

# Sensor Signal Processing

**Prof. Dr.-Ing. Andreas König**  
Lehrstuhl Integrierte Sensorsysteme



FB Elektrotechnik und Informationstechnik  
Technische Universität Kaiserslautern

Fall Semester 2005



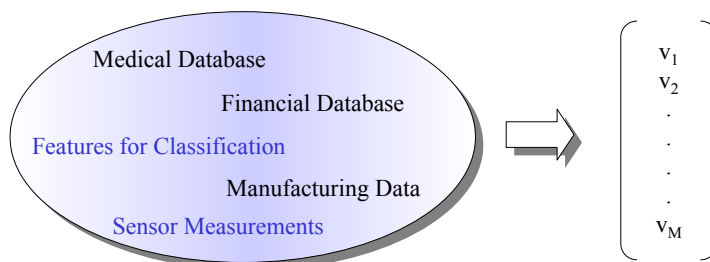
## Course Contents

1. Introduction
2. Signal Processing and Analysis
3. Feature Computation
4. Cluster Analysis
5. Dimensionality Reduction Techniques
6. Data Visualization & Analysis
7. Classification Techniques
8. Sensor Fusion
9. Systematic Design of Sensor Systems
10. Outlook

**5. Dimensionality Reduction Techniques**

- 5.1 Motivation
- 5.2 Assessment functions
- 5.3 Feature selection
- 5.4 Feature extraction
- 5.5 Accelerated methods
- 5.6 Summary

- Technical problems are often characterized by sets of high dimensional data
- Significance, correlations, redundancy, or irrelevancy of the variables  $v_i$  with regard to the given application a priori unknown

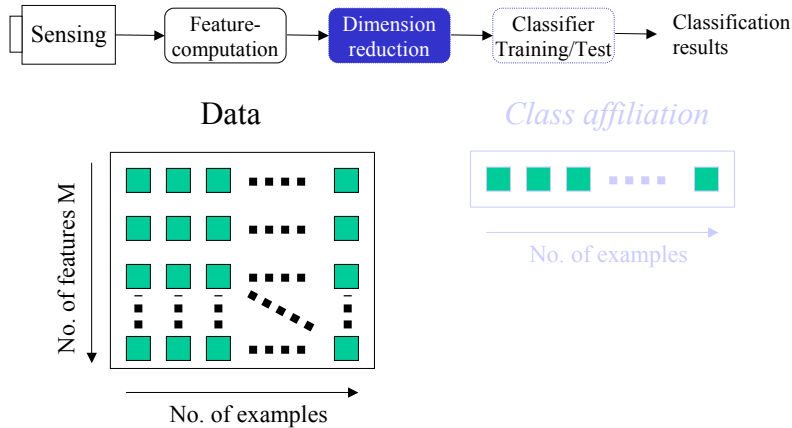


- Computational complexity and *Curse of Dimensionality* promote efficient dimensionality reduction, visualization can provide transparency & insight
- Dimensionality reduction is an ubiquitous problem in many disciplines !

## Sensor Signal Processing Dimensionality Reduction

### Motivation

- After signal processing and feature computation further condensing or compression of the data **in terms of attributes or features** is desirable
- Supervised as well as unsupervised **dimensionality reduction techniques** can serve for that purpose !

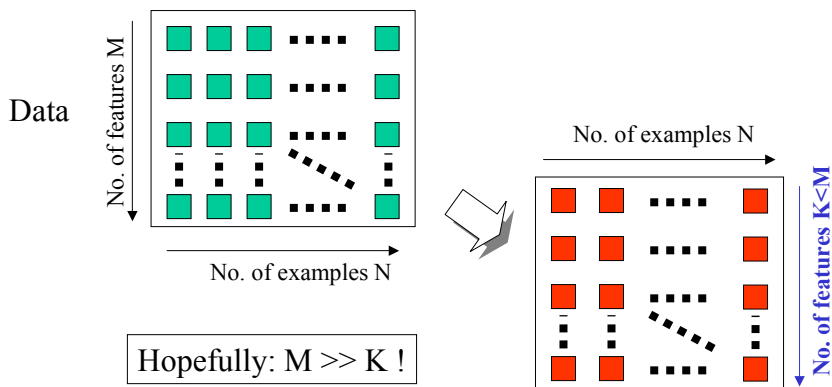


© Andreas König Slide5-5

## Sensor Signal Processing Dimensionality Reduction

### Motivation

- There is a close relation between data analysis and compression based on dimensionality reduction techniques
- The purpose is to represent the data by a reduced number of variables or attributes or features according to **appropriate assessment functions** !



© Andreas König Slide5-6

## Sensor Signal Processing Dimensionality Reduction

### Motivation

- Projection of multivariate data by DR mappings & ensuing visualization and interactive analysis is a research topic of interest for more than 3 decades
- Recently, data mining/warehouse & knowledge discovery applications give renewed strong incentive & drive to the field
- For classification, **feature extraction & selection** for vectors  $\vec{v} = [v_1, v_2, \dots, v_o]^T$  are typically defined as (see, e.g., [Kittler 86])

Feature Extraction :	$J(A) = \max_A J(A(\vec{v}))$
Feature Selection :	$J(X) = \max_X J(\chi)$

A mapping  $\Phi : \mathbb{R}^o \rightarrow \mathbb{R}^d$  is optimized with regard to assessment criterion  $J$  and with  $d < o$  and  $\vec{y} = [y_1, y_2, \dots, y_d]^T$

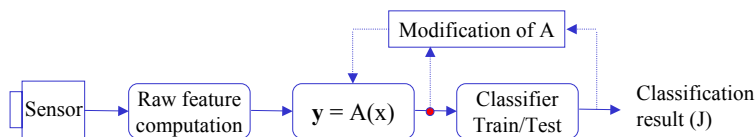
- Methods for classification may be salient for visualization and vice versa
- $\vec{y} = A(\vec{v})$  can be a linear or a nonlinear mapping

© Andreas König Slide5-7

## Sensor Signal Processing Dimensionality Reduction

### Motivation

- The dimensionality reducing mapping can be unsupervised or employ supervised information



- Optimization criteria can be, e.g., signal, topology, distance preservation or discriminance gain
- Dimensionality reduction (DR) methods can be salient both for multivariate data classification and visualization
- **Classification:** discriminance optimization with DR constraint for lean and well performing recognition system
- **Visualization:** DR fixed to dimension two or three with, e.g., structure preservation constraint for data analysis, advanced MMI & system design

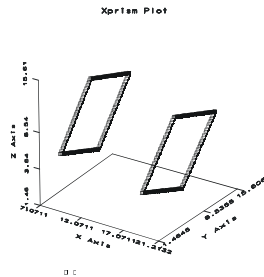
© Andreas König Slide5-8

## Sensor Signal Processing Dimensionality Reduction

### Motivation

- Benchmark Data is required for method demonstration & assessment

- **Cube** data: Artificial data, 3 dimensions, 400 vectors, 8 classes. Data points on 8 edges of two opposite sides of a cube rotated by 45°



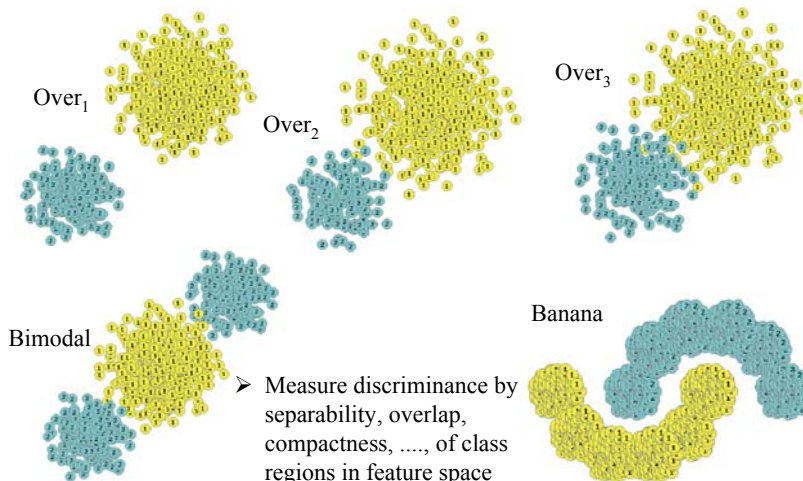
- **Iris** data: Well known Iris flower data, 4 dimensions, training & test set with 75 vectors each, 3 classes virginica, setosa, and versicolor
- **Mech<sub>x</sub>** data: Mechatronic data from turbine jet engine compressor monitoring. Five data sets with 375 24-D vectors each, 4 classes corresponding to operating regions and compressor stability
- **X-Ray** data: X-ray inspection of ball grid array packages in electronics manufacturing. Two sets, 40 10-D vectors each, 3 classes for ok & defect type

© Andreas König Slide5-9

## Sensor Signal Processing Dimensionality Reduction

### Motivation

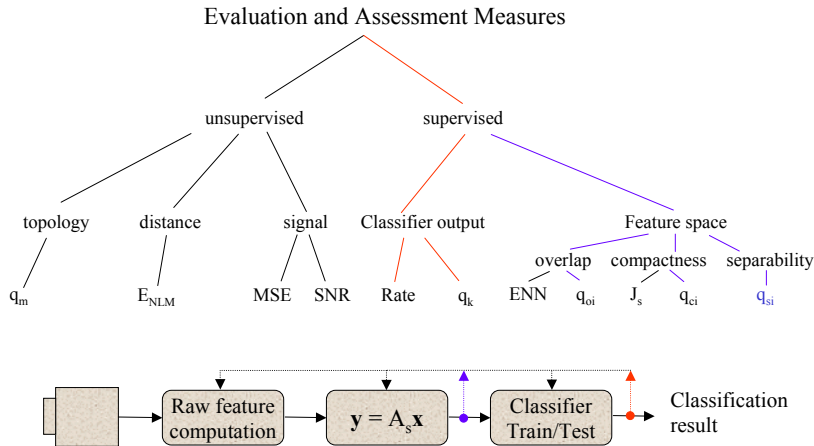
- Supervised Assessment Options of Feature Spaces



- Measure discriminance by separability, overlap, compactness, ....., of class regions in feature space

