

Linear Classification: Maximizing the Margin

Yufei Tao

Department of Computer Science and Engineering
Chinese University of Hong Kong

Recall:

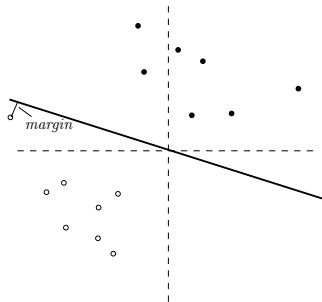
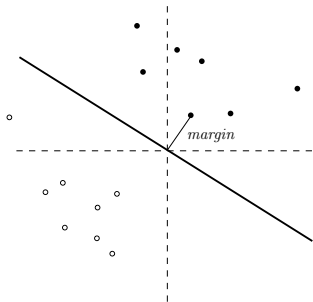
S is **linearly separable** if there is a d -dimensional vector \mathbf{w} such that for each $\mathbf{p} \in S$:

- $\mathbf{w} \cdot \mathbf{p} > 0$ 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 .

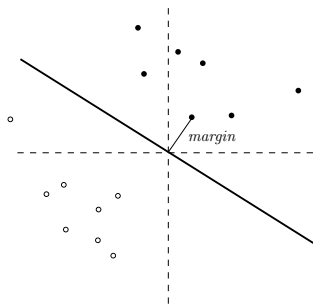
There can be many separation planes. As discussed previously, we should find the plane with the **largest margin**. In this lecture, we will discuss how to achieve the purpose.

Review: Margins



We prefer the left plane.

Let S be a linearly separable set of points in \mathbb{R}^d . In the **large margin separation problem**, we want to find a separation plane with the maximum margin.



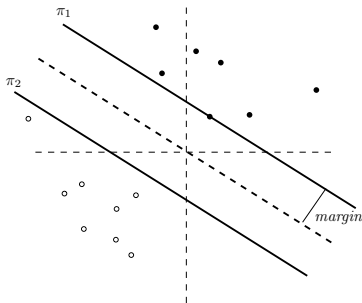
An algorithm solving this problem is called a **support vector machine**.

Next, we will discuss two methods. The first one finds the **optimal** solution but is computationally expensive. The second method is (much) faster but gives an **approximate** solution close to optimality.

Finding the Optimal Plane

We will model the problem as a **quadratic programming** problem.

Consider an arbitrary separation plane $w' \cdot x = 0$. Imagine two copies of the plane, one moving up and the other down, at the same speed. They stop as soon as a plane hits a point in S .



Now, focus on the two copies of the plane in their final positions. If one copy has equation $\mathbf{w}' \cdot \mathbf{x} = \tau$, the other copy must have equation $\mathbf{w}' \cdot \mathbf{x} = -\tau$, where $\tau > 0$.

For each point $\mathbf{p} \in S$, we must have:

- if \mathbf{p} has label 1, then $\mathbf{w}' \cdot \mathbf{p} \geq \tau$;
- if \mathbf{p} has label -1 , then $\mathbf{w}' \cdot \mathbf{p} \leq -\tau$.

By dividing τ on both sides of each inequality, we have:

- if \mathbf{p} has label 1, then $\mathbf{w} \cdot \mathbf{p} \geq 1$;
- if \mathbf{p} has label -1 , then $\mathbf{w} \cdot \mathbf{p} \leq -1$

where

$$\mathbf{w} = \frac{\mathbf{w}'}{\tau}.$$

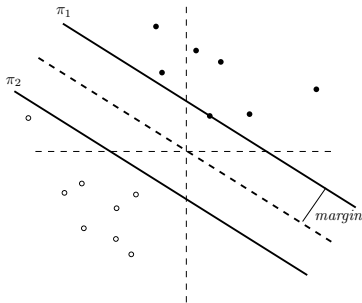
We will refer to the following plane as π_1

$$\mathbf{w} \cdot \mathbf{x} = 1$$

the following plane as π_2

$$\mathbf{w} \cdot \mathbf{x} = -1$$

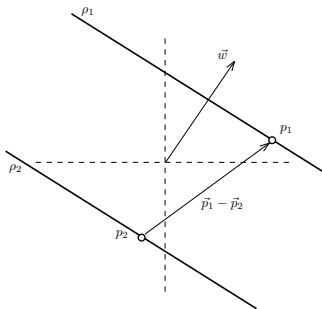
The margin of the original separation plane is exactly half of the distance between π_1 and π_2 :



Lemma: The distance between π_1 and π_2 is $\frac{2}{|\mathbf{w}|}$.

