

ROB 101 - Fall 2021

# Hyperplanes in $\mathbb{R}^n$ , Quadratic Program, and Maximum Margin Classifier

November 29, 2021



- ▶ Introduce material that is assumed in UofM Computer Science courses that have Math 214 as a prerequisite.
- ▶ Provide a resource for use after you leave ROB 101.

- ▶ Learn how to separate  $\mathbb{R}^n$  into two halves via hyperplanes.
- ▶ The notion of signed distance to a hyperplane.
- ▶ An example of a max-margin classifier, a common tool in Machine Learning.

# Separating Hyperplanes

We want to study linear structures than can be used to divide  $\mathbb{R}^n$  into pieces.

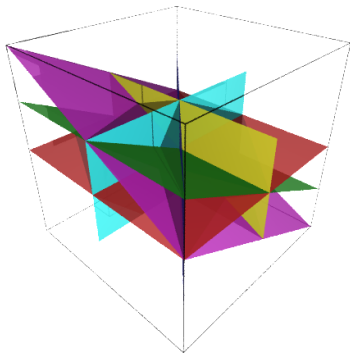
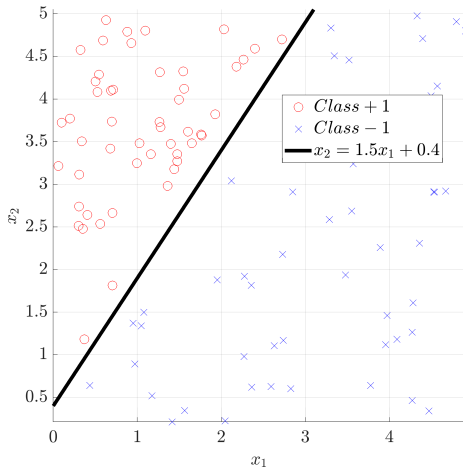
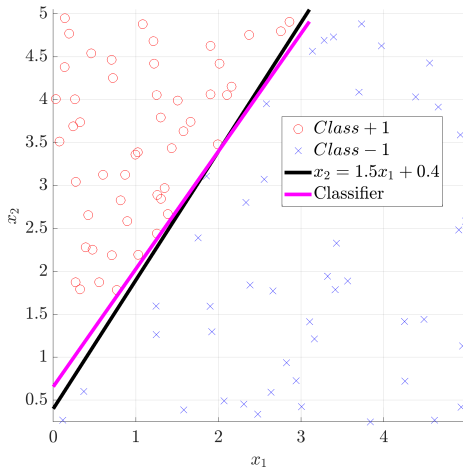


Figure: Dividing  $\mathbb{R}^3$  into disjoint regions. Image from Wikimedia Commons.

## Example: Classification



## Example: Classification

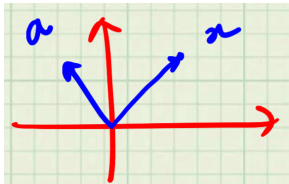


## Lines in $\mathbb{R}^2$ as Separating Hyperplanes

- Consider the set of all points,  $x \in \mathbb{R}^2$  such that

$$\langle a, x \rangle = 0, \quad a \in \mathbb{R}^2.$$

- $a = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$ ,  $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ , and  $\langle a, x \rangle = a_1x_1 + a_2x_2$ .



## Lines in $\mathbb{R}^2$ as Separating Hyperplanes

- ▶ Consider the set of all points,  $x \in \mathbb{R}^2$  such that

$$\langle a, x \rangle = 0, \quad a \in \mathbb{R}^2.$$

- ▶ Writing it in the set notation:

$$L = \{x \in \mathbb{R}^2 \mid \langle a, x \rangle = 0, \ a \in \mathbb{R}^2\}.$$

- ▶  $L$  is a line that passes through the origin.



- ▶ Consider the set of all points,  $x \in \mathbb{R}^n$  such that

$$\langle a, x \rangle = 0, \quad a \in \mathbb{R}^n.$$

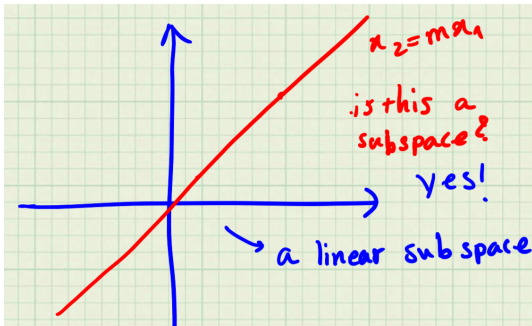
- ▶ Writing it in the set notation:

$$H = \{x \in \mathbb{R}^n \mid \langle a, x \rangle = 0, \ a \in \mathbb{R}^n\}.$$

- ▶  $H$  is a hyperplane that passes through the origin.

# Separating Hyperplanes in $\mathbb{R}^n$

- ▶  $H = \{x \in \mathbb{R}^n \mid \langle a, x \rangle = 0, a \in \mathbb{R}^n\}$ .
- ▶  $\langle a, x \rangle = 0 \iff a$  (normal vector) is orthogonal to all vectors that lie on the hyperplane.



