

# ECE595 / STAT598: Machine Learning I

## Lecture 06 Linear Separability

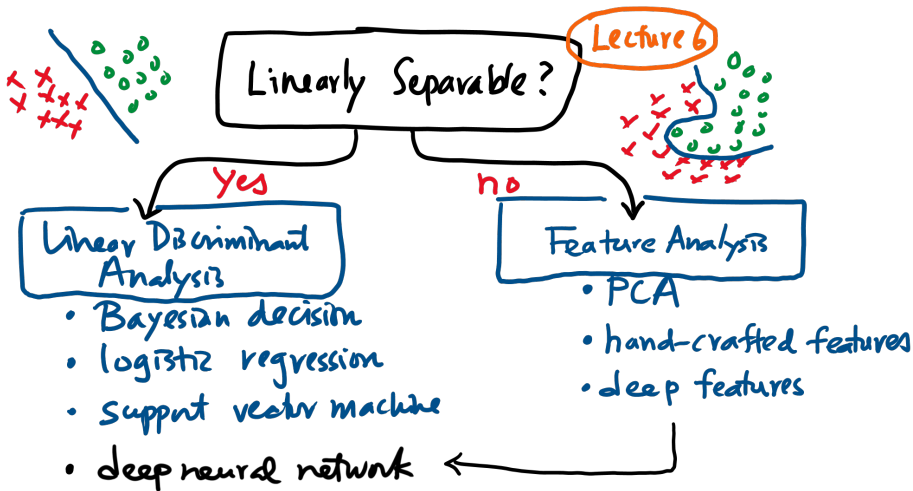
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Overview



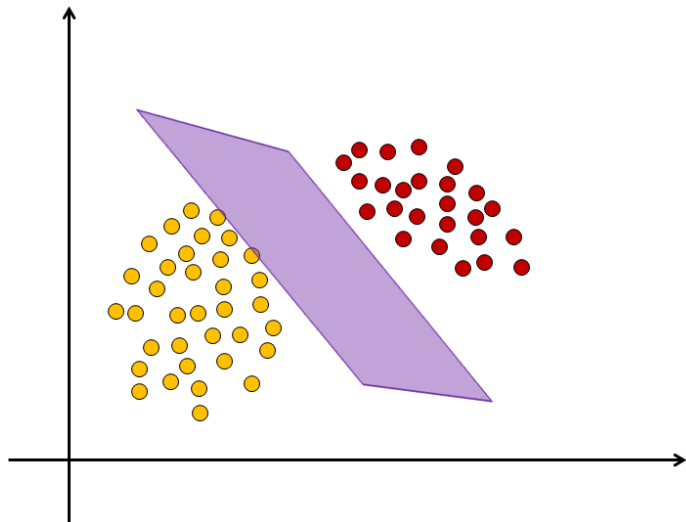
# Outline

**Goal:** Understand the geometry of linear separability.

- **Notations**
  - **Input Space, Output Space, Hypothesis**
  - **Discriminant Function**
- **Geometry of Discriminant Function**
  - Separating Hyperplane
  - Normal Vector
  - Distance from Point to Plane
- **Linear Separability**
  - Which set is linearly separable?
  - Separating Hyperplane Theorem
  - What if theorem fails?

# Supervised Classification

The goal of supervised classification is to construct a **decision boundary** such that the two classes can be (maximally) **separated**.



# Terminology

- **Input vectors:**  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ .
  - E.g., images, speech, EEG signal, rating, etc
- **Input space:**  $\mathcal{X}$ . Every  $\mathbf{x}_n \in \mathcal{X}$ .
- **Labels**  $y_1, y_2, \dots, y_N$ .
- **Label space:**  $\mathcal{Y}$ . Every  $y_n \in \mathcal{Y}$ .
  - If labels are binary, e.g.,  $y_n = \pm 1$ , then

$$\mathcal{Y} = \{+1, -1\}.$$

- Labels are arbitrary.  $\{+1, -1\}$  and  $\{0, 1\}$  has no difference.
- **Target function**  $f : \mathcal{X} \rightarrow \mathcal{Y}$ . Unknown.
  - Relationship:

$$y_n = f(\mathbf{x}_n).$$

- **Hypothesis**  $h : \mathcal{X} \rightarrow \mathcal{Y}$ . Ideally, want

$$h(\mathbf{x}) \approx f(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}.$$

