

## Machine Learning

### Introduction

Josef Kittler

email: [J.Kittler@surrey.ac.uk](mailto:J.Kittler@surrey.ac.uk)

[www.ee.surrey.ac.uk/Personal/J.Kittler](http://www.ee.surrey.ac.uk/Personal/J.Kittler)

## Course topics

- Machine learning problem
  - Examples
  - Problem formulation
  - Learning scenarios
- Basic linear machines
  - Nearest neighbour classifier
  - Perceptron
  - Sparse representation based classifier
- Nonlinear extensions
  - Kernel methods
  - Multilayer neural networks
  - Deep neural networks
- Dimensionality reduction
- Classifier design issues

- Training deep neural networks (DNN)
  - Advanced DNN
  - Recurrent neural networks
  - Deep learning libraries
  - Anomaly detection in graphs
- } Muhammad Rana
- } Fei Yan
- } Radek Marik

- **Machine Learning** is a field of study concerned with the development of algorithms that can learn from and make predictions on data
- **The aim of machine learning** is to give computers the ability to learn (find solutions to problems) without being explicitly programmed
- Applications span a vast range of problems.

Biometrics Object detect Target detect Bridge detection



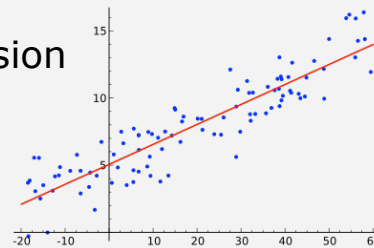
## Typical machine learning tasks

- To generate response to input data so as to achieve required functionality
- Examples include
  - Regression (predict output given input)
  - Classification (predict class membership)
  - Cluster assignment (associate input to data structure)
  - Detection (detect a specific object)
  - Anomaly detection (identification of input as an outlier)

5

## Examples of regression

- Simple linear regression



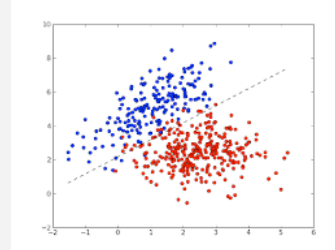
- Multivariate regression



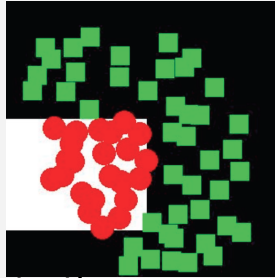
6

## Examples of classification

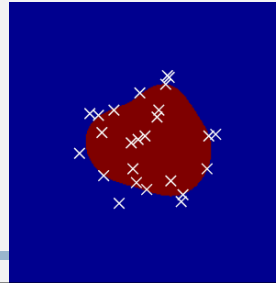
- Linear classification



- Detection



- Anomaly detection



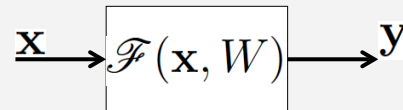
7

## General statement of the machine learning problem

- Mathematically, the machine is realising an appropriate function

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, W)$$

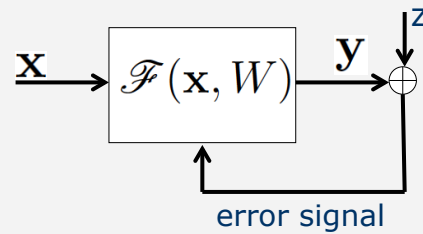
- $\mathbf{x}$ ..... D dimensional input
- $\mathbf{y}$ ..... d dimensional output
- $W$ ..... parameters



8

- Pre-requisites

- Training set  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Ground truth target values  $Z = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$
- Form of function  $\mathcal{F}$  (architecture)
- Objective function (error measure)
- Procedure for updating  $W$



