

Machine Learning (521289S)

Linear Discriminant Functions

M.Sc. Antti Isosalo

Physiological Signal Analysis Team
Center for Machine Vision and Signal Analysis (CMVS)
University of Oulu

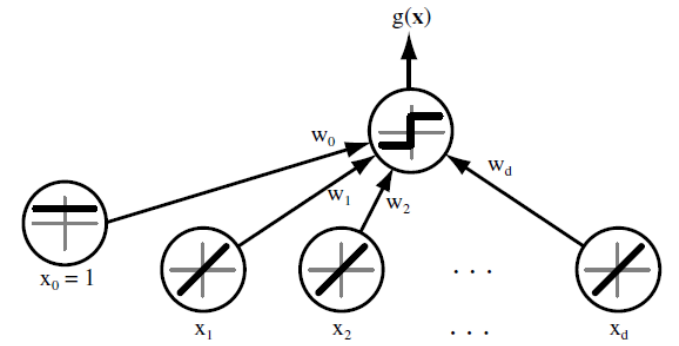
Linear Discriminant Functions

- In MLE, for example, we assumed that the forms for the underlying probability densities were known, and used the training samples to perform parameter estimation.
- ***Here we shall instead assume, that we know the proper forms for the discriminant functions,*** and use the samples to estimate the values of parameters of the classifier.
- Determining the discriminant functions does not require knowledge of the forms of underlying probability distributions.

Linear Discriminant Functions Are Attractive in Their Simplicity

- As we have seen before, linear discriminant functions can be optimal if the underlying distributions are co-operative, such as Gaussians having equal covariance, as might be obtained through an intelligent ***choice of feature detectors***.
- Even when they are not optimal, we might be willing to sacrifice some performance in order to gain the advantage of their simplicity.
- Linear discriminant functions are ***relatively easy to compute*** and in the absence of information suggesting otherwise, linear classifiers are an attractive candidates for initial, trial classifiers.
- You will come across linear discriminant functions also when we look into ***neural networks*** later on this course.

Linear Discriminant Functions



- Discriminant function, that is linear combination of the components of *feature vector* \mathbf{x} , can be written as

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

where \mathbf{w} is a *weight vector* and w_0 is a *bias constant*.

- In general, each class is assigned its own discriminant function and a *pattern* is classified into a class based on discriminant function that gives the highest value at point \mathbf{x} .
- In a two category case the discriminant functions can be combined to one

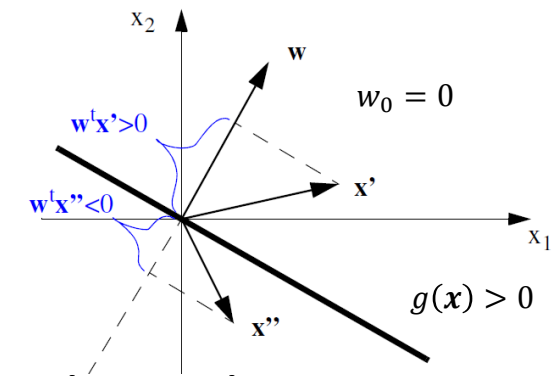
$$g_1(\mathbf{x}) > g_2(\mathbf{x}),$$

$$g_1(\mathbf{x}) - g_2(\mathbf{x}) > 0,$$

$$g(\mathbf{x}) > 0$$

- Then we can say: Decide ω_1 if $g(\mathbf{x}) > 0$ and decide ω_2 if $g(\mathbf{x}) < 0$.
 - If $g(\mathbf{x}) = 0$, \mathbf{x} can usually be assigned to either class, but in Chapter 5 the course book chooses to leave the assignment undefined.
- In other words, decide ω_1 if the inner product $\mathbf{w}^T \mathbf{x} > -w_0$, otherwise ω_2 .

Linear Discriminant Functions (cont.)



- To define a decision surface (hyperplane when $g(\mathbf{x})$ is linear) separating points assigned to ω_1 from points classified to ω_2 , we write

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = [w_1 \quad w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_0 \equiv 0.$$

- If $w_0 = 0$, then the hyperplane would be positioned to *origin*, and $\mathbf{w}^T \mathbf{x} \equiv 0$, $\mathbf{w} \cdot \mathbf{x} = 0$, $\mathbf{w} \perp \mathbf{x}$, saying \mathbf{w} is normal (perpendicular) to any vector lying in the hyperplane.
- If $w_0 > 0$, then origin is found from the positive side of the decision surface, and if $w_0 < 0$, then origin is found on the negative side of the decision surface.
- If we were to have, e.g., $w_0 > 0$, then $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$, and on the decision surface $\boxed{\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \xi \equiv 0}$ for some ξ , $\boxed{-\mathbf{w}^T \xi = w_0}$, $\mathbf{w} \perp (\mathbf{x} - (-\xi))$, resulting **a shift of origin** and allowing more flexible positioning of the hyperplane between classes.

