Bayesian N

Non-parametric Li

Linear discriminant

Sequential Stochastic

ic Conclusion

### Hidden Markov Models





Stephan Sigg | Secure communication based on noisy input data | 124 Institute of Operating Systems and Computer Networks

Introduction Recognition

Bayesian

Non-parametric

Linear discriminant

NN Sequential Stochastic

c Conclusion

# Hidden Markov Models



- At every time step t the system is in an internal state  $\omega(t)$
- Additionally, we assume that it emits a (visible) symbol v(t)
- Only access to visible symbols and not to internal states

Technische Universität Braunschweig

# Hidden Markov Models

- $V^T = \{v(1), v(2), \dots, v(T)\}$
- In any state ω(t) we have a probability of emitting a particular visible symbol v<sub>k</sub>(t)
- Probability to be in state  $\omega_j(t)$  and emit symbol  $v_k(t)$ :
  - $P(v_k(t)|\omega_j(t)) = b_{jk}$
- Transmission probabilities:  $p_{ij} = P(\omega_j(t+1)|\omega_i(t))$
- Emission probability:  $b_{jk} = P(v_k(t)|\omega_j(t))$



## Hidden Markov Models

#### Normalisation conditions

$$\sum_{j} p_{ij} = 1$$
for all  $i$   
 $\sum_{k} b_{jk} = 1$ for all  $j$ 



# Hidden Markov Models

- Central issues in hidden Markov models:
  - Evaluation problem
    - Determine the probability that a particular sequence of visible symbols  $V^{\mathcal{T}}$  was generated by a given hidden Markov model
  - Decoding problem
    - Determine the most likely sequence of hidden states  $\omega^T$  that led to a specific sequence of observations  $V^T$
  - Learning problem
    - Given a set of training observations of visible symbols, determine the parameters  $p_{ij}$  and  $b_{jk}$  for a given HMM



# Hidden Markov Models – Evaluation problem

• Probability that model produces a sequence  $V^T$ :

$$P(V^{T}) = \sum_{\overline{\omega}^{T}} P(V^{T} | \overline{\omega}^{T}) P(\overline{\omega}^{T})$$

Also:

$$egin{array}{rcl} P(\overline{\omega}^T) &=& \prod_{t=1}^T P(\omega(t)|\omega(t-1)) \ P(V^T|\overline{\omega}^T) &=& \prod_{t=1}^T P(v(t)|\omega(t)) \end{array}$$

Together:

$$P(V^T) = \sum_{\overline{\omega}^T} \prod_{t=1}^T P(v(t)|\omega(t))P(\omega(t)|\omega(t-1))$$



# Hidden Markov Models – Evaluation problem

• Probability that model produces a sequence  $V^T$ :

$$P(V^T) = \sum_{\overline{\omega}^T} \prod_{t=1}^T P(v(t)|\omega(t))P(\omega(t)|\omega(t-1))$$

- Formally complex but straightforward
- Naive computational complexity
  - $\mathcal{O}(c^T T)$



Technische

NN

# Hidden Markov Models – Evaluation problem

• Probability that model produces a sequence  $V^T$ :

$$P(V^{T}) = \sum_{\overline{\omega}^{T}} \prod_{t=1}^{T} P(v(t)|\omega(t))P(\omega(t)|\omega(t-1))$$

- Computationally less complex algorithm:
  - Calculate  $P(V^T)$  recursively
  - $P(v(t)|\omega(t))P(\omega(t)|\omega(t-1))$  involves only  $v(t), \omega(t)$  and  $\omega(t-1)$

$$\alpha_j(t) = \begin{cases} 0 & t = 0 \text{ and } j \neq \text{ initial state} \\ 1 & t = 0 \text{ and } j = \text{ initial state} \\ \left[\sum_i \alpha_i(t-1)p_{ij}\right] b_{jk} & \text{ otherwise } (b_{jk} \text{ leads to observed } v(t)) \end{cases}$$



# Hidden Markov Models – Evaluation problem

- Forward Algorithm
- Computational complexity:  $O(c^2 T)$

# Forward algorithm

1	initialise $t \leftarrow 0, p_{ij}, b_{jk}, V^T, \alpha_j(0)$
2	for $t \leftarrow t+1$
3	$j \leftarrow 0$
4	for $j \leftarrow j+1$
5	$\alpha_j(t) \leftarrow b_{jk} \sum_{i=1}^c \alpha_i(t-1) p_{ij}$
6	until $j = c$
7	until $t = T$
8	return $P(V^T) \leftarrow \alpha_j(T)$ for the final state
9	end