Hence, the margin of the separation plane $\mathbf{w} \cdot \mathbf{x} = 0$ is $\frac{1}{|\mathbf{w}|}$.

Proof: Take an arbitrary point \mathbf{p}_1 on π_1 and an arbitrary point \mathbf{p}_2 on π_2 . Hence, $\mathbf{w} \cdot \mathbf{p}_1 = 1$ and $\mathbf{w} \cdot \mathbf{p}_2 = -1$. It follows that $\mathbf{w} \cdot (\mathbf{p}_1 - \mathbf{p}_2) = 2$.



The distance between the two planes is precisely

$$\frac{\mathbf{w}}{|\mathbf{w}|} \cdot (\mathbf{p}_1 - \mathbf{p}_2) = \frac{2}{|\mathbf{w}|}.$$



In summary of the above, to solve the large margin separation problem, we want to find \mathbf{w} with the smallest $|\mathbf{w}|$, subject to:

- For each point $p \in S$ of label 1:

$$\mathbf{w} \cdot \mathbf{p} \geq 1$$

- For each point $p \in S$ of label -1 :

$$\mathbf{w} \cdot \mathbf{p} \leq -1$$

This is an instance of **quadratic programming**.

In theory, the quadratic programming instance can be solved using convex-optimization techniques whose efficiency is rather difficult to analyze. We will not discuss this direction further.

Finding an Approximate Plane

Define γ_{opt} as the maximum margin of all separation planes.

A separation plane is **c -approximate** if its margin is at least $c \cdot \gamma_{opt}$.

We will give an algorithm to find a $(1/4)$ -approximate separation plane.

Recall that a separation plane is given by $\mathbf{w} \cdot \mathbf{x} = 0$. The goal is to find a good \mathbf{w} .

Our weapon is once again Perceptron. But we will correct \mathbf{w} **not only** when a point falls on the wrong side of the plane, **but also** when the point is too close to the plane.

For now, let us assume we are given an arbitrary value $\gamma_{\text{guess}} \leq \gamma_{\text{opt}}$. A point p causes a **violation** in any of the following situations:

- Its distance to the plane $\mathbf{w} \cdot \mathbf{x} = 0$ is less than $\gamma_{\text{guess}}/2$, **regardless of its label**.
- p has label 1 but $\mathbf{w} \cdot \mathbf{p} < 0$.
- p has label -1 but $\mathbf{w} \cdot \mathbf{p} > 0$.

Margin Perceptron

The algorithm starts with $\mathbf{w} = \mathbf{0}$ and runs in **iterations**.

In each iteration, it tries to find a **violation point** $p \in S$. If found, the algorithm adjusts \mathbf{w} as follows:

- if p has label 1, $\mathbf{w} \leftarrow \mathbf{w} + p$.
- otherwise, $\mathbf{w} \leftarrow \mathbf{w} - p$.

The algorithm finishes where no violation points are found.

Define $R = \max_{\mathbf{p} \in S} \{|\mathbf{p}|\}$, i.e., the maximum distance from the origin to the points in S .

Theorem: If $\gamma_{\text{guess}} \leq \gamma_{\text{opt}}$, margin Perceptron terminates in at most

$$12R^2/\gamma_{\text{opt}}^2$$

iterations and returns a separation plane with margin at least $\gamma_{\text{guess}}/2$.

The proof can be found in the appendix.

Margin Perceptron requires a parameter $\gamma_{guess} \leq \gamma_{opt}$. By the theorem on the previous slide, a larger γ_{guess} promises a better quality guarantee.

Ideally, an ideal value for γ_{guess} is γ_{opt} , but unfortunately, we do not know γ_{opt} . Next, we present a strategy to estimate γ_{opt} up to a factor of $1/2$.

