

Linear Classification: Perceptron

Yufei Tao

Department of Computer Science and Engineering
Chinese University of Hong Kong

Today, we start a series of lectures devoted to **linear classification**, which harbors a deep theory and is one of the most important topics in machine learning.

Linear Classification

Let A_1, \dots, A_d be d **attributes**, each with a domain \mathbb{R} , i.e., $\text{dom}(A_i) = \mathbb{R}$ for each $i \in [1, d]$.

Instance space: $\mathcal{X} = \text{dom}(A_1) \times \text{dom}(A_2) \times \dots \times \text{dom}(A_d) = \mathbb{R}^d$.

Label space: $\mathcal{Y} = \{-1, 1\}$ (where -1 and 1 are **class labels**).

Instance-label pair (a.k.a. **object**): a pair (\mathbf{x}, y) in $\mathcal{X} \times \mathcal{Y}$.

- \mathbf{x} is a d -dimensional vector. Since every dimension has a real domain, we can regard \mathbf{x} as a d -dimensional point.
- We use $\mathbf{x}[i]$ to represent the i -th coordinate of point \mathbf{x} .

Linear Classification

Linear classifier: A function $h: \mathcal{X} \rightarrow \mathcal{Y}$ where h is defined by a d -dimensional **weight vector** \mathbf{w} such that

- $h(\mathbf{x}) = 1$ if $\mathbf{x} \cdot \mathbf{w} \geq 0$ (note: “ \cdot ” represents dot product);
- $h(\mathbf{x}) = -1$ otherwise.

Suppose that Alice chooses a linear classifier h^* and a distribution \mathcal{D} over \mathcal{X} (note: \mathcal{D} is defined in the instance space, not the instance-label space).

For any linear classifier h , its **error on** \mathcal{D} is defined as:

$$\text{err}_{\mathcal{D}}(h) = \Pr_{\mathbf{x} \sim \mathcal{D}}[h(\mathbf{x}) \neq h^*(\mathbf{x})].$$

Note that the error of h^* on \mathcal{D} is 0.

Linear Classification

Alice provides a **training set** S which contains objects (\mathbf{x}, y) obtained as follows:

- First, draw \mathbf{x} independently from \mathcal{X} .
- Then, set $y = h^*(\mathbf{x})$.

The goal of linear classification is to learn a classifier h from S whose error on \mathcal{D} is as low as possible.

Linear Classification

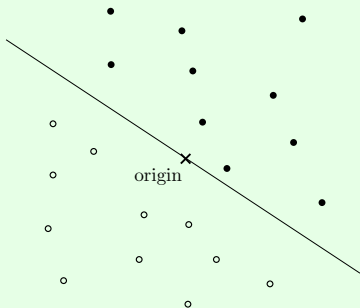
S is **linearly separable** if there is a d -dimensional vector \mathbf{w} such that for each $\mathbf{p} \in S$:

- $\mathbf{w} \cdot \mathbf{p} > 0$ if \mathbf{p} has label 1;
- $\mathbf{w} \cdot \mathbf{p} < 0$ if \mathbf{p} has label -1 .

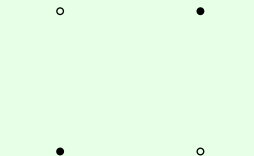
The plane $\mathbf{w} \cdot \mathbf{x} = 0$ is a **separation plane** of S .

We will discuss only the scenario where S is linearly separable.

Example:



Linearly separable



Linearly non-separable

In this lecture, we will study the following problem:

Problem (Finding a Separation Plane): Given a linearly separable set S , find a separation plane.

The separation plane gives a linear classifier h with $err_S(h) = 0$, i.e., empirical error 0.

We will solve the problem with a surprisingly simple algorithm called **perceptron**.

Perceptron

The algorithm starts with $\mathbf{w} = (0, 0, \dots, 0)$ and, then, runs in **iterations**.