© Andreas König Slide5-10

### ➤ Taxonomy of Evaluation and Assessment Measures



© Andreas König Slide5-11

- Simple Parametric Overlap Measure
- Class specific distributions modelled by Gaussian functions
- Pairwise overlap can be computed from mean values  $\mu_i$ ,  $\mu_j$  and standard deviations  $\sigma_i$ ,  $\sigma_j$  by

$$q_{v_{ij}} = \frac{|\mu_i - \mu_j|}{(N_i - 1)\sigma_i + (N_j - 1)\sigma_j} \quad (5.1)$$

- The feature merit for separating one class from all others is given by

$$q_{v_{i_i}} = \frac{1}{L-1} \sum_{j \neq i}^L q_{v_{ij}} \quad (5.2)$$

- The merit of a single feature to distinguish all classes could be computed by

$$q_{v_i} = \frac{1}{L} \sum_{i=1}^L q_{v_{i_i}} \quad (5.3)$$

- Global summation can be misleading in some cases
- Simple & efficient measure for fast *parametric first order feature selection*

© Andreas König Slide5-12

- Inspired by the work of Fukunaga et al. on scatter matrices & nonlinear mappings, a nonparametric compactness measure can be derived based on intra/inter class distances

$$q_{ci} = \frac{\frac{1}{L} \sum_{l=1}^L \frac{2}{N_l(N_l-1)} \sum_{i=1}^N \sum_{j=i+1}^N \partial(\omega_i, \omega_j) * \partial(\omega_i, l) * d_{x_{ij}}}{\frac{1}{N^B} \sum_{i=1}^N \sum_{j=i+1}^N (1 - \partial(\omega_i, \omega_j)) * d_{x_{ij}}} \quad (5.4)$$

$$\text{with } d_{x_{ij}} = \sqrt{\sum_{q=1}^M (x_{iq} - x_{jq})^2} \quad \text{and} \quad \partial(\omega_i, \omega_j) = \begin{cases} 1 & \text{if } \omega_i = \omega_j \\ 0 & \text{if } \omega_i \neq \omega_j \end{cases}$$

- The Kronecker Delta  $\partial(\omega_i, l)$  assures, that only vectors of class  $l$  are regarded
  - + Nonparametric, multivariate compactness measure
  - + The measure is free of user-definable parameters
  - Measure has  $O(N^2)$  complexity
  - *Normalization properties* only allow observation of relative changes

- Inspired by probability estimation in Edited-Nearest-Neighbor method (ENN) or Leave-One-Out (LOO) classification
- For each vector, a number of  $k$  nearest-neighbors are computed & according to the neighbors' affiliation to same/different class a *nonparametric overlap measure* is computed

$$q_{oi} = \frac{1}{L} \sum_{c=1}^L \frac{1}{N_c} \sum_{j=1}^{N_c} \frac{\sum_{i=1}^k q_{NN_{ji}} + \sum_{i=1}^k n_i}{2 * \sum_{i=1}^k n_i} \quad (5.5)$$

$$\text{with } n_i = 1 - \frac{d_{NN_{ji}}}{d_{NN_{jk}}} \quad \text{and} \quad q_{NN_{ji}} = \begin{cases} n_i & \text{if } \omega_j = \omega_i \\ -n_i & \text{if } \omega_j \neq \omega_i \end{cases}$$

- Simplification of the measure can be achieved by only computing the rank or just the number of nearest neighbors in the measure instead of the distances

## Assessment Functions

## Sensor Signal Processing Dimensionality Reduction

- Assessment of the nonparametric overlap measure:
  - +  $q_{oi}$  is well normalized in  $[0,1]$ ; 1.0 indicates no overlap of class regions
  - +  $q_{oi}$  is fine grained and thus well suited for optimization
  - The method has one required parameter  $k$  ( $k$  typically set to 5-10)
  - $q_{oi}$  has  $O(N^2)$  complexity
- Application example of  $q_o$  for synthetic & application data:

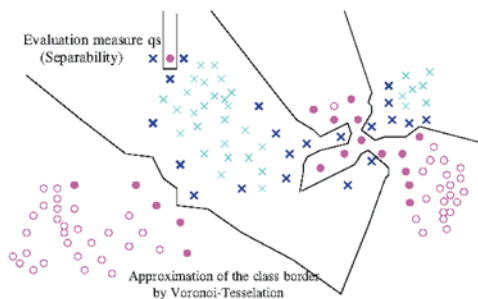
Data set	$q_o$	$q_o'$	$q_o''$
Over <sub>1</sub>	1,000000	1,000000	1,000000
Over <sub>2</sub>	0,991100	0,993100	0,992800
Over <sub>3</sub>	0,975300	0,979400	0,978200
Banana	1,000000	1,000000	1,000000
Bimodal	0,990900	0,989800	0,990000
Iris train	0,915600		
Pins	0,952700		
Mech <sub>1</sub>	0,978800		

© Andreas König Slide5-15

## Assessment Functions

## Sensor Signal Processing Dimensionality Reduction

- Nonparametric Separability Measure exploits the fast training run of a Reduced-Nearest-Neighbor-Classfier (RNN)
- Separability is proportional to the number of chosen reference vectors  $T_{RNNi}$  per class during dynamic RNN configuration



$$q_{si} = \frac{1}{L} \sum_{i=1}^L \frac{N_i - (T_{RNN_i} - 1)}{N_i} \quad (5.6)$$

- Here,  $N_i$  denotes the number of patterns per class &  $L$  the number of classes
  - + Fast due to  $O(N)$  complexity, no required user-definable parameters
  - + normed response  $[0,1]$ , where 1.0 means linear separability
  - Coarse grained & thus less well suited for optimization

© Andreas König Slide5-16



## Assessment Functions

## Sensor Signal Processing Dimensionality Reduction

- Application example of  $q_{si}$  to synthetic and application data in comparison with Fukunaga's measure  $J$  from scatter matrices

Data set	$q_s$	$J_s$
Over <sub>1</sub>	1,000000	0,030863
Over <sub>2</sub>	0,985612	0,025812
Over <sub>3</sub>	0,960884	0,026084
Banana	0,993421	0,002372
Bimodal	0,983333	0,007937
Iris train	0,906667	2,659621
Pins	0,908108	3,461134
Mech <sub>1</sub>	0,989333	4,824150

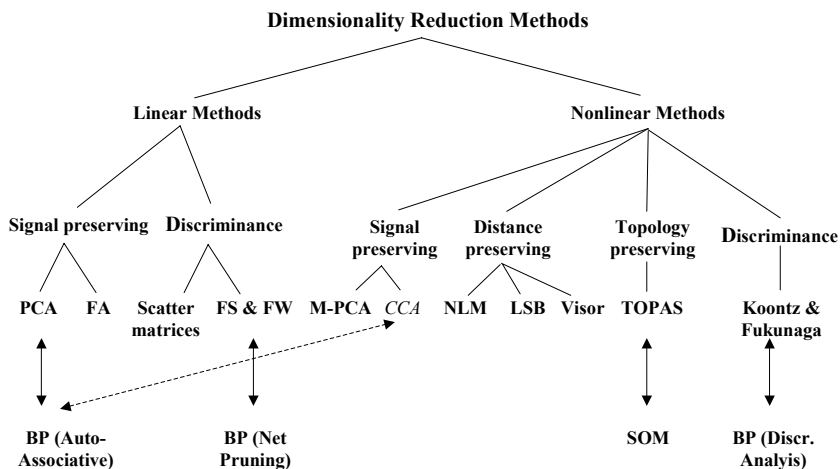
- Synthetic data assessment results show the sensitivity to underlying & gradually decreasing separability
- Identical vectors in the data set with different class affiliations return an assessment value of  $q_{si} = 0.0$

© Andreas König Slide5-17

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- Taxonomy of Dimensionality Reduction Methods:



© Andreas König Slide5-18

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- **Quest:** Find minimum feature subset with optimum discriminance
- AFS chooses for a given sample set the subset of variables or configuration  $X$  that maximizes a given cost function  $J$ :

$$J(X) = \max_{\chi} J(\chi) \quad (5.7)$$

- AFS is a linear mapping, based on the selection matrix  $A_S$  with  $\bar{y}_j = A_S \bar{x}_j$
  - **Feature selection:**  $c_i \in \{0,1\}$   
switch variables,  $2^M$  combinations
  - **Feature weighting:**  $c_i \in [0,1]$  or arbitrary real numbers
- $$A_S = \begin{pmatrix} c_1 & 0 & 0 & \dots & 0 \\ 0 & c_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & c_{M-1} & 0 \\ 0 & 0 & 0 & \dots & c_M \end{pmatrix} \quad (5.8)$$
- Cost function  $J$  can be one of the measures given in the taxonomy, e.g.,  $q_{si}$
  - Combinatorial optimization problem, **NP-complete**
  - Exhaustive search for global optimum only feasible for small  $M$
  - **Way out:** Apply **heuristics & optimization strategies** to find at least a **local optimum** with bounded time and effort

© Andreas König Slide5-19

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- **Simplification:**  $J$  is computed for each individual feature and a selected combination of classes, e.g., pairwise class discriminance
- Considerable computational savings by neglecting possible higher order correlations between feature pairs or tuples
- $q_{ij}$  could serve as simple **parametric** measure (**1st order par. selection**)
- For each class separation feature are ranked according to their  $J$  value
- Rank table example for **Iris** data:

Feature	Rank	C 1-2	Rank	C 1-3	Rank	C 2-3
$v_1$	4	1.020	3	1.482	3	0.442
$v_2$	3	1.065	4	0.890	4	0.255
$v_3$	2	4.139	<b>1</b>	<b>5.451</b>	2	1.218
$v_4$	<b>1</b>	<b>4.387</b>	2	5.180	<b>1</b>	<b>1.660</b>



- Selection takes place by choosing the features on top rank positions, e.g., features **3 & 4** for first rank position ( $C=[0, 0, 1, 1]^T$ )