## Binary Case

If we restrict ourselves to binary classifier, then

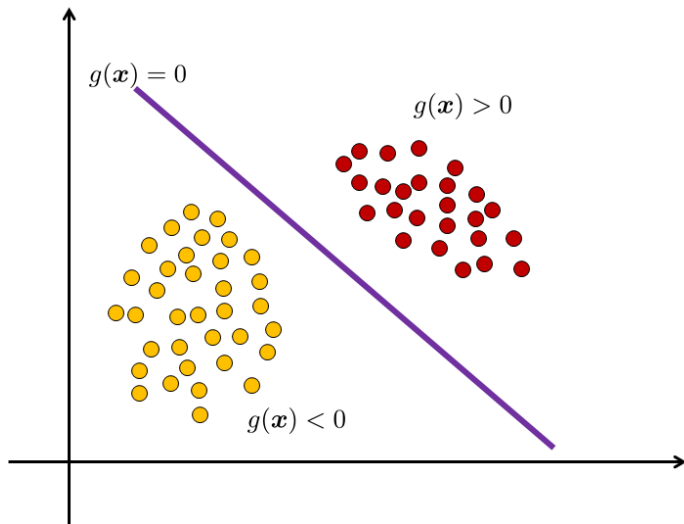
$$h(\mathbf{x}) = \begin{cases} 1, & \text{if } g(\mathbf{x}) > 0 \\ 0, & \text{if } g(\mathbf{x}) < 0 \\ \text{either,} & \text{if } g(\mathbf{x}) = 0 \end{cases}$$

- $g : \mathcal{X} \rightarrow \mathbb{R}$  is called a **discriminant function**.
- $g(\mathbf{x}) > 0$ :  $\mathbf{x}$  lives on the positive side of  $g$ .
- $g(\mathbf{x}) < 0$ :  $\mathbf{x}$  lives on the negative side of  $g$ .
- $g(\mathbf{x}) = 0$ : The decision boundary.
- You can also claim

$$h(\mathbf{x}) = \begin{cases} +1, & \text{if } g(\mathbf{x}) > 0 \\ -1, & \text{if } g(\mathbf{x}) < 0 \\ \text{either,} & \text{if } g(\mathbf{x}) = 0 \end{cases}$$

No difference as far as decision is concerned.

## Binary Case



# Linear Discriminant Function

A linear discriminant function takes the form

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

- $\mathbf{w} \in \mathbb{R}^d$ : linear coefficients
- $w_0 \in \mathbb{R}$ : bias / offset
- Define the overall **parameter**

$$\boldsymbol{\theta} = \{\mathbf{w}, w_0\} \in \mathbb{R}^{d+1}.$$

- Example:
  - If  $d = 2$ , then

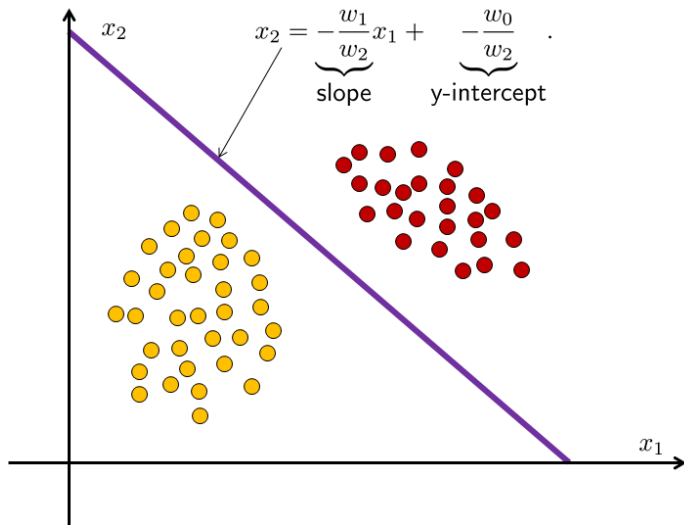
$$g(\mathbf{x}) = w_2 x_2 + w_1 x_1 + w_0.$$

- $g(\mathbf{x}) = 0$  means

$$x_2 = \underbrace{-\frac{w_1}{w_2}}_{\text{slope}} x_1 + \underbrace{-\frac{w_0}{w_2}}_{\text{y-intercept}}.$$



