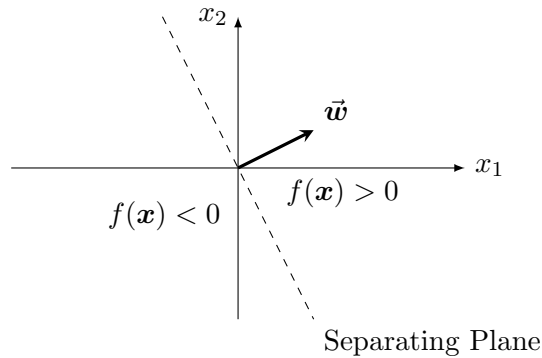# THE PERCEPTRON

The perceptron is a learning algorithm for binary classification of real valued vectors.

## 1  Introduction

The perceptron binary classifier can be thought of as the following function,

$$f(\boldsymbol{x}) = \begin{cases} 1 \text{ if } \boldsymbol{w}^T \boldsymbol{x} > 0 \\ 0 \text{ otherwise} \end{cases}$$



Separating Plane

One severe limitation of this formulation is that the separating hyperplane always passes through origin which might be undesirable in many cases. However, this limitation can be overcome in the following way,

$$f(\boldsymbol{x}) = \begin{cases} 1 \text{ if } \boldsymbol{w}^T \boldsymbol{x} + b > 0 \\ 0 \text{ otherwise} \end{cases}$$

The bias term $b$ frees the separating hyperplane from origin. Additionally, we can employ a trick to simplify the expression as follows.

$$f(\boldsymbol{x}) = \begin{cases} 1 \text{ if } \left(\boldsymbol{w}^T, b\right) \begin{pmatrix} \boldsymbol{x} \\ 1 \end{pmatrix} > 0 \\ 0 \text{ otherwise} \end{cases}$$

Hence, by increasing the dimension of input by one and defaulting the intercept on the new dimension to 1 for every input vector, any separating hyperplane in the old dimensional space becomes a hyperplane passing through origin in the new dimensional space. This simplifies the mathematics without constraining the classifier.

The vector $\vec{w}$ is learned from training data as usual. The algorithm for learning the same is presented in the next section.

## 2 Learning Algorithm

The training data consists of $N$ pairs of $d$-dimensional real input vectors $\boldsymbol{x_i}$s and the output binary labels $y_i$s. For mathematical convenience, positive output labels are denoted by +1 and negative output labels by -1. The learning algorithm makes the following assumptions:

- It is assumed that the input vectors have already been extended to account for the bias term.

- The training data is assumed to be linearly separable at origin. So, there exists a unit vector $\boldsymbol{w}^*$ such that $y_i \boldsymbol{w}^{*T} \boldsymbol{x_i} > \gamma \ \ \forall i \ \ \text{where} \ \ \gamma > 0$.

- It is assumed that all training input vectors are finite.

The algorithm is defined as follows.
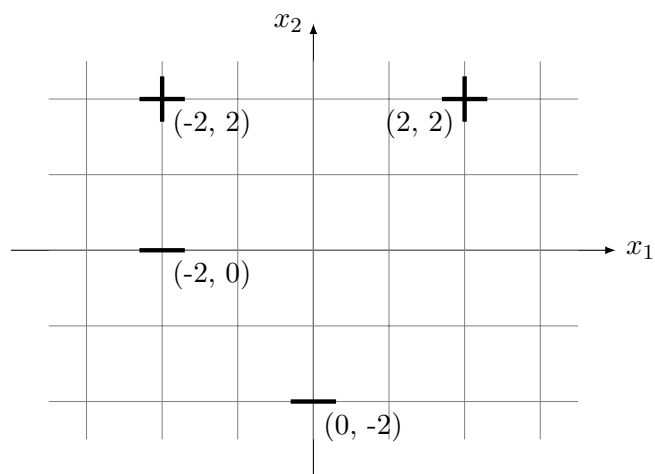
---
**Algorithm**  The Perceptron Learning Algorithm
---
$k \leftarrow 1$
$\boldsymbol{w}_k \leftarrow 0$
**while** there exists $j \in \{1, 2, .., N\}$ such that $y_j \boldsymbol{w}^{*T} \boldsymbol{x_j} <= 0$ **do**
  pick $i \in \{1, 2, .., N\}$ such that $y_i \boldsymbol{w}^{*T} \boldsymbol{x_i} <= 0$
  $\boldsymbol{w}_{k+1} \leftarrow \boldsymbol{w}_k + y_i \boldsymbol{x_i}$
  $k \leftarrow k + 1$
**end while**
**return** $\boldsymbol{w}_k$

---

## 3 Example

The algorithm is illustrated on a very simple example.