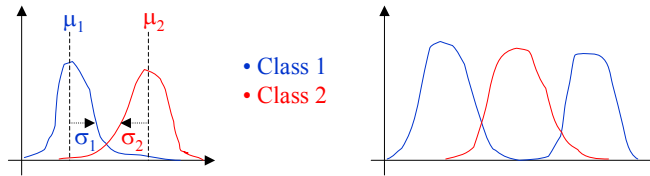
© Andreas König Slide5-20

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

### Assessment of Parametric First-Order Feature Selection (FOFS):

- The method is very fast and its complexity is  $O(M)$
- However, for pairwise class separation, the complexity grows with  $O(L)$
- The rank tables grow exponentially, visual inspection becomes infeasible !
- Only a local optimum solution can be expected
- The method can be extended to different class separation (one vs. all, all)
- Inclusion of lower ranking features can affect the solution
- **Summary:** If parametric assumption is met, method can be fast & effective
- **Problem:** Nonparametric and multimodal one-dimensional distributions:



© Andreas König Slide5-21

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- **Remedy:** A nonparametric J is computed, e.g.,  $q_{oi}$  restricted to one dimension
- Example for application data with two classes from visual inspection:

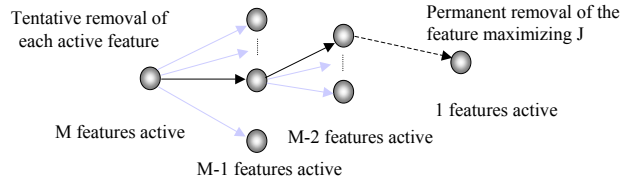
Feature	Rank Par.)	C 1-2	Rank (Nonpar.)	C 1-2
$v_1$	5	0.2917	5	0.6977
$v_2$	4	0.6489	3	0.8211
$v_3$	<b>1</b>	<b>1.1558</b>	4	0.7058
$v_4$	3	0.8547	<b>2</b>	<b>0.8808</b>
$v_5$	2	1.0047	<b>1</b>	<b>0.9270</b>

- The global solution  $C=[0, 0, 0, 1, 1]^T$  was determined by exhaustive search
- The nonparametric FOFS found the global optimum, parametric FOFS failed
- However, the number of rank positions salient for the given problem is not immediately obvious
- FOFS can serve as first step to confine search space in selection hierarchy

© Andreas König Slide5-22

## Feature Selection

- Nonparametric cost function is computed for each configuration, i.e., the currently selected subset of features
- Heuristic search strategy is applied to confine the search space
- Sequential-Backward/Forward-Selection (SBS/SFS):



- $M*(M-1)/2 + M$  combinations have to be assessed for SBS/SFS  
 $\Rightarrow O(M^2)$  complexity
- For  $M=16$ , a local optimum solution will be found in approx. 136 s (1 s per assessment assumed) in contrast to 18h for exhaustive search

© Andreas König Slide5-23

## Feature Selection

- Application example of SBS for *Iris* train:

1	2	3	4	0.90667
-	2	3	4	0.94667
-	-	3	4	0.96000
-	-	-	4	0.00000
Optimum quality:				0.96
Significant Features:				3 4

- Application example of SBS for visual inspection data:

1	2	3	4	5	0.903846
-	2	3	4	5	0.942308
-	-	3	4	5	0.961538
-	-	-	4	5	0.961538
-	-	-	-	5	0.913462
Optimum quality:					0.961538
Significant Features:					4 5



- Further applicable heuristics: Branch & bound, floating search, etc.

© Andreas König Slide5-24

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- Stochastic methods, e.g., Simulated Annealing (SA), applicable as heuristic
- Perturbation Method (PM) is a simple variant, based on random state changes
- PM application for *Iris*train:

				Proposed	Accepted
1	-	-	-	0.00000	0.000000
1	-	3	-	0.88000	0.880000
1	-	3	-	0.00000	0.880000
1	-	3	4	0.94667	0.946667
1	-	3	4	0.92000	0.946667
1	-	3	4	0.90667	0.946667
1	-	3	4	0.90667	0.946667
1	-	3	4	0.88000	0.946667
1	-	3	4	0.88000	0.946667
1	-	3	4	0.88000	0.946667
-	-	3	4	0.96000	0.960000
-	-	3	4	0.00000	0.960000
-	-	3	4	0.00000	0.960000
-	-	3	4	0.00000	0.960000
-	-	3	4	0.94667	0.960000

Best Quality: 0.96000  
Best Features: 3 4



- Random initial state
- Random selection of variable to propose state transition
- Selected variable is toggled and transition is accepted if

$$\Delta q_{si} = q_{si}^{new} - q_{si}^{old} \quad (5.9)$$

- Random tossing surprisingly effective at moderate effort
- Different results for each run !
- Several enhancements possible (multiple variable change or SA)

© Andreas König Slide5-25

## Feature Selection

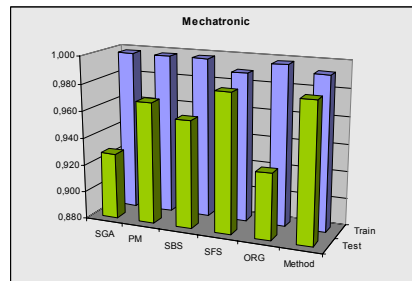
## Sensor Signal Processing Dimensionality Reduction

- Additional optimization techniques can be applied to feature selection, e.g., from evolutionary computation like genetic algorithms & swarm intelligence

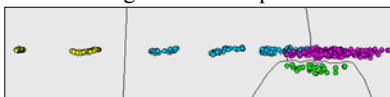
### ➤ Results for Mech<sub>1</sub> & Mech<sub>2</sub>:

- SFS and FSSPEA find same optimum solution !
- Inclusion of further constraints

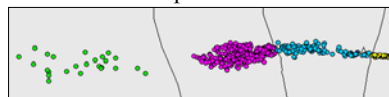
Method	Mech <sub>1</sub>	Mech <sub>2</sub>	Features
ORG	0.992500	0.981667	24
➔ SFS	0.997500	0.928333	3
SBS	0.989333	0.981667	12
PM	0.997500	0.959167	6
SGA	0.997333	0.969167	12
➔ FSSPEA	0.997500	0.928333	3



Original feature space



Feature space for FSSPEA

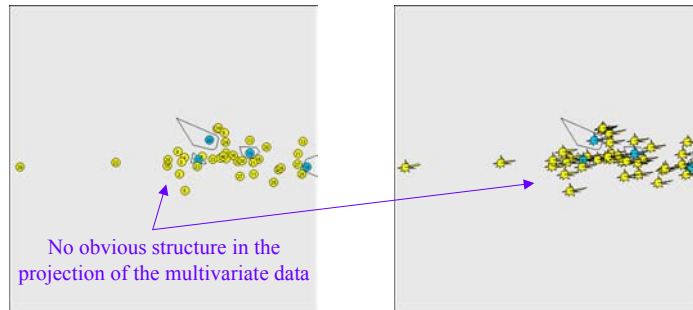


© Andreas König Slide5-26

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- Small application example from microelectronic manufacturing: Wafer Data Analysis
- Simple example of DR methodology potential for general fab data analysis
- The database was gathered from an MPC-Run with 10 Wafers and 4 Parameter Extraction Sites per Wafer, 59 parameters each:



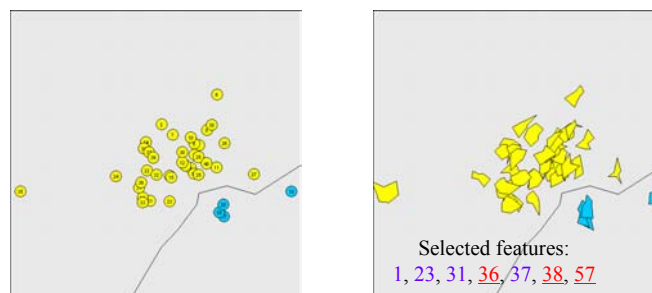
- All chips with unsatisfactory behavior on Wafer 5 (blue; 1-4, 6-10 yellow)
- **Assume Wafer 5 processing was abnormal, what makes it different ?**

© Andreas König Slide5-27

## Feature Selection

## Sensor Signal Processing Dimensionality Reduction

- The lack of obvious structure can be due to absence of abnormality or to occlusion by a majority of normal parameters and high intrinsic dimension
- Thus, AFS is employed to find parameters supporting abnormality hypothesis



- From 59 parameters AFS chose 3 for  $q_{si}/SBS$  and 7 for  $q_{oi}/SBS$
- **There is (weak) hypothesis support: Can these features be responsible ?**

© Andreas König Slide5-28

## Sensor Signal Processing

### Dimensionality Reduction

- Analyzing the AFS choice by probing the parameter meaning:

Plane	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10	PC11	PC12	PC13	PC14	PC15	PC16	PC17	PC18	PC19	PC20	PC21	PC22	PC23	PC24	PC25	PC26	PC27	PC28	PC29	PC30	PC31	PC32	PC33	PC34	PC35	PC36	PC37	PC38	PC39	PC40	PC41	PC42	PC43	PC44	PC45	PC46	PC47	PC48	PC49	PC50	PC51	PC52	PC53	PC54	PC55	PC56	PC57	PC58	PC59	PC60	PC61	PC62	PC63	PC64	PC65	PC66	PC67	PC68	PC69	PC70	PC71	PC72	PC73	PC74	PC75	PC76	PC77	PC78	PC79	PC80	PC81	PC82	PC83	PC84	PC85	PC86	PC87	PC88	PC89	PC90	PC91	PC92	PC93	PC94	PC95	PC96	PC97	PC98	PC99	PC100	PC101	PC102	PC103	PC104	PC105	PC106	PC107	PC108	PC109	PC110	PC111	PC112	PC113	PC114	PC115	PC116	PC117	PC118	PC119	PC120	PC121	PC122	PC123	PC124	PC125	PC126	PC127	PC128	PC129	PC130	PC131	PC132	PC133	PC134	PC135	PC136	PC137	PC138	PC139	PC140	PC141	PC142	PC143	PC144	PC145	PC146	PC147	PC148	PC149	PC150	PC151	PC152	PC153	PC154	PC155	PC156	PC157	PC158	PC159	PC160	PC161	PC162	PC163	PC164	PC165	PC166	PC167	PC168	PC169	PC170	PC171	PC172	PC173	PC174	PC175	PC176	PC177	PC178	PC179	PC180	PC181	PC182	PC183	PC184	PC185	PC186	PC187	PC188	PC189	PC190	PC191	PC192	PC193	PC194	PC195	PC196	PC197	PC198	PC199	PC200	PC201	PC202	PC203	PC204	PC205	PC206	PC207	PC208	PC209	PC210	PC211	PC212	PC213	PC214	PC215	PC216	PC217	PC218	PC219	PC220	PC221	PC222	PC223	PC224	PC225	PC226	PC227	PC228	PC229	PC230	PC231	PC232	PC233	PC234	PC235	PC236	PC237	PC238	PC239	PC240	PC241	PC242	PC243	PC244	PC245	PC246	PC247	PC248	PC249	PC250	PC251	PC252	PC253	PC254	PC255	PC256	PC257	PC258	PC259	PC260	PC261	PC262	PC263	PC264	PC265	PC266	PC267	PC268	PC269	PC270	PC271	PC272	PC273	PC274	PC275	PC276	PC277	PC278	PC279	PC280	PC281	PC282	PC283	PC284	PC285	PC286	PC287	PC288	PC289	PC290	PC291	PC292	PC293	PC294	PC295	PC296	PC297	PC298	PC299	PC300	PC301	PC302	PC303	PC304	PC305	PC306	PC307	PC308	PC309	PC310	PC311	PC312	PC313	PC314	PC315	PC316	PC317	PC318	PC319	PC320	PC321	PC322	PC323	PC324	PC325	PC326	PC327	PC328	PC329	PC330	PC331	PC332	PC333	PC334	PC335	PC336	PC337	PC338	PC339	PC340	PC341	PC342	PC343	PC344	PC345	PC346	PC347	PC348	PC349	PC350	PC351	PC352	PC353	PC354	PC355	PC356	PC357	PC358	PC359	PC360	PC361	PC362	PC363	PC364	PC365	PC366	PC367	PC368	PC369	PC370	PC371	PC372	PC373	PC374	PC375	PC376	PC377	PC378	PC379	PC380	PC381	PC382	PC383	PC384	PC385	PC386	PC387	PC388	PC389	PC390	PC391	PC392	PC393	PC394	PC395	PC396	PC397	PC398	PC399	PC400	PC401	PC402	PC403	PC404	PC405	PC406	PC407	PC408	PC409	PC410	PC411	PC412	PC413	PC414	PC415	PC416	PC417	PC418	PC419	PC420	PC421	PC422	PC423	PC424	PC425	PC426	PC427	PC428	PC429	PC430	PC431	PC432	PC433	PC434	PC435	PC436	PC437	PC438	PC439	PC440	PC441	PC442	PC443	PC444	PC445	PC446	PC447	PC448	PC449	PC450	PC451	PC452	PC453	PC454	PC455	PC456	PC457	PC458	PC459	PC460	PC461	PC462	PC463	PC464	PC465	PC466	PC467	PC468	PC469	PC470	PC471	PC472	PC473	PC474	PC475	PC476	PC477	PC478	PC479	PC480	PC481	PC482	PC483	PC484	PC485	PC486	PC487	PC488	PC489	PC490	PC491	PC492	PC493	PC494	PC495	PC496	PC497	PC498	PC499	PC500	PC501	PC502	PC503	PC504	PC505	PC506	PC507	PC508	PC509	PC510	PC511	PC512	PC513	PC514	PC515	PC516	PC517	PC518	PC519	PC520	PC521	PC522	PC523	PC524	PC525	PC526	PC527	PC528	PC529	PC530	PC531	PC532	PC533	PC534	PC535	PC536	PC537	PC538	PC539	PC540	PC541	PC542	PC543	PC544	PC545	PC546	PC547	PC548	PC549	PC550	PC551	PC552	PC553	PC554	PC555	PC556	PC557	PC558	PC559	PC560	PC561	PC562	PC563	PC564	PC565	PC566	PC567	PC568	PC569	PC570	PC571	PC572	PC573	PC574	PC575	PC576	PC577	PC578	PC579	PC580	PC581	PC582	PC583	PC584	PC585	PC586	PC587	PC588	PC589	PC590	PC591	PC592	PC593	PC594	PC595	PC596	PC597	PC598	PC599	PC600	PC601	PC602	PC603	PC604	PC605	PC606	PC607	PC608	PC609	PC610	PC611	PC612	PC613	PC614	PC615	PC616	PC617	PC618	PC619	PC620	PC621	PC622	PC623	PC624	PC625	PC626	PC627	PC628	PC629	PC630	PC631	PC632	PC633	PC634	PC635	PC636	PC637	PC638	PC639	PC640	PC641	PC642	PC643	PC644	PC645	PC646	PC647	PC648	PC649	PC650	PC651	PC652	PC653	PC654	PC655	PC656	PC657	PC658	PC659	PC660	PC661	PC662	PC663	PC664	PC665	PC666	PC667	PC668	PC669	PC670	PC671	PC672	PC673	PC674	PC675	PC676	PC677	PC678	PC679	PC680	PC681	PC682	PC683	PC684	PC685	PC686	PC687	PC688	PC689	PC690	PC691	PC692	PC693	PC694	PC695	PC696	PC697	PC698	PC699	PC700	PC701	PC702	PC703	PC704	PC705	PC706	PC707	PC708	PC709	PC710	PC711	PC712	PC713	PC714	PC715	PC716	PC717	PC718	PC719	PC720	PC721	PC722	PC723	PC724	PC725	PC726	PC727	PC728	PC729	PC730	PC731	PC732	PC733	PC734	PC735	PC736	PC737	PC738	PC739	PC740	PC741	PC742	PC743	PC744	PC745	PC746	PC747	PC748	PC749	PC750	PC751	PC752	PC753	PC754	PC755	PC756	PC757	PC758	PC759	PC760	PC761	PC762	PC763	PC764	PC765	PC766	PC767	PC768	PC769	PC770	PC771	PC772	PC773	PC774	PC775	PC776	PC777	PC778	PC779	PC780	PC781	PC782	PC783	PC784	PC785	PC786	PC787	PC788	PC789	PC790	PC791	PC792	PC793	PC794	PC795	PC796	PC797	PC798	PC799	PC800	PC801	PC802	PC803	PC804	PC805	PC806	PC807	PC808	PC809	PC810	PC811	PC812	PC813	PC814	PC815	PC816	PC817	PC818	PC819	PC820	PC821	PC822	PC823	PC824	PC825	PC826	PC827	PC828	PC829	PC830	PC831	PC832	PC833	PC834	PC835	PC836	PC837	PC838	PC839	PC840	PC841	PC842	PC843	PC844	PC845	PC846	PC847	PC848	PC849	PC850	PC851	PC852	PC853	PC854	PC855	PC856	PC857	PC858	PC859	PC860	PC861	PC862	PC863	PC864	PC865	PC866	PC867	PC868	PC869	PC870	PC871	PC872	PC873	PC874	PC875	PC876	PC877	PC878	PC879	PC880	PC881	PC882	PC883	PC884	PC885	PC886	PC887	PC888	PC889	PC890	PC891	PC892	PC893	PC894	PC895	PC896	PC897	PC898	PC899	PC900	PC901	PC902	PC903	PC904	PC905	PC906	PC907	PC908	PC909	PC910	PC911	PC912	PC913	PC914	PC915	PC916	PC917	PC918	PC919	PC920	PC921	PC922	PC923	PC924	PC925	PC926	PC927	PC928	PC929	PC930	PC931	PC932	PC933	PC934	PC935	PC936	PC937	PC938	PC939	PC940	PC941	PC942	PC943	PC944	PC945	PC946	PC947	PC948	PC949	PC950	PC951	PC952	PC953	PC954	PC955	PC956	PC957	PC958	PC959	PC960	PC961	PC962	PC963	PC964	PC965	PC966	PC967	PC968	PC969	PC970	PC971	PC972	PC973	PC974	PC975	PC976	PC977	PC978	PC979	PC980	PC981	PC982	PC983	PC984	PC985	PC986	PC987	PC988	PC989	PC990	PC991	PC992	PC993	PC994	PC995	PC996	PC997	PC998	PC999	PC1000
min	max	25	50	75	100	140	80	1.4	23	15	21	20	33	28	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																															

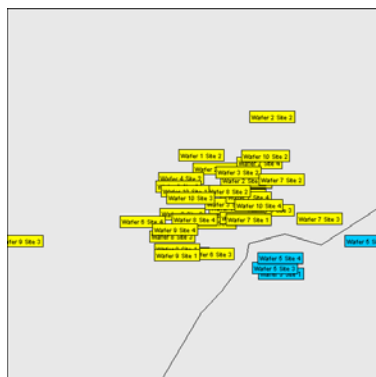
- The subcluster was found for the assumption of Wafer 5 abnormality
- The parameters must be checked for physical plausibility with regard to observed chip behavior
- The visualization & analysis result is imposed on the original database in Excel
- In the given case, the detected abnormality in the selected parameters cannot be held responsible
- Probably a design error, not a manufacturing problem !

© Andreas König Slide5-29

## Sensor Signal Processing

### Dimensionality Reduction

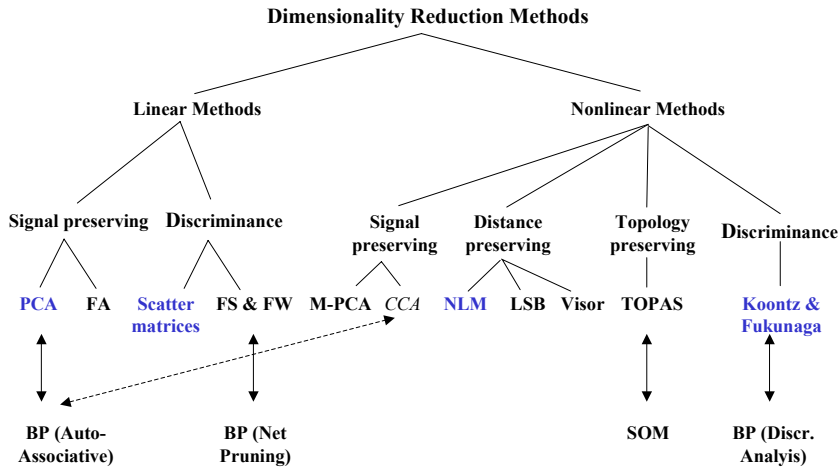
- This simple example only scratched the surface of the application potential
- Feasibility demonstration of DR methodology for microelectronic data