In each iteration, it looks for a **violation point** $\mathbf{p} \in S$:

- If \mathbf{p} has label 1, \mathbf{p} is a violation point if $\mathbf{w} \cdot \mathbf{p} \leq 0$;
- If \mathbf{p} has label -1 , \mathbf{p} is a violation point if $\mathbf{w} \cdot \mathbf{p} \geq 0$;

If \mathbf{p} exists, the algorithm adjusts \mathbf{w} as follows:

- If \mathbf{p} has label 1, then $\mathbf{w} \leftarrow \mathbf{w} + \mathbf{p}$.
- If \mathbf{p} has label -1 , then $\mathbf{w} \leftarrow \mathbf{w} - \mathbf{p}$.

The algorithm finishes when there are no more violation points.

Example: Suppose that S has points: $\mathbf{p}_1 = (1, 0)$, $\mathbf{p}_2 = (0, -1)$, $\mathbf{p}_3 = (0, 1)$, and $\mathbf{p}_4 = (-1, 0)$. Points \mathbf{p}_1 and \mathbf{p}_3 have label 1, and the other have label -1 .

The algorithm starts with $\mathbf{w} = (0, 0, \dots, 0)$.

- Iteration 1: \mathbf{p}_1 is a violation point because it has label 1 but $\mathbf{p}_1 \cdot \mathbf{w} = 0$. Hence, we update \mathbf{w} to $\mathbf{w} + \mathbf{p}_1 = (1, 0)$.
- Iteration 2: \mathbf{p}_2 is a violation point because it has label -1 but $\mathbf{p}_2 \cdot \mathbf{w} = 0$. Hence, we update \mathbf{w} to $\mathbf{w} - \mathbf{p}_2 = (1, 0) - (0, -1) = (1, 1)$.
- Iteration 3: No more violation points. The algorithm finishes with $\mathbf{w} = (1, 1)$.

We now analyze the number of iterations performed by Perceptron.

Given a vector $\mathbf{v} = (v_1, \dots, v_d)$, we define its **length** as

$$|\mathbf{v}| = \sqrt{\mathbf{v} \cdot \mathbf{v}} = \sqrt{\sum_{i=1}^d v[i]^2}.$$

For any vectors $\mathbf{v}_1, \mathbf{v}_2$, it holds that $\mathbf{v}_1 \cdot \mathbf{v}_2 \leq |\mathbf{v}_1| |\mathbf{v}_2|$.

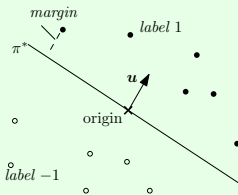
Define:

$$R = \max_{\mathbf{p} \in S} \{|\mathbf{p}|\}.$$

In other words, all the points of S fall in a ball that centers at the origin and has radius R .

Given a separation plane π , define its **margin** as the smallest distance from the points of S to π .

Example:



Denote by γ the **largest** margin of all the separation planes. Let π^* be the origin-passing plane with margin γ ; the plane has a **unit normal vector** u^* such that

- for every $p \in S$ with label 1, $u^* \cdot p > 0$;
- for every $p \in S$ with label -1 , $u^* \cdot p < 0$.

We have:

$$\gamma = \min_{p \in S} |u^* \cdot p|.$$

Theorem: Perceptron terminates after at most $(R/\gamma)^2$ adjustments of \mathbf{w} .

Proof: Let \mathbf{w}_i ($i \geq 1$) be the value of \mathbf{w} after the i -th adjustment. As a special case, define $\mathbf{w}_0 = (0, \dots, 0)$. Denote by k the total number of violations.

