

# Machine Learning: Algorithms and Applications

Philip O. Ogúnbona

Advanced Multimedia Research Lab  
University of Wollongong

Pattern Classification  
Autumn

# Outline of Topics

- 1 **Pattern Recognition/Classifier**
- 2 **Approaches to Classification**
- 3 **Elementary Decision Theory**
- 4 **Support vector machine**
- 5 **References**

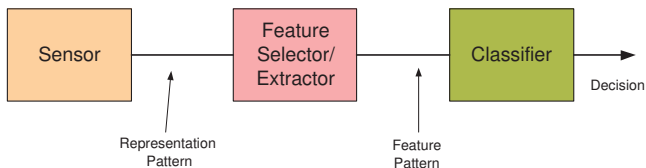
# What is pattern recognition/classifier?

- We take for granted the fact that we are able to move around in our world and recognize:
  - cars
  - people
  - animals
  - objects in general

despite the variety in their form and existence.

- What features help in these tasks?
- In **pattern recognition** we study how to design machines that can recognize and classify "things".
- We study the statistics of the features that describe "things".
- We study how to measure the performance of pattern recognition systems and select good systems.

# Basic Model - Pattern Classifier



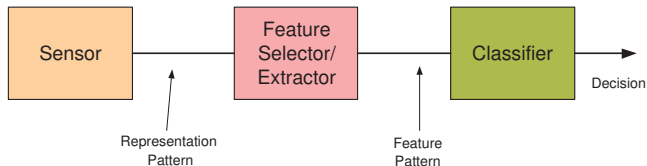
**Figure 1:** Pattern Classifier Webb (2002)

- The **pattern** is a set of numbers or values represented as a  $p$  – dimensional vector,

$$\mathbf{x} = [x_1 \quad x_2 \quad \cdots \quad x_p]^t$$

where  $t$  (or sometimes  $T$ ) denotes vector transpose

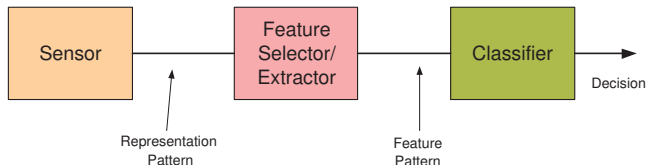
# Basic Model - Pattern Classifier



**Figure 1:** Pattern Classifier Webb (2002)

- The pattern could be:
  - pixels in an image
  - closing prices of a share on the stock market
  - recordings of a speech signal
  - measurements on weather variables
  - group of measurements about a real estate
  - group of measurements about the behaviour and life style of people
  - etc.

# Basic Model - Pattern Classifier



**Figure 1:** Pattern Classifier Webb (2002)

- We assume that there are  $C$  classes denoted by,

$$\omega_1, \dots, \omega_C$$

- There is a variable,  $z$ , that indicates which class,  $\omega_i$ , a pattern  $x$  belongs. That is,

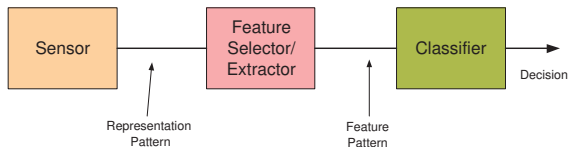
if  $z = i$ , then the pattern,  $x$ , belongs to  $\omega_i$ ,  $i \in \{1, \dots, C\}$

# Basic Model - Pattern Classifier

- The problem is how to design the pattern classifier.
- Designing a pattern classifier entails:
  - specifying the classifier model parameters
  - ensuring that response for a given pattern is optimal
- The design process assumes we have a set of patterns of known class,  $\{(\mathbf{x}_i, z_i)\}$ , called the training or design set, used to design the classifier.
- Part of the design process is to evaluate and set optimal operating parameters.
- The idea is that once we have a designed classifier we can estimate the class membership of an unknown pattern.
- There is an assumption that the samples used for training are drawn from the same probability distribution as the test samples and the operational samples.

# Basic Model - Pattern Classifier

- A closer look at the simplified classifier model:



**Figure 2:** Pattern Classifier Webb (2002)

- 1 Representation pattern is the raw data we obtain from the sensor e.g. image or video pixels, price of stock, etc.
- 2 Feature pattern is a small set of variables obtained through some transformation - **feature selection and/or extraction**
- 3 The trained classifier uses the feature pattern to make a decision regarding the pattern presented at its input.



# Basic Model - Pattern Classifier

- Further consideration about classifier design:

## Problem

*Given a training set of patterns of known class, we seek to design a classifier that is optimal for the expected operating conditions.*

- 1 The given set of training patterns is *finite*.
- 2 The classifier model cannot be too complex. In other words it cannot have too many parameters. This situation may lead to **over-fitting**.
- 3 It is not important to achieve optimal performance on the design set.
- 4 It is very important to achieve optimal **generalization** performance.
  - Expected performance on data representative of the true operating condition - the infinite set from which the design set is drawn.

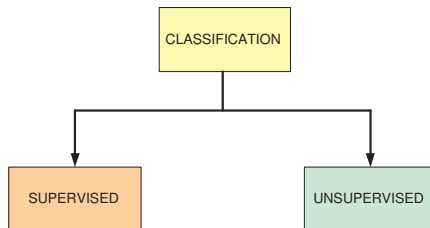
# Supervised and Unsupervised Classification

There are two main categories of classifiers :

**Supervised:** The classifier design process has a set of data samples with associated labels (class type). These are exemplars or training data.