- In this case found parameters distinguish Wafer 5 but are not relevant to observed problem
- The outlined approach offers fast, efficient, and transparent access to multivariate, complex data met in today's manufacturing processes
- Other nonobvious information can be found & employed for process optimization & centering
- Circuit design & design centering is another potential application field

© Andreas König Slide5-30

- Taxonomy of Dimensionality Reduction Methods:



- Dimensionality reduction by **Principal Component Analysis**
- Dimensionality reduction method from communication science
- Objective: Find a linear transformation with  $\bar{y}_i = W_{1..m} \bar{v}_i$  and  $\hat{\bar{v}}_i = W_{1..m}^T \bar{y}_i$  so that the reconstruction error is minimum for given m
- Visualization employs the first two (three) principal components
- Compute covariance matrix of the database:

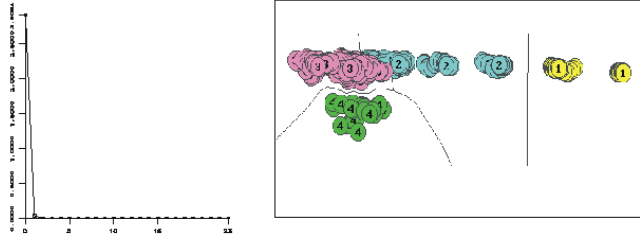
$$K = \frac{1}{N} \sum_{i=1}^N (\bar{v}_i - \bar{\mu})(\bar{v}_i - \bar{\mu})^T \quad \text{with} \quad \bar{\mu} = \frac{1}{N} \sum_{i=1}^N \bar{v}_i \quad (5.11)$$

- Compute Eigenvalues  $\lambda_i$  and Eigenvectors  $\psi_i$  of matrix K
- Data is decorrelated by applying  $\bar{y}_i = \Psi \bar{v}_i$  (5.12)
- Largest Eigenvalue corresponds to component with largest variance in the data
- Eigenvalues are sorted and the  $m$  Eigenvectors corresponding to the largest Eigenvalues are selected for projection:  $\bar{y}_{i..m} = \Psi_{1..m} \bar{v}_i$  (5.13)



## Feature Extraction

- For compression in channel coding or classification,  $m$  can be chosen according to Scree-plot or achieved classification rate  $\mathbf{R}$
- PCA is a linear, signal-preserving approach based on the assumption of unimodal Gaussian data. Example for Mech<sub>1</sub> data:

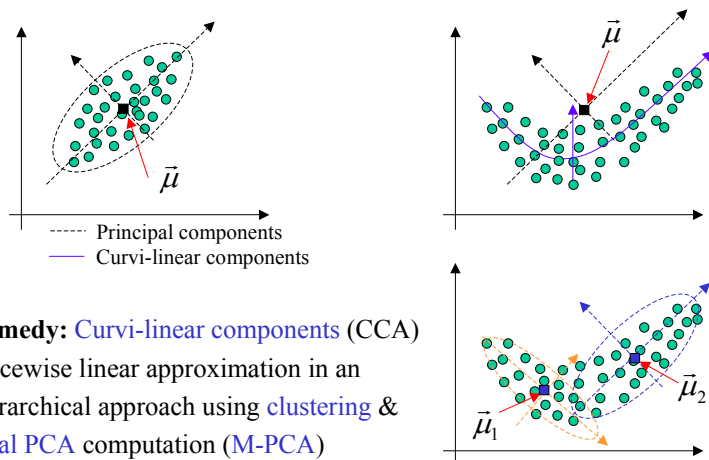


- Most variance is embodied in the first two principal components
- $\mathbf{R}$  grows with  $M^2$  for M-D data (Numerical & accuracy problems)
- Bad visualization for high *intrinsic dimensionality* or nonlinear components

© Andreas König Slide5-33

## Feature Extraction

- For *nonlinear data*, the assumptions of PCA do not hold:



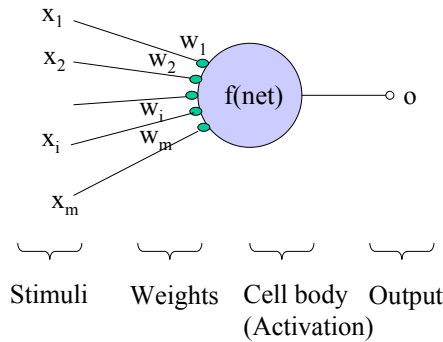
- **Remedy:** Curvi-linear components (CCA)
- Piecewise linear approximation in an hierarchical approach using *clustering* & *local PCA* computation (M-PCA)

© Andreas König Slide5-34

## Feature Extraction

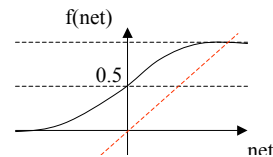
## Sensor Signal Processing Dimensionality Reduction

- Artificial neural networks are commonly applied for function approximation in classification and dimensionality reduction
- Simple artificial neuron model:



$$o = f\left(\sum_{i=1}^m x_i w_i\right) \quad (5.14)$$

$$f(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \quad (5.15)$$



- Coarse abstraction of the natural nerve cell in biological systems

© Andreas König Slide5-35

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Adaptation of a neuron weight is commonly achieved by gradient descent based on an error function:

$$E = \frac{1}{2} \sum_{k=1}^N \left( y^k - f\left(\sum_{i=1}^m x_i^k w_i\right) \right)^2 \quad (5.16)$$

- Every weight is adapted after (random) initialization according to:

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} \quad (5.17)$$

- The gradient is computed as:

$$\frac{\partial E}{\partial w_i} = -\sum_{k=1}^N \left( y^k - f\left(\sum_{i=1}^m x_i^k w_i\right) \right) \cdot f'\left(\sum_{i=1}^m x_i^k w_i\right) \cdot x_i^k \quad (5.18)$$

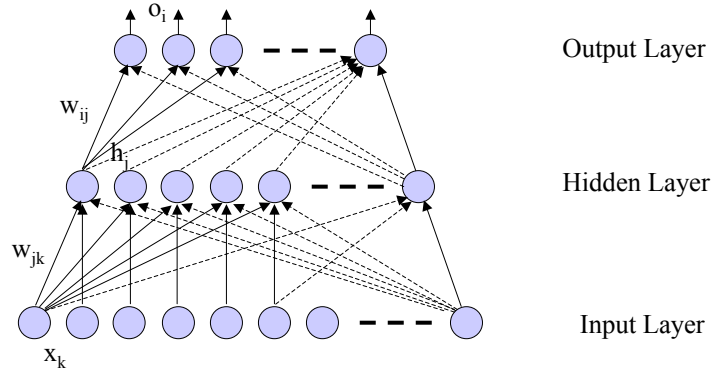
- Inserting (5.18) in (5.17) return the **batch learning rule**, reducing the batch to one returns **on-line learning rule** with immediate weight adaptation

© Andreas König Slide5-36

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Commonly, a multi-layered network of multiple of such neuron models is applied, denoted as multi-layer feedforward neural network :



- The network can be extended by more hidden layers, but already the given topology is proven to be a **universal function approximator**
- A **learning rule** is required for this network, in particular the hidden layer

© Andreas König Slide5-37

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Introducing the following abbreviations using the notation given with the network structure:

$$o_i = f\left(\sum_j h_j w_{ij}\right) ; h_j = f\left(\sum_k x_k w_{jk}\right) \quad (5.19)$$

- With (5.19) the error can be expressed as:

$$E = \frac{1}{2} \sum_{\mu=1}^N \sum_{i=1}^L \left( y_i^{\mu} - f\left(\sum_{j=1}^m h_j^{\mu} w_{ij}\right) \right)^2 \quad (5.20)$$

- Every weight is adapted after (random) initialization according to:

$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}} \quad (5.21)$$

- The gradient for the output layer weights is computed as:

$$\frac{\partial E}{\partial w_{ij}} = -\sum_{\mu=1}^N \left( y_i^{\mu} - f\left(\sum_j h_j^{\mu} w_{ij}\right) \right) \cdot f'\left(\sum_j h_j^{\mu} w_{ij}\right) \cdot h_j^{\mu} \quad (5.22)$$

© Andreas König Slide5-38

- This can be expressed employing the abbreviations of (5.19) as

$$\frac{\partial E}{\partial w_{ij}} = -\sum_{\mu=1}^N (y_i^{\mu} - o_i^{\mu}) \cdot o_i^{\mu} \cdot h_j^{\mu} \quad (5.23)$$

- Inserting in (5.21) gives the output layer batch adaptation rule

$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}} = \eta \sum_{\mu=1}^N (y_i^{\mu} - o_i^{\mu}) \cdot o_i^{\mu} \cdot h_j^{\mu} \quad (5.24)$$

- For the hidden layer adaptation rule, the error function must be expanded:

$$E = \frac{1}{2} \sum_{\mu=1}^N \sum_{i=1}^L \left( y_i^{\mu} - f \left( \sum_j w_{ij} f \left( \underbrace{\sum_k x_k^{\mu} w_{jk}}_{\text{net}_j} \right) \right) \right)^2 \quad (5.25)$$

$\underbrace{\hspace{10em}}_{\text{net}_i}$   
 $\underbrace{\hspace{10em}}_{o_i}$   
 $\underbrace{\hspace{10em}}_{e_i}$

- Every hidden weight is adapted after (random) initialization according to:

$$\Delta w_{jk} = -\eta \frac{\partial E}{\partial w_{jk}} \quad (5.26)$$

- The gradient for the hidden layer weights is computed by application of the chain rule:

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial e_i} \cdot \frac{\partial e_i}{\partial o_i} \cdot \frac{\partial o_i}{\partial \text{net}_i} \cdot \frac{\partial \text{net}_i}{\partial h_j} \cdot \frac{\partial h_j}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{jk}} \quad (5.27)$$

$$\frac{\partial E}{\partial w_{jk}} = -\sum_{\mu=1}^N \sum_{i=1}^L \left( y_i^{\mu} - f \left( \sum_j h_j^{\mu} w_{ij} \right) \right) \cdot f' \left( \sum_j h_j^{\mu} w_{ij} \right) \cdot w_{ij} \cdot f' \left( \sum_k x_k^{\mu} w_{jk} \right) \cdot x_k^{\mu} \quad (5.28)$$

