

Week 03 - Part 01

Before we start, some announcements:

1 - TA OH are added to Course Page

There's at least one OH per day.

2 - WS1 posted. Solve the questions before the tutorial session. TAs will take your questions and review WS1.

3 - A1 posted: Pocket Alg. & linear Regression
I created a Colab notebook so that you can easily do your work there.

Please download the latest version of notebook,

Visit "Assignments/A1" to get the link.

4. A1 is to be done in groups of two.
you can join a different group for each assignments

5. Are there any hidden test cases for A1?
- No. the marking scheme is provided in
the handout.

Week03 - Part 01

So far: Linear classification & linear Regression

Given : $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$

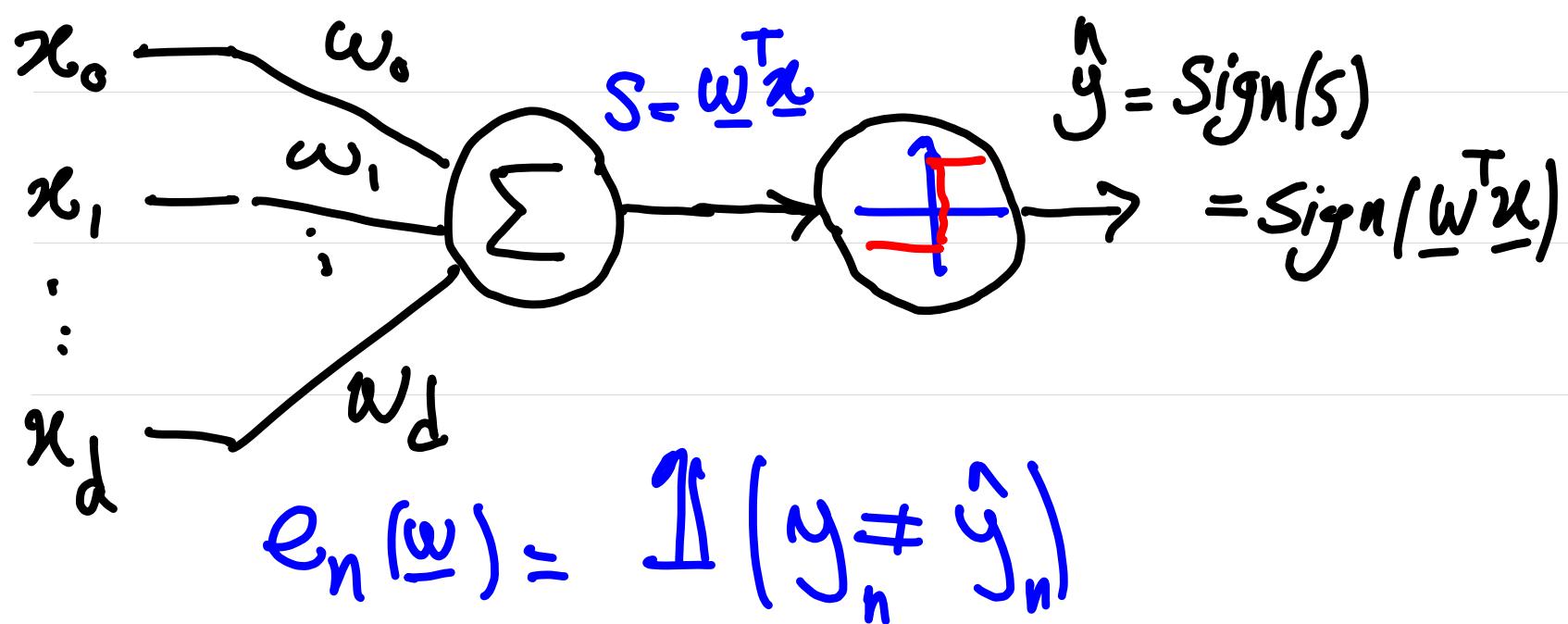
Unknown Target Function : $y = f(x)$

Hypothesis Set : $\hat{y} = h(x)$ where $h \in \mathcal{H}$

Linear Classification

$$x \in \mathbb{R}^{d+1}, y \in \{-1, +1\}, \hat{y} = \text{Sign}(w^T x)$$

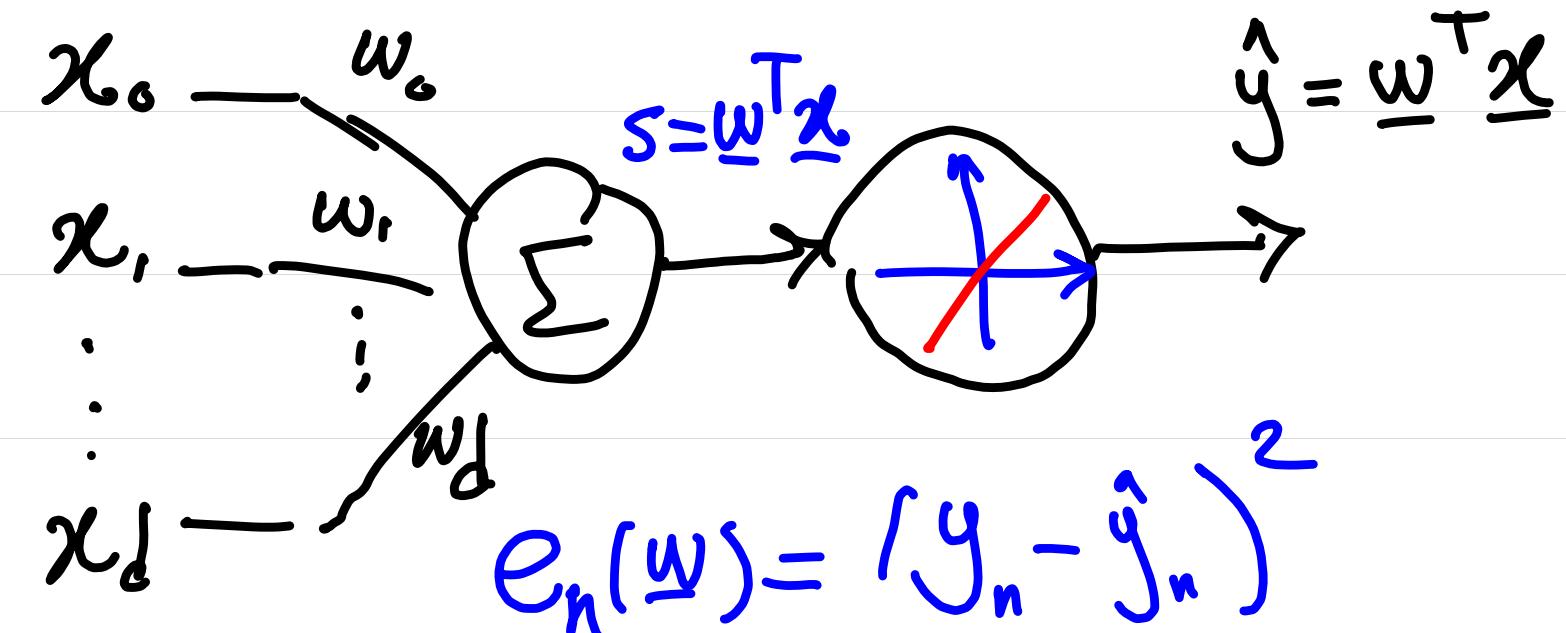
Illustrating as Neuron



Linear Regression

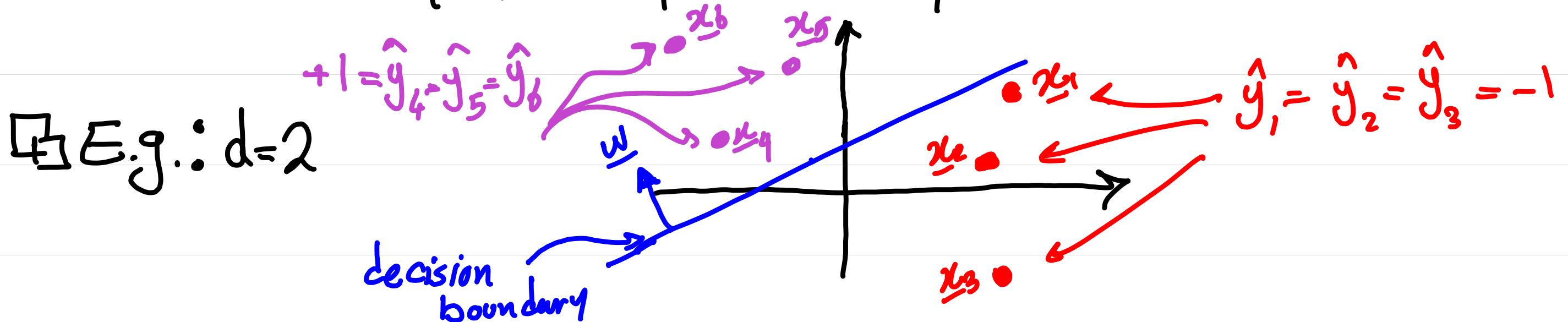
$$x \in \mathbb{R}^{d+1}, y \in \mathbb{R}, \hat{y} = w^T x$$

Illustrating as Neuron



■ So far, We have studied deterministic hypothesis.

□ Given an input, they tell you what is its label

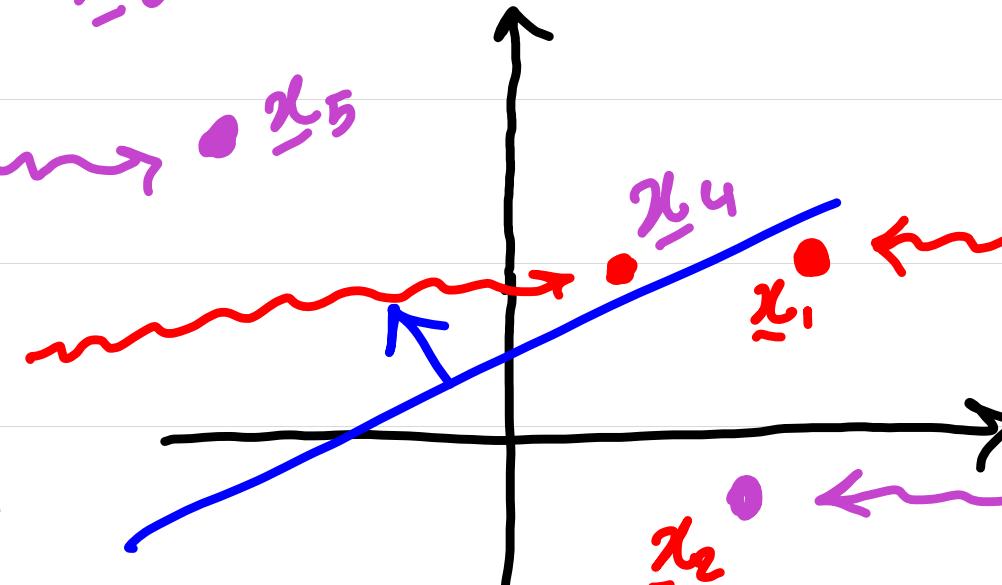


■ What if we want randomness in our prediction?

95% sure + x_6

80% sure it is + x_5

Very uncertain,
52% sure it
is +1



x_3 95% sure it is -1

■ Why would Someone Need a Predictor with randomness?

■ E.g.: Suppose we want to predict the occurrence of heart attacks based on diet.

- \underline{x} = (average sugar in the diet, average fat in the diet)

- y = heart attack

- Given $D = \{(x_n, y_n)\}_{n=1}^N$,

- For a new input \underline{x} , Predict ...

- How likely it is for the new person to have a heart attack, i.e., $\hat{P}(+ | \underline{x})$

The likelihood of having a heart attack given that the person has \underline{x} as its feature vector.