**Unsupervised:** The given data is not labelled and the idea is to find groups in the data and the features that distinguish one group from another.

- There is a third category, namely semi-supervised classifiers, in which both labelled and unlabelled data are used for the training.



**Figure 3:** Main categories of classifiers

# Supervised Classification

## Example - from Duda-Hart-Stork Duda et al. (2001)

We are required to design a classifier for a fishing company so as to automate the sorting process. The company is interested in sorting **salmon** from **bass**. The cost of selling a salmon as bass could be high when misclassified!



Salmon

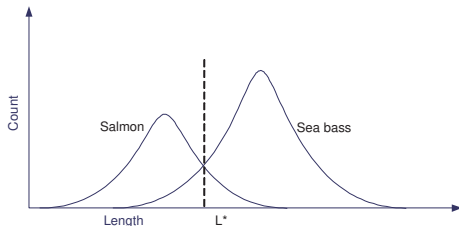


Bass

- Possible features of interest
  - length
  - width
  - number and shape of fins
  - position of mouth
  - lightness
- These features will vary because of measurement errors or conditions.

# Supervised Classification

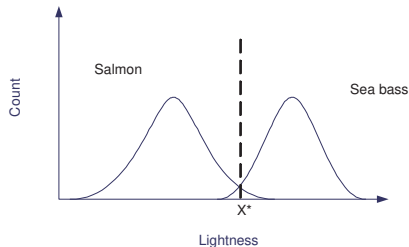
- Interesting observation may show that bass is usually longer than salmon.
- We take several samples of the two fishes and measure their lengths.
- We may represent our measurement as a histogram.
- We may ask the question, "Will this feature sufficiently classify the fishes?"



**Figure 4:** “Histogram” of fishes lengths; Length marked  $L^*$  will lead to smallest number of errors.

# Supervised Classification

- Perhaps the cost of using *length* alone to classify is too high.
- We take several samples of the two fishes and measure their *lightness*.
- We again represent our measurement as a histogram.
- The answer to the question, "Will this feature sufficiently classify the fishes?" is more satisfying.
- The  $X^*$  or  $L^*$  is a **decision threshold**.



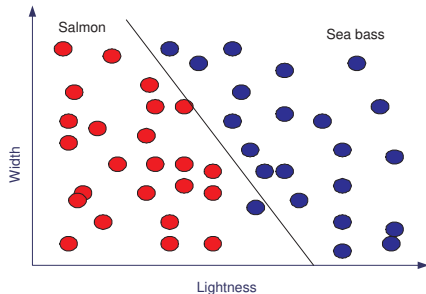
**Figure 5:** “Histogram” of lightness of fishes; Lightness marked  $X^*$  will lead to smallest number of errors.

# Supervised Classification

- Assume that we believe that we could do better at classifying the fishes by using two (2) features.
- We now have a two-dimensional feature vector,

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

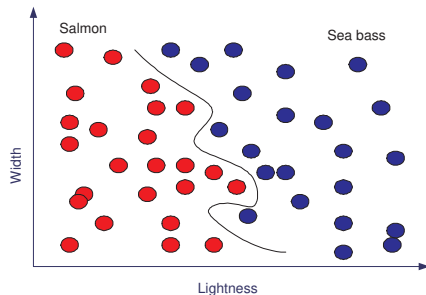
- The feature space can be visualized.
- How to obtain the “best” decision boundary is the classifier design problem.



**Figure 6:** Feature space with decision boundary of classifier.

# Supervised Classification

- As we increase the number of features there is need to deal with a high dimensional feature vector.
- The problem of “too many features” is referred to as **dimensionality curse**.
- A very complicated model may also result in over fitting - training data is separated “perfectly”; new patterns are poorly classified. Generalization problem.



**Figure 7:** Feature space with complex decision boundary of classifier.

# Supervised Classification

## Occam's Razor

The principle of using as simple as is necessary model to describe systems is captured in the so-called “Occam's razor” - favour simpler explanations over those that are needlessly complicated.

- The principle underlies the very popular method of **sparse representation**.



# Bayes decision rule - minimum error

- This approach to classification (also called discrimination) assumes that we have full knowledge of the probability density function of each class
- Let the  $C$  classes have known *a priori* probabilities,  $P(\omega_1), \dots, P(\omega_C)$
- We make use of the measurement vector  $x$  to assign  $x$  to one of the  $C$  classes
- Formulate a decision rule to **assign**  $x$  to class  $\omega_j$  if the probability of class  $\omega_j$  given the observation  $x$ , (i.e.  $P(\omega_j|x)$  - posterior probability), is the highest over all classes  $\omega_1, \dots, \omega_C$ ;

$x \in \omega_j$  if,

$$P(\omega_j|x) > P(\omega_k|x) \quad k = 1, \dots, C; \quad k \neq j$$

- Measurement space is partitioned into  $C$  regions,  $\Omega_1, \Omega_2, \dots, \Omega_C$ ;  $\mathbf{x} \in \Omega_j \Rightarrow \mathbf{x}$  is in class  $\omega_j$

# Bayes decision rule - minimum error

- We use Bayes' theorem to express the *a posteriori* probabilities  $P(\omega_j|x)$  in terms of the *a priori* probabilities and the *class-conditional density functions*  $p(x|\omega_i)$

$$P(\omega_i|x) = \frac{p(x|\omega_i)P(\omega_i)}{p(x)}$$

where

$$p(x) = \sum_{j=1}^C p(x|\omega_j)P(\omega_j)$$

- In terms of the class-conditional density we can write the decision rule as, **assign  $x$  to  $\omega_j$**  if,

$$p(x|\omega_j)P(\omega_j) > p(x|\omega_k)P(\omega_k) \quad k = 1, \dots, C \quad k \neq j$$

This is the **Bayes' rule for minimum error**.