$x_2$

$(-2, 2)$  $(2, 2)$

$x_1$

$(-2, 0)$

$(0, -2)$

## 3.1 Iteration 1

$$\boldsymbol{w}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

| $\boldsymbol{x}$ | | |
|---|---|---|
| $\begin{pmatrix} -2 \\ 0 \end{pmatrix}$ | -1 | 0 |
| $\begin{pmatrix} 0 \\ -2 \end{pmatrix}$ | -1 | 0 |
| $\begin{pmatrix} -2 \\ 2 \end{pmatrix}$ | 1 | 0 |
| $\begin{pmatrix} 2 \\ 2 \end{pmatrix}$ | 1 | 0 |

## 3.2 Iteration 2

$$\boldsymbol{w}_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} + (-1) \begin{pmatrix} -2 \\ 0 \end{pmatrix}$$
$$= \begin{pmatrix} 2 \\ 0 \end{pmatrix}$$

| $\boldsymbol{x}$ | | |
|---|---|---|
| $\begin{pmatrix} -2 \\ 0 \end{pmatrix}$ | -1 | 4 |
| $\begin{pmatrix} 0 \\ -2 \end{pmatrix}$ | -1 | 0 |
| $\begin{pmatrix} -2 \\ 2 \end{pmatrix}$ | 1 | -4 |
| $\begin{pmatrix} 2 \\ 2 \end{pmatrix}$ | 1 | 4 |

### 3.3 Iteration 3

$$\boldsymbol{w}_3 = \begin{pmatrix} 2 \\ 0 \end{pmatrix} + (-1) \begin{pmatrix} 0 \\ -2 \end{pmatrix}$$
$$= \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

| $\boldsymbol{x}$ | | |
|---|---|---|
| $\begin{pmatrix} -2 \\ 0 \end{pmatrix}$ | -1 | 4 |
| $\begin{pmatrix} 0 \\ -2 \end{pmatrix}$ | -1 | 4 |
| $\begin{pmatrix} -2 \\ 2 \end{pmatrix}$ | 1 | 0 |
| $\begin{pmatrix} 2 \\ 2 \end{pmatrix}$ | 1 | 8 |

### 3.4 Iteration 4

$$\boldsymbol{w}_4 = \begin{pmatrix} 2 \\ 2 \end{pmatrix} + (1) \begin{pmatrix} -2 \\ 2 \end{pmatrix}$$
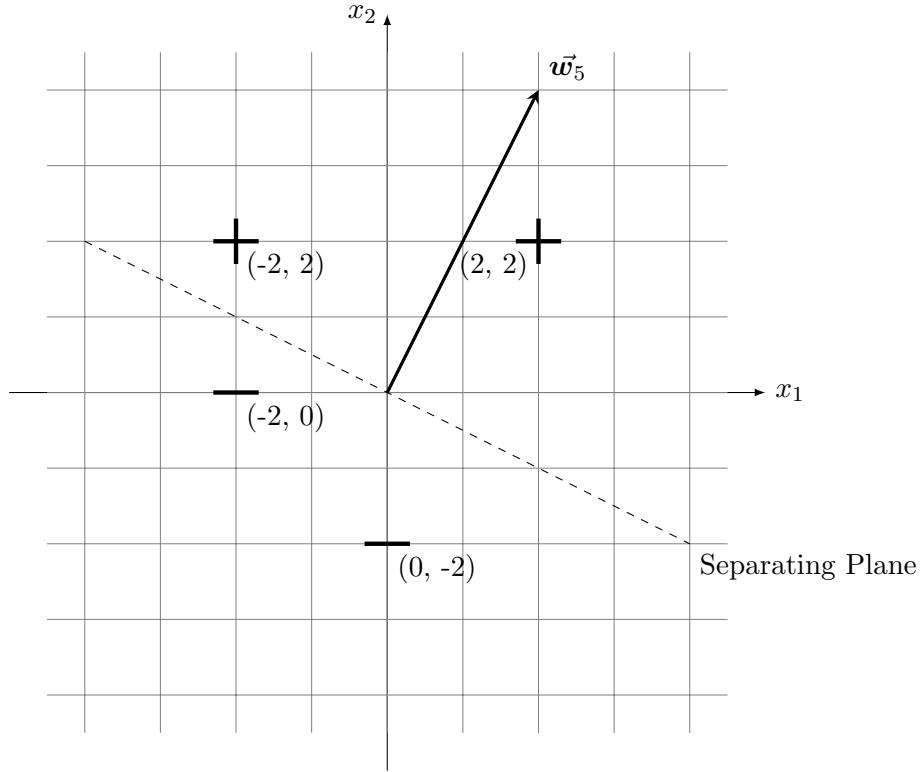$$= \begin{pmatrix} 0 \\ 4 \end{pmatrix}$$

| $\boldsymbol{x}$ | | |
|---|---|---|
| $\begin{pmatrix} -2 \\ 0 \end{pmatrix}$ | -1 | 0 |
| $\begin{pmatrix} 0 \\ -2 \end{pmatrix}$ | -1 | 8 |
| $\begin{pmatrix} -2 \\ 2 \end{pmatrix}$ | 1 | 8 |
| $\begin{pmatrix} 2 \\ 2 \end{pmatrix}$ | 1 | 8 |