Sequential Stochastic Conclusion

# Hidden Markov Models – Decoding problem

- Given a sequence  $V^T$ , find most probable sequence of hidden states
- Enumeration of every possible path will cost  $O(c^{T})$ 
  - Not feasible



# Hidden Markov Models – Decoding problem

• Given a sequence  $V^T$ , find most probable sequence of hidden states

# Decoding algorithm

1	$\texttt{initialise:}  \texttt{path} \leftarrow \{\}, t \leftarrow 0$
2	for $t \leftarrow t+1$
3	$j \leftarrow 0;$
4	for $j \leftarrow j+1$
5	$lpha_j(t) \leftarrow b_{jk} \sum_{i=1}^{c} lpha_i(t-1) p_{ij}$
6	until $j = c$
7	$j' \leftarrow \mathit{arg} \max_j lpha_j(t)$
8	append $\omega_{j'}$ to path
9	until $t = T$
10	return path

#### 11 end



Linear discriminant

Sequential Stochastic C

ic Conclusion

### Hidden Markov Models – Decoding problem



- Computational time of the decoding algorithm
  - $O(c^2T)$



Conclusion

# Hidden Markov Models – Learning problem

- Determine the model parameters p<sub>ii</sub> and b<sub>ik</sub>
  - Given: Training sample of observed values  $V^{T}$
- No method known to obtain the optimal or most likely set of parameters from the data
  - However, we can nearly always determine a good solution by the forward-backward algorithm
  - General expectation maximisation algorithm
  - Iteratively update weights in order to better explain the observed training sequences



# Hidden Markov Models – Learning problem

• Probability that the model is in state  $\omega_i(t)$  and will generate the remainder of the given target sequence:

$$\beta_i(t) = \begin{cases} 0 & t = T \text{ and } \omega_i(t) \text{ not final hidden state} \\ 1 & t = T \text{ and } \omega_i(t) \text{ final hidden state} \\ \sum_j \beta_j(t+1)p_{ij}b_{jk} & \text{otherwise } (b_{jk} \text{ leads to } v(t+1)) \end{cases}$$



Introduction Recognition

Technische

Bayesian No

Non-parametric Li

# Hidden Markov Models – Learning problem

- α<sub>i</sub>(t) and β<sub>i</sub>(t) only estimates of their true values since transition probabilities p<sub>ij</sub>, b<sub>jk</sub> unknown
- Probability of transition between  $\omega_i(t-1)$  and  $\omega_j(t)$  can be estimated
  - Provided that the model generated the entire training sequence  $V^{\, T}$  by any path

$$\gamma_{ij}(t) = rac{lpha(t-1) 
ho_{ij} b_{jk} eta_j(t)}{P(V^T | \Omega)}$$

• Probability that model generated sequence  $V^T$ :

 $P(V^T|\Omega)$ 



Technische

Braunschweig

### Hidden Markov Models – Learning problem

• Calculate improved estimate for  $p_{ij}$  and  $b_{jk}$ 

$$\overline{p_{ij}} = \frac{\sum_{t=1}^{T} \gamma_{ij}(t)}{\sum_{t=1}^{T} \sum_{k} \gamma_{ik}(t)}$$

$$\overline{b_{jk}} = \frac{\sum_{t=1,v(t)=v_k}^T \sum_l \gamma_{jl}(t)}{\sum_{t=1}^T \sum_l \gamma_{jl}(t)}$$

- Start with rough estimates of  $p_{ij}$  and  $b_{jk}$
- Calculate improved estimates
- Repeat until some convergence is reached

# Hidden Markov Models – Learning problem

### Forward-Backward algorithm



# Outline

### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



# Stochastic methods

- When problem structure is not well known, it might be hard to configure classification methods correctly
- Then, randomised search approaches are an option to find optimum parameters for a classifier
- The search process then virtually becomes a problem of iterative learning
- Search space is then spanned by possible configurations for all distinct parameters



