

Tfy-99.275 lecture 7

Pattern Recognition, Classification techniques

Problem Statement

- z In many applications we measure many variables from different subjects with the aim of building some useful input-output system.
- z We often start with whatever variables that may be expected to be useful and then try to figure out how that variable set (or part of it) can be associated with known outputs.
- z This involves issues like feature extraction, feature selection and classification performance evaluation

Feature Extraction

- z reduces the amount of raw data while maintaining (or enhancing) the vital information for the pattern recognition task at hand
- z features can range from low-level features (typically measurements, e.g., numerical values like signal amplitudes) to high-level features (e.g., 'symptom' feature for classification into 'diseases' classes)

Feature Selection

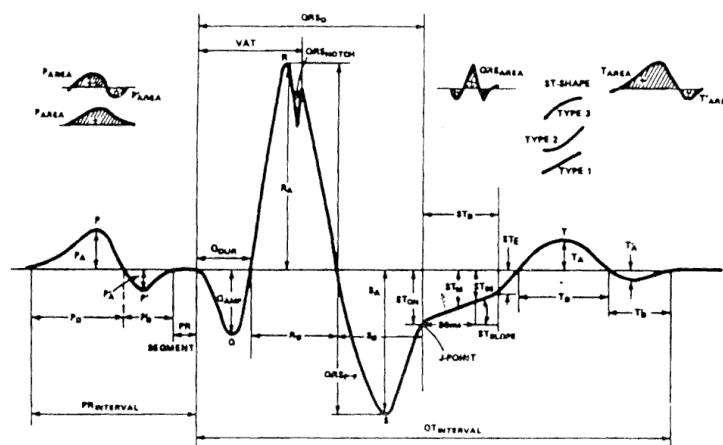
- z reduce the set of data even further; select those features that are most relevant to the pattern recognition task at hand (for example, using principal component analysis)

Some examples of features

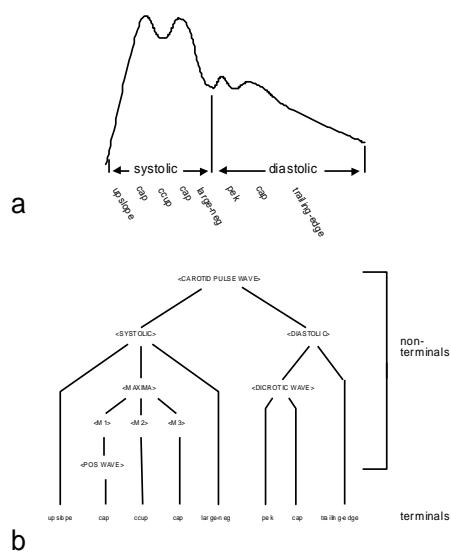
Features can be obtained from different ways of processing the signals

| measurements | results from transforms | underlying structure |
|---|--|---|
| amplitude bias phase energy moments | polynomials Fourier cosine wavelets | peaks slopes corners parametric models |

Examples of the many features that can be obtained from a 'simple' ECG curve



a structural description
of a blood pressure
waveform

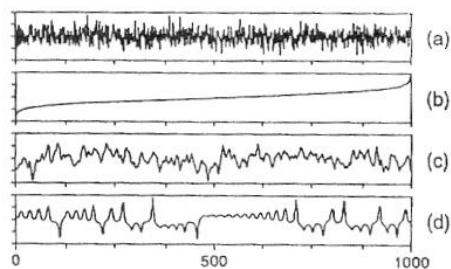


Feature Vectors

- z comprise all features used to describe a pattern; it is a reduced-dimensional representation of that pattern
- z the set of all features that could be used to describe a pattern is limited to those actually stated in the feature vector.

Watch out!

- z A feature vector usually can not replace the original signal on a one-to-one basis; keep backups of your raw data at all times!



Example of four synthesized time series with identical means, standard deviations, and ranges. Series (c) and (d) also have identical autocorrelation functions and therefore identical power spectra.

Classification (pattern recognition)

- z Identify to which class, C , an observation (pattern), x , from which the selected-feature vector was calculated, belongs.
- z Patterns that are quantitatively similar may belong to the same class, in such cases clustering techniques or matching of patterns against prototypes can be used.
- z However, also patterns that bear no obvious similarities to one another can also belong to the same class.

Pattern Recognition Techniques

- z The huge amount of pattern recognition techniques is being used in biomedical applications can be (roughly) fitted into 3 main groups of techniques

| techniques | | principle of operation |
|------------------------------|------------------|--|
| statistical | | classify upon probability distribution of characteristic measurement |
| syntactic | | use a structural description of the shape of the signals, use interrelationships between features rather than values |
| artificial intelligence (AI) | symbolic AI | use symbolic representations of class structures, manipulate information by logic statements (expert systems) |
| | connectionist AI | relate characteristic feature sets to class structures using examples (neural networks) |

Statistical Pattern Recognition Techniques

- z A set of N features extracted from a pattern can be treated as an N -dimensional vector \mathbf{x} in an N -dimensional space. Such a vector is called a *feature vector* and the N -dimensional space is called *feature space*. Usually, the feature space is \mathbf{R}^N , but it can also be a subspace of \mathbf{R}^N .
- z Classifications are based on the determination of *decision regions*. The feature space is thus divided into *decision regions* that are associated with classes.
- z The borders between decision regions are called *decision boundaries*.

Decision Regions

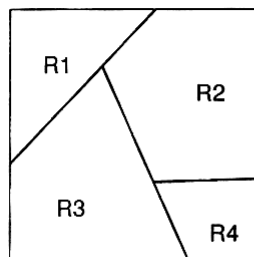
- z The classification is performed by determining the decision region in which the applied feature vector falls. The main problem in constructing a good classifier is the method of determination of the decision boundaries.