Linear Discriminant Functions: Augmented Vector Representation

- So, a linear discriminant function divides the feature space by a decision surface.
- The orientation of the surface is determined by the normal vector \mathbf{w} (weight vector), and the location of the surface is determined by the bias w_0 .
- We can also write

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$
$$g(\mathbf{x}) = w_0 \cdot \underbrace{1}_{x_0} + [w_1 \quad w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



$$g(\mathbf{x}) = \underbrace{[w_0 \quad w_1 \quad w_2]}_{\mathbf{a}} \underbrace{\begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}}_{\mathbf{y}}$$

$$g(\mathbf{x}) = \mathbf{a}^T \mathbf{y},$$

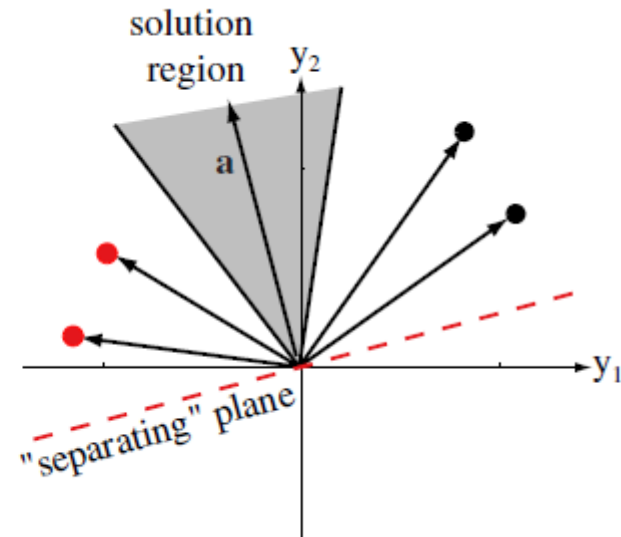
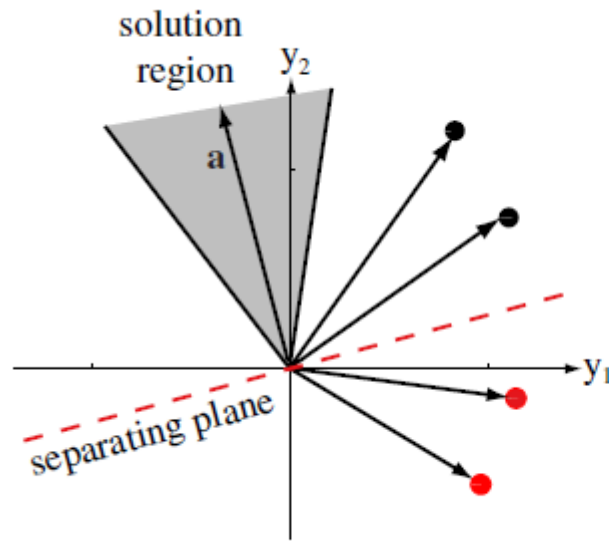
where \mathbf{a} is an **augmented weight vector** and \mathbf{y} is an **augmented feature vector**.

- You can see that the decision surface now lies in origin in \mathbf{y} space.

Linear Discriminant Functions: Linear Separation and Solution Vector

- Now we want to use *samples* in our feature vector \mathbf{y} to determine the weight vector \mathbf{a} in a linear discriminant function $g(\mathbf{x}) = \mathbf{a}^T \mathbf{y}$.
- If we are able to find vector \mathbf{a} for a discriminant function which classifies correctly all samples, we say that classes ω_1 and ω_2 are ***linearly separable***.
- Let's decide that we achieve a correct classification to ω_1 if $\mathbf{a}^T \mathbf{y}_i > 0$ and correctly to ω_2 if $\mathbf{a}^T \mathbf{y}_i < 0$.
- In order to find *solution vector* \mathbf{a} for which $\mathbf{a}^T \mathbf{y}_i > 0$ for all samples we can replace all feature vectors \mathbf{y}_i from class ω_2 with their negation $-\mathbf{y}_i$.
 - All vectors point now to the positive side of the decision boundary.


Linear Discriminant Functions: Linear Separation and Solution Vector (cont.)