9

- Supervised (target set available)
- Non supervised (target set unavailable)
  - clustering
- Semi-supervised (some training data labelled)
- Transfer learning
  - Supervised learning in the source domain
  - Target domain different from the source domain
  - Only unlabelled data available in the target domain

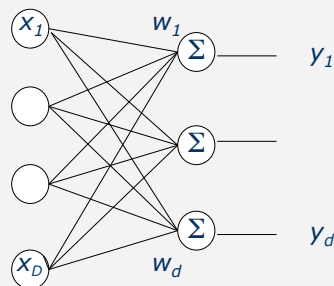
10

- Function  $\mathcal{F}$  is linear, namely

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

where  $\mathbf{W}$  is a  $d \times D$  matrix of parameters

$$\mathbf{y} = \begin{bmatrix} w_1^T \\ \cdot \\ \cdot \\ w_d^T \end{bmatrix} \mathbf{x}$$



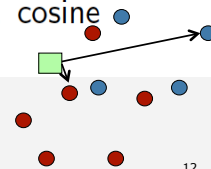
11

# Similarity based classification

- Nearest Neighbour classifier labels patterns based on similarity, gauged in terms of distance
- Squared Euclidean distance between  $\mathbf{x}$  and  $\mathbf{x}_j$  is given as

$$(\mathbf{x} - \mathbf{x}_j)^T (\mathbf{x} - \mathbf{x}_j) = \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{x}_j + \mathbf{x}_j^T \mathbf{x}_j \quad (10)$$

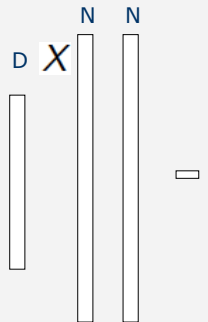
- The distance between  $\mathbf{x}$  and  $\mathbf{x}_j$  minimal when  $\mathbf{x}^T \mathbf{x}_j$  is maximal
- scalar product  $\mathbf{x}^T \mathbf{x}_j$  gauges similarity
- other notions of similarity can be defined, e.g. cosine similarity, Gaussian kernel, subjective grading



12

# NN classifier machine

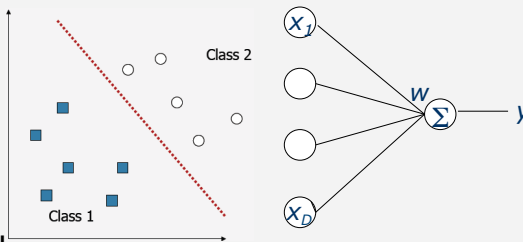
- N... number of training samples
- D... dimensionality of each sample



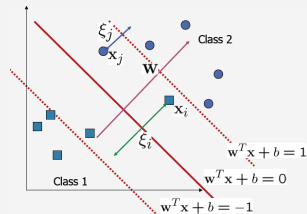
- Augment the result by input layer
- Compute scalar products for all training samples
- Min

# Learning objectives

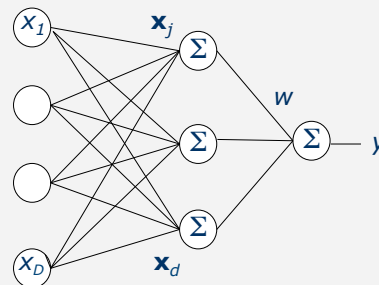
- Linear classifier



- Support vector machine



$$\mathcal{F}(\mathbf{x}) = \sum_{j(\text{support vectors})} w_j z_j (\mathbf{x}_j^T \mathbf{x}) + b$$



## Sparse representation classification (SRC)

- Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_M]$  be training data matrix
- Reconstruct an unknown sample  $\mathbf{y}$  as

$$\mathbf{y} = X\mathbf{a}$$

where  $\mathbf{a}$  is a vector of coefficients

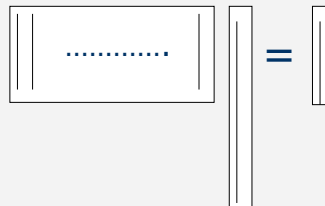
- Premise: for a sample  $\mathbf{y} \in \omega_i$ , we would expect the reconstruction to be constituted by training samples from class  $i$ , i.e. vector  $\mathbf{a}$  should be sparse (all entries for samples from other classes should be zero)
- The solution has to be regularised by imposing a minimum norm on  $\mathbf{a}$

