

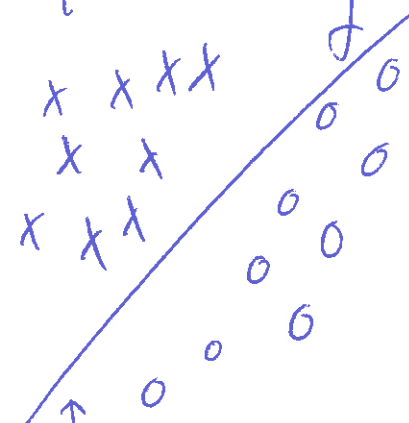
Lecture # 14

①

KNN classifier basically does no training... we only use the data when it is time to label, and even then, we only look locally.

An alternative (and a preferable one, generally), is to use all the data ahead of time to first do some work (training), that will make the subsequent labelling of points easier.

Let's try to "split" data by learning a hyperplane that separates two classes:



Can we efficiently learn this, or a similar hyperplane?

Set-up: Let $\{(x_i, y_i)\}_{i=1}^n$, $y_i \in \{1, -1\}$ (classes)
 $x_i \in \mathbb{R}^d$ (data points)

A hyperplane is a set of $x \in \mathbb{R}^d$ satisfying the equation
 $w^T x = -b$, some $w, b \in \mathbb{R}$.

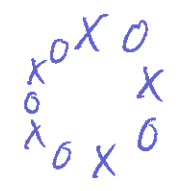
In our sketch, w^T enforces the "direction" of the dividing plane and b is the "offset" from the origin; not ~~unlike~~ unlike slope & y-intercept when discussing the equation of a line.

We want $\{x_i\}$ to be "discriminated" by w^T, b , in the following sense:

$$w^T x_i + b = \begin{cases} +, & y_i = 1 \\ -, & y_i = -1 \end{cases}$$

That is, points are on one side of the hyperplane if they have label 1, on the other side if they have label -1.

Remark: Not hard to construct datasets in which this is impossible, i.e. no separating hyperplanes:



Such examples can be dealt with using non-linear boundary fitting with kernel support vector machines.

Linear algebra set-up: Let $V_i := (\underbrace{y_i}_{\mathbb{R}^d}, \underbrace{x_i}_{\mathbb{R}^d}, y_i) \in \mathbb{R}^{d+1}$
 (from C.B. slides) $a = (w, b) \in \mathbb{R}^{d+1}$
 $\uparrow \mathbb{R}^d \quad \uparrow \mathbb{R}^d$

Choose u such that $v_i \cdot u > 0, i=1, \dots, n$.

This makes sense because $v_i \cdot u$

$$= (y_i x_i, y_i) \cdot (w, b)$$

$$= y_i (w^T x_i + b),$$

So want $y_i, w^T x_i + b$ to have the same sign.

In essence, we have n inequalities that need to hold simultaneously. This is easily addressed using linear programming.

A standard formulation to linear programming problems is:

$$\begin{aligned} & \text{minimize/maximize } c^T x \\ & \text{s.t. } Ax \leq b \\ & \quad \text{or} \\ & \quad \geq b \end{aligned}$$

for some matrix A .

Note that we don't even need to minimize/maximize here... we

just need $y_i (w^T x_i + b)$

$$= v_i \cdot u > 0 \quad \text{for every } i, \text{ which is analogous to}$$

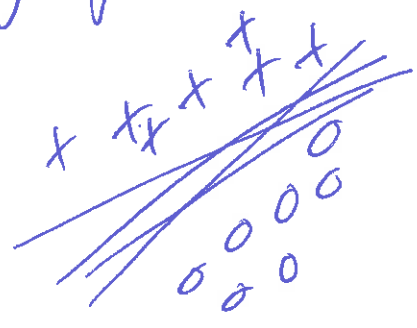
finding u s.t.

$$\begin{cases} \mathbb{R}^{n \times (d+1)} \\ \begin{pmatrix} -v_1 \\ \vdots \\ -v_n \end{pmatrix} u > (0, 0, \dots, 0) \end{cases}$$

If we require u is bounded (for example, $u_i \in [-1, 1], i=1, \dots, d+1$), this can be solved in a variety of ways; take a course on optimization for details. MATLAB has "linprog" which is a useful blackbox.