Examples of decision regions and boundaries in a 2-dimensional feature space:

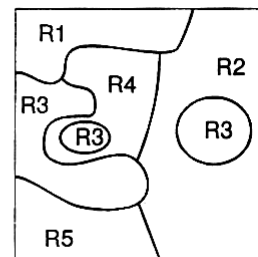
a: decision regions

separated by linear

decision boundaries;
b: a general division of
feature space into
decision regions.



a



b

Linear Separation of Decision Regions

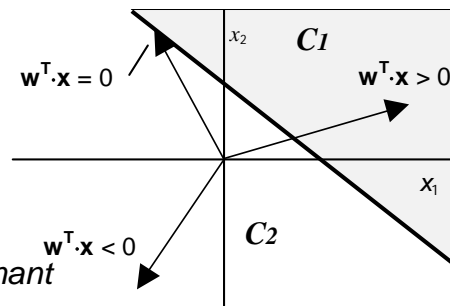
n -dimensional feature vector \mathbf{x} with x_0 set to 1 and a weight vector, \mathbf{w}

$$y = w_0 \cdot 1 + w_1 \cdot x_1 + \dots + w_n \cdot x_n$$

2-class separation problem

$$d(\mathbf{x}) = \mathbf{w}^T \cdot \mathbf{x} \begin{cases} > 0 & \text{if } \mathbf{x} \in C_1 \\ \leq 0 & \text{if } \mathbf{x} \in C_2 \end{cases}$$

For an M -class separation problem we'll need $M-1$ such *linear discriminant functions*



Decision Functions*

In the general case the decision region to which a feature vector belongs is determined by a *decision function* (which is not necessary linear). A decision function $g_i(\mathbf{x})$ is defined for each decision region i . Often, the feature vector \mathbf{x} is assigned to region m for which the relationship

$$g_m(\mathbf{x}) > g_i(\mathbf{x}) \quad \forall_{i \neq m}$$

holds.

The linear discriminant function of the previous slide is a special (but popular) choice.

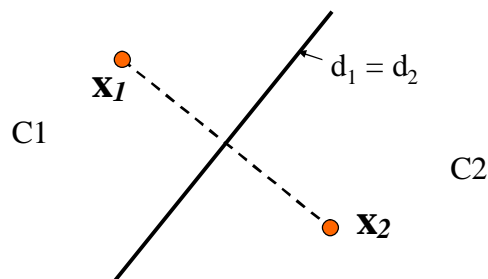
*Often, also the term 'Discriminant Functions' is used.

Minimum Distance Classifier

An example of a feature space with two decision regions their decision functions being:

$$d_1(\mathbf{x}) = |\mathbf{x} - \mathbf{x}_1| \text{ and } d_2(\mathbf{x}) = |\mathbf{x} - \mathbf{x}_2|.$$

In this case the feature vector is assigned to the decision region with the **smallest** decision function; this is a minimum distance classifier.



k-Nearest Neighbour Classifiers

Suppose we have N sample patterns $\{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_N\}$ that we know belong to one of M classes $\{\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_M\}$.

For a new, to classify, pattern \mathbf{x} , we calculate the distance $D(\mathbf{s}_i, \mathbf{x})$ and assign \mathbf{x} to the class of the sample pattern that is closest to \mathbf{x} (the 'nearest neighbour').

This is not a very stable way of classifying as the classification is made upon the basis of one single sample pattern (which may be an outlier).

A more reliable estimate may be obtained by classifying using several sample patterns; k nearest neighbours and then use a majority vote. This is the *k-nearest-neighbour (k-NN) rule*.

Estimation of Discriminant Functions

The crucial point is to create these decision regions and discriminant functions. Two possible methods;

1. convert the so-called *a priori* class probability $P(c_i)$ into an *a posteriori* (based on measured patterns) probability $P(c_i|\mathbf{x})$;
2. minimize the expected error the classifier will make.

There are many methods available to those two approaches, it is beyond the scope of this course however to provide an in-depth review of them.

Clustering Methods

- z Sometimes we don't want (or cannot) assign classes beforehand. We may then examine possible formation of inherent groups or clusters of the data in feature space and (possibly) label the clusters with class labels.
- z This may be relatively easy in 2D, but for higher-dimensional feature spaces this is a difficult task.
- z We have to group the feature vectors on the basis of similarity, dissimilarity or distance measures.
- z Examples of commonly used measures are the Euclidean distance, the Mahalanobis distance, and the normalised dot product (there exist many more)

Distance measures

- z Euclidean distance, between feature vectors \mathbf{x} and \mathbf{z}

$$D_E^2 = \|\mathbf{x} - \mathbf{z}\|^2 = (\mathbf{x} - \mathbf{z})^T (\mathbf{x} - \mathbf{z}) = \sum_{i=1}^n (x_i - z_i)^2$$

- z Normalised dot product $D_d = \frac{\mathbf{x}^T \mathbf{z}}{\|\mathbf{x}\| \|\mathbf{z}\|}$

- z Mahalanobis distance $D_M^2 = (\mathbf{x} - \mathbf{m})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{m})$

\mathbf{x} is a feature vector compared to a pattern class with \mathbf{m} as mean class vector and \mathbf{C} as covariance matrix

$$\mathbf{C} = E[(\mathbf{y} - \mathbf{m})(\mathbf{y} - \mathbf{m})^T]$$

The expectation is calculated over all vectors \mathbf{y} that belong to the class. \mathbf{C} provides covariance of all possible pairs of features. The diagonal contains variance of individual features. The matrix represents the 'scatter' of features belonging to a certain class.

Example of a clustering algorithm, the K-means algorithm

- z goal: (iteratively) minimise sum of squared distances from all points in a cluster domain to the cluster centre.
1. choose K initial cluster centres, $\mathbf{z}_1(1), \mathbf{z}_2(1), \dots, \mathbf{z}_K(1)$ (the index in parenthesis indicates iteration number)
 2. at iteration k distribute samples $\{\mathbf{x}\}$ among the K clusters using

$$\mathbf{x} \in S_j(k) \text{ if } \|\mathbf{x} - \mathbf{z}_j(k)\| < \|\mathbf{x} - \mathbf{z}_i(k)\| \quad \forall i = 1, 2, \dots, K, i \neq j$$
 3. from results of 2. calculate new cluster centres $\mathbf{z}_j(k+1)$ so that the sum of squared distances from all points in $S_j(k)$ to the new cluster centre is minimised. The $\mathbf{z}_j(k+1)$ that minimises this distance is simply the sample mean of $S_j(k)$, so

$$\mathbf{z}_j(k+1) = \frac{1}{N_j} \sum_{\mathbf{x} \in S_j(k)} \mathbf{x}, \quad j = 1, 2, \dots, K.$$
 4. if $\mathbf{z}_j(k+1) = \mathbf{z}_j(k)$ for $j=1, 2, \dots, K$ the algorithm has converged: terminate the iteration, otherwise go to step 2.