- This can again be expressed employing the abbreviations of (5.19) as

$$\frac{\partial E}{\partial w_{jk}} = -\sum_{\mu=1}^N \sum_{i=1}^L (y_i^{\mu} - o_i^{\mu}) \cdot o_i^{\mu} \cdot w_{ij} \cdot h_j^{\mu} \cdot x_k^{\mu} \quad (5.29)$$

## Feature Extraction

- Introduction of error or  $\delta$ -terms with

$$\delta_i^\mu = (y_i^\mu - o_i^\mu) \cdot o_i^\mu \quad (5.30)$$

$$\delta_j^\mu = h_j^\mu \cdot \sum_{i=1}^L w_{ij} \delta_i^\mu \quad (5.31)$$

- ... allows a compact representation of the adaptation rules

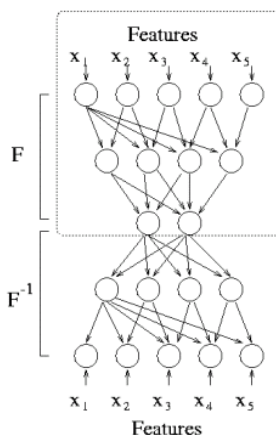
$$\Delta w_{ij} = \eta \sum_{\mu=1}^N \delta_i^\mu \cdot h_j^\mu \quad (5.32)$$

$$\Delta w_{jk} = \eta \sum_{\mu=1}^N \delta_j^\mu \cdot x_k^\mu \quad (5.33)$$

- This learning rule is denoted as **error-backpropagation learning rule**
- Numerous variants of this **vanilla approach** are in existence to improve learning behavior, e.g., introduction of a **momentum term** or **adaptive  $\eta$**

## Feature Extraction

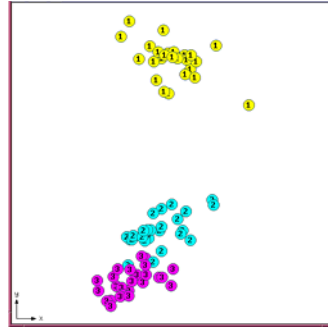
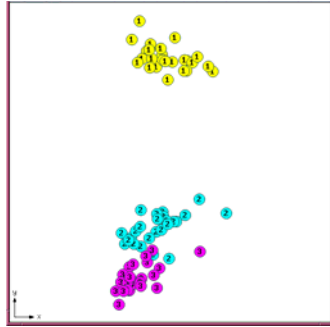
- **Nonlinear Signal Preserving Mappings**
- Application of backpropagation networks in autoassociative mode:



- Network learning tries to preserve the feature values over the bottleneck layer
- 3L networks perform similar to PCA
- Number of bottleneck layer neurons define projection dimension  $d$
- Example with 4-2-4 BP network
- 5L networks performs nonlinear compression, extracting principal curves
- Topology must be specified by user
- Hard to train and to interpret
- Example with 4-9-2-9-4 BP network

**Feature Extraction**

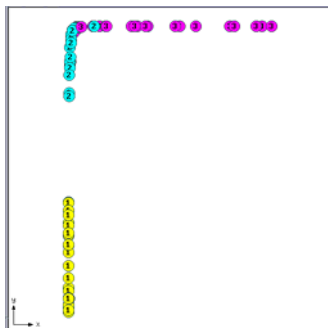
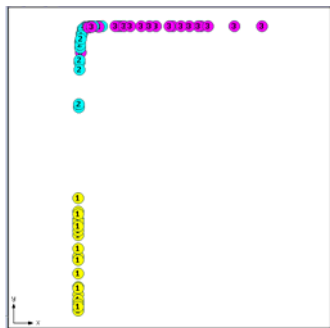
- Visualization of *Iris*train data by 3L-BP in autoassociative mode:



- Visualization of *Iris*test data by 3L-BP in autoassociative mode:

**Feature Extraction**

- Visualization of *Iris*train data by 5L-BP in autoassociative mode:



- Visualization of *Iris*test data by 3L-BP in autoassociative mode:

- Distance Preserving Nonlinear Mapping: **Sammon's NLM**
- Particular case of Multi-dimensional-scaling (MDS) approach
- Sammon's stress assesses distortion of distances while mapping the data vectors from **X** space to corresponding pivot points in **Y** space by his NLM:

$$E(m) = \frac{1}{c} \sum_{j=1}^N \sum_{i=1}^j \frac{(d_{X_{ij}} - d_{Y_{ij}}(m))^2}{d_{X_{ij}}} \quad \text{with} \quad c = \sum_{j=1}^N \sum_{i=1}^j d_{X_{ij}} \quad (5.34)$$

- Here,  $d_{X_{ij}}$  and  $d_{Y_{ij}}$  denote the interpoint distance in feature space **X** and projection space **Y**, respectively:

$$d_{Y_{ij}}(m) = \sqrt{\sum_{q=1}^d (y_{iq}(m) - y_{jq}(m))^2} \quad (5.35)$$

$$d_{X_{ij}} = \sqrt{\sum_{q=1}^M (x_{iq} - x_{jq})^2} \quad (5.36)$$

- Minimization of Sammon's stress using gradient descent is achieved by iterative computation of new coordinates for the pivot vectors in **Y** space:

$$y_{iq}(m+1) = y_{iq}(m) - MF * \Delta y_{iq}(m) \quad (5.37)$$

$$\text{with} \quad \Delta y_{iq}(m) = \frac{\partial E(m)}{\partial y_{iq}(m)} \bigg/ \left| \frac{\partial^2 E(m)}{\partial y_{iq}(m)^2} \right| \quad (5.38)$$

- The partial derivatives for gradient descent are given by

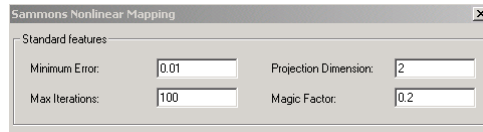
$$\frac{\partial E(m)}{\partial y_{iq}} = \frac{-2}{c} \sum_{\substack{j=1 \\ j \neq i}}^N \left[ \frac{1}{d_{Y_{ij}}} - \frac{1}{d_{X_{ij}}} \right] (y_{iq} - y_{jq}) \quad (5.39)$$

$$\frac{\partial^2 E(m)}{\partial y_{iq}^2} = \frac{-2}{c} \sum_{\substack{j=1 \\ j \neq i}}^N \left[ \frac{1}{d_{Y_{ij}}} - \frac{1}{d_{X_{ij}}} - \frac{(y_{iq} - y_{jq})^2}{d_{Y_{ij}}^3} \right] \quad (5.40)$$

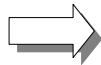
## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- NLM maps each point and preserves distances & structure in data
- Salient visualization method for direct data projection & visual analysis
- Mapping control:



- Error  $E(m) > 0.1$  indicates unacceptable mapping result [Sammon 69]
- Computation of NLM has  $O(N^2)$  complexity
- Method becomes infeasible for large databases
- Gradient descent optimization is not capable to exactly preserve all distances
- Inevitable mapping error is introduced



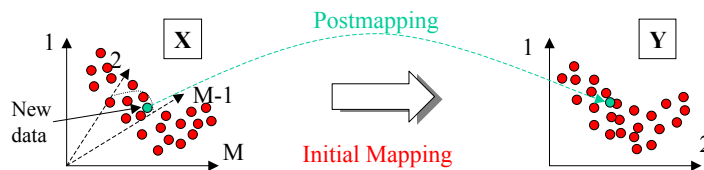
**Acceleration for approximated mapping by heuristics !**

© Andreas König Slide5-47

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- After initial dimensionality reducing mapping computation for a given data set, single or groups of data points shall be subject to the same mapping, too
- In **classification**, test & application data shall be mapped in **recall**
- In **visualization**, **new data** points shall be **mapped & displayed**



- In both cases, the mapping of additional points shall take place without recomputation of the entire mapping
- PCA & ANN's satisfy this requirement, NLM originally does not
- **Remedy 1:** Train an ANN with NLM mapping data & use it for recall
- **Remedy 2:** Introduce NLM recall mapping (NLMR)

© Andreas König Slide5-48



## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Introduction of a NLM recall method (NLMR) basing on the previously computed NLM
- New data vectors  $\vec{v}_i^r$  subject to recall are placed according to their distances
- $d_{\hat{x}_{ij}}$  to the training data points  $\vec{v}_j^t$
- Mutual distances between test vectors are neglected in NLMR
- Thus the cost function is modified to:

$$\hat{E}_i(m) = \frac{1}{\hat{c}} \sum_{j=1}^K \left( \frac{(d_{\hat{x}_{ij}} - d_{\hat{y}_{ij}}(m))^2}{d_{\hat{x}_{ij}}} \right) \quad (5.41)$$

$$\text{with } d_{\hat{x}_{ij}} = \sqrt{\sum_{q=1}^M (\vec{v}_{iq}^r - \vec{v}_{jq}^t)^2} \quad \text{and} \quad \hat{c} = \sum_{j=1}^K d_{\hat{x}_{ij}}$$

- K corresponds to the number of initially mapped training samples

© Andreas König Slide5-49

## Feature Extraction

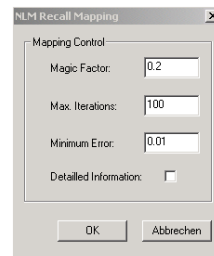
## Sensor Signal Processing Dimensionality Reduction

- Each point is randomly initialized in  $\mathbf{Y}$  space
- Iterative adjustments by gradient descent take place according to:

$$\hat{y}_{iq}(m+1) = \hat{y}_{iq}(m) - MF * \Delta \hat{y}_{iq}(m) \quad (5.42)$$

$$\text{with } \Delta \hat{y}_{iq}(m) = \frac{\partial \hat{E}_i(m)}{\partial \hat{y}_{iq}(m)} / \left| \frac{\partial^2 \hat{E}_i(m)}{\partial \hat{y}_{iq}(m)^2} \right| \quad \text{and} \quad 0 < MF \leq 1$$

- Mapping of test data for recall or postmapping is now feasible
- The number of iterations required per data vector mapping vary considerably
- Computational savings can be achieved for given Minimum Error threshold