# Linear Discriminant Function



# Outline

**Goal:** Understand the geometry of linear separability.

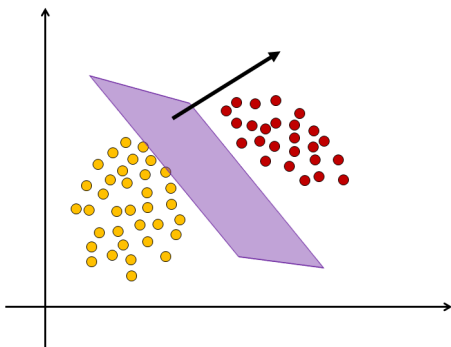
- Notations
  - Input Space, Output Space, Hypothesis
  - Discriminant Function
- Geometry of Discriminant Function
  - Separating Hyperplane
  - Normal Vector
  - Distance from Point to Plane
- Linear Separability
  - Which set is linearly separable?
  - Separating Hyperplane Theorem
  - What if theorem fails?

# Linear Discriminant Function

- In high-dimension,

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

is a hyperplane.

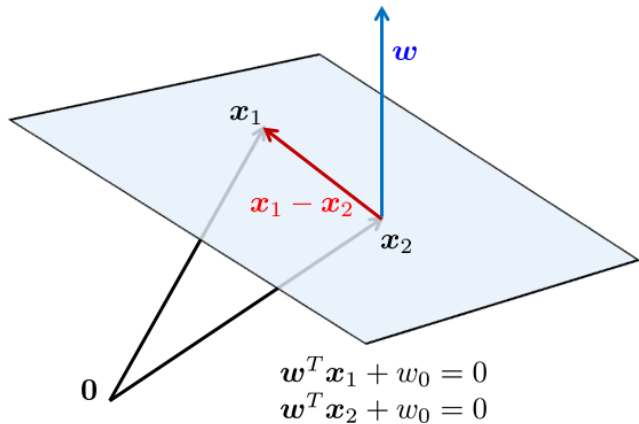


- **Separating Hyperplane:**

$$\begin{aligned}\mathcal{H} &= \{\mathbf{x} \mid g(\mathbf{x}) = 0\} \\ &= \{\mathbf{x} \mid \mathbf{w}^T \mathbf{x} + w_0 = 0\}\end{aligned}$$

- $\mathbf{x} \in \mathcal{H}$  means  $\mathbf{x}$  is on the decision boundary.
- $\mathbf{w} / \|\mathbf{w}\|_2$  is the **normal vector** of  $\mathcal{H}$ .

## Why is $w$ the Normal Vector?



## Why is $\mathbf{w}$ the Normal Vector?

- Pick  $\mathbf{x}_1$  and  $\mathbf{x}_2$  from  $\mathcal{H}$ .
- So  $g(\mathbf{x}_1) = 0$  and  $g(\mathbf{x}_2) = 0$ .
- This means:

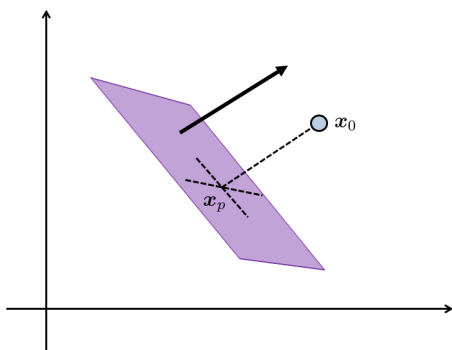
$$\mathbf{w}^T \mathbf{x}_1 + w_0 = 0, \quad \text{and} \quad \mathbf{w}^T \mathbf{x}_2 + w_0 = 0.$$

- Consider the difference vector  $\mathbf{x}_1 - \mathbf{x}_2$ .
- $\mathbf{x}_1 - \mathbf{x}_2$  is the tangent vector on the surface of  $\mathcal{H}$ .
- Check

$$\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = (\mathbf{w}^T \mathbf{x}_1 + w_0) - (\mathbf{w}^T \mathbf{x}_2 + w_0) = 0.$$

- So  $\mathbf{w}$  is perpendicular to  $\mathbf{x}_1 - \mathbf{x}_2$ , hence it is the normal.
- Normalize  $\mathbf{w} / \|\mathbf{w}\|_2$  so that it has unit norm.