Note: one has to define the number of clusters and initial cluster centres beforehand.

Syntactic Pattern Recognition Techniques

- These techniques concentrate on the structure of the presented patterns rather than on assumed distributions of feature sets.
- Basic tenet: patterns belonging to the same class can be characterised by a common structural description.
- Descriptions are made using hierarchies of features. Low-level features are called *primitives* (e.g., arcs, corners, straight lines)
- Signals are decomposed into primitives by using *grammars*. A mechanism (the *parser*) matches the decomposed structural description against known structures and recognises the patterns.

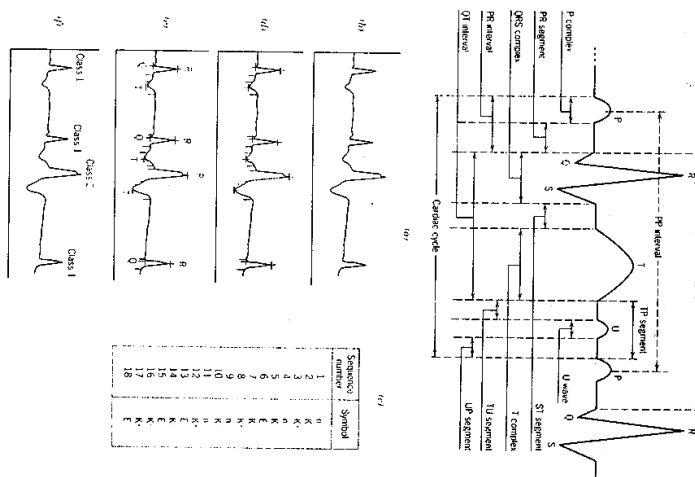


Figure 13: Use of SyntPR for ECG waveform description and classification

(from Tripathi/Skordalakis 1990).

- (a) A cardiac cycle and its constituent patterns.
- (b) Input waveform.
- (c) Extracted primitives from waveform of (b).
- (d) Graphical display of extracted primitives from (c).
- (e) Recognized ECG cycles and constituent patterns from (d).
- (f) Classification of QRS complexes from ECG cycles of (e).

Syntactic Pattern Recognition

- z The usefulness of these syntactic techniques is sometimes believed to be limited for practical applications; parsing difficulties occur especially for noisy signals.
- z There are some fields where they seem to be more effective than statistical approaches though. Example application areas include ECG and EP processing.

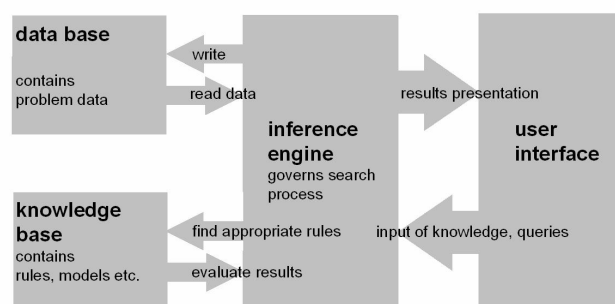
Artificial Intelligence (AI) Techniques

- z These techniques concentrate mainly on two concepts:
 - y *knowledge presentation* - the choice of a method to represent facts;
 - y *learning* - the AI-based system should be able to increase its efficiency and effectivity when a similar or identical task is processed repeatedly.

Symbolic AI Techniques (Expert Systems)

- z *symbolic representations* (or simply: symbols) of class structures are handled in such a way that the system is able to describe and identify the structures and states referred to by the symbols.
- z For example, a symbol can be the word *square* to represent a square-shaped part of an image, or the words *low temperature* to represent a temperature in the range of -5 ... +5 °C. Usually, these symbols are manipulated by using logical statements (e.g. an *if...then...else...* construction).

General Structure of an Expert System



The *inference engine* searches the *knowledge base* (which contains the domain-specific knowledge) during a problem-solving task and also evaluates the various rules in the knowledge base to add results to the *data base*. The data base contains the current problem data.

Expert Systems

- z The inference engine can use different approaches in its search process, for example *backward chaining* (start with the goal of the problem and repeatedly infer subgoals on the way to the final solution of the problem), or *forward chaining* (start with inferring the immediate consequences of the data and use consequences to come to further consequences, repeat until the final goal of the problem has been reached).
- z *Expert systems* use knowledge of the problem domain to interpret data and to provide feedback to the user of the system, i.e. they can 'explain' how they interpreted the data.
- z This domain knowledge is put into the expert system as a result of the interaction between a domain expert and a knowledge engineer: an AI specialist skilled in analysing the expert's problem-solving processes and encoding them in a computer system.

- z Expert systems operate in a transparent manner, i.e., the path to their conclusions can be traced. The systems can explain or justify their conclusions, and they also enable developers to check on the operation of the systems themselves.
- z One essential step in the construction of an expert system is the storage of knowledge in the knowledge base. To this end, this knowledge has to be present in a structured form. The transition of the domain expert's knowledge to rules cannot always be made; this is especially the case when this knowledge is partially based on experience. Another problem may occur when the domain expert's opinion or method of acting changes over time; it is then often difficult to change the knowledge base in order to reflect the altered situation.