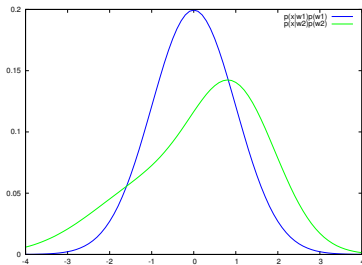
# Bayes decision rule - minimum error

- In a two-class case we can write the Bayes' minimum error rule in terms of **likelihood ratio**,  $L_r(x)$ , for  $x \in \text{class } \omega_1$ ,

$$L_r(x) = \frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}$$

- Take as an example a two-class discrimination problem, with class  $\omega_1$  normally distributed as,  $p(x|\omega_1) = N(x|0, 1)$  and class  $\omega_2$  as a normal mixture with  $p(x|\omega_2) = 0.6N(x|1, 1) + 0.4N(x|-1, 2)$

- The plots of  $p(x|\omega_i)p(\omega_i)$ ,  $i = 1, 2$  with  $P(\omega_1) = P(\omega_2) = 0.5$  are shown below;



# Bayes decision rule - minimum error

Plots of the likelihood ratio  $L_r(x)$  and threshold,  $P(\omega_2)/P(\omega_1)$  are shown below:

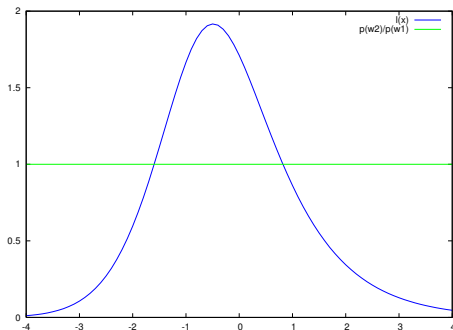


Figure 8: Likelihood function

If  $L_r(x) > \frac{P(\omega_2)}{P(\omega_1)}$ , the observed sample is classified as  $\omega_1$ ,

# Bayes decision rule - minimum risk

- This decision rule minimizes an expected loss or risk.
- Define a loss matrix,  $\Lambda$ , with components,

$\lambda_{ji}$  = cost of assigning a pattern  $x$  to  $\omega_i$  when  $x \in \omega_j$

- The **conditional risk** of assigning a pattern  $x$  to class  $\omega_i$  is defined as

$$l_i(x) = \sum_{j=1}^C \lambda_{ji} P(\omega_j|x)$$

- The average risk over decision region  $\Omega_i$  is

$$\begin{aligned} r_i &= \int_{\Omega_i} l_i(x) p(x) dx \\ &= \int_{\Omega_i} \sum_{j=1}^C \lambda_{ji} P(\omega_j|x) p(x) dx \end{aligned}$$

# Bayes decision rule - minimum risk

- The overall expected cost or risk is obtained by summing the risks associated with all the classes

$$r = \sum_{i=1}^C r_i = \sum_{i=1}^C \int_{\Omega_i} \sum_{j=1}^C \lambda_{ji} P(\omega_j|x) p(x) dx$$

- The risk is minimized if the regions  $\Omega_i$  are chosen such that if

$$\sum_{j=1}^C \lambda_{ji} P(\omega_j|x) p(x) \leq \sum_{j=1}^C \lambda_{jk} P(\omega_j|x) p(x) \quad k = 1, \dots, C$$

then  $x \in \Omega_i$ .

- This is the **Bayes decision rule for minimum risk**
- The Bayes risk,  $r^*$ , is

$$r^* = \int_x \min_{i=1, \dots, C} \sum_{j=1}^C \lambda_{ji} P(\omega_j|x) p(x) dx$$

# Bayes decision rule - minimum risk

- For a two-category classification problem we can write the conditional risks as:

$$l_i(x) = \sum_{j=1}^C \lambda_{ji} P(\omega_j|x)$$

$$l_1(x) = \lambda_{11}P(\omega_1|x) + \lambda_{21}P(\omega_2|x)$$

$$l_2(x) = \lambda_{12}P(\omega_1|x) + \lambda_{22}P(\omega_2|x)$$

- The minimum risk decision rule is simply to decide  $\omega_1$  if  $l_1(x) < l_2(x)$ .
- This can be expressed in terms of posterior probabilities as: Decide  $\omega_1$  if

$$(\lambda_{11} - \lambda_{12})P(\omega_1|x) < (\lambda_{22} - \lambda_{21})P(\omega_2|x)$$

- In terms of the prior probabilities and conditional densities we decide  $\omega_1$  if,

$$(\lambda_{11} - \lambda_{12})p(x|\omega_1)P(\omega_1) < (\lambda_{22} - \lambda_{21})p(x|\omega_2)P(\omega_2)$$

# Bayes decision rule - minimum risk

- If we consider the special case of equal cost (also called symmetrical or zero-one) loss matrix,  $\Lambda$ , in which,

$$\lambda_{ij} = \begin{cases} 1, & i \neq j; \\ 0, & i = j \end{cases}$$

a substitution of this condition into the Bayes decision rule for minimum risk, gives,

$$\sum_{j=1}^C P(\omega_j|x)p(x) - P(\omega_i|x)p(x) \leq \sum_{j=1}^C P(\omega_j|x)p(x) - P(\omega_k|x)p(x)$$

for  $k = 1, \dots, C$ . This is easily simplified as,

$$p(x|\omega_i)p(\omega_i) \geq p(x|\omega_k)p(\omega_k), \quad k = 1, \dots, C$$

when  $x \in$  class  $\omega_i$ .

