

Pattern Classification

Lecture 06: Linear Discriminant Functions

Kundan Kumar

<https://github.com/erkundanec/PatternClassification>

Linear Discriminant Functions

Introduction

- In parametric estimation, we assumed that the forms for the underlying **probability densities were known**, and used the training samples to estimate the values of their parameters.
- Instead, assume that the proper forms for the discriminant functions is known, and use the samples to estimate the values of parameters of the classifier.
- None of the various procedures for determining discriminant functions require knowledge of the forms of underlying probability distributions so called **nonparametric** approach.
- Linear discriminant functions are relatively **easy to compute** and estimate the form using training samples.

Linear discriminant functions and decisions surfaces

- A **discriminant function** is a linear combination of the components of \mathbf{x} can be written as

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

where \mathbf{w} is the **weight vector** and w_0 the **bias** or **threshold** weight.

- The equation $g(\mathbf{x}) = 0$ defines the **decision surface** that separates points from different classes.
- Linear discriminant functions are going to be studied for
 - two-category case,
 - multi-category case, and
 - general case

For the general case there will be c such discriminant functions, one for each of c categories.

Two-Category Case

- A two-category classifier with a discriminant function of the form $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ uses the following rule:

$$\text{Decide } \begin{cases} \omega_1 & \text{if } g(\mathbf{x}) > 0 \\ \omega_2 & \text{otherwise} \end{cases}$$

- Thus, \mathbf{x} is assigned to ω_1 if the **inner product** $\mathbf{w}^T \mathbf{x}$ exceeds the threshold $-w_0$ and to ω_2 otherwise.
- If $g(\mathbf{x}) = 0$, \mathbf{x} can ordinarily be assigned to either class, or can be left undefined.

A simple linear classifier

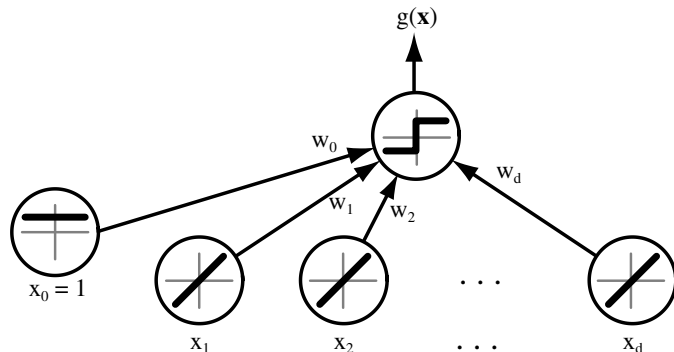
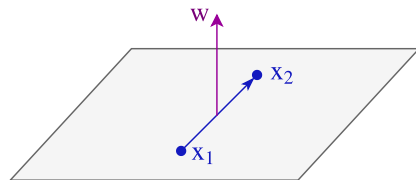


Figure: A simple linear classifier having d input units, each corresponding to the values of the components of an input vector. Each input feature value x_i is multiplied by its corresponding weight w_i ; the output unit sums all these products and emits $+1$ if $\mathbf{w}^T \mathbf{x} + w_0 > 0$ or -1 otherwise

Two-Category Case

- The equation $g(x) = 0$ defines the decision surface that separates points assigned to the category ω_1 from points assigned to the category ω_2
- When $g(x)$ is linear, the decision surface is a hyperplane.
- If x_1 and x_2 are both on the decision surface, then

$$\begin{aligned}w^T x_1 + w_0 &= w^T x_2 + w_0 \\ \Rightarrow w^T (x_1 - x_2) &= 0\end{aligned}$$



- This shows that w is normal to any vector lying in the hyperplane.

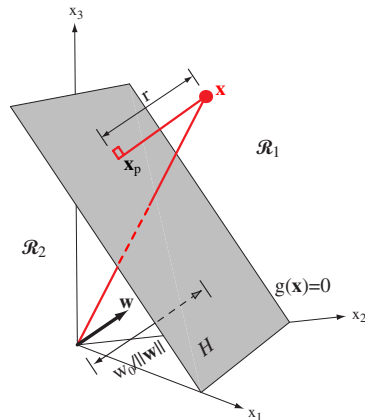
Two-Category Case

- The discriminant function $g(x)$ gives an algebraic measure of the distance from x to the hyperplane. The easiest way to see this is to express x as

$$x = x_p + r \frac{w}{\|w\|}$$

- where x_p is the normal projection of x onto H , and r is the desired algebraic distance which is positive if x is on the positive side and negative if x is on the negative side.
- Because, $g(x_p) = 0$

$$r = \frac{g(x)}{\|w\|}$$



Two-Category Case

- The distance from the origin to H is given by $\frac{w_0}{\|w\|}$.
- If $w_0 > 0$, the origin is on the positive side of H , and if $w_0 < 0$, it is on the negative side.
- If $w_0 = 0$, then $g(x)$ has the homogeneous form $w^T x$, and the hyperplane passes through the origin.

