

COMS 4771 Lecture 5

1. Linear classifiers.
2. Linearly separable instances.

ANNOUNCEMENTS

1. Homework 1 due Monday Feb. 9 @ 1 PM.
 - ▶ Write-up, as a PDF file.
 - ▶ All MATLAB function files and scripts.
 - ▶ Everything in a ZIP file; submit on Courseworks.
2. TA office hours:
 - ▶ Thursday Feb. 5 @ 6–7 PM cancelled (sorry).
 - ▶ Additional TA office hours on Tuesday Feb. 10 @ 6–7 PM.
 - ▶ Need help with MATLAB? Your TAs can help! Lots of tips and tricks to be learned ...

LINEAR CLASSIFIERS

AXIS-ALIGNED THRESHOLD FUNCTIONS

Decision tree learning

Basic step in greedy decision tree learning (with axis-aligned splits in $\mathcal{X} = \mathbb{R}^d$):

$$\arg \min_h \frac{|S_{h,0}|}{|S|} u(S_{h,0}) + \frac{|S_{h,1}|}{|S|} u(S_{h,1})$$

where u is some uncertainty measure,

$$S_{h,0} = \{(x, y) \in S : h(x) = 0\}, \quad S_{h,1} = \{(x, y) \in S : h(x) = 1\},$$

and the minimization is over splitting rules of the form

$$h(\mathbf{x}) = \mathbb{1}\{x_i > t\}, \quad i \in [d], t \in \mathbb{R}.$$

AXIS-ALIGNED THRESHOLD FUNCTIONS

Decision tree learning

Basic step in greedy decision tree learning (with axis-aligned splits in $\mathcal{X} = \mathbb{R}^d$):

$$\arg \min_h \frac{|S_{h,0}|}{|S|} u(S_{h,0}) + \frac{|S_{h,1}|}{|S|} u(S_{h,1})$$

where u is some uncertainty measure,

$$S_{h,0} = \{(x, y) \in S : h(x) = 0\}, \quad S_{h,1} = \{(x, y) \in S : h(x) = 1\},$$

and the minimization is over splitting rules of the form

$$h(\mathbf{x}) = \mathbb{1}\{x_i > t\}, \quad i \in [d], t \in \mathbb{R}.$$

Axis-aligned threshold functions

When u is classification error and $\mathcal{Y} = \{-1, +1\}$, we are equivalently doing the following:

$$\arg \min_{i \in [d], v \in \{-1, +1\}, t \in \mathbb{R}} \sum_{(x, y) \in S} \mathbb{1}\{\text{sign}(vx_i - t) \neq y\}.$$

AXIS-ALIGNED THRESHOLD FUNCTIONS

Decision tree learning

Basic step in greedy decision tree learning (with axis-aligned splits in $\mathcal{X} = \mathbb{R}^d$):

$$\arg \min_h \frac{|S_{h,0}|}{|S|} u(S_{h,0}) + \frac{|S_{h,1}|}{|S|} u(S_{h,1})$$

where u is some uncertainty measure,

$$S_{h,0} = \{(x, y) \in S : h(x) = 0\}, \quad S_{h,1} = \{(x, y) \in S : h(x) = 1\},$$

and the minimization is over splitting rules of the form

$$h(\mathbf{x}) = \mathbb{1}\{x_i > t\}, \quad i \in [d], t \in \mathbb{R}.$$

Axis-aligned threshold functions

When u is classification error and $\mathcal{Y} = \{-1, +1\}$, we are equivalently doing the following:

$$\arg \min_{i \in [d], v \in \{-1, +1\}, t \in \mathbb{R}} \sum_{(x, y) \in S} \mathbb{1}\{\text{sign}(vx_i - t) \neq y\}.$$

i.e., looking at classifiers of the form $f_{i,v,t}(\mathbf{x}) = \text{sign}(vx_i - t)$.