An Incremental Algorithm

- 1 $R \leftarrow$ the maximum distance from the origin to the points in S
- 2 $\gamma_{guess} \leftarrow R$
- 3 Run margin Perceptron with parameter γ_{guess} .
 - **[Self-Termination]**
If the algorithm terminates with a plane π , return π as the final answer.
 - **[Forced-Termination]**
If the algorithm has not terminated after $\frac{12R^2}{\gamma_{guess}^2}$ iterations:
 - Stop the algorithm manually.
 - Set $\gamma_{guess} \leftarrow \gamma_{guess}/2$.
 - Repeat Line 3.

Theorem: Our incremental algorithm returns a separation plane with margin at least $\gamma_{opt}/4$. Furthermore, it performs $O(R^2/\gamma_{opt}^2)$ iterations in total (including all the repeats at Line 3).

Proof: Suppose that we repeat Line 3 in total h times. For each $i \in [1, h]$, denote by γ_i the value of γ_{guess} at the i -th time we execute Line 3.

By the fact that the $(i-1)$ -th repeat required a forced termination, we know that $\gamma_{h-1} > \gamma_{opt}$. Hence, $\gamma_h = \gamma_{h-1}/2 > \gamma_{opt}/2$. It thus follows that the plane we return must have a margin at least $\gamma_h/2 > \gamma_{opt}/4$.

The total number of iterations performed is

$$\begin{aligned} O\left(\sum_{i=1}^h \frac{R^2}{\gamma_i^2}\right) &= O\left(\frac{R^2}{\gamma_h^2} + \frac{R^2}{4\gamma_h^2} + \frac{R^2}{4^2\gamma_h^2} + \dots\right) \\ &= O(R^2/\gamma_h^2) = O(R^2/\gamma_{opt}^2). \end{aligned}$$

□

Appendix: Proof of the Theorem on Slide 18.

Let π^* be the the optimal plane with margin γ_{opt} .

Define \mathbf{u} as the unit normal vector of π^* pointing to the positive side of π^* ; in other words, we have:

- $|\mathbf{u}| = 1$.
- For every point $p \in S$ of label 1, $\mathbf{p} \cdot \mathbf{u} > 0$.
- For every point $p \in S$ label -1 , $\mathbf{p} \cdot \mathbf{u} < 0$.
- $\gamma_{opt} = \min_{p \in S} \{|\mathbf{p} \cdot \mathbf{u}|\}$.

Recall that the perceptron algorithm adjusts \mathbf{w} in each iteration. Let k be the total number of adjustments. Denote by \mathbf{w}_i ($i \geq 1$) the value of \mathbf{w} after the i -th adjustment; and define $\mathbf{w}_0 = (0, \dots, 0)$.

Claim 1: $|\mathbf{w}_k| \geq \mathbf{w}_k \cdot \mathbf{u} \geq k\gamma_{opt}$.

Proof: We will first prove: for any $i \geq 0$, it holds that.

$$\mathbf{w}_{i+1} \cdot \mathbf{u} \geq \mathbf{w}_i \cdot \mathbf{u} + \gamma_{opt}. \quad (1)$$

Due to symmetry, we will prove the above only for the case where \mathbf{w}_{i+1} is adjusted from \mathbf{w}_i due to a violation point \mathbf{p} of label 1. In this case, $\mathbf{w}_{i+1} = \mathbf{w}_i + \mathbf{p}$; and hence, $\mathbf{w}_{i+1} \cdot \mathbf{u} = \mathbf{w}_i \cdot \mathbf{u} + \mathbf{p} \cdot \mathbf{u}$. From the definition of γ_{opt} , we know that $\mathbf{p} \cdot \mathbf{u} \geq \gamma_{opt}$, which gives (1).

It then follows from (1) that

$$\begin{aligned} |\mathbf{w}_k| &\geq \mathbf{w}_k \cdot \mathbf{u} \\ &\geq \mathbf{w}_{k-1} \cdot \mathbf{u} + \gamma_{opt} \\ &\geq \mathbf{w}_{k-2} \cdot \mathbf{u} + 2\gamma_{opt} \\ &\dots \\ &\geq \mathbf{w}_0 + k\gamma_{opt} = k\gamma_{opt}. \end{aligned}$$

□

Claim 2: $|\mathbf{w}_{i+1}| \leq |\mathbf{w}_i| + R$.