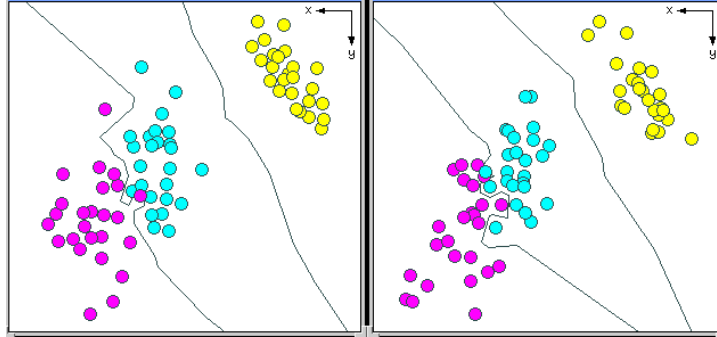


© Andreas König Slide5-50

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- NLMR Application example for *Iris*train and *Iris*test data:



- NLM(R) can now serve as nonlinear alternative to PCA. Quantitative discriminance comparison showed NLM(R) to be superior for given dimension

© Andreas König Slide5-51

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Quantitative Assessment of selected unsupervised Mapping Techniques
- Compression from M- to 2D data for linear/nonlinear unsupervised methods
- Evaluation by overlap  $q_o$  and separability  $q_s$  measures for training & test sets

Method	Dim.	Train	$q_o$	$q_s$	Test	$q_o$	$q_s$
Original	4	Iris <sub>train</sub>	0.95503	0.90666	Iris <sub>test</sub>	0.91683	0.88000
PCA	2	Iris <sub>train</sub>	0.91329	0.86667	Iris <sub>test</sub>	0.91970	0.89333
NLM	2	Iris <sub>train</sub>	0.94250	0.89333	Iris <sub>test</sub>	0.93128	0.89333
BP	2	Iris <sub>train</sub>	0.93008	0.90666	Iris <sub>test</sub>	0.91630	0.90667
Original	24	Mech <sub>1</sub>	1.000	0.98933	Mech <sub>2</sub>	0.99799	0.96308
PCA	2	Mech <sub>1</sub>	0.98959	0.97067	Mech <sub>2</sub>	0.94889	0.91384
NLM	2	Mech <sub>1</sub>	0.97898	0.94400	Mech <sub>2</sub>	0.95730	0.91384
BP	2	Mech <sub>1</sub>	0.97861	0.94400	Mech <sub>2</sub>	0.93619	0.89538

- BP has 4-9-2-9-4 network topology, 1000 epochs of quickprop learning
- PCA is simple and excellent when applicable to the given data
- NLM(R) outperforms PCA for nonlinear data & is easy to use
- BP results are very hard to obtain & approximately as those of NLM(R)
- This comparison assessed **discriminance not structure preservation** !

© Andreas König Slide5-52

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Supervised linear technique for dimensionality reduction [Fukunaga 90]
- **Objective:** Find a linear transformation that minimizes the within or intraclass distances between data points and maximizes the between or interclass distances, i.e., class regions shall become compact & well separated
- Scatter matrices are computed for this aim. In the parametric approach interclass  $S_b$  and intraclass  $S_w$  scatter matrices are computed by:

$$S_w = \frac{1}{N} \sum_{i=1}^L \sum_{j=1}^{N_i} (\vec{v}_j^{\omega_i} - \vec{\mu}_j^{\omega_i})(\vec{v}_j^{\omega_i} - \vec{\mu}_j^{\omega_i})^T \quad \text{with} \quad \vec{\mu}^{\omega_i} = \frac{1}{N_i} \sum_{j=1}^{N_i} \vec{v}_j^{\omega_i} \quad (5.43)$$

$$S_b = \frac{1}{N} \sum_{i=1}^L \frac{N_i}{N} (\vec{\mu}_j^{\omega_i} - \vec{\mu})(\vec{\mu}_j^{\omega_i} - \vec{\mu})^T \quad \text{with} \quad \vec{\mu} = \frac{1}{N} \sum_{j=1}^N \vec{v}_j = \sum_{j=1}^N \frac{N_i}{N} \vec{\mu}^{\omega_i} \quad (5.44)$$

© Andreas König Slide5-53

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- One possible measure of the underlying separability is given by  $J_s$  with:

$$J_s = \text{Tr}(S_w^{-1} S_b) \quad (5.45)$$

- Computation of Eigenvalues  $\lambda_i$  and Eigenvectors  $\psi_i$  of matrix  $(S_w^{-1} S_b)$
- Largest Eigenvalues correspond to components with largest **discriminance**
- Eigenvalues are sorted and the  $m$  Eigenvectors corresponding to the largest Eigenvalues are selected for dimensionality reduction:

$$\vec{y}_{i..m} = \Psi_{1..m} \vec{v}_i \quad (5.46)$$

- Dimensionality is reduced and discriminance improved in one step
- For classification,  $m$  is chosen according to Scree plot of  $\lambda_i$  or one of the measures given in the taxonomy of section 2
- Visualization by 2 (3)  $\psi_i$ , for  $m > 2$  combination with unsupervised methods
- The parametric approach is limited in scope and ability !

© Andreas König Slide5-54

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Extension of scatter matrices for nonparametric or even multimodal data [Fukunaga 90]
- Interclass scatter is computed based on k-nearest-neighbor technique for nonparametric scatter matrix  $\mathbf{S}_b$  by:

$$S_b = \frac{1}{N} \sum_{i=1}^L \sum_{j=1}^{N_i} g_i (\vec{v}_j^{\omega_i} - \vec{\mu}_{jNN}^{\omega_i}) (\vec{v}_j^{\omega_i} - \vec{\mu}_{jNN}^{\omega_i})^T \quad (5.47)$$

- Here,  $\vec{\mu}_{jNN}^{\omega_i} = \frac{1}{k} \sum_{l=1}^k \vec{v}_{lNN}^{\omega_i}$  denotes the mean vector of the k-nearest-neighbors of  $\vec{v}_j^{\omega_i}$  with different class affiliations
- The weighting factor for off-class-borders vectors  $g_j$  is given by:

$$g_j = \frac{\min\{d^\zeta(v_j^{\omega_i}, \mu_{jNN}^{\omega_i}), d^\zeta(v_j^{\omega_i}, \mu_{jNN}^{\omega_i})\}}{d^\zeta(v_j^{\omega_i}, \mu_{jNN}^{\omega_i}) + d^\zeta(v_j^{\omega_i}, \mu_{jNN}^{\omega_i})} \quad \text{with weight decay factor} \quad (5.48)$$

$$\zeta = [0, \infty[$$

© Andreas König Slide5-55

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- The parametric intraclass scatter matrix  $S_w$  is retained
- Decorrelation with regard to  $S_w$ , i.e.,  $S_w$  in space  $\mathbf{Y}$  will be  $S_w = \mathbf{I}$
- Pattern data is transformed by

$$\vec{y}_j = (\Psi_{1..M}^w (\Lambda^w)^{-\frac{1}{2}})^T \vec{v}_j \quad (5.49)$$

- In space  $\mathbf{Y}$ , now  $Tr(S_w^{-1} \mathbf{S}_b) = Tr(\mathbf{S}_b)$  holds (whitening according to  $S_w$ )
- Computation of Eigenvalues  $\lambda_i$  and Eigenvectors  $\psi_i$  of matrix  $\mathbf{S}_b$
- Eigenvalues are sorted and the  $m$  Eigenvectors corresponding to the largest Eigenvalues are selected for dimensionality reduction:

$$\vec{y}_{j1..m}^* = (\Psi_{1..m}^b)^T \vec{y}_j \quad (5.50)$$

- Finally, the dimensionality reducing projection is computed by:

$$\vec{y}_{j1..m}^* = (\Psi_{1..m}^b)^T \vec{y}_j = (\Psi_{1..m}^b)^T (\Psi_{1..M}^w (\Lambda^w)^{-\frac{1}{2}})^T \vec{v}_j \quad (5.51)$$

© Andreas König Slide5-56

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Visualization example of nonparametric scatter matrix (NPSCM) application for *Mech<sub>1</sub>* data:



- Numerical assessment of the method with regard to discriminance will follow after presentation of nonlinear supervised mapping methods

© Andreas König Slide5-57

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Koontz & Fukunaga extended Sammon's stress in the NLM by an additional term dedicated to assessment of intraclass distances
- The KFM mapping pursues intraclass reduction & structure preservation and thus nonlinear discriminance analysis by the following cost function:

$$E(m) = \frac{1}{c} \sum_{j=1}^N \sum_{i=1}^j \frac{\delta(\omega_i, \omega_j) d_{y_{ij}}^2(m) + \hat{\lambda} (d_{x_{ij}} - d_{y_{ij}}(m))^2}{d_{x_{ij}}} \quad (5.52)$$

- Here,  $c$ ,  $d_{x_{ij}}$  and  $d_{y_{ij}}$  are the same terms defined for Sammon's stress in the NLM and

$$\delta(\omega_i, \omega_j) = \begin{cases} \hat{\alpha} & : \omega_i = \omega_j \\ 0 & : \omega_i \neq \omega_j \end{cases} \quad \text{with} \quad \hat{\alpha} = 1 \quad (5.53)$$

© Andreas König Slide5-58

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Minimization of KFM cost function using gradient descent is achieved by iterative computation of new coordinates for the pivot vectors in  $\mathbf{Y}$  space:

$$y_{iq}(m+1) = y_{iq}(m) - MF * \Delta y_{iq}(m) \quad (5.54)$$

$$\text{with} \quad \Delta y_{iq}(m) = \frac{\partial E(m)}{\partial y_{iq}(m)} / \left| \frac{\partial^2 E(m)}{\partial y_{iq}(m)^2} \right|$$

- The partial derivatives for gradient descent are given by

$$\frac{\partial E(m)}{\partial y_{iq}} = \frac{2}{c} \sum_{\substack{j=1 \\ j \neq i}}^N \left[ \frac{\delta(\omega_i, \omega_j)}{d_{X_{ij}}} - \hat{\lambda} \left( \frac{1}{d_{Y_{ij}}} - \frac{1}{d_{X_{ij}}} \right) \right] (y_{iq} - y_{jq}) \quad (5.55)$$