Connectionist AI techniques

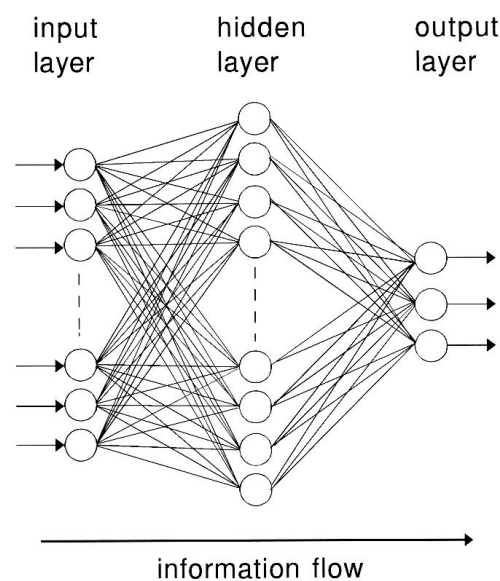
- z This mainly concerns the use of artificial neural networks (ANNs).
- z They are especially useful in finding solutions for problems where no, or very little, explicit knowledge is available and where sufficient representative specimen cases from which heuristic knowledge can be derived are available.

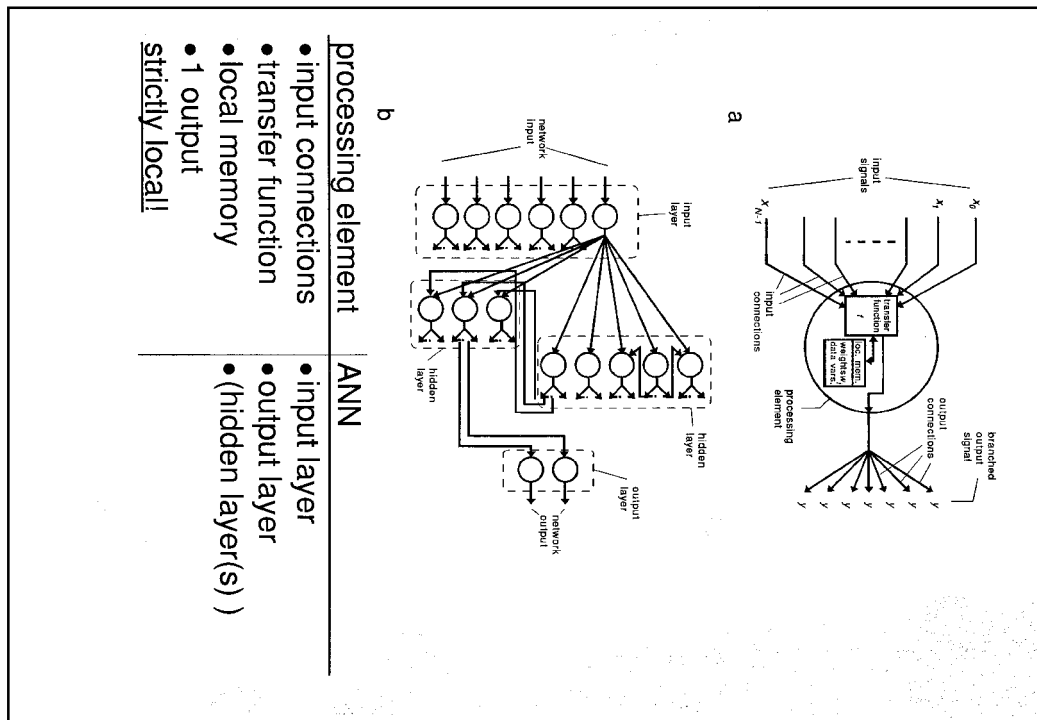
Programmed vs. Neurocomputing

- z programmed computing: success dependent on the possibility to convert the solution to a given problem into a set of procedures or rules
- z neurocomputing: success dependent on the existence of enough usable 'examples' of aspects of the problem

Artificial Neural Networks

- z a *large* number of *simple (non-linear)* units that are *densely interconnected*
- z information is stored in *weights* that are associated with the *interconnections*
- z the network 'learns' by *adapting* the network *weights* in response to information present in its environment

