

Linear Classification

1. Recalling Bayes Classifier
2. What is Linear Classification
3. Why is Linear Classification required?
4. Linear Classifier Equation
5. Perceptron Algorithm

Our major concern in Bayes Classifier was to design classifiers based on probability density or probability functions. Today, we will focus on the design of linear classifiers.

The major advantage of linear classifiers is their simplicity and computational attractiveness.

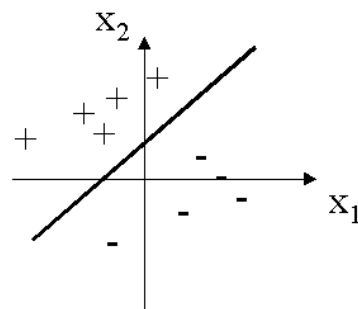
The lecture starts with the assumption that **all feature vectors from the available classes can be classified correctly using a linear classifier**

- We will learn about linear classification.
- We will develop techniques for the computation of the corresponding linear functions.
- In the sequel, **we will focus on to train the perceptron to classify inputs correctly by adjusting the connecting weights.**

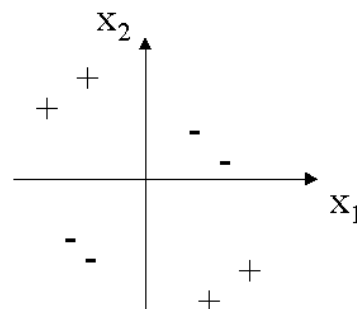
Linear Classification:

The goal of statistical classification is to use an object's characteristics to identify which class (or group) it belongs to. A **linear classifier** achieves this by making a classification decision based on the value of a linear combination of the characteristics. An object's characteristics are also known as **feature values** and are typically presented to the machine in a vector called a **feature vector**.

We can use linear classifier where things are linearly separable.



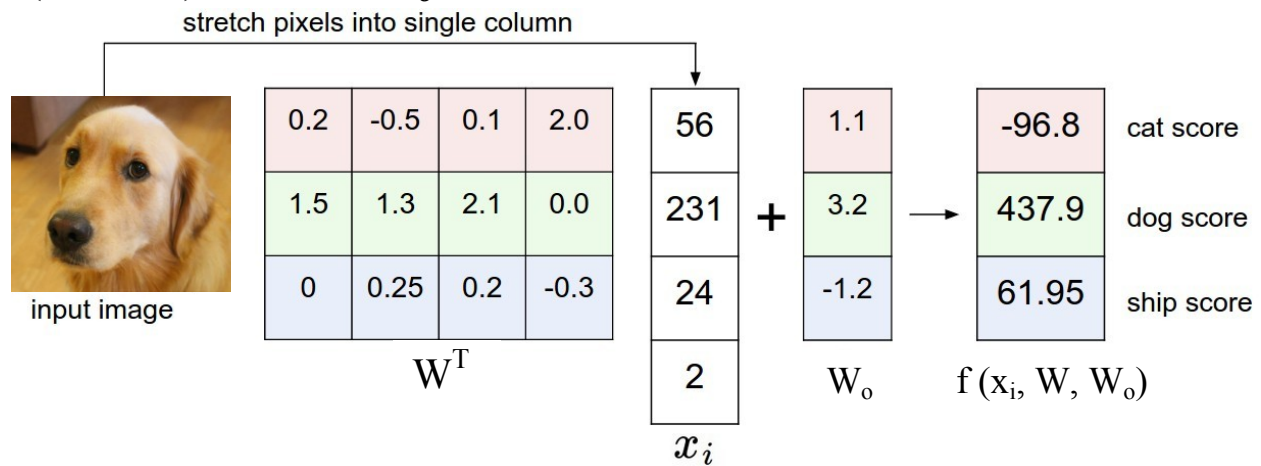
Linearly Separable



Not Linearly Separable

The simplest possible function, a linear mapping:

$$f(x, W, W_o) = W^T x + W_o$$



Two class discriminator:

Let assume, we have two classes C_1 and C_2 .

$$f(x, W, W_o) = y(x) = W^T x + W_o = \sum_{i=1}^L W_i x_i + W_o$$

$y(x) \geq 0 \rightarrow x$ assigned to C_1

$y(x) < 0 \rightarrow x$ assigned to C_2

Consider the things as point, so we have two features: x value and y value.

$$X = [1, X_1, X_2]$$

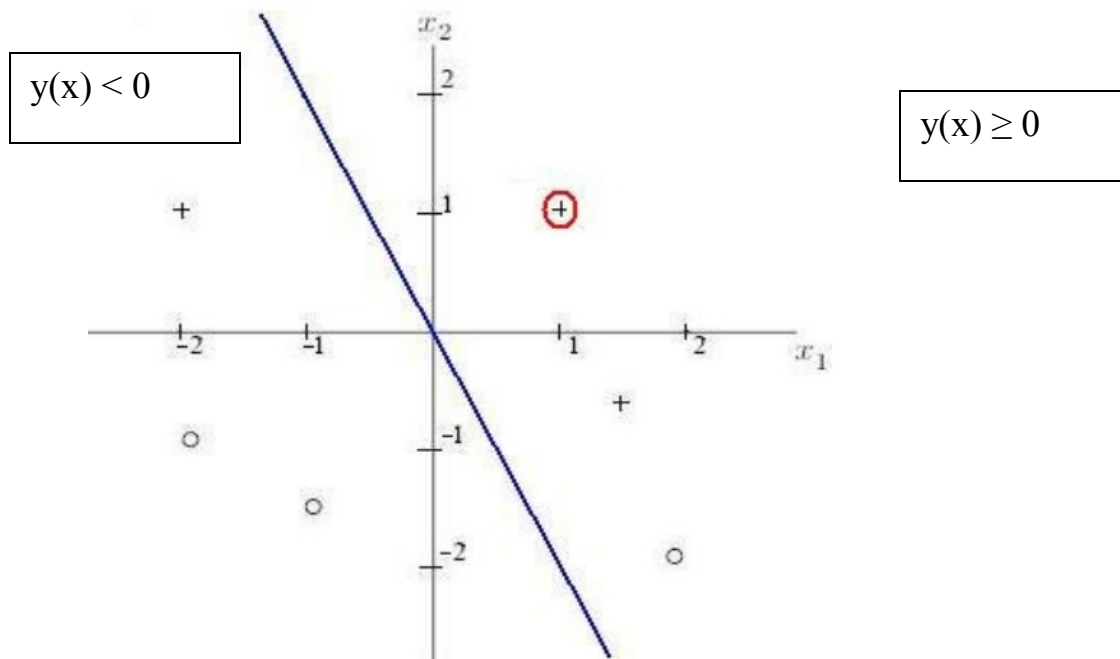
$$W = [W_o, W_1, W_2] = [0, 1, 0.5]$$

The linear equation will be:

$$W_1 x_1 + W_2 x_2 + W_o = 0$$

$$\text{Or, } x_1 + 0.5 x_2 + 0 = 0$$

$$\text{Or, } x_2 = -2x_1$$



$$x_1 = 1$$

$$x_2 = 1$$

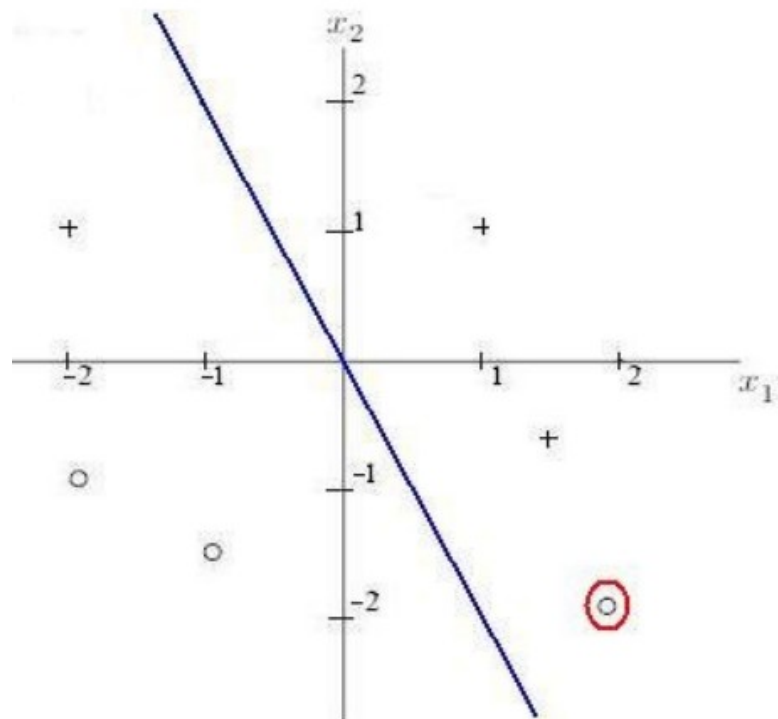
From $W^T x + W_0$

$1.5 \geq 0$, So, this will be classified in class C_1

Do the same thing for the following point:

$$x_1 = 2$$

$$x_2 = -2$$



Perceptron Learning Algorithm

In this section, we assume that the two classes C_1, C_2 are linearly separable. In other words, we assume that there exists a hyperplane, defined by $\mathbf{W}^T \mathbf{x} = 0$ such that

$$\mathbf{W}^T \mathbf{x} > 0 \in C_1$$

$$\mathbf{W}^T \mathbf{x} < 0 \in C_2$$

The formulation above also covers the case of a hyperplane not crossing the origin, that is, $\mathbf{W}^T \mathbf{x} + \mathbf{W}_0 = 0$, since this can be brought into the previous formulation by defining the extended $(l+1)$ -dimensional vectors $\mathbf{x}' \equiv [\mathbf{x}^T, 1]^T$, $\mathbf{W}' \equiv [\mathbf{W}^T, \mathbf{W}_0]^T$

$$\text{So, } \mathbf{W}^T \mathbf{x} + \mathbf{W}_0 = \mathbf{W}'^T \mathbf{x}'$$

Thus we need to adopt

- (a) an appropriate cost function and
- (b) an algorithmic scheme to optimize it.

To this end, we choose the perceptron cost defined as,

$$J(\mathbf{w}) = \sum_{\mathbf{x} \in Y} (\delta_{\mathbf{x}} \mathbf{w}^T \mathbf{x})$$

where Y is the subset of the training vectors, which are misclassified by the hyperplane defined by the weight vector \mathbf{w} . The variable $\delta_{\mathbf{x}}$ is chosen so that $\delta_{\mathbf{x}} = -1$ if $\mathbf{x} \in C_1$ and $\delta_{\mathbf{x}} = +1$ if $\mathbf{x} \in C_2$, positive, and it becomes zero when Y becomes the empty set, that is, if there are not misclassified vectors \mathbf{x} . Indeed, if $\mathbf{x} \in \omega_1$ and it is misclassified, then $\mathbf{w}^T \mathbf{x} < 0$ and $\delta_{\mathbf{x}} < 0$, and the product is positive.

The perceptron cost function is continuous and piecewise linear. Indeed, if we change the weight vector smoothly, the cost $J(\mathbf{w})$ changes linearly until the point at which there is a change in the number of misclassified vectors.

We will adopt an iterative scheme in the spirit of the gradient descent method,

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \rho_t \left. \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}(t)}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_{\mathbf{x} \in Y} \delta_{\mathbf{x}} \mathbf{x}$$

Substituting,

$$\mathbf{w}(t + 1) = \mathbf{w}(t) - \rho_t \sum_{\mathbf{x} \in Y} \delta_{\mathbf{x}} \mathbf{x}$$

The algorithm is known as the perceptron algorithm and is quite simple in its structure. The algorithm is initialized from an arbitrary weight vector $\mathbf{w}(0)$, and the correction vector is formed using the misclassified features. The weight vector is then corrected according to the preceding rule. This is repeated until the algorithm converges to a solution, that is, all features are correctly classified.