$$\hat{P}(+ | \underline{x})$$

the event of having a heart attack

Conditioned on the observed input feature vector

Lecture Outline

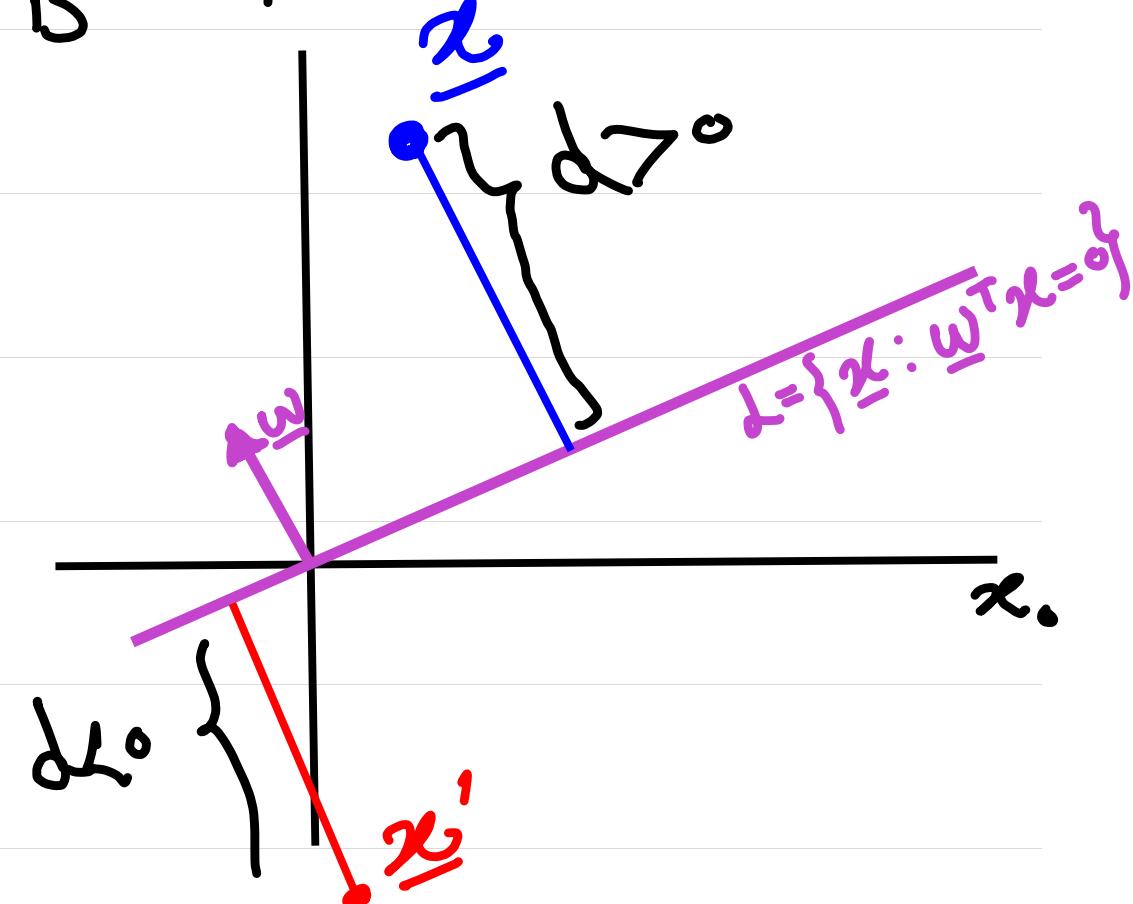
we want to
predict with randomness
 $\hat{P}_{\text{wo}}(y=+1 | \underline{x})$