We first show that $\mathbf{w}_{i+1} \cdot \mathbf{u}^* \geq \mathbf{w}_i \cdot \mathbf{u}^* + \gamma$ for any $i \geq 0$. Consider the violation point \mathbf{p} used to change \mathbf{w} from \mathbf{w}_i to \mathbf{w}_{i+1} :

- Case 1: \mathbf{p} has label 1. Thus, $\mathbf{p} \cdot \mathbf{w}_i < 0$ and $\mathbf{w}_{i+1} = \mathbf{w}_i + \mathbf{p}$. Hence, $\mathbf{w}_{i+1} \cdot \mathbf{u}^* = \mathbf{w}_i \cdot \mathbf{u}^* + \mathbf{p} \cdot \mathbf{u}^*$. From the definition of γ , we know that $\mathbf{p} \cdot \mathbf{u}^* \geq \gamma$. This gives $\mathbf{w}_{i+1} \cdot \mathbf{u}^* \geq \mathbf{w}_i \cdot \mathbf{u}^* + \gamma$.
- Case 2: \mathbf{p} has label -1 . The proof is similar and left to you.

Therefore:

$$\begin{aligned} \mathbf{w}_k \cdot \mathbf{u}^* &\geq \mathbf{w}_{k-1} \cdot \mathbf{u}^* + \gamma \\ &\geq \mathbf{w}_{k-2} \cdot \mathbf{u}^* + 2\gamma \\ &\dots \\ &\geq \mathbf{w}_0 \cdot \mathbf{u}^* + k\gamma \\ &= k\gamma. \end{aligned} \tag{1}$$

Next, we show that $|\mathbf{w}_{i+1}|^2 \leq |\mathbf{w}_i|^2 + R^2$ for any $i \geq 0$. Consider the violation point \mathbf{p} used to change \mathbf{w} from \mathbf{w}_i to \mathbf{w}_{i+1} :

- Case 1: \mathbf{p} has label 1. Thus, $\mathbf{p} \cdot \mathbf{w}_i < 0$ and $\mathbf{w}_{i+1} = \mathbf{w}_i + \mathbf{p}$. Hence:

$$\begin{aligned} |\mathbf{w}_{i+1}|^2 &= \mathbf{w}_{i+1} \cdot \mathbf{w}_{i+1} &= (\mathbf{w}_i + \mathbf{p}) \cdot (\mathbf{w}_i + \mathbf{p}) \\ &= \mathbf{w}_i \cdot \mathbf{w}_i + 2\mathbf{w}_i \cdot \mathbf{p} + |\mathbf{p}|^2 \\ &\text{(by def. of } R) \leq |\mathbf{w}_i|^2 + 2\mathbf{w}_i \cdot \mathbf{p} + R^2 \\ &\leq |\mathbf{w}_i|^2 + R^2 \end{aligned}$$

where the last step used the fact that $\mathbf{p} \cdot \mathbf{w}_i < 0$.

- Case 2: \mathbf{p} has label -1 . The proof is similar and left to you.

Therefore:

$$|\mathbf{w}_k|^2 \leq |\mathbf{w}_{k-1}|^2 + R^2 \leq |\mathbf{w}_{k-2}|^2 + 2R^2 \dots \leq |\mathbf{w}_0|^2 + kR^2 = kR^2. \quad (2)$$

From (1), we know:

$$|\mathbf{w}_k| = |\mathbf{w}_k| |\mathbf{u}^*| \geq \mathbf{w}_k \cdot \mathbf{u}^* \geq k\gamma.$$

Therefore, $|\mathbf{w}_k|^2 \geq k^2\gamma^2$. Comparing this to (2) gives:

$$\begin{aligned} kR^2 &\geq k^2\gamma^2 \Rightarrow \\ k &\leq \frac{R^2}{\gamma^2} \end{aligned}$$



We have learned how to obtain a linear classifier h with 0 empirical error on S . Does h have a small generalization error $err_{\mathcal{D}}(h)$? The answer is yes, but this does not follow from the generalization theorem we currently have (**think**: why not?). In the next lecture, we will discuss a more powerful generalization theorem that will allow us to bound $err_{\mathcal{D}}(h)$.