LINEAR CLASSIFIERS (FOR BINARY CLASSIFICATION)

A natural generalization of axis-aligned threshold functions

$$f_{i,v,t}(\mathbf{x}) = \text{sign}(vx_i - t), \quad i \in [d], v \in \{-1, +1\}, t \in \mathbb{R},$$

are **linear threshold functions** (or **linear classifiers**):

$$f_{\mathbf{w},t}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle - t), \quad \mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}.$$

LINEAR CLASSIFIERS (FOR BINARY CLASSIFICATION)

A natural generalization of axis-aligned threshold functions

$$f_{i,v,t}(\mathbf{x}) = \text{sign}(vx_i - t), \quad i \in [d], v \in \{-1, +1\}, t \in \mathbb{R},$$

are **linear threshold functions** (or **linear classifiers**):

$$f_{\mathbf{w},t}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle - t), \quad \mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}.$$

Interpretation: does a weighted linear combination of input features exceed a threshold?

$$\langle \mathbf{w}, \mathbf{x} \rangle = \sum_{i=1}^d w_i x_i \stackrel{?}{>} t$$

LINEAR CLASSIFIERS (FOR BINARY CLASSIFICATION)

A natural generalization of axis-aligned threshold functions

$$f_{i,v,t}(\mathbf{x}) = \text{sign}(vx_i - t), \quad i \in [d], v \in \{-1, +1\}, t \in \mathbb{R},$$

are **linear threshold functions** (or **linear classifiers**):

$$f_{\mathbf{w},t}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle - t), \quad \mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}.$$

Interpretation: does a weighted linear combination of input features exceed a threshold?

$$\langle \mathbf{w}, \mathbf{x} \rangle = \sum_{i=1}^d w_i x_i \stackrel{?}{>} t$$

Tie at zero can go either way; we'll use the following convention:

$$\text{sign}(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0. \end{cases}$$

LINEAR CLASSIFIERS (FOR BINARY CLASSIFICATION)

A natural generalization of axis-aligned threshold functions

$$f_{i,v,t}(\mathbf{x}) = \text{sign}(vx_i - t), \quad i \in [d], v \in \{-1, +1\}, t \in \mathbb{R},$$

are **linear threshold functions** (or **linear classifiers**):

$$f_{\mathbf{w},t}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle - t), \quad \mathbf{w} \in \mathbb{R}^d, t \in \mathbb{R}.$$

Interpretation: does a weighted linear combination of input features exceed a threshold?

$$\langle \mathbf{w}, \mathbf{x} \rangle = \sum_{i=1}^d w_i x_i \stackrel{?}{>} t$$

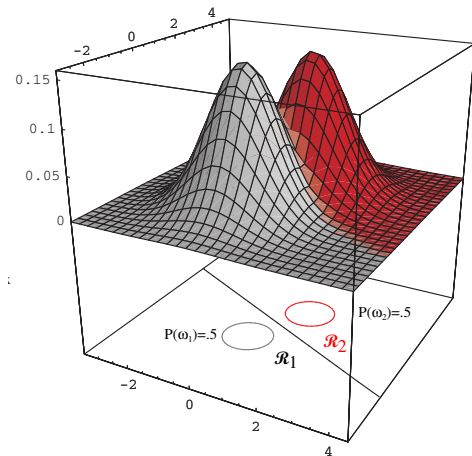
Tie at zero can go either way; we'll use the following convention:

$$\text{sign}(z) = \begin{cases} +1 & \text{if } z > 0 \\ -1 & \text{if } z \leq 0. \end{cases}$$

For now, only considering **binary classification**, where $\mathcal{Y} = \{-1, +1\}$.

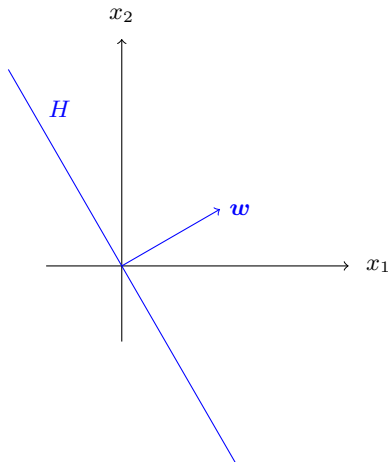
LINEAR CLASSIFIERS (FOR BINARY CLASSIFICATION)

We've seen these before: the (binary) Bayes classifier when class conditional densities are multivariate Gaussians with the same covariance.



HYPERPLANES

Geometric interpretation of linear classifiers



A **hyperplane** in \mathbb{R}^d is a linear subspace of dimension $(d - 1)$.