This is the same as Bayes rule for minimum error.



# Bayes decision rule - minimum risk

- If we consider the special case of zero-one loss matrix,  $\Lambda$ , in which,

$$\lambda_{ij} = \begin{cases} 1, & i \neq j; \\ 0, & i = j \end{cases}$$

and a two-category classification the Bayes decision rule for minimum risk, gives,

$$p(x|\omega_1)p(\omega_1) \geq p(x|\omega_2)p(\omega_2),$$

when  $x \in$  class  $\omega_1$ .

This is the same as Bayes rule for minimum error in the two-category case.

- The corresponding risks in the case of the zero-one loss matrix are

$$\begin{aligned} I_i(x) &= \sum_{j=1}^C \lambda_{ij} P(\omega_j|x) \\ &= \sum_{j \neq i} P(\omega_j|x) \\ &= 1 - P(\omega_i|x) \end{aligned}$$

# Discriminant Functions

- Bayes decision rules requires knowledge of prior class probabilities and class conditional densities which are often not available in practice and must be estimated from data
- The class of techniques being introduced makes no assumption about  $p(x|\omega_i)$  but rather assumes a form of the **discriminant** functions
- A **discriminant** function is a function of the feature vector  $x$  that leads to a classification rule
- Consider a two-class problem, a discriminant function  $h(x)$  is such that

$$\begin{aligned}h(x) > k &\Rightarrow x \in \omega_1 \\ < k &\Rightarrow x \in \omega_2\end{aligned}$$

for a constant  $k$

# Discriminant Functions

- Discriminant functions are not unique. If  $f(\cdot)$  is a monotonic function, then

$$\begin{aligned}g(x) = f(h(x)) > k' &\Rightarrow x \in \omega_1 \\ &< k' \Rightarrow x \in \omega_2\end{aligned}$$

where  $k' = f(k)$ , gives the same decision as  $h(x)$ .

- For classification problem with  $C$  classes we define  $C$  discriminant functions,  $g_i(x)$  such that,

$$g_i(x) > g_j(x) \Rightarrow x \in \omega_i \quad j = 1, \dots, C; j \neq i$$

This implies that a feature vector is assigned to the class with the largest discriminant.

- The **discriminant** techniques rely on the **form of the function being specified** and not on the underlying distribution
- Parameters of the functional form are adjusted by a **training procedure**

# Linear discriminant functions

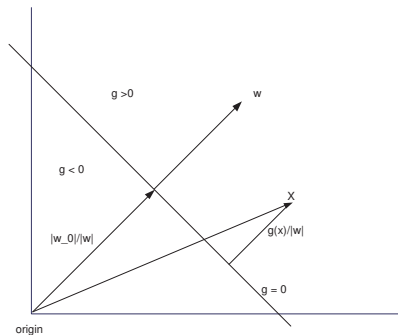
- Linear discriminant functions are a linear combination of the components of the measurement (or feature) vector,  $\mathbf{x} = [x_1, x_2, \dots, x_p]^t$ , such that,

$$g(\mathbf{x}) = \omega^t \mathbf{x} + \omega_0 = \sum_{i=1}^p \omega_i x_i + \omega_0$$

where we need to specify the *weight vector*  $\omega$  and *threshold weight*  $\omega_0$

- The equation describes a *hyperplane* with unit normal in the direction of  $\omega$  and a perpendicular distance,  $|\omega_0|/|\omega|$  from origin.

# Linear discriminant functions



**Figure 9:** Geometry of linear discriminant function

- The value of the discriminant function for a pattern  $x$  is the perpendicular distance from the hyperplane

# Linear discriminant functions

- Classifiers that use a linear discriminant function are called **linear machine**.
- The **minimum-distance classifier** is an example. It uses the nearest-neighbour decision rule.
- Let the prototype points of the classifier be  $p_1, \dots, p_C$ . Each point represents a class,  $\omega_j$ . The minimum distance classifier assigns  $x$  to the  $\omega_i$  with nearest point  $p_i$

$$\|x - p_i\|^2 = x^t x - 2x^t p_i + p_i^t p_i$$

The class assigned to  $x$  is

$$\omega_i = \max_i (x^t p_i - \frac{1}{2} p_i^t p_i)$$

# Linear discriminant functions

- We can relate this assignment to the linear discriminant function

$$g_i(x) = \omega_i^t x + \omega_{i0}$$

where

$$\omega_i = p_i \tag{1}$$

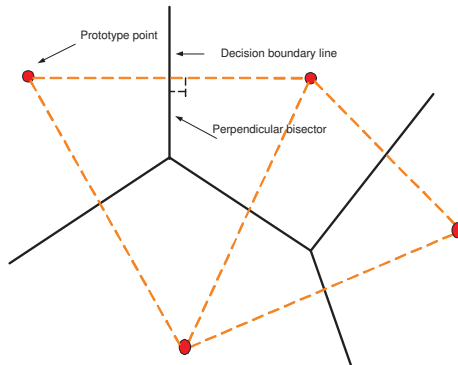
$$\omega_{i0} = -\frac{1}{2} \|p_i\|^2 \tag{2}$$

to show that it is indeed a linear machine.