15

## SRC

- By imposing sparsity on the reconstruction solution, we should be able to identify the class of  $\mathbf{y}$
- This can be achieved by using  $l_1$  and solving

$$\operatorname{argmin} \|\mathbf{a}\|_1 \quad \text{s.t. } \mathbf{y} = X\mathbf{a} \quad (7)$$



$$X \mathbf{a} = \mathbf{y}$$

16



## SRC algorithm

1 Solve

$$\operatorname{argmin} \|\mathbf{a}\|_1 \quad \text{s.t.} \|\mathbf{y} - \mathbf{X}\mathbf{a}\| < \epsilon \quad (8)$$

2 Let  $\mathbf{a}_i$  be vector  $\mathbf{a}$  with all entries associated with samples from class  $j \neq i$  set to zero, and compute the residual

$$r_i(\mathbf{y}) = \|\mathbf{y} - \mathbf{X}\mathbf{a}_i\|_2 \quad (9)$$

3 assign  $\mathbf{y} \rightarrow \omega_i$  if  $r_i(\mathbf{y}) = \min_j r_j(\mathbf{y})$

- Relationship to the k-NN classifier

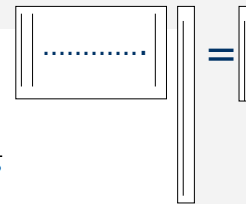
- k-NN classifier minimises the distances to  $\mathbf{y}$
- in addition, SRC involves pairwise interactions of residual error vectors

- Solution to be found for every test sample

17

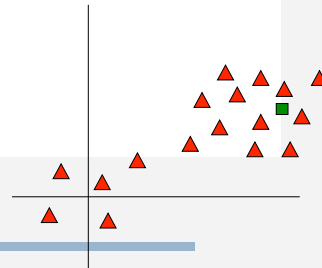
## The effect of norm

- $\ell_p$  norm of  $\mathbf{a} = [a_1, \dots, a_d]^T$



$$\ell_p = \left[ \sum_{j=1}^d |a_j|^p \right]^{\frac{1}{p}}$$

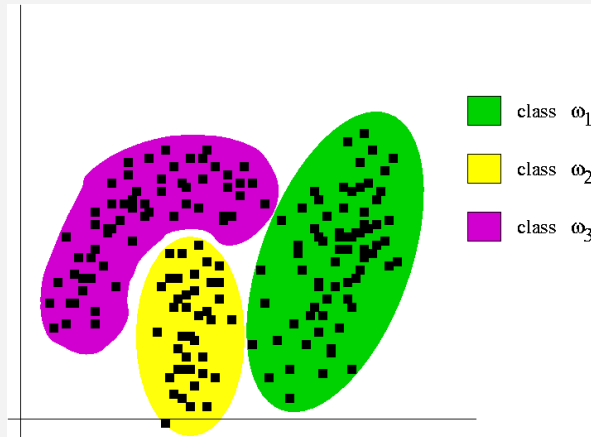
- $\ell_0$  ..... counts the number of nonzero elements
- $\ell_1$  ..... induces sparsity
- $\ell_2$  ..... length of vector  $\mathbf{a}$
- $\ell_\infty$  ..... selects  $\arg \max_j a_j$