what hypothesis Set
 \mathcal{H} should we use

Bingo! \mathcal{H} is

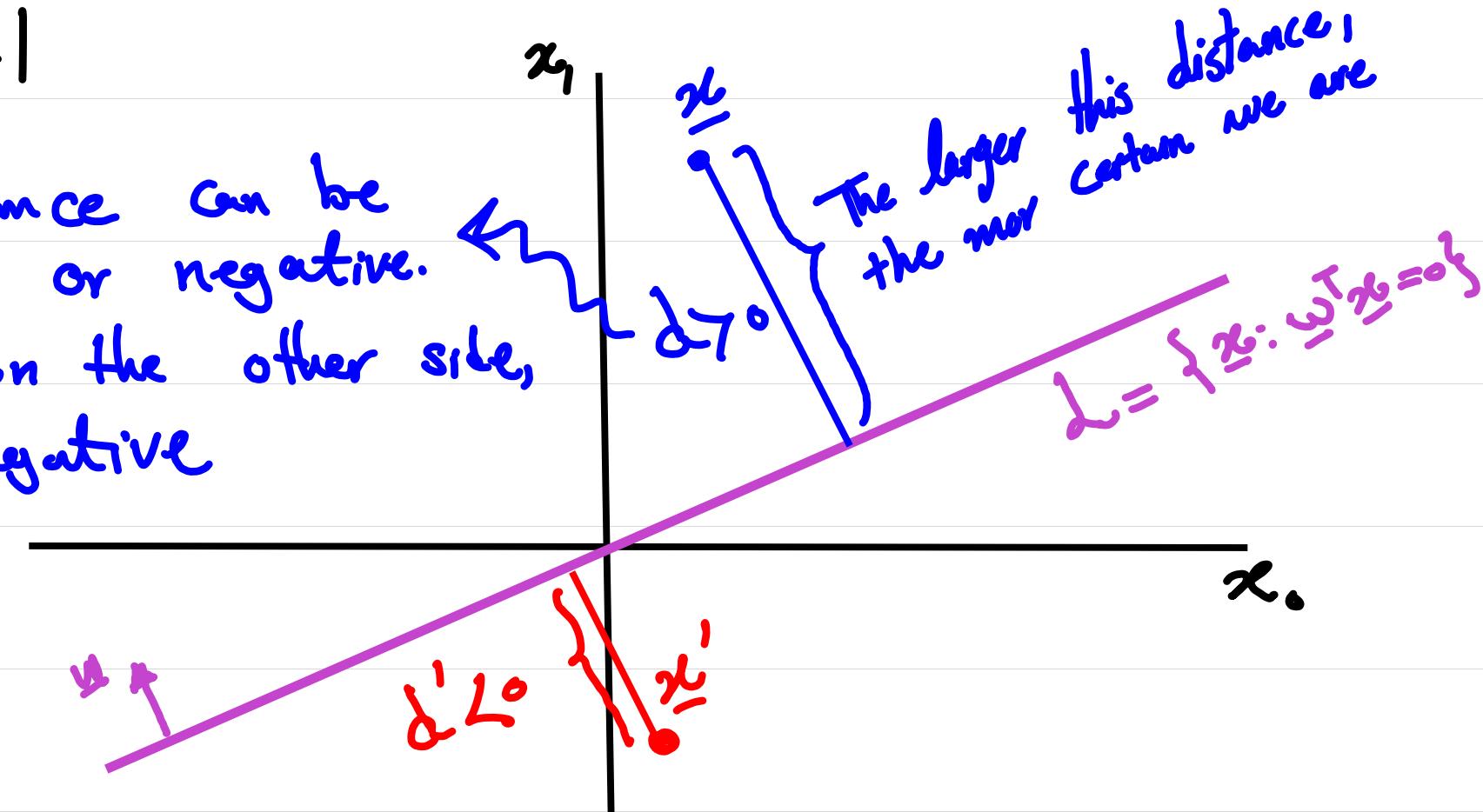
Today: Logistic Regression

- Predicts the probability of label: $\hat{P}_w(+1|\underline{x})$
- Let's see what could be a reasonable formulation for logistic regression model to output probability of label, $\hat{P}_w(+1|\underline{x})$
 - We are looking for a decision rule that its prediction depends on
 - which side of the hyperplane the datapoint is
 - How far away the datapoint is from the hyperplane.
 - The larger the distance b/w the datapoint and the hyperplane the more certain we should be



◻ E.g.: $d=1$

This distance can be positive or negative.
If it's on the other side, it is negative



◻ The distance b/w the new point \underline{x}_c and the decision boundary indicates the probability of our Prediction.

- As $d \rightarrow +\infty$, $\hat{P}_{w^T}(+1 | \underline{x}) \rightarrow 1$
- As $d \rightarrow -\infty$, $\hat{P}_{w^T}(+1 | \underline{x}) \rightarrow 0$
- When $d \approx 0$, $\hat{P}_{w^T}(+1 | \underline{x}) \approx 0.5$

Lecture Outline

We want to predict with randomness
 $\hat{P}_{\omega}(y=+1 | \underline{x})$

What hypothesis Set \mathcal{H} should we use

Let's use hyperplane $\underline{w}^T \underline{x} = 0$, and make prediction based on the distance b/w the point and the hyperplane

Bingo! \mathcal{H} is

what is the distance?