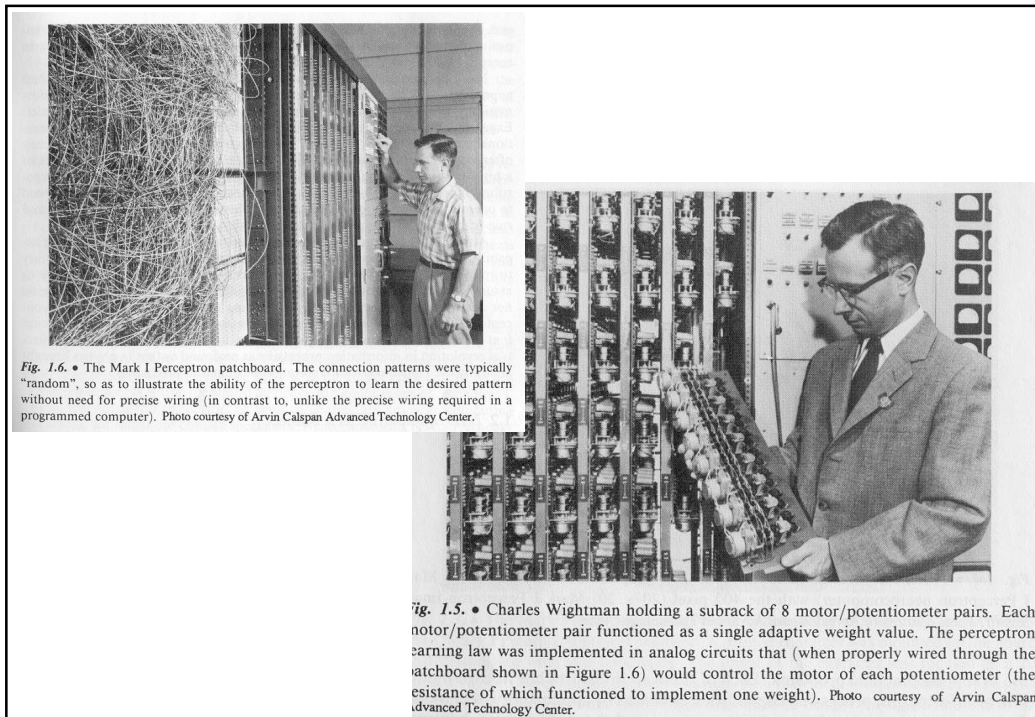
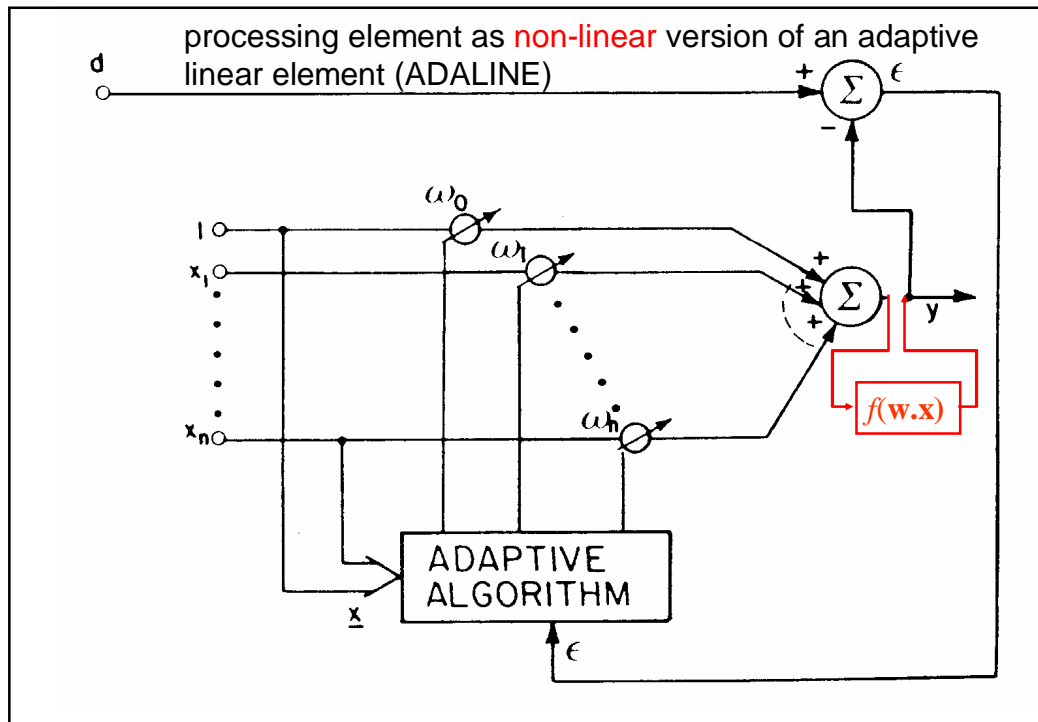




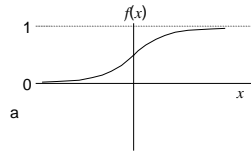
Some terms

- z inputs of a processing element are described by an *input vector*, \mathbf{x}
- z weights associated with input connections are described by a *weight vector*, \mathbf{w}
- z *output signal*, y
- z *transfer function*, f

for example:
$$y = f(\mathbf{w}, \mathbf{x}) = \sum_i w_i \cdot x_i = \mathbf{w} \cdot \mathbf{x}$$

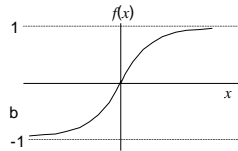


often used transfer functions



a sigmoid function

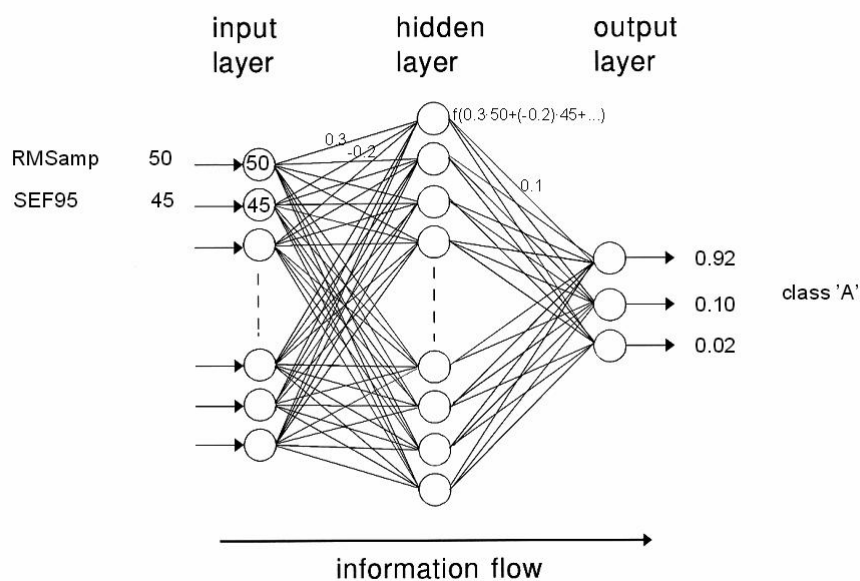
$$f(x) = \frac{1}{1 + e^{-x}}$$



a hyperbolic tangent

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

An extra nonlinearity parameter, β , may be used to adjust the 'steepness' of the function, in which case the functions do not operate upon x , but on βx .