18

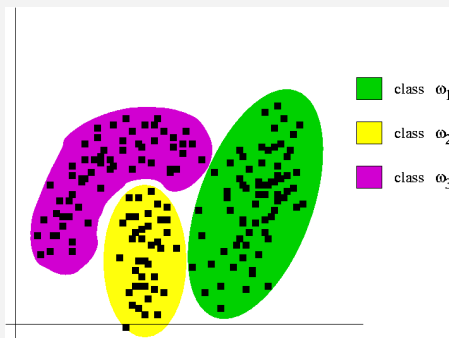
# Nonlinear separation

Geometric viewpoint of the pattern recognition problem

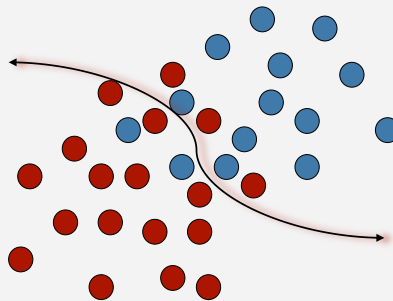


# Overlapping classes

- Class boundaries may be nonlinear



- Classes may be overlapping



# Kernel SVM

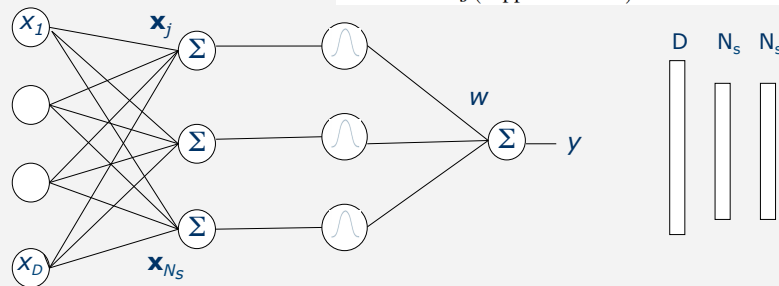
- Scalar product replaced by a kernel function

$$k(\mathbf{x}_j, \mathbf{x})$$

for example, a radial basis function

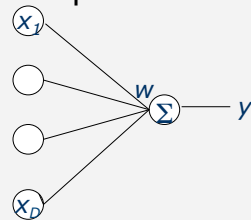


$$\mathcal{F}(\mathbf{x}) = \sum_{j(\text{support vectors})} w_j z_j k(\mathbf{x}_j, \mathbf{x}) + b$$

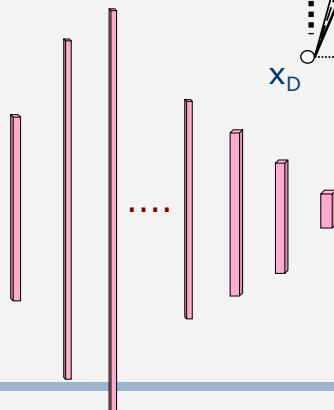


# Neural networks

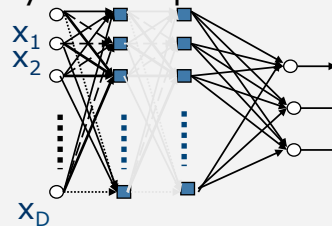
- Perceptron



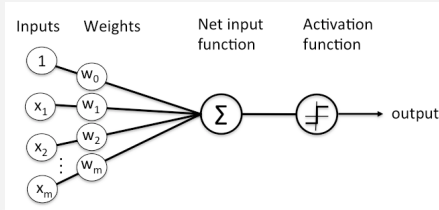
- Deep neural network



- Multilayer Perceptron

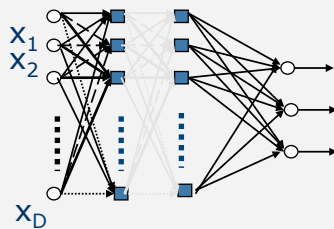
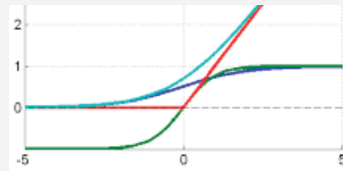