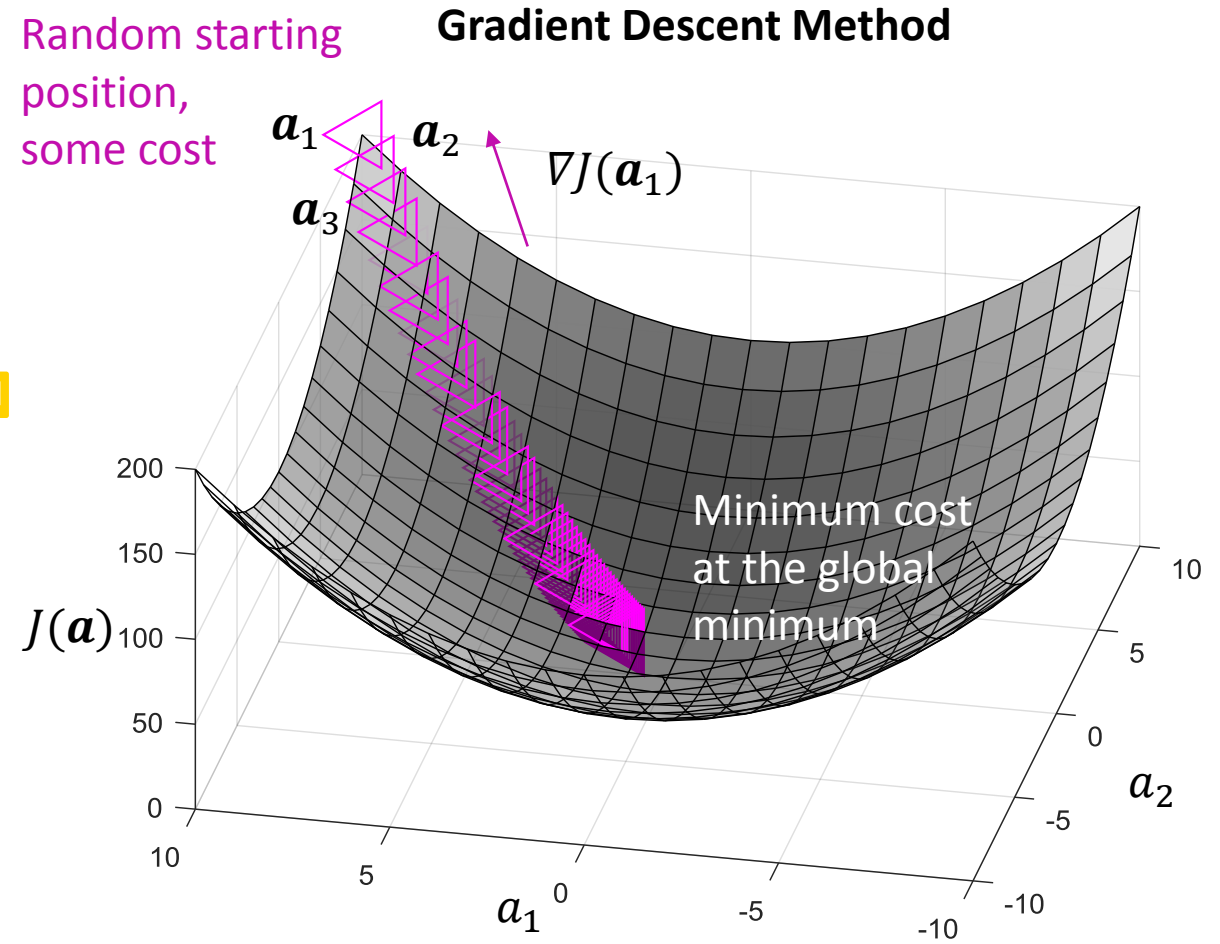


Gradient Search Methods for Finding the Weight Vector

- There are different ways to find a weight vector \mathbf{a} such that $\mathbf{a}^T \mathbf{y}_i > 0$ for all the samples (assuming one exists).
- So, we find a solution for linear systems of inequalities $\mathbf{a}^T \mathbf{y}_i > 0$ by using a certain criterion function $J(\mathbf{a})$ which minimizes when \mathbf{a} is our desired solution vector.
- So, *we move to solve a minimization problem* for scalar valued function, and that can be accomplished using **gradient search** (e.g. gradient decent procedure).
 - Gradient descent is an optimization algorithm used to find the values of parameters (coefficients) of a function that minimizes a criterion function (cost).

Gradient Decent Procedure

- In basic gradient descent method, we start with some arbitrarily chosen weight vector \mathbf{a}_1 and compute the gradient vector $\nabla J(\mathbf{a}_1)$.
- The next value \mathbf{a}_2 is obtained by moving some distance from \mathbf{a}_1 in the direction of steepest descent, i.e., along the negative of the gradient. 
- In general, \mathbf{a}_{k+1} is obtained from \mathbf{a}_k by the equation $\mathbf{a}_{k+1} = \mathbf{a}_k - \eta(k)\nabla J(\mathbf{a}_k)$, where $\eta(k)$ is a **learning rate** that sets the step size.
- We hope that our sequence of weight vectors will converge to a solution minimizing the criterion function $J(\mathbf{a})$.
- If $\eta(k)$ is too small, convergence is needlessly slow.
- If $\eta(k)$ is too large, the correction process will overshoot and can even diverge.



Here the learning rate was fixed to 0.02