## Distance from $\mathbf{x}_0$ to $g(\mathbf{x}) = 0$



Therefore, we can show that

$$\begin{aligned}g(\mathbf{x}_0) &= \mathbf{w}^T \mathbf{x}_0 + w_0 \\ &= \mathbf{w}^T \left( \mathbf{x}_p + \eta \frac{\mathbf{w}}{\|\mathbf{w}\|_2} \right) + w_0 \\ &= g(\mathbf{x}_p) + \eta \|\mathbf{w}\|_2 = \eta \|\mathbf{w}\|_2.\end{aligned}$$

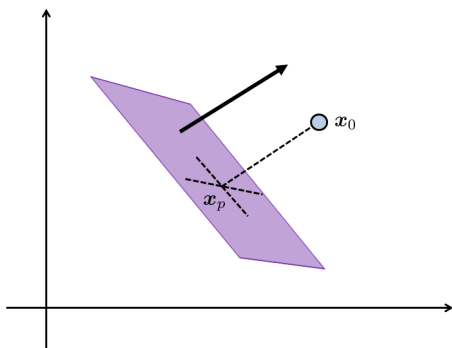
- Pick a point  $\mathbf{x}_p$  on  $\mathcal{H}$
- $\mathbf{x}_p$  is the closest point to  $\mathbf{x}_0$
- $\mathbf{x}_0 - \mathbf{x}_p$  is the normal direction
- So, for some scalar  $\eta > 0$ ,

$$\mathbf{x}_0 - \mathbf{x}_p = \eta \frac{\mathbf{w}}{\|\mathbf{w}\|_2}$$

- $\mathbf{x}_p$  is on  $\mathcal{H}$ . So

$$g(\mathbf{x}_p) = \mathbf{w}^T \mathbf{x}_p + w_0 = 0$$

## Distance from $x_0$ to $g(x) = 0$



- So distance is

$$\eta = \frac{g(x_0)}{\|w\|_2}$$

- The closest point  $x_p$  is

$$\begin{aligned}x_p &= x_0 - \eta \frac{w}{\|w\|_2} \\ &= x_0 - \frac{g(x_0)}{\|w\|_2} \cdot \frac{w}{\|w\|_2}.\end{aligned}$$

**Conclusion:**

$$x_p = x_0 - \underbrace{\frac{g(x_0)}{\|w\|_2}}_{\text{distance}} \cdot \underbrace{\frac{w}{\|w\|_2}}_{\text{normal vector}}$$

## Distance from $\mathbf{x}_0$ to $g(\mathbf{x}) = 0$

### Alternative Solution:

We can also obtain the same result by solving the optimization:

$$\mathbf{x}_p = \underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 \quad \text{subject to} \quad \mathbf{w}^T \mathbf{x} + w_0 = 0.$$

- Let Lagrangian

$$\mathcal{L}(\mathbf{x}, \lambda) = \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 - \lambda(\mathbf{w}^T \mathbf{x} + w_0)$$

- Stationarity condition implies

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda) &= (\mathbf{x} - \mathbf{x}_0) - \lambda \mathbf{w} = \mathbf{0}, \\ \nabla_{\lambda} \mathcal{L}(\mathbf{x}, \lambda) &= \mathbf{w}^T \mathbf{x} + w_0 = 0. \end{aligned}$$



## Distance from $\mathbf{x}_0$ to $g(\mathbf{x}) = 0$

Let us do some derivation:

$$\begin{aligned}\nabla_{\mathbf{x}}\mathcal{L}(\mathbf{x}, \lambda) &= (\mathbf{x} - \mathbf{x}_0) - \lambda \mathbf{w} = \mathbf{0}, \\ \nabla_{\lambda}\mathcal{L}(\mathbf{x}, \lambda) &= \mathbf{w}^T \mathbf{x} + w_0 = 0.\end{aligned}$$

- This gives

$$\begin{aligned}\mathbf{x} &= \mathbf{x}_0 + \lambda \mathbf{w} \\ \Rightarrow \mathbf{w}^T \mathbf{x} + w_0 &= \mathbf{w}^T (\mathbf{x}_0 + \lambda \mathbf{w}) + w_0 \\ \Rightarrow 0 &= \mathbf{w}^T \mathbf{x}_0 + \lambda \|\mathbf{w}\|^2 + w_0 \\ \Rightarrow 0 &= g(\mathbf{x}_0) + \lambda \|\mathbf{w}\|^2 \\ \Rightarrow \lambda &= -\frac{g(\mathbf{x}_0)}{\|\mathbf{w}\|^2} \\ \Rightarrow \mathbf{x} &= \mathbf{x}_0 + \left(-\frac{g(\mathbf{x}_0)}{\|\mathbf{w}\|^2}\right) \mathbf{w}.\end{aligned}$$