## Lines in $\mathbb{R}^2$ as Separating Hyperplanes

- ▶ We can divide  $\mathbb{R}^2$  into two halves.
- ▶ Indeed, we define the following half-planes.

$$H^+ := \{x \in \mathbb{R}^2 \mid \langle a, x \rangle > 0\},$$

$$H^- := \{x \in \mathbb{R}^2 \mid \langle a, x \rangle < 0\}.$$

### Remark

*Using the angle between  $a$  and  $x$ ,  $\cos \theta = \frac{\langle a, x \rangle}{\|a\| \cdot \|x\|}$ , we see that  $\langle a, x \rangle > 0$  is the side ( $H^+$ ) where the angle is less than 90 deg, and  $\langle a, x \rangle < 0$  is the side ( $H^-$ ) where the angle is greater than 90 deg.*

The same observation holds in  $\mathbb{R}^n$ .

- ▶ We can divide  $\mathbb{R}^n$  into two halves.
- ▶ Indeed, we define the following half-planes.

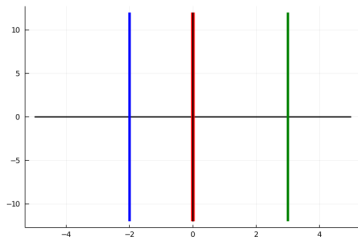
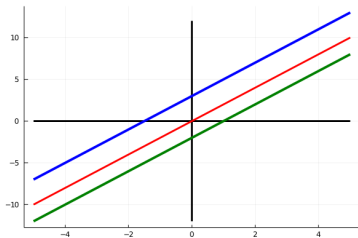
$$H^+ := \{x \in \mathbb{R}^n \mid \langle a, x \rangle > 0\},$$

$$H^- := \{x \in \mathbb{R}^n \mid \langle a, x \rangle < 0\}.$$

# Affine Subspace (Linear Variety)

If we take a subspace and translate it by  $x_c$ , then we get an affine subspace.

►  $M = x_c + H = \{x \in \mathbb{R}^n \mid \langle a, x - x_c \rangle = 0, a, x_c \in \mathbb{R}^n\}.$



If we take a subspace and translate it by  $x_c$ , then we get an affine subspace.

►  $M = x_c + H = \{x \in \mathbb{R}^n \mid \langle a, x - x_c \rangle = 0, a, x_c \in \mathbb{R}^n\}.$

### Remark

*Previously, we had  $x_c = 0$ . If  $x_c \neq 0$ , then the vector  $x - x_c$  lie on the hyperplane. Hence,  $\langle a, x - x_c \rangle = 0$ .*

We can look into the translated subspace as follows.

$$\langle a, x - x_c \rangle = a^\top (x - x_c) = 0$$

$$a^\top x = a^\top x_c =: d, \quad d \in \mathbb{R}.$$

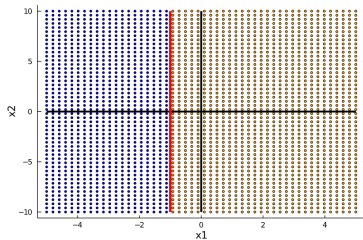
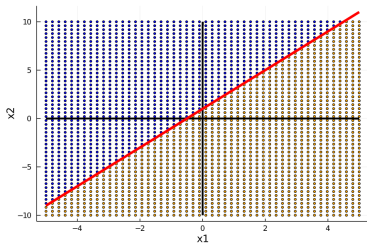
### Remark

*If  $a$  is normalized, then  $d = a^\top x_c$  is the distance from the origin to the hyperplane, i.e., the length of a vector parallel to  $a$  that starts from the origin and ends at the hyperplane.*

We can divide  $\mathbb{R}^n$  into two halves.

$$H^+ := \{x \in \mathbb{R}^n \mid \langle a, x - x_c \rangle > 0\},$$

$$H^- := \{x \in \mathbb{R}^n \mid \langle a, x - x_c \rangle < 0\}.$$





## Signed Distance to a Hyperplane

- ▶ For any point that does not lie on the hyperplane, we have  $\langle a, x - x_c \rangle \neq 0$  or  $a^\top x \neq a^\top x_c$ .
- ▶ We define the signed distance of a point to the hyperplane by the amount of deviation from the hyperplane equation.

$$y(x) = \frac{\langle a, x - x_c \rangle}{\|a\|}$$

- ▶ We normalize by  $\|a\|$  to avoid scaling the space.

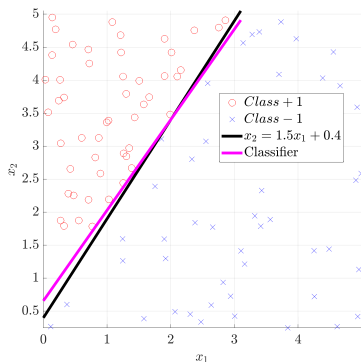
Supervised machine learning can be divided into two categories:

- 1 Regression; in this case, the outputs (also called target values) are continuous (real numbers).
- 2 Classification; in this case, the outputs (targets) are discrete categories (called labels).