# Randomised search approaches

- Most frequently applied in engineering design
  - Assumptions to compute extrema are not fulfilled (e.g. unfriendly/unknown conditions)
  - Difficulties to carry out necessary differentiations
  - Solution to the equations describing all conditions does not always lead to optimum point in the search space
  - Equations to describe conditions are not immediately solvable



## Randomised search approaches

- Local random search
- Metropolis algorithm
- Simulated annealing
- Evolutionary algorithms


NN

# Local random search heuristics

- Local random search
  - Intuitive way to climb a mountain (by a sightless climber)

# Local random search

For every point x in a search space S, a non-empty neighbourhood  $N(x) \subseteq S$  is defined. The local random search approach iteratively draws one sample  $x' \in N(x)$ . When the fitness of the new value is better than the old one (F(x) < F(x')), the new value is utilised as the new best search point. Otherwise it is discarded.



Recognition

Bayesian Non-parametric Linear discriminant NN

Sequential Stochastic

# Local random search heuristics



- In principle, N(x) = x or N(x) = S is valid, but the original idea is that N(x) is a relatively small set of search points.
- The points  $x' \in N(x)$  are expected to be nearer to x than those points  $x'' \notin N(x)$
- Typically,  $x \in N(x)$



Technische Universität

Braunschweig

Stephan Sigg | Secure communication based on noisy input data | 146 Institute of Operating Systems and Computer Networks

Recognition Bayesian Non-parametric Linear discriminant NN

Sequential Stochastic

Conclusion

# Local random search heuristics

Complexity reduction by restriction of the search space size

	$d \leq 1$	$d \leq 2$	$d \leq 3$
	1010	1010	1010
Example: $S = \{0, 1\}^n$ and $N_1(x)$ are	0010	0010	1110
Example: $\mathbf{J} = \{0, 1\}$ and $\mathcal{M}_{d}(\mathbf{x})$ are	1110	1000	1000
all points y with Hamming distance	1011	1011	1011
smaller than $d(H(x,y) \leq d)$	1011	0110	0110
smaller than $u(n(x, y) \leq u)$		0000	0000
• For constant d we obtain:		0011	0011
		1100	1100
$ N_d(x)  = \Theta(n^d) \ll  S  = 2^n$		$1 \ 1 \ 1 \ 1 \ 1$	$1 \ 1 \ 1 \ 1 \ 1$
		$1 \ 0 \ 0 \ 1$	1001
			0100
			0111
			1101
			1001
(n) $(n)$ $(n)$	1)	(n)	
$ N_d(x)  = (1, 1) + (1, 1) + \dots + (1, 1)$	()+		
(d)	L / É	$\left( 0 \right)$	
	Example: $S = \{0, 1\}^n$ and $N_d(x)$ are all points $y$ with Hamming distance smaller than $d$ ( $H(x, y) \le d$ ) • For constant $d$ we obtain: $ N_d(x)  = \Theta(n^d) \ll  S  = 2^n$ $ N_d(x)  = \begin{pmatrix} n \\ d \end{pmatrix} + \begin{pmatrix} n \\ d-1 \end{pmatrix} + \dots + \begin{pmatrix} n \\ 1 \end{bmatrix}$	Example: $S = \{0,1\}^n$ and $N_d(x)$ are all points $y$ with Hamming distance smaller than $d$ ( $H(x,y) \le d$ ) • For constant $d$ we obtain: $ N_d(x)  = \Theta(n^d) \ll  S  = 2^n$ $ N_d(x)  = \begin{pmatrix} n \\ d \end{pmatrix} + \begin{pmatrix} n \\ d-1 \end{pmatrix} + \dots + \begin{pmatrix} n \\ 1 \end{pmatrix} +$	Example: $S = \{0,1\}^n$ and $N_d(x)$ are all points $y$ with Hamming distance smaller than $d$ $(H(x,y) \le d)$ • For constant $d$ we obtain: $ N_d(x)  = \Theta(n^d) \ll  S  = 2^n$ $ N_d(x)  = \begin{pmatrix} n \\ d \end{pmatrix} + \begin{pmatrix} n \\ d-1 \end{pmatrix} + \dots + \begin{pmatrix} n \\ 1 \end{pmatrix} + \begin{pmatrix} n \\ 0 \end{pmatrix}$