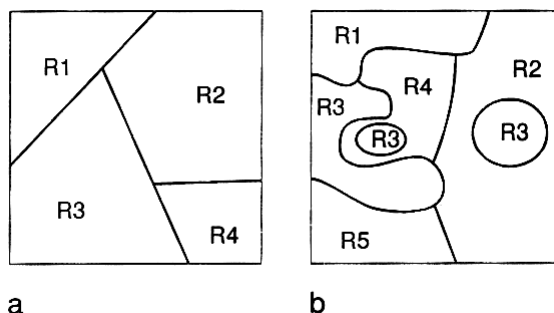


- z often, one input is kept at a value of 1, the *bias*
- z *local memory* contains *data variables* such as learning rate and momentum
- z a *learning rule* influences the behaviour of the ANN by *adapting* the network *weights*

Use for Recognition and Classification

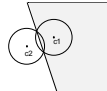
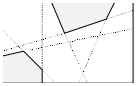

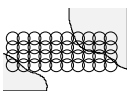


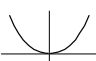
- z generate suitable decision regions in feature space by looking for the right weight configuration

a: decision regions realised by a *linear* classifier
b: decision regions realised by a *non-linear* classifier



ANN-based classifiers

- z Many artificial neural networks operate similarly to methods used in the statistical approach, i.e., they try to estimate decision regions in feature space. We can distinguish:
 - y probabilistic classifiers that use *a priori* information about distributions and use supervised learning to estimate their parameters
 - y hyperplane classifiers that construct hyperplane boundaries by calculating a weighted sum of inputs and passing this through a non-linear function
 - y kernel (receptive-field) classifiers that establish receptive fields across feature space; if an input is near the centre of a receptive field of a processing unit, the output of that unit will be high
 - y exemplar classifiers that use the distance between input pattern and training exemplars to come to classifications

| group | estimation of decision regions | processing unit | examples of classifiers |
|---------------|---|--|---------------------------------|
| probabilistic |  | distribution dependent | Gaussian mixture |
| hyperplane |  | sigmoid  | multi-layer perceptron networks |
| kernel |  | kernel  | radial basis networks |
| exemplar |  | (Euclid.)norm  | LVQ, ART, Kohonen networks |

The first column gives the name of the group of ANN classifiers, the second column shows schematically for each group how the decision regions in feature space are estimated. The third column shows the kind of computations performed in the processing units and the fourth column contains examples of ANNs that belong to each group

Learning in ANNs

z Different groups of training/learning methods

- y supervised learning
- y unsupervised learning
- y reinforcement learning

z each group contains many methods

z examples discussed here:

- y *(generalised) delta rule* (supervised learning)
- y *Kohonen learning rule* (unsupervised learning)

Delta Rule

(aka Widrow, Widrow/Hoff, or LMS learning rule)

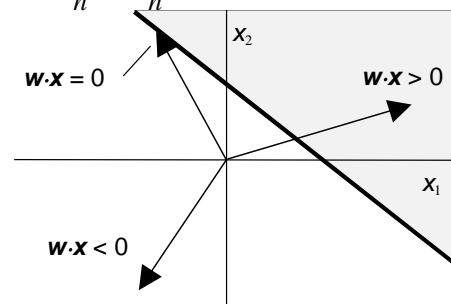
z originally used in the 60's in ADALINEs (ADAPtive LINEar Elements)

z in its generalised form very popular (backpropagation)

n -dimensional input vector \mathbf{x}

$$y = w_0 \cdot 1 + w_1 \cdot x_1 + \dots + w_n \cdot x_n$$

hyperplane classifier



- input/desired output pair k is (\mathbf{x}_k, y_k^*)
- define cost function, G , as expectation of the squared error

$$G(\mathbf{w}) = E[(y_k^* - y_k)^2] \text{ or}$$

$$G(\mathbf{w}) = \lim_{N \rightarrow \infty} \left(\frac{1}{N} \right) \sum_{k=1}^N (y_k^* - y_k)^2$$

parabolic surface; 'slide' towards minimum by using $-\nabla_{\mathbf{w}} G$
 (note: $\nabla_{\mathbf{w}} (-y_k) = -\mathbf{x}_k$)

$$\nabla_{\mathbf{w}} G(\mathbf{w}) = \lim_{N \rightarrow \infty} \left(\frac{1}{N} \right) \sum_{k=1}^N 2 \cdot (y_k^* - y_k) \cdot (-\mathbf{x}_k)$$

$$\nabla_{\mathbf{w}} G(\mathbf{w}) = -2 E[\delta_k \mathbf{x}_k]$$

δ_k is the error for input k

- this implies: average a large number of $\delta_k \mathbf{x}_k$ vectors, multiply by -2 and move weights in that direction
- Widrow & Hoff: update weights after every input presentation → **delta rule**:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \delta_k \mathbf{x}_k$$

with learning rate α

alternatives:

- update only after a number of input presentations (**batched version**)
- use most recent weight update also in current weight update (**momentum version**):

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha \delta_k \mathbf{x}_k + \mu (\mathbf{w}_k - \mathbf{w}_{k-1})$$

with momentum μ

limited performance due to its **linearity**

ANN architectures

often designed to

- **store and retrieve information**

associative networks (e.g., Hopfield networks)

process $(\mathbf{x}_i, \mathbf{y}_i^*)$ examples \rightarrow

in response to $\mathbf{x}_i + \epsilon$: produce \mathbf{y}_i^*

- **approximate mathematical functions**

mapping networks (e.g., 'backpropagation' networks)

process $(\mathbf{x}_i, \mathbf{y}_i^*)$ examples \rightarrow

approximate $f(\mathbf{x}) = \mathbf{y}^*$ for $\mathbf{x} \neq \mathbf{x}_i$

map vectors from \mathbf{R}^n to \mathbf{R}^m

Kolmogorov's theorem

any continuous function $f: [0,1]^n \rightarrow \mathbf{R}^m$, $f(\mathbf{x}) = \mathbf{y}$ can be approximated exactly by a three-layer ANN
(but how to create such an ANN???)

Multi-Layer Perceptron (MLP) trained with the backpropagation algorithm

⚡ Perceptron: 'non-linear ADALINE'

$$y = f\left(\sum_{i=0}^n w_i \cdot x_i\right)$$

if the argument of f is larger than or equal to 0, f is +1, else it is -1
learning rule similar to delta rule



decision regions built of 'half planes'

more complex decision regions by using multiple layers of perceptrons

Q: how to train such a configuration?