- The prototype points could be chosen as the mean of each class and we have a **nearest class mean classifier**.

# Linear discriminant functions

- Each boundary is the perpendicular bisector of the lines joining the prototype points of regions that are contiguous.
- Note also that the decision regions of a linear machine are always convex.

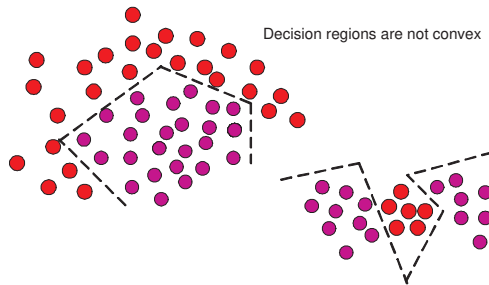


**Figure 10:** Decision regions of minimum distance classifier



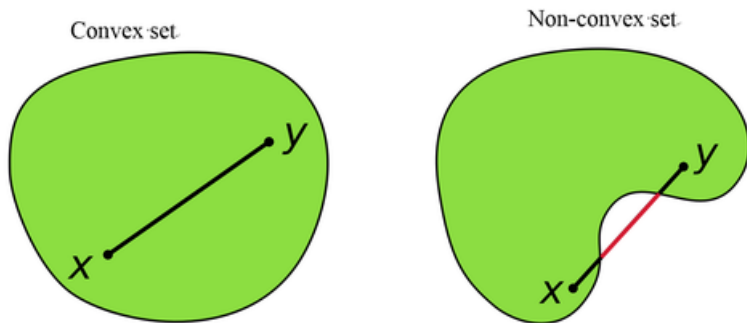
# Piecewise linear discriminant functions

- The linear machine has a simple form but suffers the limitation of not being able to separate situations where the decision regions have to be non-convex.
- The examples below show two-class problems where a linear discriminant will fail to separate. They require **piece-wise linear discriminant functions**.



**Figure 11:** Groups not separable by linear discriminant functions

# Piecewise linear discriminant functions



**Figure 12:** Quick illustration of convex and non-convex regions

# Piecewise linear discriminant functions

- We may solve the previous two-class problem by using piece-wise linear discriminant function to generalize the minimum-distance classifier.
- We allow **more than one prototype** for each class.
- Suppose there are  $n_i$  prototypes in class  $\omega_i$ ,  $p_i^1, \dots, p_i^{n_i}$ ,  $i = 1, \dots, C$ .
- The discriminant function which assigns pattern  $x$  to class,  $\omega_i$ , is defined as

$$g_i(x) = \max_{j=1, \dots, n_i} g_i^j(x)$$

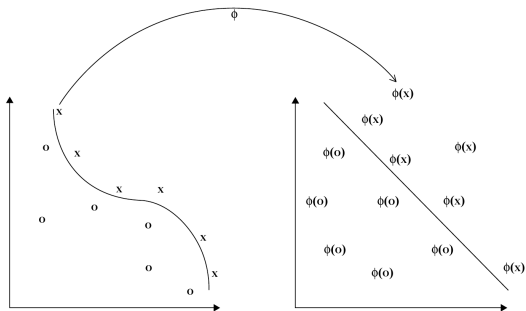
where  $g_i^j$  is a linear subsidiary discriminant function, given by

$$g_i^j(x) = x^t p_i^j - \frac{1}{2} p_i^{j^t}, \quad j = 1, \dots, n_i; \quad i = 1, \dots, C$$

## Main idea

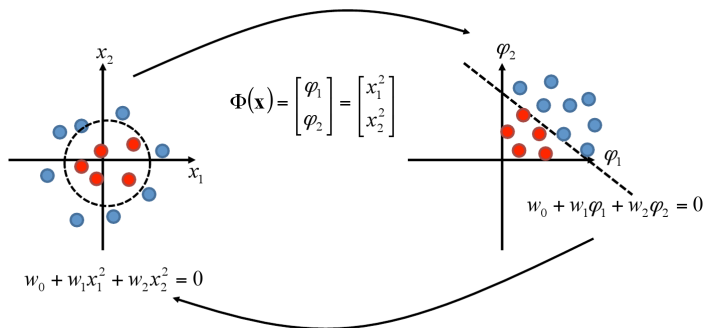
- Embed given data into a space where the patterns can be discovered as linear relations
- Two steps:
  - Mapping is defined implicitly by a so-called kernel function (depends on domain knowledge regarding pattern in data source)
  - Use general purpose algorithm that is robust
- Algorithm is efficient, requiring computational resource that is polynomial in the size and number of data items; dimension of the embedding space grows exponentially and does not affect computational burden

# Kernel methods



**Figure 13:** Function  $\phi$  embeds data into a feature space; nonlinear pattern now appears linear

# Kernel methods



**Figure 13:** Example of nonlinear mapping in classification problem

# Ridge regression revisited

- Formulation is somewhat similar to linear regression;
- This model gives better theoretical guarantees and improved performance partly because we constrain the norm of the weight vector. The optimization problem is written compactly as:

$$\min_{\mathbf{W}} F(\mathbf{W}) = \lambda \|\mathbf{W}\|^2 + \|\mathbf{X}^T \mathbf{W} - \mathbf{Y}\|^2 \quad (3)$$

where  $\lambda$  is a positive parameter that determines the trade-off between the regularization term  $\|\mathbf{W}\|^2$  and the empirical mean squared error;  $\mathbf{X} \in \mathbb{R}^{N \times m}$  is the matrix of feature vectors,  $\mathbf{X} = [x_1, \dots, x_m]$  and  $\mathbf{W}$  and  $\mathbf{Y}$  are as defined previously (see lecture slides on regression).