Stephan Sigg | Secure communication based on noisy input data | 147 Institute of Operating Systems

and Computer Networks

Introduction

Recognition Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

c Conclusion

# Local random search heuristics



- Small neighbourhood: Easy conversion to local optima
- Huge neighbourhood: Similar to random search
- Change neighbourhood radius during optimisation
  - Initially, big neighbourhood to allow huge steps, then decrease size
  - Challenging: Not to decrease neighbourhood size too fast

Introduction Recognition

Bayesian I

Non-parametric Lir

Linear discriminant

NN Sequential

Stochastic Conclusion

# Local random search heuristics

- Alternative to avoid local optima: Multistart strategies
  - Local search approach applied t times on the problem domain
  - Probability amplification results in respectable search result also when single success probability is low.
    - Assume a success probability of  $\delta > 0$  for one iteration of the algorithm
    - When the algorithm is applied t times, the overall probability of success is  $1-(1-\delta)^t$
    - Small polynomial success probabilities are enough for the multistart strategy to obtain very good overall success probabilities





Technische Universität Braunschweig

Stephan Sigg | Secure communication based on noisy input data | 149 Institute of Operating Systems and Computer Networks

ntroduction Recognition

Bayesian Non-parametric

Linear discriminant

N Sequential Stochastic

Conclusion

# Metropolis algorithms

- For local random search, only multistart strategies can avoid the termination in local optima.
- A Metropolis approach allows to accept also new search points that decrease the fitness value
- If F(x') < F(x) the search point x' is discarded only with probability

$$1-\frac{1}{e^{(F(x)-F(x'))/T}}$$





# Metropolis algorithms

- Probability to accept search points with decreasing fitness value dependent on degree by which fitness decreased
- $\bullet~$  For  $~\mathcal{T} \rightarrow 0$  the Metropolis approach becomes a random search
- $\bullet\,$  For  $\,\mathcal{T} \to \infty$  the Metropolis approach becomes an uncontrolled local search
- Choice of *T* impacts the performance
- $\bullet\,$  Knowledge on the problem or the fitness function might impact the choice of  ${\cal T}\,$



# Simulated annealing

- Choice of optimal T not easy: Change parameter in the pace of the optimisation
- Initially: T should allow to 'jump' to other regions of the search space with increased fitness value
- Finally: Process should gradually 'freeze' until local search approach propagates the local optimum in the neighbourhood.
- Name chosen in analogy to natural cooling processes in the creation of crystals
  - In this process, the temperature is gradually decreased so that Molecules that could move freely at the beginning are slowly put into their place



# Simulated annealing

- Optimal choice of the cooling schedule for T?
- Non-Adaptive approaches
  - Fixed temperature function T(t)
  - $\bullet\,$  Every few steps the original value is multiplied with a factor  $\alpha < 1$
- Adaptive approaches
  - React on the optimisation process
  - Probably dependent on the frequency of accepted iterations.



Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion

### Simulated annealing

- Problem: No natural problem known for which it has been proved that Simulated Annealing is sufficiently more effective than the Metropolis algorithm with optimum stationary temperature.
- However, artificially constructed problems exist, for which it could be shown that Simulated Annealing is superior to the Metropolis algorithm



Introduction Recognition

Bayesian

Non-parametric Linear discriminant

IN Sequential

Stochastic Conclusion

# Evolutionary algorithms



- Several researchers have studied the use of evolutionary approaches for optimisation purposes
- To-date, evolutionary algorithms combine these different approaches so that no clear distinction can be made



Universität Braunschweig

Technische

Introduction Recognition

Bayesian

Non-parametric

Linear discriminant

NN

Sequential Stochastic

Conclusion

# Evolutionary algorithms





Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential **Stochastic** Conclusion

# Evolutionary algorithms

#### Genetic algorithms

23

Technische

Braunschweig

- Proposed by John Holland <sup>23</sup>
- Binary discrete search spaces:  $\{0,1\}^n$
- Fitnessproportional selection
  - For *m* individuals  $x_1, \ldots, x_m$  the probability to choose  $x_i$  is  $\frac{f(x_i)}{f(x_1)+\cdots+f(x_m)}$ .
- Main evolution operator is crossover
  - Originally One-point crossover
- The main goal was not optimisation but the adaptation of an environment

daptation in Natural and Artificial Systems, University of Michigan Press, 1975.

