$$P = \begin{pmatrix} v_1^T \\ v_2^T \\ v_m^T \end{pmatrix} \text{ where } m < n \text{ is the number of components we keep}$$

Given an *n*-dimensional input pattern x we can construct an *m*-dimensional pattern z  $z = P(x - \mu)$ 

Then use *z* are our new pattern

#### Reconstruction k)

 $\hat{x}_k = \mu + \sum_{i=1}^m z_i^k v_i$  Where  $z_i^k = v_i^T (x_k - \mu)$  and eigenvectors are normalised

Reconstruction error  $||x_k - \hat{x}_k||^2$  can be used to measure how much information has been lost PCA finds a subspace with the smallest reconstruction error

## 7 PCA in practice

#### a) Images

When processing real images there are too many features to look at so it creates a very large covariance matrix

Mega pixel image  $\rightarrow$  million wide matrix This will be much smaller that the *NxN* space described by the images with trindard PCA The Dual Matrix 'trick' can be used to over come this **b)** Dual Matrix Typical covariance matrix  $D = XX^T$  consider the matrix  $D = X^T X$ Suppose we an effect vector of D $Dv = \lambda v$ 

Suppose  $v_{i}$  as the vector of  $Dv = \lambda v$  $X^T X v = \lambda v$  $XX^T X v = \lambda X v$  $(XX^T)Xv = \lambda Xv$  $CXv = \lambda Xv$  let u = Xv

 $Cu = \lambda u$ *u* is an eigen vector of C

## Single Value Decomposition C)

Any NxP matrix X can be written as  $X = USV^T$ Where: U is an NxN orthogonal matrix V is a PxP orthogonal matrix S is an NxP diagonal matrix of singular values A matrix M can also be described in terms of its eigenvectors and eigenvalues  $M = V \Lambda V^T$ 

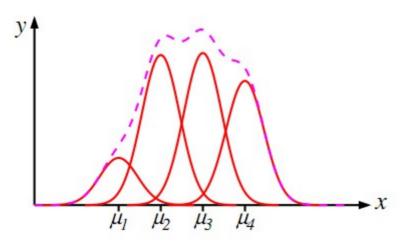
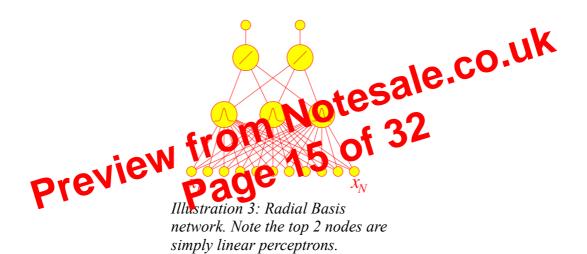


Illustration 2: Example output of a radial basis function. Note the height of the functions is the weight.

This can be combined further to produce a network of nodes and functions with multiple outputs.



# d) Finding centers

Radial basis functions depend on the distance to centres Finding these centres and how many are needed can be hard They are normally found in clusters

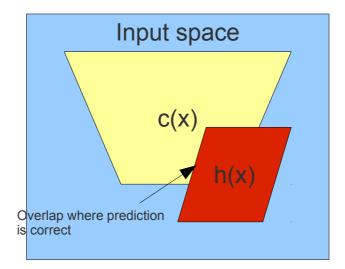
# e) k-Means Clustering

Find k clusters in the data Can be tricky to find the right kEasy algorithm and popular

Takes a set of unlabelled data  $D = \{x_j\}_{1}^{N}$ 

Algorithm:

- 1. Choose K
- 2. Randomly assign each data point to a cluster  $C_i$



#### d) PAC-Learnable

A concept class C is PAC learn-able if for all concepts in it  $c \in C$  and a probability distribution p a learner choosing a hypothesis  $h \in H$  can:

achieve a generalisation performance of  $E(h|p) < \epsilon$  for  $0 < \epsilon < \frac{1}{2}$ with a probability of  $1-\delta$  and  $0 < \delta < \frac{1}{2}$ 

using *P* training examples

where *P* is a polynomial in  $\frac{1}{\epsilon}$  and  $\frac{1}{\delta}$ 

iesale.co.uk The idea is to restrict the number of hypothesis's h can be learn by eliminating those with a to the sectainty mat the machine hasn't learnt a generalisation error greater than hypothesis with an error start than

## Finite Consistent Hypothese Spaces e)

Suppose the hypothesis space is finite and contains the true answer

Suppose our learning machine is able to correctly categorise all the training examples

$$E(h|p)=0$$

Consider examples  $D = \{(x_k, c(x_k))\} \stackrel{P}{k=1}$ 

Consider a hypothesis with an error  $E(h|p) > \epsilon$  this means that the probability of it getting all the training examples correct is:  $(1-\epsilon)^{P}$ 

There will be k < |H| hypothesises with  $E(h|p) > \epsilon$ 

The probability of any of these being correct means marginalising over all the k hypothesis  $k(1-\epsilon)^{P} < |H|(1-\epsilon)^{P} \le |H|e^{-\epsilon P}$ 