A pseudocode for the perceptron algorithm is given below.

The Perceptron Algorithm

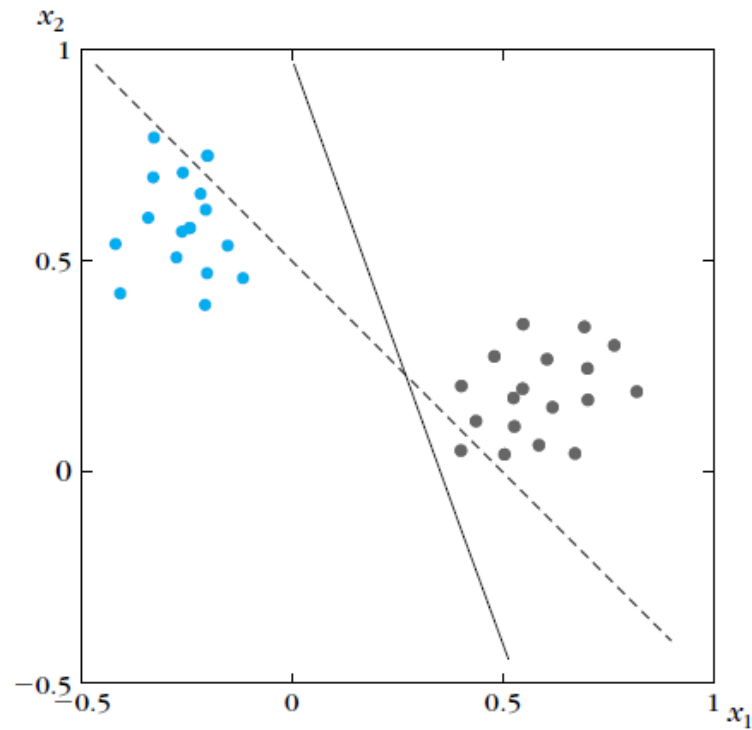
- Choose $\mathbf{w}(0)$ randomly
- Choose ρ_0
- $t = 0$
- Repeat
 - $Y = \emptyset$
 - For $i = 1$ to N
 - If $\delta_{x_i} \mathbf{w}(t)^T \mathbf{x}_i \geq 0$ then $Y = Y \cup \{\mathbf{x}_i\}$
 - End {For}
 - $\mathbf{w}(t + 1) = \mathbf{w}(t) - \rho_t \sum_{\mathbf{x} \in Y} \delta_{\mathbf{x}} \mathbf{x}$
 - Adjust ρ_t
 - $t = t + 1$
- Until $Y = \emptyset$

Example 3.1

Figure 3.3 shows the dashed line

$$x_1 + x_2 - 0.5 = 0$$

corresponding to the weight vector $[1, 1, -0.5]^T$, which has been computed from the latest iteration step of the perceptron algorithm (3.9), with $\rho_t = \rho = 0.7$. The line classifies correctly



all the vectors except $[0.4, 0.05]^T$ and $[-0.20, 0.75]^T$. According to the algorithm, the next weight vector will be

$$w(t+1) = \begin{bmatrix} 1 \\ 1 \\ -0.5 \end{bmatrix} - 0.7(-1) \begin{bmatrix} 0.4 \\ 0.05 \\ 1 \end{bmatrix} - 0.7(+1) \begin{bmatrix} -0.2 \\ 0.75 \\ 1 \end{bmatrix}$$

or

$$w(t+1) = \begin{bmatrix} 1.42 \\ 0.51 \\ -0.5 \end{bmatrix}$$

The resulting new (solid) line $1.42x_1 + 0.51x_2 - 0.5 = 0$ classifies all vectors correctly, and the algorithm is terminated.