- Therefore, we arrive at the same result:

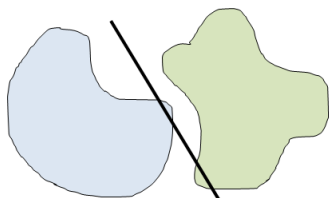
$$\mathbf{x}_p = \mathbf{x}_0 - \underbrace{\frac{g(\mathbf{x}_0)}{\|\mathbf{w}\|_2}}_{\text{distance}} \cdot \underbrace{\frac{\mathbf{w}}{\|\mathbf{w}\|_2}}_{\text{normal vector}}$$

# Outline

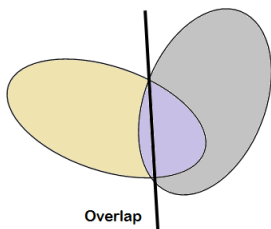
**Goal:** Understand the geometry of linear separability.

- Notations
  - Input Space, Output Space, Hypothesis
  - Discriminant Function
- Geometry of Discriminant Function
  - Separating Hyperplane
  - Normal Vector
  - Distance from Point to Plane
- Linear Separability
  - Which set is linearly separable?
  - Separating Hyperplane Theorem
  - What if theorem fails?

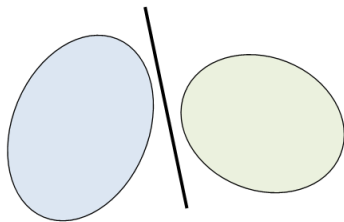
# Which one is Linearly Separable? Which one is Not?



Non-convex



Overlap



Linearly Separable

# Separating Hyperplane Theorem

Can we always find a separating hyperplane?

- No.
- Unless the classes are linearly separable.
- If convex and not overlapping, then yes.

## Theorem (Separating Hyperplane Theorem)

Let  $\mathcal{C}_1$  and  $\mathcal{C}_2$  be two *closed convex sets* such that  $\mathcal{C}_1 \cap \mathcal{C}_2 = \emptyset$ . Then, there exists a linear function

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

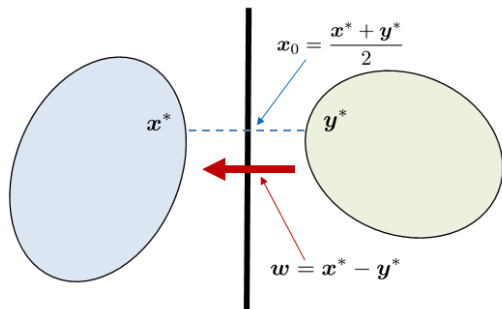
such that  $g(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \mathcal{C}_1$  and  $g(\mathbf{x}) < 0$  for all  $\mathbf{x} \in \mathcal{C}_2$ .

**Remark:** The theorem above provides sufficiency but not necessity for linearly separability.

# Separating Hyperplane Theorem

Pictorial “proof”:

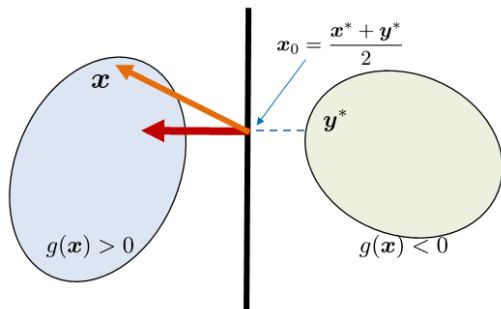
- Pick two points  $\mathbf{x}^*$  and  $\mathbf{y}^*$  s.t. the distance between the sets is minimized.
- Define the mid-point as  $\mathbf{x}_0 = (\mathbf{x}^* + \mathbf{y}^*)/2$ .
- Draw the separating hyperplane with normal  $\mathbf{w} = \mathbf{x}^* - \mathbf{y}^*$
- Convexity implies any inner product must be positive.