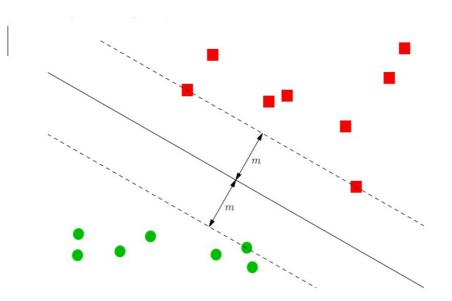
Using  $\delta \ge |H|e^{-\epsilon P}$  the number of training examples needed to get a generalisation error better than  $\in$  with a probability of  $1-\delta$  is given by:

$$\delta \ge |H| e^{-\epsilon P}$$
  

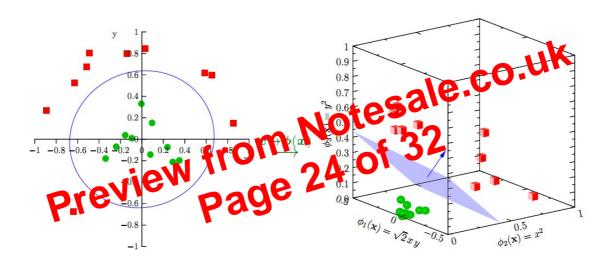
$$\ln(\delta) - \ln(|H|) \ge -\epsilon P$$
  

$$\ln(|H|) - \ln(\delta) \le \epsilon P$$
  

$$\frac{1}{\epsilon} \ln(|H|) - \ln(\delta) \le P$$



Can increase the likelihood of finding a separating plane by projecting into high dimensional space using the Kernel trick (see later)



SVMs have the ability to separate almost any data set However, by choosing the largest margin it selects the most robust solution Thus giving a good generalisation performance despite having a large capacity to learn data

This is termed structural risk minimisation

## b) Maximal Margins

Consider a linearly separable data set

 $\{(x_k, y_k)\}_{k=1}^{p} \text{ and } y_k \in \{-1, 1\}$ 

Task is to find the vector of weights w and bias b such that

 $y_k(\frac{w^T x_k}{\|w\|} - b) \ge m$  where *m* is the margin

In the image below, we want to maximise the distance from the central separating plane to the nearest data points

The vector from the line to these points is known as *Support Vectors* and these need to have their length maximise

 $\lambda_i \ge 0 \forall i$  otherwise the 'distance' between points in the extended space can be negative  $K(\mathbf{x}, \mathbf{y}) = \sum_{i} \phi_i(\mathbf{x}) \phi_i(\mathbf{y})$ 

This is the same as saying:

 $K(\boldsymbol{x}, \boldsymbol{y}) = \boldsymbol{\phi}(\boldsymbol{x})^T \boldsymbol{\phi}(\boldsymbol{y})$ 

So if you like a kernel is the sum of multiple eigenfunctions and their eigenvalues For SVMs to work the kernels must be projecting into euclidean space so that distances are always positive

Therefore we must use Positive Semi-definite Kernels (  $\lambda_i \ge 0 \forall i$  ) which can always be decomposed into as sum of positive functions

 $K(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i} \phi_{i}(\boldsymbol{x}) \phi_{i}(\boldsymbol{y})$ 

## d) Integratal of Kernels

$$K(\mathbf{x}, \mathbf{y}) = \sum_{i} \phi_{i}(\mathbf{x}) \phi_{i}(\mathbf{y}) \text{ And is positive semi-definite}$$

$$\int \int f(\mathbf{x}) K(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \ge 0$$

$$\int \int f(\mathbf{x}) \sum_{i} \phi_{i}(\mathbf{x}) \phi_{i}(\mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y}$$

$$\sum_{i} \int f(\mathbf{x}) \phi_{i}(\mathbf{x}) d\mathbf{x} \int \phi_{i}(\mathbf{y}) f(\mathbf{y}) d\mathbf{y}$$

$$\int f(\mathbf{x}) \phi_{i}(\mathbf{x}) d\mathbf{x} = \int \phi_{i}(\mathbf{y}) f(\mathbf{y}) d\mathbf{y} \text{ therefore}$$

$$\sum_{i} (\int f(\mathbf{x}) \phi_{i}(\mathbf{x}) d\mathbf{x})^{2} = \sum_{i} (\int \phi_{i}(\mathbf{y}) f(\mathbf{y}) d\mathbf{y})^{2} \text{ NOTESALE. CO.UK}$$
e) Combining Kerect from 27 of 32  
Adding Kerect from 27 of 32  
Adding Karect from 27 of 32  
Adding Karect from 27 of 32  
Multiplied  

$$K_{3}(x, y) = K_{1}(x, y) + K_{2}(x, y) \text{ is valid}$$

$$K_{2}(x, y) = K_{1}(x, y)^{n} \text{ is valid too}$$
Exponential  

$$K_{2}(x, y) = e^{K_{1}(x, y)} \text{ is valid}$$
An important hereal is :

An important kernel is :  $\frac{-||x-y||^2}{2}$ 

$$K(x, y) = e^{\frac{\pi}{2}}$$

As it a an infinite number of eigenvalues so if you like its projecting the data into an infinite dimensional space. This is because e is an irrational number which can be estimated by an infinite sequence of sums