- ▶ A \mathbb{R}^2 -hyperplane is a line.
- ▶ A \mathbb{R}^3 -hyperplane is a plane.
- ▶ As a linear subspace, a hyperplane always contains the origin.

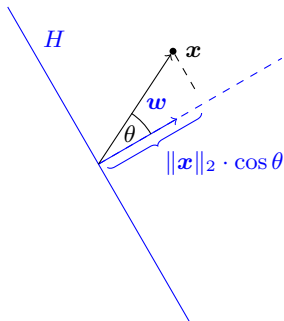
A hyperplane H can be specified by a (non-zero) **normal vector**.

The hyperplane with normal vector $w \in \mathbb{R}^d$ is the set

$$H = \{x \in \mathbb{R}^d : \langle w, x \rangle = 0\}.$$

It becomes **oriented** if we pick a *particular* normal vector $w \in \mathbb{R}^d$.

WHICH SIDE OF THE HYPERPLANE ARE WE ON?



Distance from the hyperplane

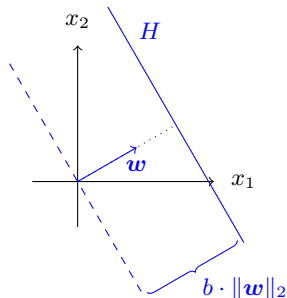
- ▶ The projection of x onto the direction of w has length $\frac{|\langle w, x \rangle|}{\|w\|_2}$.
- ▶ Cosine rule: $\cos \theta = \frac{\langle w, x \rangle}{\|w\|_2 \|x\|_2}$.
- ▶ The distance of x from the hyperplane is given by $\frac{|\langle w, x \rangle|}{\|w\|_2} = \|x\|_2 \cdot |\cos \theta|$.

Which side of the hyperplane?

- ▶ The cosine satisfies $\cos \theta > 0$ iff $\theta \in (-\frac{\pi}{2}, \frac{\pi}{2})$.
- ▶ We can determine which side of the hyperplane H that x is on, using

$$\text{sign}(\cos \theta) = \text{sign}(\langle w, x \rangle).$$

AFFINE HYPERPLANES



- ▶ An **affine hyperplane** H is a hyperplane translated (shifted) by a vector b : i.e., $H = b + H'$ for some hyperplane H' .
Without loss of generality, $H = bw + H'$, for some hyperplane H , $b \geq 0$, and normal vector w for H' .
- ▶ If $b > 0$, naturally oriented by which side contains the origin.

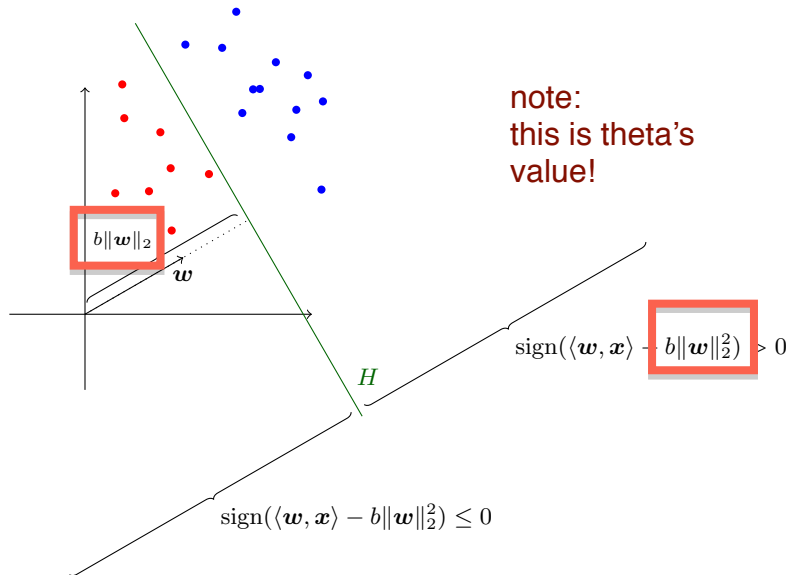
Which side of the affine hyperplane?

- ▶ We can determine which side of the affine hyperplane H that x is on using

$$\text{sign}(\langle x, w \rangle - b\|w\|_2^2).$$

side of affine hyperplane that x is on \equiv linear classification of x

LINEAR CLASSIFIERS



LEARNING LINEAR CLASSIFIERS

Even if the Bayes classifier is not a linear classifier, we can hope that it has a good linear approximation.