Stephan Sigg | Secure communication based on noisy input data | 157 Institute of Operating Systems and Computer Networks Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential Stochastic Conclusion

# **Evolutionary algorithms**

#### **Evolution strategies**

- Proposed by Bienert, Rechenberg and Schwefel<sup>24</sup> <sup>25</sup>
- At first only steady search spaces as  $\mathbb{R}^n$

Evolution and optimum seeking, 1993

- No Crossover
- Only mutation
  - First mutation operator: Each component  $x_i$  is replaced by  $x_i + \sigma Z_i$ ( $Z_i$  normally distributed,  $\sigma^2$  Variance)

<sup>24</sup> I. Rechenberg, Evolutionsstrategie: Optimierung technischer Systeme nach Prinzipien der biologischen Evolution, 1973.



Stephan Sigg | Secure communication based on noisy input data | 158 Institute of Operating Systems

# **Evolutionary algorithms**

#### Evolutionary programming

- The approach was proposed by Lawrence J. Fogel<sup>2627</sup>
- Various similarities to evolution strategies
- Search Space: Space of deterministic finite automata that well adapt to their environment.



Technische

Braunschweig

and Computer Networks

<sup>&</sup>lt;sup>26</sup>L.J. Fogel, *Autonomous automata*, Industrial Research, Vol. 4, 1962.

Diptechnology: Concepts and Applications, Prentice-Hall, 1963

Introduction Recognition Bayesian Non-parametric Linear discriminant NN Sequential

al Stochastic Conclusion

### **Evolutionary algorithms**

#### Genetic programming

- Proposed by John Koza<sup>28</sup>
- Search space: Syntactically correct programs
- Crossover more important than mutation

<sup>&</sup>lt;sup>28</sup> John Koza Genetic Programming: On the Programming of Computers by Means of Natural Selection, MIT Press,



Sequential Stochastic

c Conclusion

- In the early days of evolutionary algorithms it has been argued that
  - Problem specific algorithms are better than evolutionary algorithms on a very small subset of problems
  - Evolutionary algorithms perform better on average over all problems
- Therefore, evolutionary algorithms have been proposed as a good choice for a general purpose optimisation scheme



Introduction Recognition

Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

c Conclusion





Recognition

Bavesian

Non-parametric

Linear discriminant

NN

Sequential Stochastic

Conclusion

- Can an algorithm be suited for 'all' problems?
  - Distinct coding of the search space
  - Various fitness functions
- What does 'all problems' mean?
  - For all possible representations and sizes of the search space
  - All possible fitness functions on the feature space
  - For a given search space and feature space, all possible fitness functions
  - Every single point in the search space is the optimum point in several of these problems
- Can one algorithm be better on average than another algorithm on 'all' problems?



# Restrictions of evolutionary approaches

- To understand this scenario, Wolpert and Macready formalised the assertion<sup>29</sup>
- Assumptions:
  - The set of all functions  $f: S \rightarrow W$  considered is given by F
  - S and W are finite (as every computation on physical computers can only have finite resources)
  - The fitness function is evaluated only once for each search point
  - A(f) is the number of search points requested until the optimum is found

<sup>&</sup>lt;sup>29</sup>D.H. Wolpert and W.G. Macready, *No Free Lunch Theorems for Optimisation*, IEEE Transactions on Evolutionary



Technische

Braunschweig

# Restrictions of evolutionary approaches

### No free lunch theorem

Assume that the average performance of an algorithm in the No Free Lunch Scenario for S and W is given by  $A_{S,W}$ , the average over all  $A(f), f \in F$ . Given two algorithms A and A', we obtain  $A_{S,W} = A'_{S,W}$ 

• This means that two arbitrary algorithms perform equally well on average on all problems



Linear discriminant

Sequential Stochastic

Conclusion

# Restrictions of evolutionary approaches

### Proof of the No Free Lunch Theorem

W.I.o.g.:  $W = \{1, \dots, N\}$ 

We consider sets  $F_{s,i,N}$  of all functions f on a search space of non-visited search points of size s with at least one x with f(x) > iObserve that for every function f and every permutation  $\pi$  also  $f_{\pi}$ belongs to  $F_{s,i,N}$ 



Introduction Recognition

Bayesian I

Non-parametric

Linear discriminant

NN

Sequential Stochastic

Conclusion

# Restrictions of evolutionary approaches

#### Proof of the No Free Lunch Theorem

Proof by induction over s := |S|.