# Separating Hyperplane Theorem

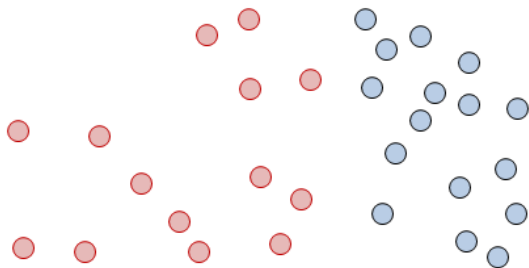
Pictorial “proof”:

- Pick two points  $\mathbf{x}^*$  and  $\mathbf{y}^*$  s.t. the distance between the sets is minimized.
- Define the mid-point as  $\mathbf{x}_0 = (\mathbf{x}^* + \mathbf{y}^*)/2$ .
- Draw the separating hyperplane with normal  $\mathbf{w} = \mathbf{x}^* - \mathbf{y}^*$
- Convexity implies any inner product must be positive.



## Linearly Separable?

- I have data  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .
- Closed. Convex. Non-overlapping.
- Separating hyperplane theorem: I can find a line.
- Victory?
- Not quite.



## When Theory Fails

### Theorem (Separating Hyperplane Theorem)

Let  $\mathcal{C}_1$  and  $\mathcal{C}_2$  be two closed convex sets such that  $\mathcal{C}_1 \cap \mathcal{C}_2 = \emptyset$ . Then, there exists a linear function

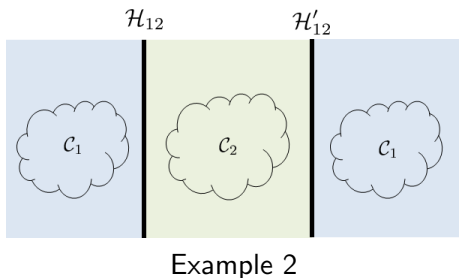
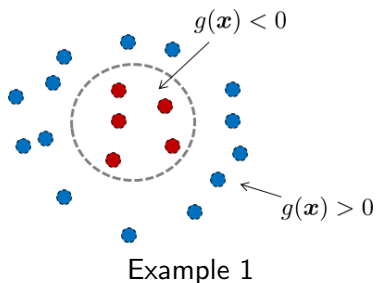
$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0,$$

such that  $g(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \mathcal{C}_1$  and  $g(\mathbf{x}) < 0$  for all  $\mathbf{x} \in \mathcal{C}_2$ .

- Finding a separating hyperplane for **training set** does not imply it will work for the **testing set**.
- Separating hyperplane theorem is more often used in **theoretical analysis** by assuming properties of the testing set.
- If a dataset is linearly separable, then you are guaranteed to find a perfect classifier. Then you can say how good is the classifier you designed compared to the perfect one.



# Linear Classifiers Do Not Work



- **Intrinsic geometry** of the two classes could be bad.
- The training set could be **lack of training samples**.
- Solution 1: Use non-linear classifiers, e.g.,  
$$g(\mathbf{x}) = \mathbf{x}^T \mathbf{W} \mathbf{x} + \mathbf{w}^T \mathbf{x} + \omega_0.$$
- Solution 2: Kernel method, e.g., Radial basis function.
- Solution 3: Extract features, e.g.,  $g(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$ .

# Reading List

## Separating Hyperplane:

- Duda, Hart and Stork's *Pattern Classification*, Chapter 5.1 and 5.2.
- Princeton ORFE-523, Lecture 5 on Separating hyperplane  
[http://www.princeton.edu/~amirali/Public/Teaching/ORF523/S16/ORF523\\_S16\\_Lec5\\_gh.pdf](http://www.princeton.edu/~amirali/Public/Teaching/ORF523/S16/ORF523_S16_Lec5_gh.pdf)
- Cornell ORIE-6300, Lecture 6 on Separating hyperplane  
<https://people.orie.cornell.edu/dpw/orie6300/fall2008/Lectures/lec06.pdf>
- Caltech, Lecture Note <http://www.its.caltech.edu/~kcborder/Notes/SeparatingHyperplane.pdf>