LEARNING LINEAR CLASSIFIERS

Even if the Bayes classifier is not a linear classifier, we can hope that it has a good linear approximation.

Goal: learning algorithm for linear classifiers with low **excess error**:

$$\underbrace{\mathbb{E}[\text{err}(f_{\hat{w}, \hat{t}})]}_{\text{expected error of your classifier}} - \underbrace{\min_{w, t} \text{err}(f_{w, t})}_{\text{error of best linear classifier}}$$

where $f_{\hat{w}, \hat{t}}$ is the linear classifier picked by the learning algorithm on the basis of an i.i.d. sample S from P . (Expectation is over S .)

LEARNING LINEAR CLASSIFIERS

Even if the Bayes classifier is not a linear classifier, we can hope that it has a good linear approximation.

Goal: learning algorithm for linear classifiers with low **excess error**:

$$\underbrace{\mathbb{E}[\text{err}(f_{\hat{\mathbf{w}}, \hat{t}})]}_{\text{expected error of your classifier}} - \underbrace{\min_{\mathbf{w}, t} \text{err}(f_{\mathbf{w}, t})}_{\text{error of best linear classifier}}$$

where $f_{\hat{\mathbf{w}}, \hat{t}}$ is the linear classifier picked by the learning algorithm on the basis of an i.i.d. sample S from P . (Expectation is over S .)

A natural approach is “**empirical risk minimization**” (ERM): find a linear classifier $f_{\mathbf{w}, t}$ with low training error (or **empirical risk**):

$$\begin{aligned} \arg \min_{\mathbf{w}, t} \text{err}(f_{\mathbf{w}, t}, S) &= \arg \min_{\mathbf{w}, t} \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \mathbb{1}\{\text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle - t) \neq y\} \\ &= \arg \min_{\mathbf{w}, t} \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} \mathbb{1}\{y(\langle \mathbf{w}, \mathbf{x} \rangle - t) \leq 0\}. \end{aligned}$$

EMPIRICAL RISK MINIMIZATION

Unfortunately, this is not possible in general.

- ▶ The following problem is NP-hard:

Given a set of labeled examples S in $\mathbb{R}^d \times \{\pm 1\}$ with the **promise** that there is a linear classifier with **training error 0.01**, find a linear classifier with **training error ≤ 0.49** .

EMPIRICAL RISK MINIMIZATION

Unfortunately, this is not possible in general.

- ▶ The following problem is NP-hard:

Given a set of labeled examples S in $\mathbb{R}^d \times \{\pm 1\}$ with the **promise** that there is a linear classifier with **training error 0.01**, find a linear classifier with **training error ≤ 0.49** .

Potential saving grace:

- ▶ Real-world problems we need to solve do not look like reductions from difficult SATISFIABILITY instances.

EMPIRICAL RISK MINIMIZATION

Unfortunately, this is not possible in general.

- ▶ The following problem is NP-hard:

Given a set of labeled examples S in $\mathbb{R}^d \times \{\pm 1\}$ with the **promise** that there is a linear classifier with **training error 0.01**, find a linear classifier with **training error ≤ 0.49** .

Potential saving grace:

- ▶ Real-world problems we need to solve do not look like reductions from difficult SATISFIABILITY instances.

Plan:

1. Study the **linearly separable** instances: where there is a linear classifier with **zero training error**.
2. Study **convex loss functions**, which can be efficiently minimized, and how they are related to classification error.

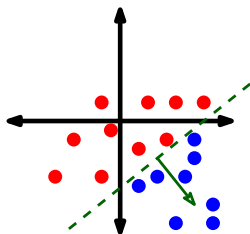
LINEARLY SEPARABLE INSTANCES

EASY CASE: LINEARLY SEPARABLE DATA

Suppose there is a linear classifier with **zero training error** on S : for some $\mathbf{w}_\star \in \mathbb{R}^d$ and $\theta \in \mathbb{R}$,

$$y(\langle \mathbf{w}_\star, \mathbf{x} \rangle - \theta) > 0, \quad \text{for all } (x, y) \in S.$$

In this case, we say the training data is **linearly separable**.