A: backpropagation

Backpropagation

- ⌚ Minimise a cost function, G_p , defined for an input pattern p and a network with m outputs

$$G_p = \frac{1}{2} \sum_{j=1}^m (y_{pj}^* - y_{pj})^2$$

With y_{pj} the output of output element j when pattern p is presented.
and y_{pj}^* the desired output-

Again, we use the gradient descent to change the weights - the elements w_{ij} of weight vector \mathbf{w}_j change according to:

$$\Delta_p w_{ij} \propto -\frac{\partial G_p}{\partial w_{ij}} = -\frac{\partial G_p}{\partial I_{pj}} \cdot \frac{\partial I_{pj}}{\partial w_{ij}}$$

I_{pj} is the input applied to processing element j as a result of presenting p

for output elements the weight change can be calculated straightforwardly using

$$I_{pj} = \sum_{i=0}^k w_{ij} \cdot y_{pi}$$

We would like to use the following weight update rule:

$$\Delta_p w_{ij} = \alpha \cdot \delta_{pj} \cdot y_{pi}$$

this is called the **generalised delta rule**

with learning rate α (usually decreasing with time), and the error made by processing element j as a result of the application of pattern p denoted as δ_{pj}

For output elements, we can easily calculate δ_{pj} since we know what their desired and actual output (and thus, error) is.

$$\Delta_p w_{ij} = -\frac{\partial G_p}{\partial I_{pj}} \frac{\partial \left(\sum_i w_{ij} \cdot y_{pi} \right)}{\partial w_{ij}} = -\frac{\partial G_p}{\partial I_{pj}} \cdot y_{pi}$$

thus, if we would like to have a rule of the form

$\Delta_p w_{ij} = \alpha \cdot \delta_{pj} \cdot y_{pi}$, we need to estimate δ_{pj} as :

$$\begin{aligned} \delta_{pj} &\propto -\frac{\partial G_p}{\partial I_{pj}} = -\frac{\frac{1}{2} \partial \left(\sum_j (y_{pj}^* - y_{pj})^2 \right)}{\partial I_{pj}} \\ &= -\frac{1}{2} \cdot 2 \cdot (y_{pj}^* - y_{pj}) \cdot -\frac{\partial y_{pj}}{\partial I_{pj}} \\ &= (y_{pj}^* - y_{pj}) \cdot f'_j(I_{pj}) \end{aligned}$$

Eventually we get for the output elements:

$$\delta_{pj} = (y_{pj}^* - y_{pj}) \cdot f'_j(I_{pj})$$

For 'hidden' nodes however, this is a bit more complex - we have to apply the chain rule for differentiation and use the error of the output elements. We thus 'backpropagate' the error through the network to calculate all weight updates.

For the elements in the hidden layer(s) we get

$$\delta_{pj} = \sum_{k=1}^m \delta_{pk} \cdot w_{jk} \cdot f'_j(I_{pj})$$

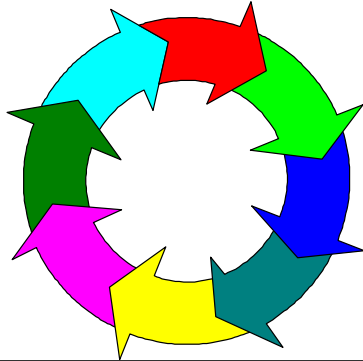
f'_j is the derivative of the transfer function of processing element j

Note: since the derivative of the transfer function is used, this function must be differentiable for every input value.

Backpropagation training

z with the weight update rule and the expressions for δ_{pj} we have the tools to calculate weight changes for every processing element after each pattern presentation p .

z Iterative training process



- present a new input pattern
- calculate network output using summations and transfer functions
- calculate errors of output elements
- calculate errors of hidden elements (backpropagation)
- update weights using generalised delta rule

Backpropagation process

- z this iterative process will eventually lead us to a weight configuration that is associated with the global minimum in the, very-many-dimensional, cost-function surface (or, at least we hope so)
- z There is no guarantee that we will actually reach the global minimum - the process can get stuck in local minima as well
- z The training process can be very time-consuming - often training is stopped when a predefined number of iterations has been reached, the error on the training set or the 'training evaluation set' drops below a certain threshold, or the weights do not significantly change over a long time.
- z Again, variations like 'batched weight updating' and use of momentum terms may be used to speed up the process.

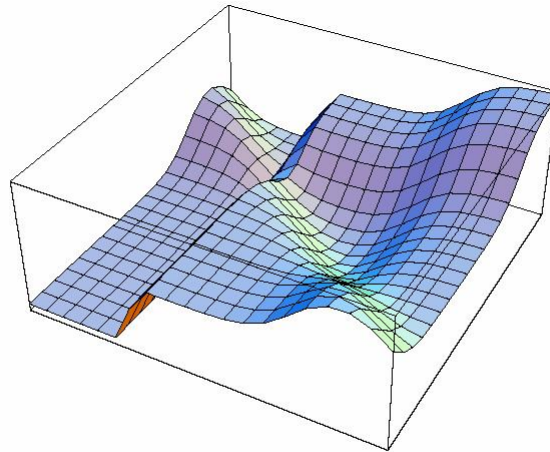
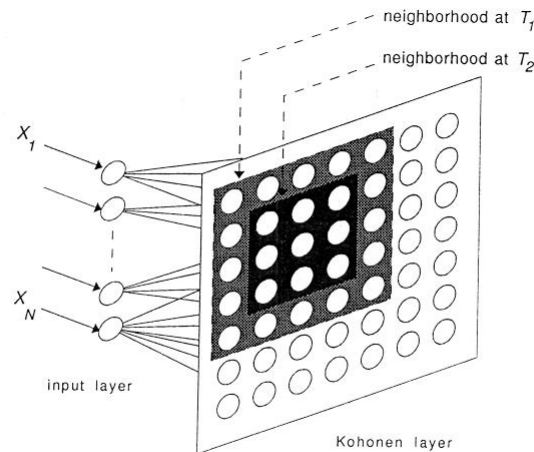


Fig. 7.5. A local minimum of the error function

Kohonen learning rule

- z unsupervised, i.e., no desired outputs
- z allows creation of an estimation of the probability density function (PDF) of applied inputs
- z two layers;
 - y input layer (with N elements)
 - y output (Kohonen) layer (M elements)

