

Linear separation and feature maps

DATA 607 — Session 6 — 11/03/2020

Lines in \mathbb{R}^2 : Direction vectors

A line in \mathbb{R}^2 is a set of points of the form

$$L_{\mathbf{u},\mathbf{v}} := \{\mathbf{u} + t\mathbf{v} : t \in \mathbb{R}\},$$

where

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \in \mathbb{R}^2, \quad \mathbf{v} \neq \mathbf{0} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Note that $\mathbf{u} = \mathbf{u} + 0\mathbf{v} \in L_{\mathbf{u},\mathbf{v}}$.

$L_{\mathbf{u},\mathbf{v}}$ is called the line through \mathbf{u} with direction vector \mathbf{v} .

- $\mathbf{u}' \in L_{\mathbf{u},\mathbf{v}} \implies L_{\mathbf{u}',\mathbf{v}} = L_{\mathbf{u},\mathbf{v}}$
- $\mathbf{v}' = c\mathbf{v}, c \in \mathbb{R}, c \neq 0 \implies L_{\mathbf{u},\mathbf{v}'} = L_{\mathbf{u},\mathbf{v}}$

Lines in \mathbb{R}^2 : Half-planes; sides; normal vectors

Dot product:

$$\mathbf{u} \cdot \mathbf{v} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \cdot \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = u_1 v_1 + u_2 v_2 \in \mathbb{R}$$

\mathbf{u} and \mathbf{v} are orthogonal or perpendicular.

\mathbf{w} is a normal vector to $L_{\mathbf{u},\mathbf{v}}$ if $\mathbf{v} \cdot \mathbf{w} = 0$.

Lines in \mathbb{R}^2 : Half-planes and sides

If you delete a line, L , from \mathbb{R}^2 , you're left with two half-planes called the sides of L .

If \mathbf{w} is a nonzero normal vector to $L_{\mathbf{u},\mathbf{v}}$, then the sets

$$H_{\mathbf{u},\mathbf{w}}^- = \{\mathbf{x} : \mathbf{w} \cdot (\mathbf{x} - \mathbf{u}) < 0\} \quad \text{and} \quad H_{\mathbf{u},\mathbf{w}}^+ = \{\mathbf{x} : \mathbf{w} \cdot (\mathbf{x} - \mathbf{u}) > 0\}$$

are the sides of $L_{\mathbf{u},\mathbf{v}}$.

The normal vector \mathbf{w} , plotted with its tail on $L_{\mathbf{u},\mathbf{v}}$, points into $H_{\mathbf{u},\mathbf{w}}^+$.

Linear separation

Let

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

be a dataset, where $\mathbf{x}_i \in \mathbb{R}^2$ and $y_i \in \{-1, 1\}$.

Let

$$D^- = \{(\mathbf{x}_i, y_i) \in D : y_i = -1\}, \quad D^+ = \{(\mathbf{x}_i, y_i) \in D : y_i = +1\}$$

Let L be a line in \mathbb{R}^2 . We say that L separates D if D^- and D^+ are contained in opposite sides of L .

We say that D is linearly separable if there is a line L that separates D .

Not all datasets are linearly separable: The dataset

$$D := \left\{ ((0, 0), 0), ((1, 0), 1), ((1, 1), 0), ((0, 1), 1) \right\}$$

is not linearly separable.

Linear separators need not be unique: The linear separators of

$$D := \left\{ ((0, -1), 0), ((0, 1), 1) \right\}$$

are precisely the lines

$$y = mx + b, \quad b \in (-1, 1).$$

Finding linear separators

Problem: Given a dataset, D , find a vector \mathbf{u} and a nonzero vector \mathbf{w} such that

$$D \cap H_{\mathbf{u}, \mathbf{w}}^+ = D^+ \quad \text{and} \quad D \cap H_{\mathbf{u}, \mathbf{w}}^- = D^-$$

or show that no such \mathbf{u} and \mathbf{w} exist.

We'll begin by analyzing a special case:

Special case: Given a dataset D , find a nonzero vector \mathbf{w} such that

$$D \cap H_{0, \mathbf{w}}^+ = D^+ \quad \text{and} \quad D \cap H_{0, \mathbf{w}}^- = D^-$$

$L_{0, \mathbf{w}}$ does **not** separate D if and only if

$$D^- \cap H_{0, \mathbf{w}}^+ \neq \emptyset \quad \text{or} \quad D^+ \cap H_{0, \mathbf{w}}^- \neq \emptyset,$$

...