HOMOGENEOUS LINEAR CLASSIFICATION

Linear classifiers $f_{w,\theta}$ with $\theta = 0$ are called **homogeneous linear classifiers**.

HOMOGENEOUS LINEAR CLASSIFICATION

Linear classifiers $f_{\mathbf{w},\theta}$ with $\theta = 0$ are called **homogeneous linear classifiers**.

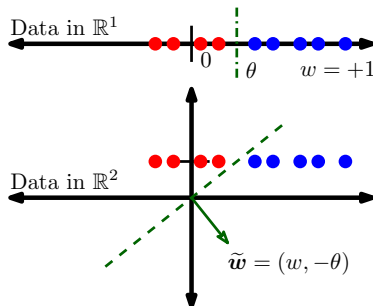
Claim: There is a mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d+1}$ with the following property. For any linear classifier $f_{\mathbf{w},\theta}: \mathbb{R}^d \rightarrow \{\pm 1\}$, there is a homogeneous linear classifier $f_{\tilde{\mathbf{w}},0}: \mathbb{R}^{d+1} \rightarrow \{\pm 1\}$ such that $f_{\mathbf{w},\theta}(\mathbf{x}) = f_{\tilde{\mathbf{w}},0}(\phi(\mathbf{x}))$ for all $\mathbf{x} \in \mathbb{R}^d$.

HOMOGENEOUS LINEAR CLASSIFICATION

Linear classifiers $f_{w,\theta}$ with $\theta = 0$ are called **homogeneous linear classifiers**.

Claim: There is a mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d+1}$ with the following property. For any linear classifier $f_{w,\theta}: \mathbb{R}^d \rightarrow \{\pm 1\}$, there is a homogeneous linear classifier $f_{\tilde{w},0}: \mathbb{R}^{d+1} \rightarrow \{\pm 1\}$ such that $f_{w,\theta}(x) = f_{\tilde{w},0}(\phi(x))$ for all $x \in \mathbb{R}^d$.

Proof: Let $\phi(x) := (x, 1)$ —i.e., add a $(d+1)$ -th coordinate that always takes value 1. For any $w \in \mathbb{R}^d$ and $\theta \in \mathbb{R}$, let $\tilde{w} := (w, -\theta)$.



FINDING A HOMOGENEOUS LINEAR SEPARATOR

Problem: given training data S in $\mathbb{R}^d \times \{\pm 1\}$, determine whether or not there exists $\mathbf{w} \in \mathbb{R}^d$

$$y\langle \mathbf{w}, \mathbf{x} \rangle > 0, \quad \text{for all } (\mathbf{x}, y) \in S;$$

(and find such a vector if one exists).

FINDING A HOMOGENEOUS LINEAR SEPARATOR

Problem: given training data S in $\mathbb{R}^d \times \{\pm 1\}$, determine whether or not there exists $\mathbf{w} \in \mathbb{R}^d$

$$y\langle \mathbf{w}, \mathbf{x} \rangle > 0, \quad \text{for all } (\mathbf{x}, y) \in S;$$

(and find such a vector if one exists).

This is a system of $|S|$ linear inequalities over d variables, and hence can be solved in polynomial time using algorithms for **linear programming** (e.g., ellipsoid algorithm, interior point).

FINDING A HOMOGENEOUS LINEAR SEPARATOR

Problem: given training data S in $\mathbb{R}^d \times \{\pm 1\}$, determine whether or not there exists $\mathbf{w} \in \mathbb{R}^d$

$$y\langle \mathbf{w}, \mathbf{x} \rangle > 0, \quad \text{for all } (\mathbf{x}, y) \in S;$$

(and find such a vector if one exists).

larger than 0 means
correct predication.

This is a system of $|S|$ linear inequalities over d variables, and hence can be solved in polynomial time using algorithms for **linear programming** (e.g., ellipsoid algorithm, interior point).

If one exists, and the inequalities in fact hold with some non-negligible

“margin” $\gamma > 0$:

$$y\langle \mathbf{w}, \mathbf{x} \rangle \geq \gamma, \quad \text{for all } (\mathbf{x}, y) \in S;$$

then there is a very **simple** algorithm that finds a solution: **Perceptron**.

the margin is great!!!