Kohonen network



output y_i of Kohonen element i :

$$y_i = D(\mathbf{w}_i, \mathbf{x})$$

depends on distance between \mathbf{w}_i and \mathbf{x}

'competition' \rightarrow one 'winner':

element with smallest distance: $y_c = 1$

all others: $y_j = 0$

kohonen learning rule:

$$\mathbf{w}_j(t+1) = \mathbf{w}_j(t) + h_{cj}(t) \cdot (\mathbf{x}_j(t) - \mathbf{w}_j(t))$$

with neighborhood kernel h_{cj}

$$h_{cj}(t) = h(\|\mathbf{r}_c - \mathbf{r}_j\|, t)$$

with \mathbf{r}_c the location of the winner in the Kohonen layer

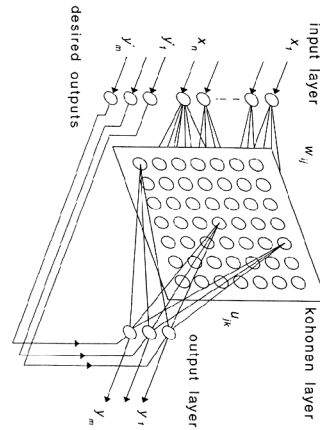
often:

$$h_{cj}(t) = \alpha(t)$$

with $\alpha(t) > 0$ for nodes in the 'neighborhood' of c , and 0 for all others
both the neighborhood size and α decrease as function of time

Counterpropagation Network

(Kohonen network extended with an extra output layer)



processing element k in the extra output layer is fully connected to all M Kohonen elements:

$$y_{pk} = \sum_{j=1}^M u_{jk} \cdot z_{pj}$$

z_{pj} is the output of Kohonen element j in reaction to pattern p , it can be 0 or 1. u_{jk} is the weight between Kohonen element j and output element k

its weights are updated using:

$$\Delta u_{jk} = \alpha \cdot (y_{pk}^* - y_{pk}) \cdot z_{pj}$$

only weights connected to the winning Kohonen element are updated.

Each weight will eventually converge to the average desired output associated with inputs that cause a certain Kohonen element to win the competition

Example: following patient state in ICU

\mathbf{z} from the routinely monitored data;

\mathbf{y} SAPs,m,d PAPs,m,d

\mathbf{y} CVPm, PCWP

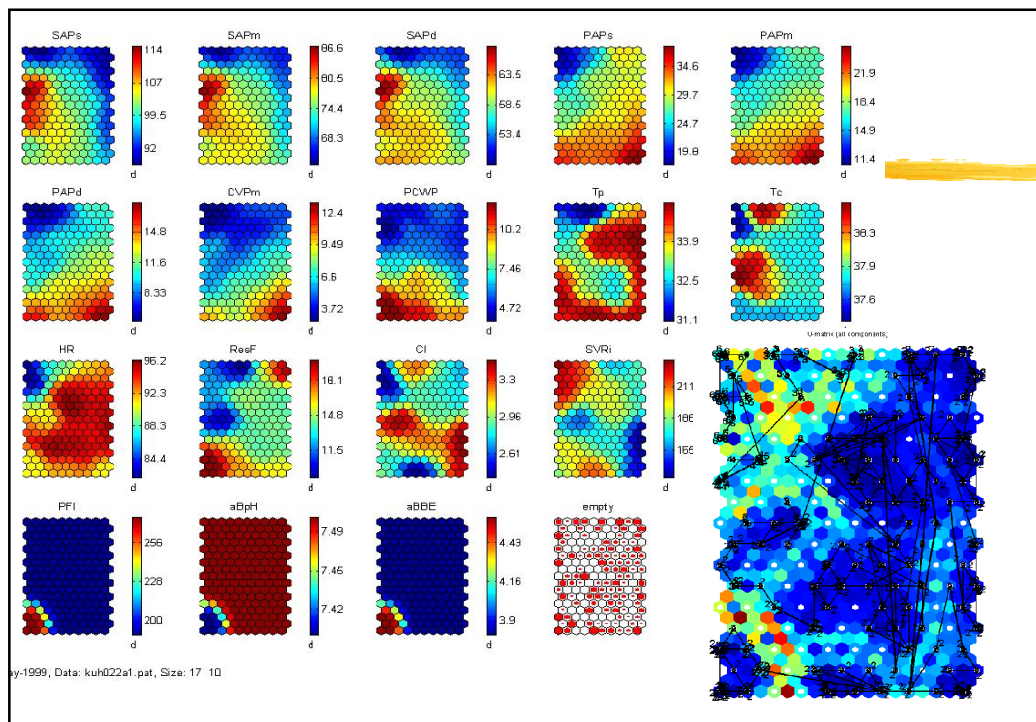
\mathbf{y} Tp, Tc

\mathbf{y} HR ResF CI SVRi

\mathbf{y} SpO2 PFI aBpH aBBE

18-dim vectors measured at 2 min.intervals

each vector labelled with the most recent Ramsay score known at that moment



Other types of networks

- z The backpropagation network and the Kohonen and counterpropagation networks are only a few examples.
- z A wide range of ANN paradigms exists, some examples:
 - y stochastic networks (e.g., try to get out of local minima using a technique called simulated annealing)
 - y spatiotemporal networks (include time aspects, e.g., for speech recognition)
 - y multi-modular networks (divide a huge problem into simpler subproblems and solve each of those with an 'expert' network)
 - y combinations with other techniques, such as expert systems or fuzzy logic ('hybrid networks')

Usage

- z In non-trivial applications ANNs are typically used as a module co-operating with other modules that may use other techniques (rules, algorithms, ...) in so-called hybrid systems. Rarely an ANN can be used to solve 'everything'.

Sufficient usable data is crucial

Large data sets are needed for training
note:

- z A data set with 999.998 occurrences of disease A and only 2 occurrences of disease B may be a small data set (quite often we would like the ANN to help us exactly with identifying those rare disease B cases)
- z Missing data often forms a serious problem in clinical applications