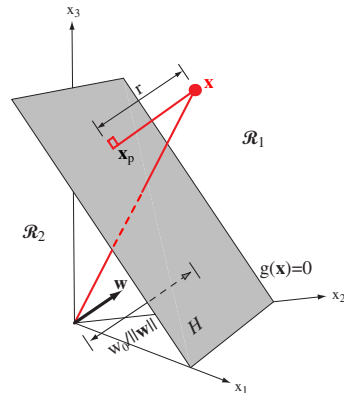


Figure: The linear decision boundary H , where $g(x) = w^T x + w_0$, separates the feature space into two half-spaces \mathcal{R}_1 (where $g(x) > 0$) and \mathcal{R}_2 (where $g(x) < 0$)

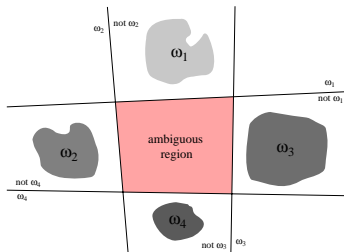
Two-Category Case

- In conclusion, a linear discriminant function divides the feature space by a hyperplane decision surface.
- The orientation of the surface is determined by the normal vector w and the location of the surface is determined by the bias w_0 .
- The discriminant function $g(x)$ is proportional to the signed distance from x to the hyperplane, with $g(x) > 0$ when x is on the positive side, and $g(x) < 0$ when x is on the negative side.

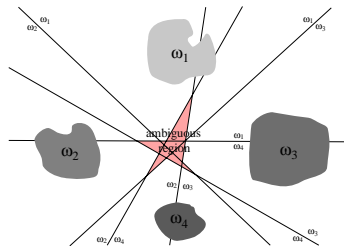
Multi-category case

- There is more than one way to devise multi-category classifiers employing linear discriminant functions.

- c two-class problem
(one-vs-rest)



- $c(c-1)/2$ linear discriminants, one for every pair of classes (one-vs-one).



- Pink regions have ambiguous category assignment.

Multi-category case

- More effective way is to define c linear discriminant functions

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0} \quad i = 1, 2, \dots, c$$

and assign \mathbf{x} to ω_i if $g_i(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq i$; in case of ties, the classification is undefined

- In this case, resulting classifier is a “*linear machine*”.
- A linear machine divides the feature space into c decision regions, with $g_i(\mathbf{x})$ being the largest discriminant if \mathbf{x} is in the region \mathcal{R}_i .
- For a two contiguous regions \mathcal{R}_i and \mathcal{R}_j ; the boundary that separates them is a portion of hyperplane H_{ij} defined by:

$$g_i(\mathbf{x}) = g_j(\mathbf{x}) \quad \text{or} \quad (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) = 0$$

Multi-category case

- It follows at once that $\mathbf{w}_i - \mathbf{w}_j$ is normal to H_{ij} , and the signed distance from \mathbf{x} to H_{ij} is given by

$$r = \frac{(g_i(\mathbf{x}) - g_j(\mathbf{x}))}{\|\mathbf{w}_i - \mathbf{w}_j\|}$$

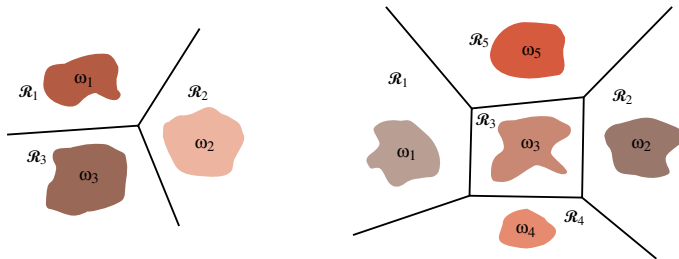


Figure: Decision boundaries produced by a linear machine for a three-class problem and a five-class problem

Generalized Linear Discriminant Functions

Generalized Linear Discriminant Functions

- The linear discriminant function $g(\mathbf{x})$ is defined as

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 \quad (1)$$

$$= w_0 + \sum_{i=1}^d w_i x_i \quad (2)$$

where $\mathbf{w} = [w_1, \dots, w_d]^T$, and $\mathbf{x} = [x_1, x_2, \dots, x_d]^T$

- We can obtain the *quadratic discriminant function* by adding second-order terms as

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j \quad (3)$$

Because $x_i x_j = x_j x_i$, we can assume that $w_{ij} = w_{ji}$ with no loss in generality. Which result in more complicated decision boundaries. (*hyperquadrics*)

Generalized Linear Discriminant Functions

- The quadratic discriminant function has an additional $d(d+1)/2$ coefficients at its disposal with which to produce more complicated separating surfaces.
- The separating surface defined by $g(\mathbf{x}) = 0$ is a second-degree or hyperquadric surface.
- If the symmetric matrix, $W = [w_{ij}]$, is nonsingular, the linear term in $g(\mathbf{x})$ can be eliminated by translating the axes.