# Appendix

# Proof of Separating Hyperplane Theorem

- Conjecture: Let's see if this is the correct hyperplane

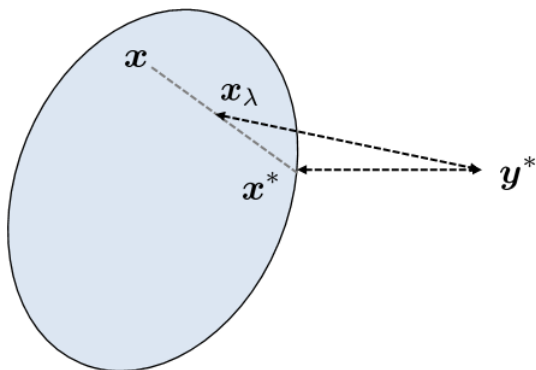
$$\begin{aligned}g(\mathbf{x}) &= \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0) \\&= (\mathbf{x}^* - \mathbf{y}^*)^T \left( \mathbf{x} - \frac{\mathbf{x}^* + \mathbf{y}^*}{2} \right) \\&= (\mathbf{x}^* - \mathbf{y}^*)^T \mathbf{x} - \frac{\|\mathbf{x}^*\|^2 - \|\mathbf{y}^*\|^2}{2}\end{aligned}$$

- According to picture, we want  $g(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \mathcal{C}_1$ .
- Suppose not. Assume

$$g(\mathbf{x}) = (\mathbf{x}^* - \mathbf{y}^*)^T \mathbf{x} - \frac{\|\mathbf{x}^*\|^2 - \|\mathbf{y}^*\|^2}{2} < 0.$$

See if we can find a contradiction.

## Proof of Separating Hyperplane Theorem



- $\mathcal{C}_1$  is convex.
- Pick  $\mathbf{x} \in \mathcal{C}_1$
- Pick  $\mathbf{x}^* \in \mathcal{C}_1$
- Let  $0 \leq \lambda \leq 1$
- Construct a point

$$\mathbf{x}_\lambda = (1 - \lambda)\mathbf{x}^* + \lambda\mathbf{x}.$$

- Convex means

$$\mathbf{x}_\lambda \in \mathcal{C}_1$$

So we must have

$$\|\mathbf{x}_\lambda - \mathbf{y}^*\| \geq \|\mathbf{x}^* - \mathbf{y}^*\|$$

# Proof of Separating Hyperplane Theorem

- Pick an arbitrary point  $\mathbf{x} \in \mathcal{C}_1$ .
- $\mathbf{x}^*$  is fixed already.
- Pick  $\mathbf{x}_\lambda$  along the line connecting  $\mathbf{x}$  and  $\mathbf{x}^*$ .
- Convexity implies  $\mathbf{x}_\lambda \in \mathcal{C}_1$ .
- So  $\|\mathbf{x}_\lambda - \mathbf{y}^*\| \geq \|\mathbf{x}^* - \mathbf{y}^*\|$ . If not, something is wrong.
- Let us do some algebra:

$$\begin{aligned}\|\mathbf{x}_\lambda - \mathbf{y}^*\|^2 &= \|(1 - \lambda)\mathbf{x}^* + \lambda\mathbf{x} - \mathbf{y}^*\|^2 \\ &= \|\mathbf{x}^* - \mathbf{y}^* + \lambda(\mathbf{x} - \mathbf{x}^*)\|^2 \\ &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 + 2\lambda(\mathbf{x}^* - \mathbf{y}^*)^T(\mathbf{x} - \mathbf{x}^*) + \lambda^2\|\mathbf{x} - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 + 2\lambda\mathbf{w}^T(\mathbf{x} - \mathbf{x}^*) + \lambda^2\|\mathbf{x} - \mathbf{x}^*\|^2.\end{aligned}$$

- Remember:  $\mathbf{w}^T(\mathbf{x} - \mathbf{x}_0) < 0$ .

## Proof of Separating Hyperplane Theorem

$$\begin{aligned}\|\mathbf{x}_\lambda - \mathbf{y}^*\|^2 &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 + 2\lambda \mathbf{w}^T (\mathbf{x} - \mathbf{x}^*) + \lambda^2 \|\mathbf{x} - \mathbf{x}^*\|^2 \\ &< \|\mathbf{x}^* - \mathbf{y}^*\|^2 + 2\lambda (\mathbf{w}^T \mathbf{x}_0 - \mathbf{w}^T \mathbf{x}^*) + \lambda^2 \|\mathbf{x} - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 + 2\lambda \left[ \left( \frac{\|\mathbf{x}^*\|^2 - \|\mathbf{y}^*\|^2}{2} \right) - \mathbf{w}^T \mathbf{x}^* \right] \\ &\quad + \lambda^2 \|\mathbf{x} - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 - \underbrace{\lambda \|\mathbf{x}^* - \mathbf{y}^*\|^2}_{=A} + \lambda^2 \underbrace{\|\mathbf{x} - \mathbf{x}^*\|^2}_{=B} \\ &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 - \lambda A + \lambda^2 B \\ &= \|\mathbf{x}^* - \mathbf{y}^*\|^2 - \lambda(A - \lambda B).\end{aligned}$$