Induction start: s = 1

Every algorithm has to choose the single optimum search point with its first request.



Linear discriminant

NN

Sequential Stochastic

c Conclusion

# Restrictions of evolutionary approaches

# Proof of the No Free Lunch Theorem

Induction:  $s - 1 \rightarrow s$ 

We define a function  $a: S \to \mathbb{N}$  so that for every  $x \in S$  the share of functions with f(x) = j is exactly a(j).

This is independent of x, since all permutations  $f_{\pi}$  of a function f also belong to  $F_{s,i,N}$ ,

a(j) is therefore the probability to choose a search point with fitness value j – Independent of the concrete algorithm A



Linear discriminant

NN

Sequential Stochastic

c Conclusion

# Restrictions of evolutionary approaches

### Proof of the No Free Lunch Theorem

Induction:  $s - 1 \rightarrow s$ 

With probability a(j) an algorithm A finds a search point with fitness value j.

Count of functions f(x) = j is equal to the number of functions

 $f_{\pi}(y) = j$ , since all permutations of f are also in  $F_{s,i,N}$ .

The probability to achieve a fitness value j > i is therefore independent of the algorithm.



NN

Sequential Stochastic

Conclusion

# Restrictions of evolutionary approaches

### Proof of the No Free Lunch Theorem

Induction:  $s - 1 \rightarrow s$ 

With probability a(i) an algorithm A finds a search point with fitness value *j*.

If  $j \leq i, x$  is not optimal in scenario  $F_{s,i,N}$  and the new scenario is  $F_{s-1,i,N}$ 



Linear discriminant

NN

Sequential Stochastic

Conclusion

# Restrictions of evolutionary approaches

### Proof of the No Free Lunch Theorem

Summary – in other words:

For any two algorithms we can state a suitable permutation of the Problem-function for one problem (i.e. state another problem), so that both algorithms in each iteration request identical search points.

 Especially, since every search point could be optimal, there are always algorithms that request the optimal search point right from the start.



NN

- The NFL is possible, since ALL algorithms and ALL problems are considered
- It is a reasonable question if an NFL is also valid in smaller, more realistic scenarios.
- In <sup>30</sup> is was proved, that a similar theorem can be stated also for more realistic problem scenarios.

<sup>&</sup>lt;sup>30</sup>S. Droste, T. Jansen and I. Wegener, *Perhaps not a free lunch but at least a free appetizer*, Proceedings of the 1st structure of the property of the property of the structure of the property of the structure of the property of the structure of the structu



Linear discriminant

NN

Sequential Stochastic

Conclusion

# Design aspects of evolutionary algorithms

Selection principles

Uniform selection Individuals chosen uniformly at random

#### Deterministic selection

Deterministically choose the highest rated individuals for the selection

### Threshold selection

Candidates for offspring population drawn uniformly at random from the t highest rated individuals



NN

Sequential Stochastic

c Conclusion

# Design aspects of evolutionary algorithms

#### Selection principles

- Fitnessproportional selection
  - For population  $x_i, \ldots, x_n$  individual  $x_i$  chosen with

$$p(x_i) = \frac{f(x_i)}{f(x_1) + \cdots + f(x_n)}$$

• Draw random variable u from [0, 1] and consider  $x_i$  if

$$p(x_1) + \cdots + p(x_{i-1}) < u \leq p(x_1) + \cdots + p(x_i)$$

#### Frequently applied for evolutionary approaches



NN

Sequential Stochastic

Conclusion

# Design aspects of evolutionary algorithms

#### Selection principles

- Problems with Fitnessproportional selection
  - Linear modification of the fitness function  $(f \rightarrow f + c)$  results in different behaviour
  - When fitness values sufficiently separated, selection is nearly deterministic
  - When deviation in fitness values is small relative to absolute values. similar to uniform selection



Introduction Recognition

Bayesian

Non-parametric

Linear discriminant

Sequential Stochastic

Conclusion

# Design aspects of evolutionary algorithms

#### Variation – Mutation



- Mutation creates one offspring individual from one individual
- Operators are designed for specific search spaces
- Shall apply only few modifications of individuals on average
- Distant individuals have smaller probability



Technische Universität Braunschweig

# Evolutionary algorithms

Mutation operators for individuals from  $\mathbb{B}^n$  (similar operators for other search spaces):