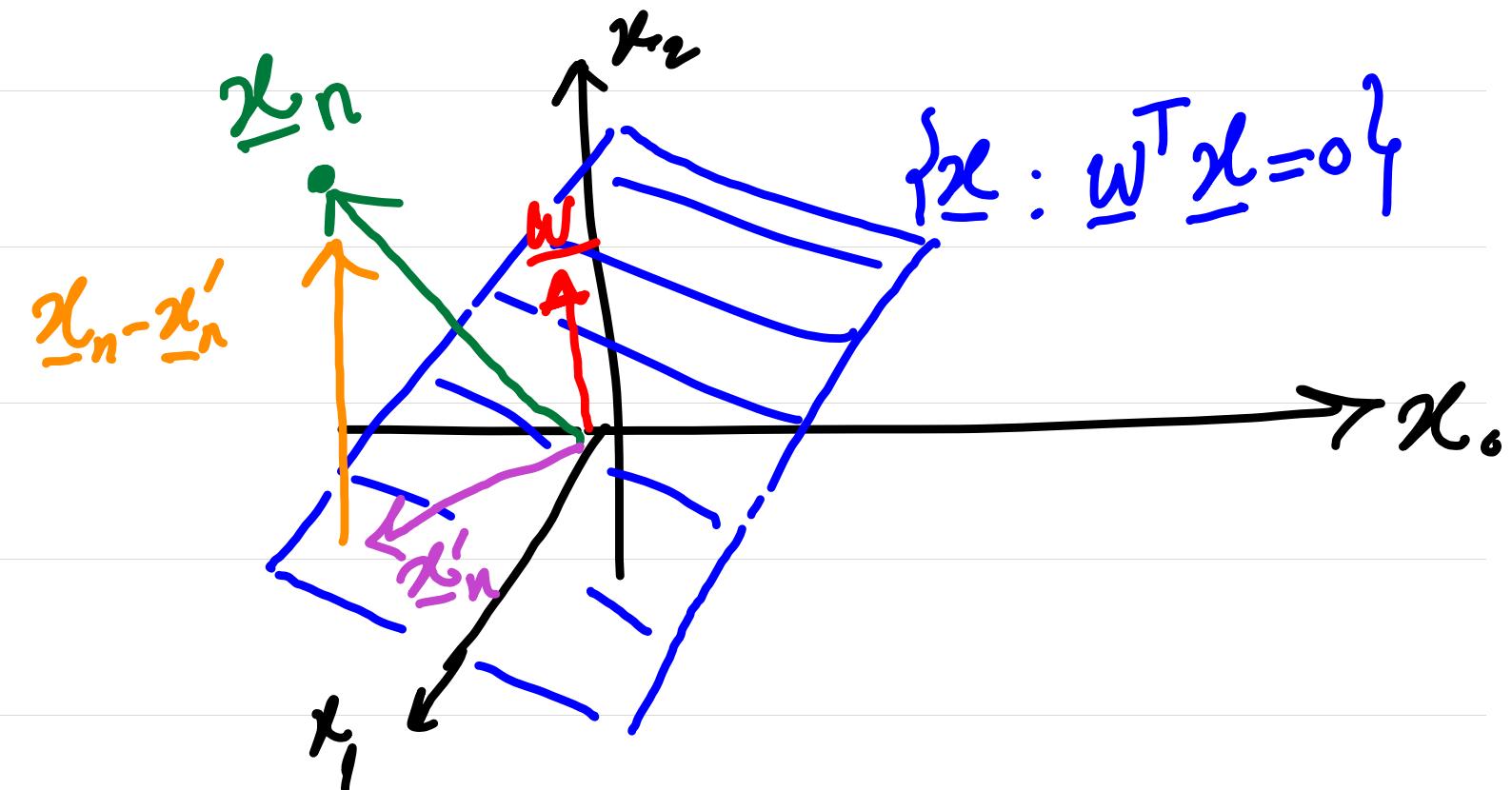
Let's find the distance b/w a point and a hyperplane

■ Assume $\{\underline{x} : \underline{w}^T \underline{x} = 0\}$ is the decision boundary (hyperplane.)