Node expansion motivation  
 ◆ Ensemble  
 ◆ Generate different terms



## ■ Activation function

- sigmoid
- rectified linear
- tanh
- softplus
- radial basis



23

- Defining the structure of the approximation of function  $\mathcal{F}$
- Making the computation of  $\mathcal{F}$  robust
- Finding the values of the unknown parameters to achieve the desired objectives (system behaviour)
  - Avoiding local optima
  - Promoting the ability of the solution to generalise to unseen data
  - Key measure
    - Dimensionality reduction
    - Training set management
    - Ensembles (dropouts)

24

## Classifier design issues

- System architecture
  - # of layers
  - connectivity
  - role of layer
  - dropouts
- Activation function
- Objective function
  - primary objective
  - constraints
  - parameters to be learned
  - regularisation of solution
  - metric
  - fusion
- Training data usage
  - training/testing
  - evaluation
  - augmentation
    - Mirroring
    - Model based
    - Random sampling
    - Perturbation
- Learning process and its parameters
  - Epoch
  - Learning rate
  - Greedy learning

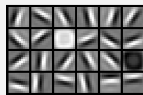
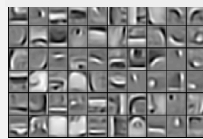
25

## Primary objective function

- Choice very important, as it induces different learning properties
- Examples
  - Classification error
  - Mean squared error  $[\mathbf{y} - \mathbf{z}]^T [\mathbf{y} - \mathbf{z}]$
  - Cross entropy  $-\sum_{j=1}^m y_j \log z_j$

26

## Dimensionality reduction

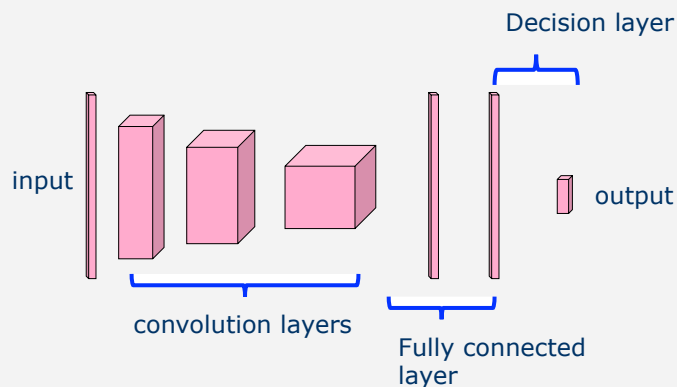


- Motivation – reduce over fitting
- Scope
  - Redundancy
  - Irrelevant content
  - Hierarchical relations – from local to global
- Implications on architecture
  - # of relevant features is low
  - At lower layers only local filters are needed: restricted connectivity
  - translation invariance -

27

## Convolutional neural network

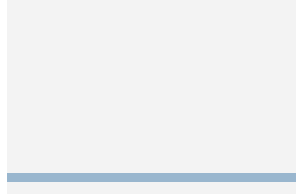
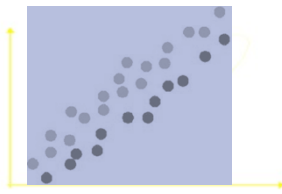
- Basic CNN architecture



28

## Dimensionality reduction - PCA

- Principal component analysis (PCA): an orthogonal basis transformation
- Transform correlated variables into uncorrelated ones (principal components)
- Can be used for dimensionality reduction
- The underlying aim is to find bases  $V$  and transformed variables  $\mathbf{y}$  so that in expectation  $\|\mathbf{x} - V\mathbf{y}\|_2$  is minimal subject to  $V^T V = I$
- Retains as much variance as possible when reducing dimensionality