...or, equivalently, if and only if

$$\left(y_i = -1 \quad \text{and} \quad \mathbf{w} \cdot \mathbf{x}_i > 0 \right) \quad \text{or} \quad \left(y_i = +1 \quad \text{and} \quad \mathbf{w} \cdot \mathbf{x}_i < 0 \right)$$

for some i , or, equivalently, if and only if

$$y_i(\mathbf{w} \cdot \mathbf{x}_i) < 0$$

for some i , or, equivalently, if and only if

$$\min(y_i(\mathbf{w} \cdot \mathbf{x}_i), 0) < 0$$

for some i , or, equivalently,

$$\sum_i \min(y_i(\mathbf{w} \cdot \mathbf{x}_i), 0) < 0,$$

or, equivalently, if and only if

$$\sum_i \max(-y_i(\mathbf{w} \cdot \mathbf{x}_i), 0) > 0.$$

View the term

$$L(\mathbf{w}, \mathbf{x}_i, y_i) := \max(-y_i(\mathbf{w} \cdot \mathbf{x}_i), 0)$$

as a **penalty** or **loss** for \mathbf{x}_i being misclassified by \mathbf{w} , i.e., lying on the wrong side of the line through 0 normal to \mathbf{w} .

Assume \mathbf{x}_i is correctly classified, then

$$\left(y_i = -1 \quad \text{and} \quad \mathbf{w} \cdot \mathbf{x}_i < 0\right) \quad \text{or} \quad \left(y_i = +1 \quad \text{and} \quad \mathbf{w} \cdot \mathbf{x}_i > 0\right),$$

in which case $y_i(\mathbf{w} \cdot \mathbf{x}_i) > 0$ and

$$-y_i(\mathbf{w} \cdot \mathbf{x}_i) < 0.$$

Thus, the penalty assessed for \mathbf{x}_i being misclassified by \mathbf{w} is

$$L(\mathbf{w}, \mathbf{x}_i, y_i) = \max(-y_i(\mathbf{w} \cdot \mathbf{x}_i), 0) = 0,$$

appropriate since, by hypothesis, \mathbf{x}_i is classified correctly!

Conversely, Assume \mathbf{x}_i is correctly classified, then

$$\left(y_i = -1 \quad \text{and} \quad \mathbf{w} \cdot \mathbf{x}_i > 0\right) \quad \text{or} \quad \left(y_i = +1 \quad \text{and} \quad \mathbf{w} \cdot \mathbf{x}_i < 0\right),$$

in which case $y_i(\mathbf{w} \cdot \mathbf{x}_i) < 0$ and

$$-y_i(\mathbf{w} \cdot \mathbf{x}_i) > 0.$$

Thus, the penalty assessed for \mathbf{x}_i being misclassified by \mathbf{w} is strictly positive:

$$L(\mathbf{w}, \mathbf{x}_i, y_i) = \max(-y_i(\mathbf{w} \cdot \mathbf{x}_i), 0) > -y_i(\mathbf{w} \cdot \mathbf{x}_i).$$

Define the **cost** associated with \mathbf{w} by

$$C(D, \mathbf{w}) = \sum_i L(\mathbf{w}, \mathbf{x}_i, y_i) = \sum_i \max(-y_i(\mathbf{w} \cdot \mathbf{x}_i), 0).$$

Then $L(\mathbf{w})$ separates D if and only if

$$C(D, \mathbf{w}) = 0.$$

$$L(\mathbf{u}, \mathbf{w}, \mathbf{x}_i, y_i) = \sum_i \max(-y_i(\mathbf{w} \cdot (\mathbf{x}_i - \mathbf{u})), 0)$$

$$\begin{aligned} C(D, \mathbf{u}, \mathbf{w}) &= \sum_i L(\mathbf{u}, \mathbf{w}, \mathbf{x}_i, y_i) \\ &= \sum_i \max(-y_i(\mathbf{w} \cdot (\mathbf{x}_i - \mathbf{u})), 0) \end{aligned}$$

Find

$$\operatorname{argmin}_{\mathbf{u}, \mathbf{w}} C(D, \mathbf{u}, \mathbf{w}).$$

The Perceptron Algorithm

$$D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$$

Suppose that D is linearly separable, i.e., that there is unit vector a $\mathbf{w} \in \mathbb{R}^n$ such that

$$\text{sign}(\mathbf{w} \cdot \mathbf{x}_i) = y_i$$

for all i , or, equivalently, that

$$y_i(\mathbf{w} \cdot \mathbf{x}_i) > 0$$

for all i .