### Standard bit mutation

- Offspring individual created bit-wise from parent individual
- Every bit 'flipped' with probability  $p_m$
- Common choice:  $p_m = \frac{1}{n}$

# 1 bit mutation

Technische

Braunschweig

- Offspring individual identical in all but one bit.
- This bit chosen uniformly at random from all *n* bits

NN

# Design aspects of evolutionary algorithms

#### Variation – Crossover

- Crossover typically takes two individuals and results in one or two offspring individuals
  - Also crossover of more than two individuals possible
  - Often generalisations of the two-individual case
- Distinct crossover methods for various search spaces
- Crossover parameter  $p_c$  specifies the probability with which crossover (and not mutation) is applied for one selected individual
- In some cases (e.g. binary coded numbers) not all positions in the individual string are allowed to apply crossover on



Linear discriminant

Stochastic Sequential

Conclusion

## Design aspects of evolutionary algorithms

### Typical crossover variants

- One-point crossover
- k-point crossover
- Uniform crossover





## Evolutionary algorithms

#### Crossover operators for $\mathbb{B}^n$

(Operators for other search spaced similar)

One-point crossover: Individual x'' from two individuals x and x' according to uniformly determined crossover position:

$$x_j'' = \begin{cases} x_j & \text{if } j \le i \\ x_j' & \text{if } j > i \end{cases}$$


# Evolutionary algorithms

#### Crossover operators for $\mathbb{B}^n$

*k*-point crossover: Choose  $k \le n$  positions uniformly at random:

$$\begin{array}{rcl} x_1 = & x_{11}, x_{1,2}, \dots, x_{1,k_1} | x_{1k_1+1}, \dots, x_{1k_2} | x_{1k_2+1}, \dots, x_{1n} \\ x_2 = & x_{21}, x_{2,2}, \dots, x_{2,k_1} | x_{2k_1+1}, \dots, x_{2k_2} | x_{2k_2+1}, \dots, x_{2n} \\ \hline y_1 = & x_{11}, x_{1,2}, \dots, x_{1,k_1} | x_{2k_1+1}, \dots, x_{2k_2} | x_{1k_2+1}, \dots, x_{1n} \\ y_2 = & x_{21}, x_{2,2}, \dots, x_{2,k_1} | x_{1k_1+1}, \dots, x_{1k_2} | x_{2k_2+1}, \dots, x_{2n} \end{array}$$



Introduction Recognition

Bayesian Non-

Non-parametric

Linear discriminant

NN

Sequential Stochastic

c Conclusion

# Evolutionary algorithms

#### Crossover operators for $\mathbb{B}^n$

# Uniform crossover: Each bit chosen with uniform probability from one of the parent individuals



NN

Sequential Stochastic Con

#### c Conclusion

# Design aspects of evolutionary algorithms

#### Discussion

- Evolutionary algorithms are easy to implement when compared to some complex specialised approaches
- However, Evolutionary algorithms are computationally complex
- It is therefore beneficial to implement efficient variants to the distinct methods



Sequential Stochastic Co

#### ic Conclusion

# Design aspects of evolutionary algorithms

### Discussion

- Generation of pseudo random bits is important for many of the theoretic results for evolutionary algorithms to hold
- It is, however possible to reduce the number of random experiments
  - It is more efficient to calculate the next flipping bit in a mutation instead of doing the calculation for every bit independently



NN

Sequential Stochastic Conclusion

# Design aspects of evolutionary algorithms

#### Discussion

- Most of the computational time is typically consumed by the fitness calculation
- One approach to reduce complexity is to prevent re-calculation of fitness for individuals
  - Dynamic data structures that support search and insert



# Outline

#### Introduction

- Recognition of patterns
- Bayesian decision theory
- Non-parametric techniques
- Linear discriminant functions
- Neural networks
- Sequential data
- Stochastic methods
- Conclusion



# **Questions?**

Stephan Sigg sigg@ibr.cs.tu-bs.de



### Literature

- C.M. Bishop: Pattern recognition and machine learning, Springer, 2007.
- P. Tulys, B. Skoric, T. Kevenaar: Security with Noisy Data On private biometrics, secure key storage and anti-counterfeiting, Springer, 2007.
- R.O. Duda, P.E. Hart, D.G. Stork: Pattern Classification, Wiley, 2001