## How PCA works

- Given  $m$  centred vectors:  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$ 
  - $X$ :  $D \times N$  data matrix
- The squared fitting error averaged over the training set
 
$$\text{tr}\{[X - VV^T X]^T [X - VV^T X]\}$$
- This can be rearranged as
 
$$\text{tr}\{[X^T X - X^T VV^T X]\}$$
- The solution to the constrained optimisation problem is a system of eigenvectors and eigenvalues satisfying
 
$$XX^T V - V\Lambda = 0$$
- Note  $C = XX^T$  is the  $D \times D$  covariance matrix
- Diagonal matrix  $\Lambda$ : eigenvalues

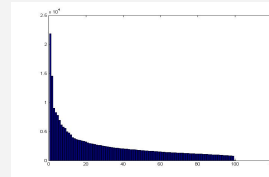
## Mapping data to low dimensional feature space

- Given  $V=[v_1, v_2, \dots, v_d]$ , the  $d$ -dimensional feature vector in the PCA space is given by

$$y = V^T x$$

- dimensionality  $d$  chosen to retain a certain fraction of variance
- $d < D$
- $d \leq N$

Distribution of eigenvalues



Example of eigenvectors



## Kernelising PCA

- Premultiplying eq (1) by data matrix  $X^T$  gives

$$X^T X X^T V - X^T V \Lambda = 0 \quad (2)$$

- an eigenvalue problem with eigenvectors  $U = X^T V$  and the same eigenvalues
- matrix  $X^T X$  is  $N \times N$
- if  $N \ll D$ , then solving eq (2) is more efficient
- matrix  $V$  can be obtained as  $V = X U \Lambda^{-1}$
- matrix  $K = X^T X$  is referred to as kernel matrix
- its element  $k_{ij} = k(x_i, x_j)$  measures "similarity" of vectors  $x_i, x_j$
- $k_{ij} = x_i^T x_j$



- the kernel formulation of PCA offers a scope for using different notions of similarity
- for instance,  $k_{ij} = \exp\{-\|\mathbf{x}_i - \mathbf{x}_j\|^2 \sigma^{-2}\}$
- other notions of similarity remap  $\mathbf{x}$  in a nonlinear way into  $\phi(\mathbf{x})$
- $\phi(\mathbf{x})$  may be of infinite dimensionality
- $\phi(\mathbf{x})$  is defined implicitly, is unknown, but satisfies  $k_{ij} = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$
- working with  $K$  is like working with data matrix  $X = [\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_N)]$
- We have kernelised PCA

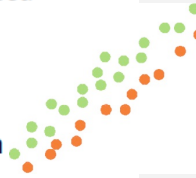
- Note
  - $V = XU\Lambda^{-1}$  is not explicitly available:  $U$  and  $\Lambda$  are, but  $X$  is not
- However... we are interested in projection onto basis  $V$ , not the basis itself
- Projection onto  $V$ :  $X^T V = X^T XU\Lambda^{-1} = KU\Lambda^{-1}$
- All  $K$  and  $U$  and  $\Lambda$  are available
- $\Lambda$  purely rescales the data and can be omitted



- Kernel Fisher discriminant analysis: another supervised learning technique
- Focusing on discrimination, rather than faithful representation
- Seeking the projection  $\mathbf{w}$  maximising Fisher criterion

$$\max_{\mathbf{w}} \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T (S_T + \lambda I) \mathbf{w}} \quad (3)$$

- $S_B$  and  $S_T$ : between class and total scatters
- $\lambda$ : regularisation parameter
- The total scatter matrix equals mixture covariance matrix  $S_T = XX^T$
- Between class scatter  $S_B$  can be expressed as  $S_B = X\Delta X^T$
- Block diagonal matrix  $\Delta$  contains a constant in block  $i$ , proportional to the number of samples from class  $i$



35



- Expressing  $\mathbf{w} = X\alpha$ , and substituting for  $S_B$ , we can kernelise Fisher criterion as

$$\max_{\alpha} \frac{\alpha^T X^T X \Delta X^T X \alpha}{\alpha^T X^T (XX^T + \lambda I) X \alpha} = \max_{\alpha} \frac{\alpha^T K \Delta K \alpha}{\alpha^T (K^2 + \lambda K) \alpha} \quad (4)$$