Define a sequence $\mathbf{w}_0, \mathbf{w}_1, \dots$ of vectors in \mathbb{R}^n by:

Algorithm The Perceptron Algorithm

```
1:  $k \leftarrow 0$ 
2:  $\mathbf{w}_0 \leftarrow 0$ 
3: while  $y_i(\mathbf{w}_k \cdot \mathbf{x}_i) \leq 0$  for some  $i$  do
4:    $k \leftarrow k + 1$ 
5:    $i_k \leftarrow \min\{i : y_i(\mathbf{w}_{k-1} \cdot \mathbf{x}_i) \leq 0\}$ 
6:    $\mathbf{w}_k \leftarrow \mathbf{w}_{k-1} + y_{i_k} \mathbf{x}_{i_k}$ 
7: end while
```

Note that $i_1 = 1$.

Theorem: (Rosenblatt, 1957) The perceptron algorithm terminates after $k < R^2/r^2$ steps.

Proof: Let k be such that $y_{i_k}(\mathbf{w}_{k-1} \cdot \mathbf{x}_{i_k}) \leq 0$.

Lower bound on $\|\mathbf{w}_k\|^2$:

Set

$$R = \max_i \|\mathbf{x}_i\| > 0.$$

$$\begin{aligned}\|\mathbf{w}_k\|^2 &= \|\mathbf{w}_{k-1} + y_{i_k} \mathbf{x}_{i_k}\|^2 \\&= \|\mathbf{w}_{k-1}\|^2 + \|\mathbf{x}_{i_k}\|^2 + 2y_{i_k}(\mathbf{w}_{k-1} \cdot \mathbf{x}_{i_k}) \\&\geq \|\mathbf{w}_{k-1}\|^2 + \|\mathbf{x}_{i_k}\|^2 \\&\geq \|\mathbf{w}_{k-1}\|^2 + R^2 \\&\geq \|\mathbf{w}_{k-2}\|^2 + 2R^2 \\&\vdots \\&\geq \|\mathbf{w}_0\|^2 + kR^2 \\&= kR^2\end{aligned}$$

Upper bound on $\|\mathbf{w}_k\|^2$:

Set

$$r = \min_i |\mathbf{w} \cdot \mathbf{x}_i| > 0.$$

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}_k &= \mathbf{u} \cdot (\mathbf{w}_{k-1} + y_{i_k} \mathbf{x}_{i_k}) \\ &= \mathbf{u} \cdot \mathbf{w}_{k-1} + y_{i_k} (\mathbf{u} \cdot \mathbf{x}_{i_k}) \\ &\geq \mathbf{u} \cdot \mathbf{w}_{k-1} + r \\ &\geq \mathbf{u} \cdot \mathbf{w}_{k-2} + 2r \\ &\vdots \\ &\geq \mathbf{u} \cdot \mathbf{w}_0 + kr \\ &= kr \end{aligned}$$

$$kr \leq \mathbf{u} \cdot \mathbf{w}_k \leq \|\mathbf{u}\| \|\mathbf{w}_k\| = \|\mathbf{w}_k\|$$

$$k^2 r^2 \leq \|\mathbf{w}_k\|^2$$

$$k^2 r^2 \leq \|\mathbf{w}_k\|^2 \leq kR^2$$

$$kr^2 \leq \|\mathbf{w}_k\|^2 \leq R^2$$

$$k \leq \frac{R^2}{r^2}$$

Feature maps

What do we do if our data isn't linearly separable?

Embed your data in a higher dimensional space in which it is linearly separable.

The higher dimensional space is called **feature space** and the function mapping **data space** — the ambient space of our data — into this feature space is called a **feature map**.

Consider the linearly inseparable dataset

$$D := \left\{ ((0, 0), 0), ((1, 0), 1), ((1, 1), 0), ((0, 1), 1) \right\}$$

Define a **feature map** $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ by

$$\phi(x_1, x_2) = (x_1, x_2, (x_1 - x_2)^2)$$

Then

$$\phi(D) = \left\{ ((0, 0, 0), 0), ((1, 0, 1), 1), ((1, 1, 0), 0), ((0, 1, 1), 1) \right\}$$

The points in with class label 0 and 1 have third coordinates 0 and 1, respectively.

The plane

$$x_3 = \frac{1}{2}$$

separates the classes.

The points in data space mapped by ϕ into this separating plane are those that satisfy

$$(x_1 - x_2)^2 = \frac{1}{2}$$

$$x_1 - x_2 = \pm \frac{1}{2}$$

This pair of lines separates the original dataset D .

More examples in the Jupyter notebook.

Where do features come from?

Many classification techniques require features to be hand-crafted for a given application.

A more robust approach is to, as much as possible, *learn* the features. Neural networks use this approach.