Proof: We will prove only the case where \mathbf{w}_{i+1} is adjusted from \mathbf{w}_i using a violation point \mathbf{p} of label 1. In this case:

$$|\mathbf{w}_{i+1}| = |\mathbf{w}_i + \mathbf{p}| \leq |\mathbf{w}_i| + |\mathbf{p}| \leq |\mathbf{w}_i| + R.$$



Claim 3: $|\mathbf{w}_{i+1}| \leq |\mathbf{w}_i| + \frac{R^2}{2|\mathbf{w}_i|} + \frac{\gamma_{opt}}{2}$.

Proof: We will prove only the case where \mathbf{w}_{i+1} is adjusted from \mathbf{w}_i using a violation point \mathbf{p} of label 1. In other words, $\mathbf{w}_{i+1} = \mathbf{w}_i + \mathbf{p}$. Hence:

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

Since \mathbf{p} is a violation point, it must hold that $\frac{\mathbf{w}_i}{|\mathbf{w}_i|} \cdot \mathbf{p} < \gamma_{guess}/2 \leq \gamma_{opt}/2$. Furthermore, obviously, $|\mathbf{p}|^2 \leq R^2$. We thus have:

$$|\mathbf{w}_{i+1}|^2 \leq |\mathbf{w}_i|^2 + \gamma_{opt}|\mathbf{w}_i| + R^2 \leq \left(|\mathbf{w}_i| + \frac{R^2}{2|\mathbf{w}_i|} + \frac{\gamma_{opt}}{2} \right)^2.$$

The claim then follows. □

Claim 4: When $|\mathbf{w}_i| \geq \frac{2R^2}{\gamma_{opt}}$, $|\mathbf{w}_{i+1}| \leq |\mathbf{w}_i| + (3/4)\gamma_{opt}$.

Proof: Directly follows from Claim 3. □

Claim 5: $|\mathbf{w}_k| \leq \frac{2R^2}{\gamma_{opt}} + \frac{3k\gamma_{opt}}{4} + R.$

Proof: Let j be the largest i satisfying $|\mathbf{w}_i| < \frac{2R^2}{\gamma_{opt}}$. If $j = k$, then $|\mathbf{w}_k| < \frac{2R^2}{\gamma_{opt}}$, and we are done. Next, we focus on the case $j < k$; note that this means $|\mathbf{w}_{j+1}|, |\mathbf{w}_{j+2}|, \dots, |\mathbf{w}_k|$ are all at least $2R^2/\gamma_{opt}$.

$$\begin{aligned} |\mathbf{w}_k| &\leq |\mathbf{w}_{k-1}| + (3/4)\gamma_{opt} && \text{(Claim 4)} \\ &\leq |\mathbf{w}_{k-2}| + 2 \cdot (3/4)\gamma_{opt} && \text{(Claim 4)} \\ &\dots \\ &\leq |\mathbf{w}_{j+1}| + (k-j-1)(3/4)\gamma_{opt} && \text{(Claim 4)} \\ &\leq |\mathbf{w}_{j+1}| + (3k/4)\gamma_{opt} \\ &\leq |\mathbf{w}_j| + R + (3k/4)\gamma_{opt} && \text{(Claim 2)} \\ &\leq \frac{2R^2}{\gamma_{opt}} + R + (3k/4)\gamma_{opt}. \end{aligned}$$



Combining Claims 1 and 5 gives:

$$\begin{aligned}k\gamma_{\text{opt}} &\leq \frac{2R^2}{\gamma_{\text{opt}}} + \frac{3k\gamma_{\text{opt}}}{4} + R \Rightarrow \\k &\leq \frac{8R^2}{\gamma_{\text{opt}}^2} + \frac{4R}{\gamma_{\text{opt}}} \\(\text{by } R \geq \gamma_{\text{opt}}) &\leq \frac{8R^2}{\gamma_{\text{opt}}^2} + \frac{4R^2}{\gamma_{\text{opt}}^2} \\&\leq \frac{12R^2}{\gamma_{\text{opt}}^2}.\end{aligned}$$

This completes the proof of the theorem.