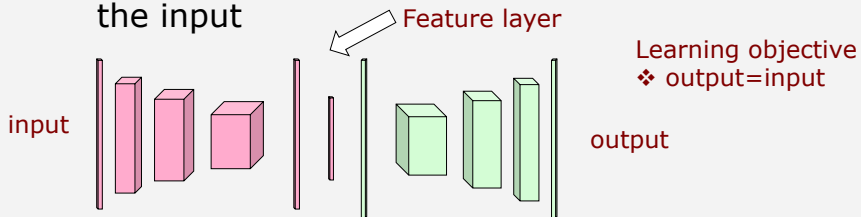
- The maximisation leads to an eigenvalue problem

$$\Delta K \alpha - \kappa (K + \lambda I) \alpha = 0 \quad (5)$$



## Examples of DNN feature extraction approaches

- Auto-encoder: find features that can regenerate the input



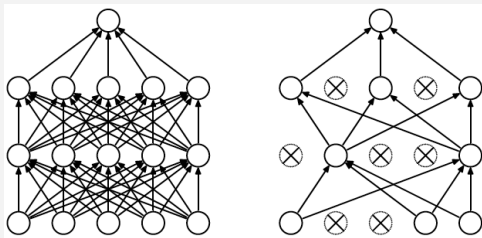
- Extracting representation that outputs similar/dissimilar features for similar/dissimilar inputs



40

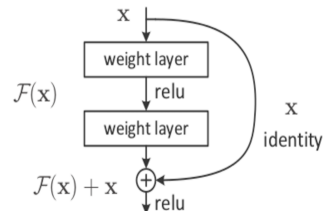
## Learning and generalisation enhancements

- Node dropout



N Srivastava et al., Dropout: A simple way to prevent neural networks from overfitting, JMLR 2014

- Residual network – vanishing gradient problem



41

- The underlying assumption in pattern recognition is that similar patterns belong to the same class.
- This is the premise behind Nearest Neighbour methods
- Also Gaussian classifier assign patterns to their classes based on similarity (distance from the class mean)
- However, the choice of metric used for measuring distances is critical

### ■ Desired properties of metric

- In general, we would like our metric  $d_A(\mathbf{x}_i, \mathbf{x}_j)$  to reflect the notion of similarity, i.e.

$$\begin{aligned} d_A(\mathbf{x}_i, \mathbf{x}_j) = \text{small} & \leq t_s & \text{similar pair} \\ d_A(\mathbf{x}_i, \mathbf{x}_j) = \text{large} & \geq t_l & \text{dissimilar pair} \end{aligned} \quad (13)$$

where threshold  $t_s$  is small, and  $t_l$  is large

- Let metric  $d_A(\mathbf{x}_i, \mathbf{x}_j)$  be parameterised as

$$d_A(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i - \mathbf{x}_j)^T A (\mathbf{x}_i - \mathbf{x}_j) \quad (14)$$

- Thus the metric learning problem can be formulated as

$$\begin{aligned} \min_A \quad & [tr(AA_0^{-1}) - \log |AA_0^{-1}|] \\ \text{s.t.} \quad & tr[A(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T] \leq t_s \quad (i, j) \in S \\ & tr[A(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T] \geq t_l \quad (i, j) \in D \end{aligned} \quad (17)$$

- As there may not exist a feasible solution, the optimisation problem can be relaxed using slack variables  $\xi_{c(i,j)}$ . These replace the constraint  $t_x$  for the  $(i, j)$  pair indexed by  $c(i, j)$ , where  $1 \leq c \leq n_S$  for the similar pairs, and  $n_S + 1 \leq c \leq n$  for the dissimilar pair.