## 3.5   Iteration 5

$$\boldsymbol{w}_5 = \begin{pmatrix} 0 \\ 4 \end{pmatrix} + (-1) \begin{pmatrix} -2 \\ 0 \end{pmatrix}$$
$$= \begin{pmatrix} 2 \\ 4 \end{pmatrix}$$

| $\boldsymbol{x}$ | | |
|---|---|---|
| $\begin{pmatrix} -2 \\ 0 \end{pmatrix}$ | -1 | 4 |
| $\begin{pmatrix} 0 \\ -2 \end{pmatrix}$ | -1 | 8 |
| $\begin{pmatrix} -2 \\ 2 \end{pmatrix}$ | 1 | 4 |
| $\begin{pmatrix} 2 \\ 2 \end{pmatrix}$ | 1 | 12 |

All the points are correctly classified, hence the algorithm terminates. The algorithm converged in 5 iterations. In fact, as will be proved later, the algorithm is always guaranteed to converge in finite iterations.The resulting separating hyperplanes is shown below.

# 4    Proof of Convergence

The intention is to show that the perceptron learning algorithm always converges to a perfect classifier in finite iterations given that all the assumptions that the algorithm makes are met.

Assume that after k+1 iterations, $\vec{w_{k+1}}$ makes an angle $\theta_{k+1}$ with $\vec{w^*}$.

$$cos(\theta_{k+1}) = \frac{\vec{w_{k+1}}.\vec{w^*}}{|\vec{w_{k+1}}| \; |\vec{w^*}|}$$

Since $\vec{w^*}$ is a unit vector,

$$cos(\theta_{k+1}) = \frac{\vec{w_{k+1}}.\vec{w^*}}{|\vec{w_{k+1}}|}$$

Now, bounds for numerator and denominator are established to proceed ahead.

## 4.1    Bound for numerator

The update rule for $\vec{w_{k+1}}$ is,

$$\vec{w_{k+1}} = \vec{w_k} + y_i\vec{x_i}$$

Taking dot product with $\vec{w^*}$ on both sides.

$$\vec{w^*}.\vec{w_{k+1}} = \vec{w^*}.\vec{w_k} + y_i\vec{w^*}.\vec{x_i}$$

Since $y_i\vec{w^*}.\vec{x_i} > \gamma$ as per the algorithm assumption.

$$\vec{w^*}.\vec{w_{k+1}} > \vec{w^*}.\vec{w_k} + \gamma$$
$$\vec{w^*}.\vec{w_{k+1}} > \vec{w^*}.\vec{w_{k-1}} + 2\gamma$$
$$\vec{w^*}.\vec{w_{k+1}} > \vec{w^*}.\vec{w_1} + k\gamma$$

since $\vec{w^*}.\vec{w_1} = 0$.

$$\vec{w^*}.\vec{w_{k+1}} > k\gamma$$

## 4.2    Bound for denominator

Again starting with the update rule for $\vec{w_{k+1}}$,

$$\vec{w_{k+1}} = \vec{w_k} + y_i\vec{x_i}$$

Taking modulus on both sides,

$$|\vec{w_{k+1}}|^2 = |\vec{w_k}|^2 + |y_i\vec{x_i}|^2 + 2y_i\vec{w_k}.\vec{x_i}$$

Since $y_i\vec{w_k}.\vec{x_i} <= 0$ according to the algorithm specification.

$$|\vec{w_{k+1}}|^2 \leq |\vec{w_k}|^2 + |y_i\vec{x_i}|^2$$

Assuming that the largest input vector has length $R$.

$$|\vec{w_{k+1}}|^2 \leq |\vec{w_k}|^2 + R^2$$
$$|\vec{w_{k+1}}|^2 \leq |\vec{w_{k-1}}|^2 + 2R^2$$
$$|\vec{w_{k+1}}|^2 \leq kR^2$$
$$|\vec{w_{k+1}}| \leq \sqrt{k}R$$

### 4.3    Bound for Number of Iterations

Using both the numerator and denominator bounds,

$$cos(\theta_{k+1}) > \frac{k\gamma}{\sqrt{k}R}$$
$$cos(\theta_{k+1}) > \sqrt{k}\frac{\gamma}{R}$$

This shows that after each iteration, the lower bound of the angle between $\vec{w}$ and $\vec{w^*}$ increases. And since $cos(\theta_{k+1})$ can never be greater than one, $\sqrt{k}\frac{\gamma}{R} \leq 1$ must hold true yielding $k \leq \frac{R^2}{\gamma^2}$.

## 5    Conclusion

The perceptron was one of the first algorithms indicating that computers could learn from data. And today, after half a century, it continues to be the basis of Artificial Neural Networks which are at the forefront of Artificial Intelligence.