You have seen regression problems in ROB 101. In this example, we will formulate a classification problem.

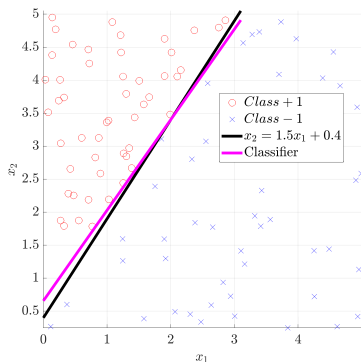
# Maximum Margin Classifier

- ▶ We wish to find a classifier (here a hyperplane) that separates  $\times$  and  $\circ$  categories.
- ▶ Furthermore, we want to predict the label for a new input (called query or test point).



# Maximum Margin Classifier

- ▶ We wish to find a classifier (here a hyperplane) that separates  $\times$  and  $\circ$  categories.
- ▶ Furthermore, we want to predict the label for a new input (called query or test point).



- ▶ We are given a data set  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ , where the inputs are  $x_i \in \mathbb{R}^2$  and targets are  $y_i \in \{+1, -1\}$ .
- ▶ Our model is a hyperplane  $a^\top x + a_0 = 0$ .  $a_0$  is called the bias term.
- ▶ Define  $w := \begin{bmatrix} a \\ a_0 \end{bmatrix}$  and  $\bar{x} := \begin{bmatrix} x \\ 1 \end{bmatrix}$ . Then  $w^\top \bar{x} = 0$ .

We define the following *hard margins*.

- ▶  $w^T \bar{x} = 1$ , anything on or above this boundary belongs to class  $+1$ .
- ▶  $w^T \bar{x} = -1$ , anything on or below this boundary belongs to class  $-1$ .
- ▶ We get the following constraints:

$$\begin{aligned}w^T \bar{x}_i &\geq 1, & \text{if } y_i = 1, \\w^T \bar{x}_i &\leq -1, & \text{if } y_i = -1.\end{aligned}$$

- ▶ We get the following constraints:

$$w^T \bar{x}_i \geq 1, \quad \text{if } y_i = 1,$$

$$w^T \bar{x}_i \leq -1, \quad \text{if } y_i = -1.$$

- ▶ We can combine both constraints into one as

$$y_i \cdot w^T \bar{x}_i \geq 1, \quad \text{for } i = 1, \dots, n.$$

We now formulate the following *constrained* optimization problem.

$$\begin{array}{ll} \min_{w \in \mathbb{R}^3} & \frac{1}{2} w^\top w \\ \text{subject to} & y_i \cdot w^\top \bar{x}_i \geq 1, \quad \text{for } i = 1, \dots, n. \end{array}$$



### Remark

*Training is the process of finding (called estimating or learning depending on the context) “optimal”  $w^*$ .*

### Remark

*Testing is the process of evaluating the trained model for a new input (an example that was not seen before).*

Given a query point (new input)  $x_*$ , we can evaluate the signed distance and pass it through the sign function. This is called the decision or response function.

$$y_* = \text{sgn}(w^{\star\top} \bar{x}_*),$$

where  $\text{sgn}$  is the sign function.

$$\text{sgn}(x) = \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases}$$

A *Quadratic Program* (QP) is a special kind of optimization problem with *constraints*. The cost to be minimized is supposed to be quadratic, meaning that  $f : \mathbb{R}^m \rightarrow \mathbb{R}$  has the form

$$f(x) = \frac{1}{2}x^\top Qx + qx,$$

where  $Q$  is an  $m \times m$  symmetric matrix, meaning that  $Q^\top = Q$ , and where  $q$  is a  $1 \times m$  row vector.

We consider the QP

$$\begin{aligned} x^* = \quad & \arg \min_{x \in \mathbb{R}^m} \quad \frac{1}{2}x^\top Qx + qx \\ & A_{in}x \preceq b_{in} \\ & A_{eq}x = b_{eq} \\ & lb \preceq x \preceq ub \end{aligned}$$

and assume that  $Q$  is symmetric ( $Q^\top = Q$ ) and *positive definite* ( $x \neq 0 \implies x^\top Qx > 0$ ), and that the subset of  $\mathbb{R}^m$  defined by the constraints is non empty, that is

$$C := \{x \in \mathbb{R}^m \mid A_{in}x \preceq b_{in}, A_{eq}x = b_{eq}, lb \preceq x \preceq ub\} \neq \emptyset.$$

Then  $x^*$  exists and is unique.

Let's switch to the Julia notebook.

- ▶ Soft Margin Classifier, Gaussian Support Vector Machine
- ▶ Read Chapter 13 of ROB 101 Book. QP is in Chapter 12.