Remark: Except in very pathological cases in which some points are "on" the plane, if ~~the~~ a separating hyperplane is slightly perturbed, it remains

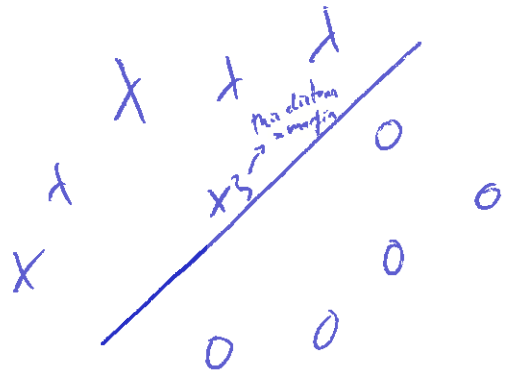
a separating hyperplane:



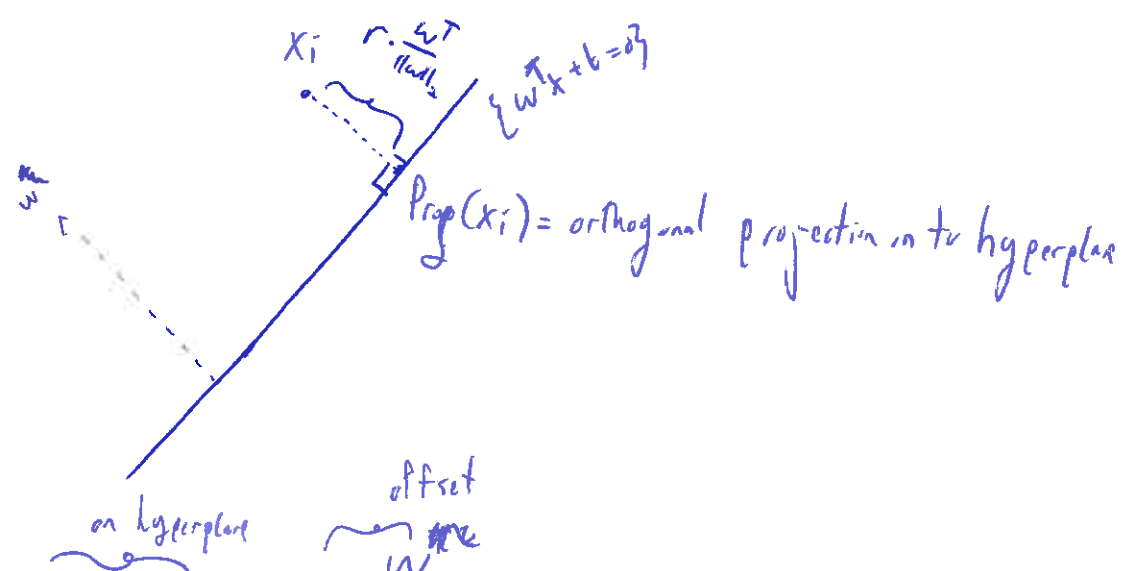
Which is "best"?

• Need an additional criterion to select "best" hyperplane. Pick the one that keeps points farthest away from the plane \rightarrow "maximize the margin."

• The margin of a hyperplane $\{w^T x + b = 0\}$ separating points $\{x_i\}_{i=1}^n \subset \mathbb{R}^d$ is the minimum distance between the hyperplane and any points:



Q: How to compute the margin? Well, how to compute the distance between a point and a hyperplane?



Set $X_i = \underbrace{\text{Proj}(x_i)}_{\text{on hyperplane}} + r \cdot \underbrace{\frac{w}{\|w\|_2}}_{\text{offset}}$; r is the "signed distance" of X_i to the hyperplane.

Since $\text{Proj}(x_i)$ is ~~orthogonal~~ on the hyperplane, we can compute:

$$\begin{aligned}
 w^T X_i + b &= w^T \left(\text{Proj}(x_i) + r \frac{w}{\|w\|_2} \right) + b \\
 &= \underbrace{w^T \text{Proj}(x_i) + b}_0 + r \frac{w^T w}{\|w\|} \\
 &= r \cdot \|w\|_2
 \end{aligned}$$

$$\Rightarrow r = \frac{w^T X_i + b}{\|w\|_2}$$

Note that if $w \mapsto \epsilon w$
 $b \mapsto \epsilon b,$

r stays the same. This gives us a choice in how to represent the hyperplane.

Let $w^T x + b$ be a separating hyperplane. It is canonical if $|w^T x_i + b| = 1$ for all x_i ~~minimizing~~ ^{achieving} the margin; these $\{x_i\}$ are called support vectors.

So, the problem of finding a separating hyperplane that maximizes the margin is the problem of

minimize $\|w\|_2$ \leftarrow maximize margin \leftrightarrow minimize $\|w\|$

s.t. $y_i (w^T x_i + b) \geq 1 \quad i = 1, \dots, n$
 \uparrow separating

Next time: What to do in the (common) case there are no separating hyperplanes?