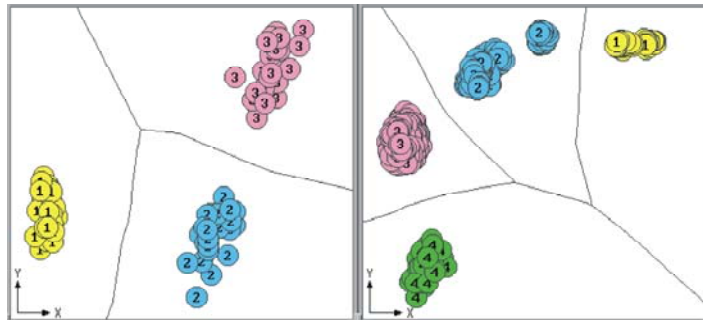
$$\frac{\partial^2 E(m)}{\partial y_{iq}^2} = \frac{2}{c} \sum_{\substack{j=1 \\ j \neq i}}^N \left[ \frac{\delta(\omega_i, \omega_j)}{d_{X_{ij}}} - \hat{\lambda} \left( \frac{1}{d_{Y_{ij}}} - \frac{1}{d_{X_{ij}}} - \frac{(y_{iq} - y_{jq})^2}{d_{Y_{ij}}^3} \right) \right] \quad (5.56)$$

© Andreas König Slide5-59

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Modification of the original approach by setting  $\hat{\alpha} = 1 - \hat{\lambda}$  reduces parameters and gives convenient control from pure structure preservation ( $\hat{\lambda} = 0.0$ ) to pure separability achievement ( $\hat{\lambda} = 1.0$ )
- Visualization of **Iris**train and **Mech**<sub>1</sub> with  $\hat{\lambda} = 0.3$  after 100 iterations

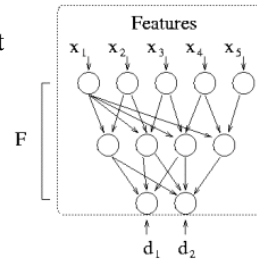


© Andreas König Slide5-60

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- As for NLM, a recall procedure is not straight forward available for KFM
- Koontz & Fukunaga presented a (restricted) method using polynomial distance approximation and a pivot point approach for mapping
- Employment of neural network, e.g., BP, as universal function approximator is an interesting alternative
- Network is supplied with original feature data at the input & KFM data at the output
- Training phase aims on mapping error minimization with generalization
- Test data is mapped later for classification or visualization employing the trained net
- **Problem:** Net configuration & convergence

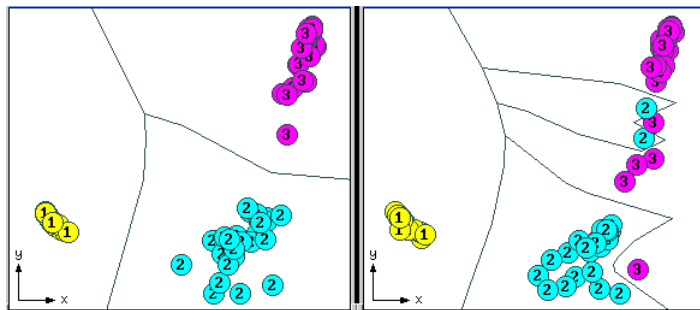


© Andreas König Slide5-61

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Visualization of KFM applied to *Iris*train and *Iris*test data with a 4-9-5-2 BP network (training & test data mapped by BP network !)



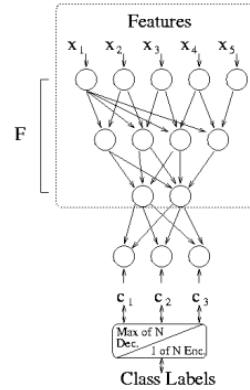
- Achieving network convergence & good approximation properties is not straightforward for combined KFM/BP nonlinear mapping

© Andreas König Slide5-62

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Employment of BP neural network with a bottleneck topology for nonlinear mapping or nonlinear discriminant analysis (NDA) is an interesting alternative:
- Network is supplied with original feature data at the input & class data at the output
- Training phase aims on classification error minimization with given bottleneck
- The number of neurons in the bottleneck layer determine projection dimension  $d$
- The trained net is cut after bottleneck
- Test data is mapped later for classification or visualization employing the cut net
- **Problem:** Net configuration & convergence

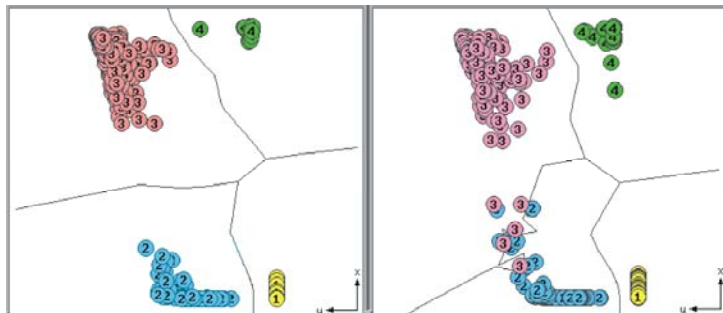


© Andreas König Slide5-63

## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Result of BP\_NDA for  $Mech_1$  and  $Mech_2$  data employing a 24-9-2-4 BP neural network with a bottleneck topology ( $d=2$ ):



- Fast and rather easy training, only 130-150 epochs were required
- Much easier & faster than KFM/BP, but no structure preservation

© Andreas König Slide5-64



## Feature Extraction

## Sensor Signal Processing Dimensionality Reduction

- Compression from M- to 2D data for linear/nonlinear supervised methods
- Evaluation by overlap  $q_o$  and separability  $q_s$  measures for training & test sets

Method	Dim.	Train	$q_o$	$q_s$	Test	$q_o$	$q_s$
Original	4	Iristrain	0.95503	0.90666	Iristest	0.91683	0.88000
NPSCM	2	Iristrain	0.98224	0.94666	Iristest	0.95788	0.90666
BP	2	Iristrain	0.97536	1.000	Iristest	0.95295	0.94667
KFM/BP	2	Iristrain	1.000	1.000	Iristest	0.94023	0.92000
Original	24	Mech <sub>1</sub>	1.000	0.98933	Mech <sub>2</sub>	0.99799	0.96308
NPSCM	2	Mech <sub>1</sub>	0.99908	0.98200	Mech <sub>2</sub>	0.98716	0.97536
BP	2	Mech <sub>1</sub>	0.99988	0.99733	Mech <sub>2</sub>	0.99373	0.99385

- KFM with trained BP network (only for *Iris* data, hard to train for *Mech* data)
- BP with 24-9-2-4 network topology for Mech<sub>1</sub> fast to train & good results
- Linear NPSCM is simple and excellent when applicable to the given data
- Nonlinear BP, KFM/BP outperform NPSCM. Problem: overfitting/overlap
- BP & KFM/BP similar but KFM/BP results are very hard to obtain

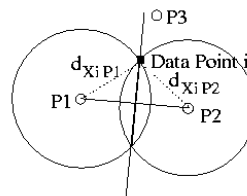
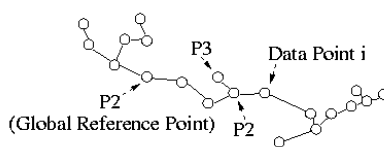
© Andreas König Slide5-65

## Accelerated Methods

## Sensor Signal Processing Dimensionality Reduction

### Lee, Slaggle & Blum's Triangulation Mapping:

- Compute *Minimal-Spanning-Tree* (MST) of the data ( $O(N^2)$  complexity !)
- Select initial data points from the MST and map to 2D space
- Step through MST & map each of the remaining points by triangulation



- Only 2N-3 distances are exactly preserved
- The algorithm tends to unfold circular structures, misleading in analysis !
- Option: Choose a fixed, global reference point for focusing the mapping (ROI)

© Andreas König Slide5-66

## Accelerated Methods

## Sensor Signal Processing Dimensionality Reduction

### Visor Mapping [König et al. 94]:

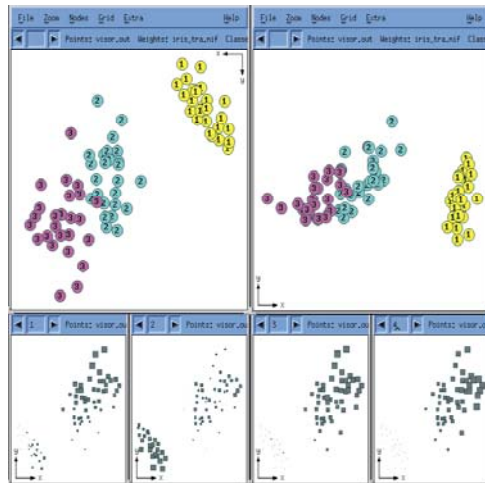
- Fast mapping with only  $O(N)$  complexity as data previewer
- Bases on the assumption to find three suited global pivot points for ensuing triangulation mapping of the remaining  $N-3$  data points
- Mapping steps (no parameters):
  1. Compute centroid of the data set
  2. Determine data point most distant from centroid
  3. Determine two more data points with maximum from centroid & each other
  4. Place these pivot points in the 2D-plane, exactly preserving their mutual distances
  5. Place the remaining  $N-3$  points by triangulation based on pivot points
- The concept is amenable to enhancement by reference or multiple pivot points

© Andreas König Slide5-67

## Accelerated Methods ,,

## Sensor Signal Processing Dimensionality Reduction

- NLM and Visor mapping results are displayed for *Irisstrain* data
- Mappings can be subject to rotation & mirroring
- Basic global data structure is quite similar
- The relevant information, e.g., for classification system design can be extracted by fast preview
- Quantitative analysis of mapping error required !



© Andreas König Slide5-68

Summary

- The chapter addressed the issue of decreasing the dimensionality of data sets to achieve **visualization** and/or improved **classification**
- A **survey** of the **basic idea** and selected **conventional** as well as **neural supervised and unsupervised methods** was provided
- In particular, potential cost functions for **assessment** and **optimization** of dimensionality reducing mappings have been regarded and applied
- Special attention was given to **automatic feature selection** (AFS)
- AFS allows the **selection** of a subset of **salient attributes** or **features** from an initially large set and the **pruning** of the architecture to a **lean system**
- Another focus was on **unsupervised methods**, e.g., **NLM**, to project arbitrary, preferably to two dimensions in the wake of **system design**
- Numerous more advanced and more recent methods can be found today
- **Application** is in **system design** and **system operation**