Generalized Linear Discriminant Functions

- The basic character of the separating surface can be described in terms of scaled matrix

$$\bar{W} = \frac{W}{w^T W^{-1} w - 4w_0}$$

where $w = (w_1, \dots, w_d)^T$ and $W = [w_{ij}]$

- The types of quadratic separating surfaces that arise in the general multivariate Gaussian case are as follows
 1. If \bar{W} is a positive multiple of the identity matrix, the separating surface is a *hypersphere* such that $\bar{W} = kI$.
 2. If \bar{W} is positive definite, the separating surfaces is a *hyperellipsoid* whose axes are in the direction of the eigenvectors of \bar{W} .
 3. If none of the above cases holds, that is, some of the eigenvalues of are positive and others are negative, the surface is one of the varieties of types of *hyperhyperboloids*.

Generalized Linear Discriminant Functions

- By continuing to add terms such as $w_{ijk}x_ix_jx_k$, we can obtain the class of *polynomial discriminant functions*. These can be thought of as truncated series expansions of some arbitrary $g(\mathbf{x})$, and this in turn suggest the *generalized linear discriminant function*.

$$g(\mathbf{x}) = \sum_{i=1}^{\hat{d}} a_i y_i(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$$

where \mathbf{a} is a \hat{d} -dimensional weight vector and \hat{d} functions $y_i(\mathbf{x})$ are arbitrary functions of \mathbf{x} .

- The physical interpretation is that the functions $y_i(\mathbf{x})$ map points \mathbf{x} from d -dimensional space to point \mathbf{y} in \hat{d} -dimensional space.
- The resulting discriminant function is not linear in \mathbf{x} , but it is linear in \mathbf{y} .

Generalized Linear Discriminant Functions

- Then, the discriminant $g(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$ separates points in the transformed space using a hyperplane passing through the origin.
- The mapping to a higher dimensional space may increase the complexity of the learning algorithms.
- However, certain assumptions can make the problem tractable.
- Let the quadratic discriminant function be

$$g(\mathbf{x}) = a_1 + a_2x + a_3x^2$$

- So that the three-dimensional vector \mathbf{y} is given by

$$\mathbf{y} = [1 \quad x \quad x^2]^T$$

Generalized Linear Discriminant Functions

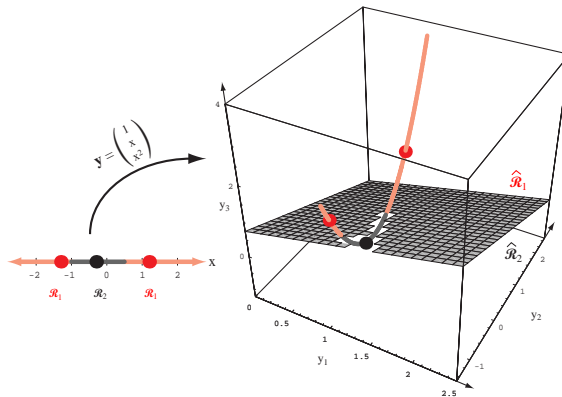


Figure: The mapping $y = (1 \ x \ x^2)^T$ takes a line and transforms it to a parabola in three dimensions. A plane splits the resulting y space into regions corresponding to two categories, and this in turn gives a non-simply connected decision region in the one-dimensional x space.

Generalized Linear Discriminant Functions

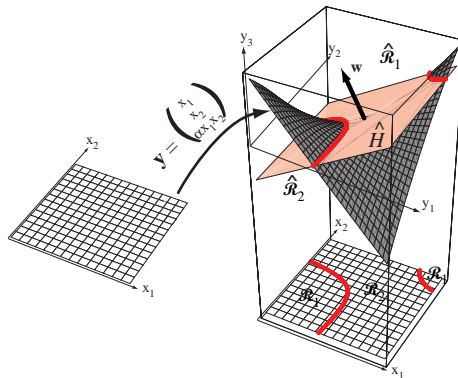


Figure: The two-dimensional input space x is mapped through a polynomial function f to y . Here the mapping is $y_1 = x_1$, $y_2 = x_2$ and $y_3 \propto x_1 x_2$. A linear discriminant in this transformed space is a hyperplane, which cuts the surface. Points to the positive side of the hyperplane \hat{H} correspond to category ω_1 , and those beneath it ω_2 . Here, in terms of the x space, \mathcal{R}_1 is not simply connected.

Problem to be solved

Question:

The following three decision functions are given for a three-class problem.

$$g_1(\mathbf{x}) = 10x_1 - x_2 - 10 = 0$$

$$g_2(\mathbf{x}) = x_1 + 2x_2 - 10 = 0$$

$$g_3(\mathbf{x}) = x_1 - 2x_2 - 10 = 0$$

- Sketch the decision boundary and regions for each pattern class.
- Assuming that each pattern class is pairwise linearly separable from every other class by a distinct decision surface and letting

$$g_{12}(\mathbf{x}) = g_1(\mathbf{x})$$

$$g_{13}(\mathbf{x}) = g_2(\mathbf{x})$$

$$g_{23}(\mathbf{x}) = g_3(\mathbf{x})$$

as listed above, sketch the decision boundary and regions for each pattern class.

References

- [1] Richard O Duda, Peter E Hart, and David G Stork. *Pattern classification*. John Wiley & Sons, 2012.



Thank you!