- Optimization problem of Equation is convex and differentiable with a global minimum if and only if

$$\nabla F(\mathbf{W}) = 0 \Leftrightarrow (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})\mathbf{W} = \mathbf{X}\mathbf{Y} \Leftrightarrow \mathbf{W} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{Y} \quad (4)$$

$\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}$  is always invertible<sup>1</sup>

- The form of the solution weight vector

$$\mathbf{W} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{Y}$$

is the primal solution

---

<sup>1</sup>because its eigenvalues are sum of non-negative eigenvalues of positive semi-definite matrix  $\mathbf{X}\mathbf{X}^T$  and  $\lambda > 0$

# Ridge regression revisited

- The solution,

$$\mathbf{W} = (\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I})^{-1}\mathbf{X}\mathbf{Y}$$

can be written as:

$$\mathbf{W} = \mathbf{X}\boldsymbol{\alpha} \quad (5)$$

where  $\boldsymbol{\alpha}$  is

$$(\mathbf{G} + \lambda\mathbf{I})^{-1}\mathbf{Y} \quad (6)$$

and  $\mathbf{G} = \mathbf{X}^T\mathbf{X}$  is called the **Gram** matrix. Each component of  $\mathbf{G}$  is an inner product,  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$ . Thus

$$\mathbf{G}_{i,j} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

- The prediction function is

$$g(\mathbf{x}) = \langle \mathbf{W}, \mathbf{x} \rangle = \left\langle \sum_{i=1}^m \alpha_i \mathbf{x}_i, \mathbf{x} \right\rangle = \sum_{i=1}^m \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle = \mathbf{Y}^T (\mathbf{G} + \lambda\mathbf{I})^{-1} \mathbf{k}$$

where  $k_i = \langle \mathbf{x}_i, \mathbf{x} \rangle$



# Ridge regression revisited

- We have two forms of solution for the ridge regression (similarly for other regression):

$$\begin{aligned} \mathbf{W} &= (\mathbf{X}\mathbf{X}^T + \lambda\mathbf{I})^{-1} \mathbf{X}\mathbf{Y} : && \text{Primal form} \\ \mathbf{W} &= \mathbf{X}\boldsymbol{\alpha} : && \text{Dual form} \end{aligned} \tag{7}$$

Primal form computes explicitly while the dual expresses solution as linear combination of training samples.

- In primal form we solve an  $(N \times N)$  system of equations while in the dual form we solve an  $(m \times m)$  system.
- If dimension of features is  $N \gg m$  (the number of samples) the computational advantage is obvious
- Key observation: Ridge regression algorithm can be solved in a form that only requires inner products between sample points

# Kernel-defined nonlinear mapping

- Consider an embedding map

$$\phi : \mathbf{x} \in \mathbb{R}^N \mapsto \phi(\mathbf{x}) \in \mathbf{F} \subseteq \mathbb{R}^N \quad (8)$$

- Choose the map  $\phi$  so that it aims to convert the nonlinear relations into linear ones.
- Map  $\phi$  recodes the given dataset,  $S$ , into  $\{(\phi(\mathbf{x}_i), y_i), \dots (\phi(\mathbf{x}_m), y_m)\}$  for the  $m$  samples in the dataset.
- Recall that the efficient **dual** form of the solution to ridge regression entails the **Gram** matrix which is made up inner products

$$G_{i,j} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (9)$$

- Computational cost of  $\alpha$  is  $\mathcal{O}(m^3 + m^2N)$  and that of evaluating the predictor on a new sample is  $\mathcal{O}(mN)$
- It turns out that the inner product can be computed **directly** in the input space rather than first computing  $\phi(\mathbf{x})$  using a **kernel function**

# Kernel-defined nonlinear mapping

## Definition:

A **Kernel** is a function,  $\kappa$ , that for all  $\mathbf{x}, \mathbf{z} \in S$  satisfies

$$\kappa(\mathbf{x}, \mathbf{z}) = \langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle ,$$

where  $\phi$  is a mapping from  $S$  to an inner product feature space  $F$

$$\phi : \mathbf{x} \mapsto \phi(\mathbf{x}) \in F$$

# Kernel-defined nonlinear mapping

## Example

Consider a two-dimensional input space  $\mathbf{x} \subseteq \mathbb{R}^2$  together with the feature map

$$\phi : \mathbf{x} = (x_1, x_2) \mapsto \phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \in F = \mathbb{R}^3.$$

Hypothesis of linear functions in  $F$  are of the form

$$g(\mathbf{x}) = w_{11}x_1^2 + w_{22}x_2^2 + w_{12}\sqrt{2}x_1x_2$$

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle &= \langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (z_1^2, z_2^2, \sqrt{2}z_1z_2) \rangle \\ &= x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (x_1z_1 + x_2z_2)^2 = \langle \mathbf{x}, \mathbf{z} \rangle^2\end{aligned}$$

- Hence  $\kappa(\mathbf{x}, \mathbf{z}) = \langle \mathbf{x}, \mathbf{z} \rangle^2$  is a kernel function with  $F$  its corresponding feature space

# Kernel-defined nonlinear mapping

## Example

Consider a two-dimensional input space  $\mathbf{x} \subseteq \mathbb{R}^2$  together with the feature map

$$\phi : \mathbf{x} = (x_1, x_2) \mapsto \phi(\mathbf{x}) = (x_1^2, x_2^2, x_1 x_2, x_2 x_1) \in F = \mathbb{R}^4.$$