Now, pick an  $\mathbf{x}$  such that  $A - \lambda B > 0$ . Then  $-\lambda(A - \lambda B) < 0$ .

$$\lambda < \frac{A}{B} = \frac{\|\mathbf{x}^* - \mathbf{y}^*\|^2}{\|\mathbf{x} - \mathbf{x}^*\|^2}.$$

## Proof of Separating Hyperplane Theorem

Therefore, if we choose  $\lambda$  such that  $A - \lambda B > 0$ , i.e.,

$$\lambda < \frac{A}{B} = \frac{\|\mathbf{x}^* - \mathbf{y}^*\|^2}{\|\mathbf{x} - \mathbf{x}^*\|^2},$$

then  $-\lambda(A - \lambda B) < 0$ , and so

$$\begin{aligned}\|\mathbf{x}_\lambda - \mathbf{y}^*\|^2 &< \|\mathbf{x}^* - \mathbf{y}^*\|^2 - \lambda(A - \lambda B) \\ &< \|\mathbf{x}^* - \mathbf{y}^*\|^2\end{aligned}$$

Contradiction, because  $\|\mathbf{x}^* - \mathbf{y}^*\|^2$  should be the smallest!

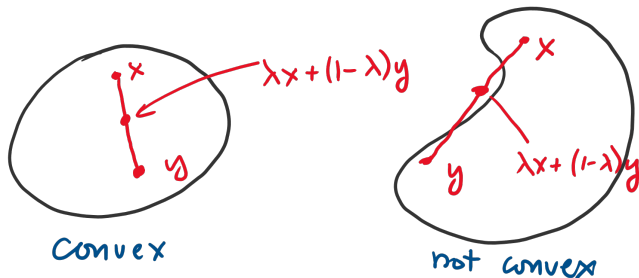
### Conclusion:

- If  $\mathbf{x} \in \mathcal{C}_1$ , then  $g(\mathbf{x}) > 0$ .
- By symmetry, if  $\mathbf{x} \in \mathcal{C}_2$ , then  $g(\mathbf{x}) < 0$ .
- And we have found the separating hyperplane  $(\mathbf{w}, w_0)$ .



## Q&A 1: What is a convex set?

- A set  $C$  is convex if the following condition is met.
- Pick  $\mathbf{x} \in C$  and  $\mathbf{y} \in C$ , and let  $0 < \lambda < 1$ . If  $\lambda\mathbf{x} + (1 - \lambda)\mathbf{y}$  is also in  $C$  for any  $\mathbf{x}$ ,  $\mathbf{y}$  and  $\lambda$ , then  $C$  is convex.
- Basically, it says that you can pick two points and draw a line. If the line is also in the set, then the set is convex.



## Q&A 2: Is there a way to check whether two sets are linearly separable?

- No, at least I do not know.
- The best you can do is to check whether a **training set** is linearly separable.
- To do so, solve the **hard SVM**. If you can solve it with zero training error, then you have found one. If the hard SVM does not have a solution, then the training set is not separable.
- Checking the **testing set** is impossible unless you know the distributions of the samples. But if you know the distributions, you can derive formula to check linear separability.
- For example, Gaussians are not linearly separable because no matter how unlikely you can always find a sample that lives in the wrong side. Uniform distributions are linearly separable.
- Bottom line: Linear separability, in my opinion, is more of a theoretical tool to describe the **intrinsic property** of the problem. It is not for computational purposes.

## Q&A 3: If two sets are not convex, how do I know if it is linearly separable?

- You can look at the convex hull.
- A convex hull is the smallest convex set that contains the original set.
- If the convex hulls are not overlapping, then linearly separable.
- For additional information about convex sets, convex hulls, you can check Chapter 2 of

<https://web.stanford.edu/class/ee364a/lectures.html>