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle &= \langle (x_1^2, x_2^2, x_1 x_2, x_2 x_1), (z_1^2, z_2^2, z_1 z_2, z_2 z_1) \rangle \\ &= x_1^2 z_1^2 + x_2^2 z_2^2 + x_1 x_2 z_1 z_2 + x_2 x_1 z_2 z_1 \\ &= (x_1 z_1 + x_2 z_2)^2 = \langle \mathbf{x}, \mathbf{z} \rangle^2\end{aligned}$$

- The same kernel computes the inner product of this feature space
- Hence feature space is not uniquely determined by the kernel function

# Kernel-defined nonlinear mapping

**Table 1:** Commonly used kernels ( $\gamma$ ,  $r$ , and  $d$  are parameters)

Nonlinearity	Mathematical form: $\kappa(\mathbf{x}_i, \mathbf{x}_j) =$
Linear	$\langle \mathbf{x}_i, \mathbf{x}_j \rangle$
Polynomial	$(\gamma \langle \mathbf{x}_i, \mathbf{x}_j \rangle + r)^d; \quad \gamma > 0$
Gaussian (Radial Basis Function - RBF)	$\exp(-\gamma \ \mathbf{x}_i - \mathbf{x}_j\ ^2); \quad \gamma > 0$
Sigmoid	$\tanh(\gamma \langle \mathbf{x}_i, \mathbf{x}_j \rangle + r)$

## Example

Let  $\mathbf{x} = [1 \quad 4 \quad 6]^T$  and  $\mathbf{z} = [3 \quad 5 \quad 2]^T$  be two feature vectors that we need to map through the mapping function  $\phi(\cdot)$  to some feature space  $\mathcal{F}$ . Further, let the kernel associated with the feature space be the RBF with parameter  $\gamma = 0.2222$ . The value of the inner product of the two vectors,  $\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle$ , in the feature space is easily computed as  $\kappa(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|^2)$

$$\begin{aligned}\langle \phi(\mathbf{x}), \phi(\mathbf{z}) \rangle &= \kappa(\mathbf{x}, \mathbf{z}) \\ &= \exp(-0.2222 \times \|\mathbf{x} - \mathbf{z}\|^2) \\ &= \exp(-4.667) = 9.4 \times 10^{-3}\end{aligned}$$

Try this simple example yourself to ensure you understand the concept

# Kernel-defined nonlinear mapping

- Recall Equations (5) and (6) and the comments that followed; each entry of the **Gram** matrix that appeared in the solution of the ridge regression is an inner product of the data in the input space.
- A mapping,  $\phi(\cdot)$ , into a high dimensional feature space  $\mathcal{F}$ , implies that each entry of the **Gram** matrix can be computed with the appropriate kernel of the feature space

$$\begin{aligned}\mathbf{G}_{i,j} &= \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \\ &= \kappa(\mathbf{x}_i, \mathbf{x}_j)\end{aligned}\tag{10}$$

- Kernalization provides a way of dealing with nonlinear relationships that may exist in the problem (e.g. regression, classification, dimensionality reduction, etc.)

# Support vector machine

- Consider a binary classification task with data points  $\mathbf{x}_i (i = 1, \dots, m)$  having corresponding labels  $y_i = \pm 1$  and decision function

$$g(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b) \quad (11)$$

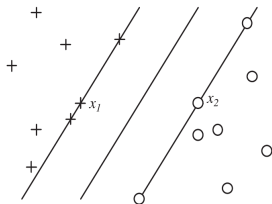
- For a separable dataset all data will be correctly classified if  $y_i(\langle \mathbf{w}, \mathbf{x} \rangle + b) > 0; \forall i$
- Define **canonical hyperplane** such that  $\langle \mathbf{w}, \mathbf{x} \rangle + b = 1$  for closest points on one side of separating plane and  $\langle \mathbf{w}, \mathbf{x} \rangle + b = -1$  for closest points on the other
- Separating plane :  $\langle \mathbf{w}, \mathbf{x} \rangle + b = 0$  and normal vector is  $\frac{\mathbf{w}}{\|\mathbf{w}\|}$
- Margin is given by projection of  $\mathbf{x}_1 - \mathbf{x}_2$  onto the separating plane
- $\langle \mathbf{w}, \mathbf{x}_1 \rangle + b = 1$  and  $\langle \mathbf{w}, \mathbf{x}_2 \rangle + b = -1$ ; margin is  $\gamma = 1/\|\mathbf{w}\|$  (see Fig 14 on next slide)



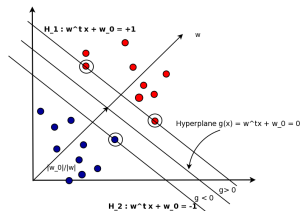
# Support vector machine

Maximize the margin  $\gamma = 1/||\mathbf{w}||$ , by

$$\begin{aligned} & \min \left[ \frac{1}{2} ||\mathbf{w}'||^2 \right] \\ & \text{subject to } y_i(\langle \mathbf{w}', \mathbf{x} \rangle + b) \geq 1; \forall i. \end{aligned}$$



**Figure 14:** SVM separating hyperplane



**Figure 15:** SVM separating hyperplane in detail

# Support vector machine

- Learning task reduces to minimizing the primal objective function:

$$\mathcal{L} = \frac{1}{2}(\langle \mathbf{w}, \mathbf{w} \rangle) - \sum_{i=1}^m \alpha_i (y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1) \quad (12)$$

$\alpha_i$  are Lagrange multipliers and  $\alpha_i \geq 0$

- After taking derivatives with respect to  $b$  and  $\mathbf{w}$  and appropriate substitution into Equ. 12, we obtain the dual objective function

$$\mathcal{W}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad (13)$$

to be maximized with respect to the  $\alpha_i$  subject to

$$\alpha_i \geq 0 \quad \sum_{i=1}^m \alpha_i y_i = 0 \quad (14)$$

- The quadratic program represented by Equ.13 gives optimal separating hyperplane with maximal margin for separable data

# Support vector machine

- Equ.13 indicates how we can incorporate the **kernel** that computes inner product in the feature space after applying a mapping  $\phi(\cdot)$ ; case of inseparable data
- Mapping achieved by

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle \mapsto \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle \quad (15)$$

- Note that the functional form of  $\phi(\mathbf{x}_i)$  need not be known since choice of kernel implicitly defines it:

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$

- Kernelized version of Equ. 13 becomes

$$\mathcal{W}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \quad (16)$$

to be maximized with respect to the  $\alpha_i$  subject to

$$\alpha_i \geq 0 \quad \sum_{i=1}^m \alpha_i y_i = 0$$

- For example the kernel could be chosen as  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$
- Decision function for a new test data  $\mathbf{z}$  is

$$f(\mathbf{z}) = \text{sign} \left( \sum_{i=1}^m y_i \alpha_i \kappa(\mathbf{x}_i, \mathbf{z}) + b \right) \quad (17)$$

# Support vector machine with soft margin

- Noisy data and outliers can lead to poor generalization
- Reduce the effect by introducing **soft margin**
- Recall (Equ. 16) that for a given kernel ( $\kappa(\cdot, \cdot)$ ), the learning task is:

$$\begin{aligned} \underset{\alpha}{\text{maximize}} \quad \mathcal{W}(\alpha) &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \kappa(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \quad \alpha_i &\geq 0 \quad \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned} \tag{18}$$

- Two ways of accounting for noisy data and outlier with
  - Use  $L_1$  norm error and introduce box constrain  $0 \leq \alpha_i \leq C$  in Equ (18)
  - Use  $L_2$  norm error and add a small positive constant to the leading diagonal of kernel matrix  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$  becomes  $\kappa(\mathbf{x}_i, \mathbf{x}_j) + \lambda$  in Equ (18)
- Parameters  $C$  and  $\lambda$  are chosen to trade-off between **training error** and **generalization** ability; achieve this with validation set; well-known library *libsvm* exposes interface to determine  $C$

# Using SVM

Beginners are advised to use the following procedure (see Hsu et al. (2003-2016)):

- Transform data to the format of an SVM package
- Conduct simple scaling on the data
- Consider the RBF kernel  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$  (see Table 42)
- Use cross-validation to find the best parameter  $C$  and  $\gamma$
- Use the best parameter  $C$  and  $\gamma$  to train the whole training set
- Test

The guide given by Hsu et al. (2003-2016) is concise and straightforward, and students are strongly encouraged to consult the publication for details.

# Bibliography

- Alpaydin, E. (2010), *Introduction to Machine Learning*, second edn, The MIT Press, Cambridge Massachusetts.
- Giampiconi, L., Elwood, A., Leonardi, M., Mohamed, A. & Rozza, A. (2023), 'A survey and taxonomy of loss functions in machine learning', online.
- Duda, R. O., Hart, P. E. & Stork, D. G. (2001), *Pattern Classification*, Second edn, John Wiley and Sons.
- Gönen, M. & Alpaydin, E. (2011), 'Multiple kernel learning algorithms', *Journal of Machine Learning Research* **12**(2011), 2211–2268.
- Goodfellow, I., Bengio, Y. & Courville, A. (2016a), *Deep Learning*, MIT Press.
- Goodfellow, I., Bengio, Y. & Courville, A. (2016b), *Deep Learning*, MIT Press.
- Hastie, T., Tibshirani, R. & Friedman, J. (2001), *The Elements of Statistical Learning - Data Mining, Inference and Prediction*, Springer Science+Business Media LLC.
- Hsu, C.-W., Chang, C.-C. & Lin, C.-J. (2003-2016), Practical guide to support vector classification, Technical report, Department of Computer Science, National Taiwan University, Taipei 106, Taiwan.
- Izenman, A. J. (2008), *Modern Multivariate Statistical Techniques - Regression, Classification and Manifold Learning*, Springer Science+Business Media LLC.
- Kelleher, J. D., Namee, B. M. & D'Arcy, A. (2015), *Fundamentals of Machine Learning for Predictive Data Analytics - Algorithms, Worked Examples and Case Studies*, The MIT Press, Cambridge Massachusetts.
- Mitchell, T. M. (1997), *Machine Learning*, WCB McGraw-Hill.
- Mohri, M., Rostamizadeh, A. & Talwalkar, A. (2012), *Foundations of Machine Learning*, MIT Press.
- Parnami, A. & Lee, M. (2022), 'Learning from few examples: A summary of approaches to few-shot learning', online.
- Webb, A. (2002), *Statistical Pattern Recognition*, Second edn, John Wiley and Sons.
- Weinberger, K. Q., Shi, F. & Saul, L. K. (2004), Learning a kernel matrix for nonlinear dimensionality reduction, in 'Proceedings, 21st International Conference on Machine Learning', Banff, Canada.