□ WLG, assume that $\|\underline{w}\|_2 = 1$.

□ Recall: \underline{w} is orthogonal

to the hyperplane $\{\underline{x} : \underline{w}^T \underline{x} = 0\}$



■ To find the distance b/w \underline{x}_n and the hyperplane, ...

□ Find the closest point on the hyperplane to \underline{x}_n , i.e.

project \underline{x}_n onto the hyperplane

□ length ($\underline{x}_n - \underline{x}'_n$) is what we are looking for.

■ Note that $\underline{x}_n - \underline{x}'_n$ is orthogonal to the hyperplane.

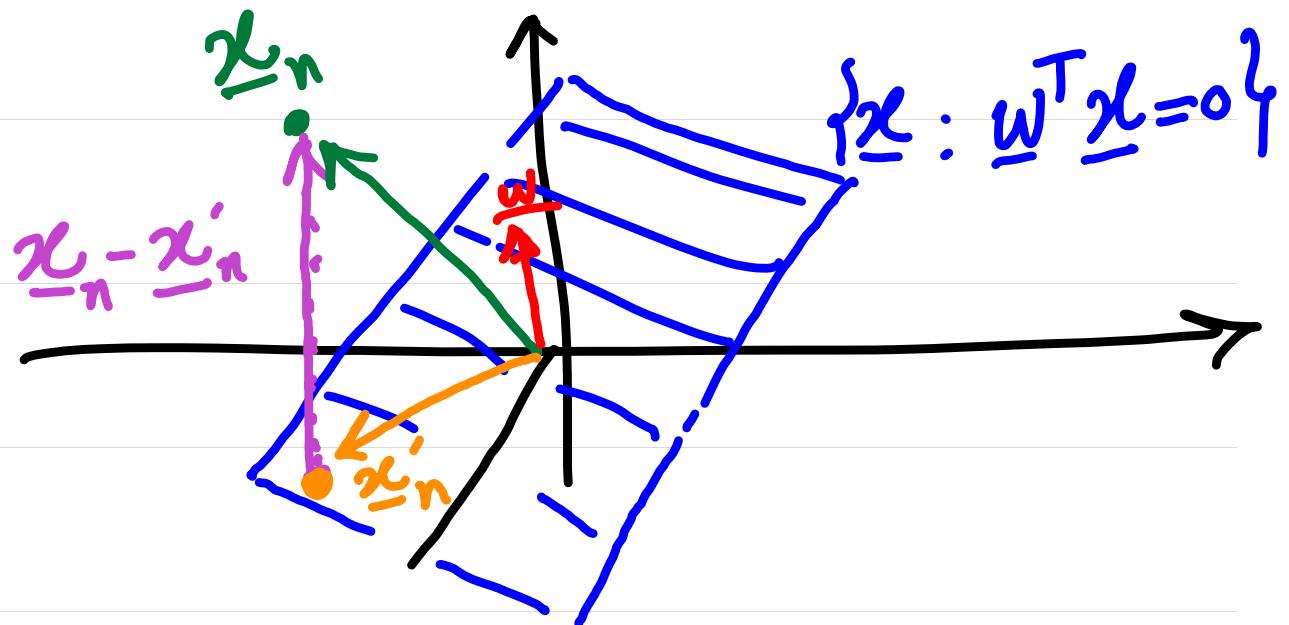
■ Thus, $\underline{x}_n - \underline{x}'_n = d \underline{w}$

■ We must find d .

$$\text{From } \underline{x}_n - \underline{x}'_n = d \underline{w} \Rightarrow \underline{x}'_n = \underline{x}_n - d \underline{w} \quad ①$$

Since, \underline{x}'_n is on the hyperplane $\underline{w}^T \underline{x}'_n = 0$

$$\Rightarrow \underline{w}^T \underline{x}_n - d \underline{w}^T \underline{w} = 0 \Rightarrow \underline{w}^T \underline{x}_n = d \|\underline{w}\|_2^2 = 1.$$



Lecture Outline

we want to predict with randomness
 $\hat{P}_{\omega}(y=+1 | \underline{x})$

what hypothesis set \mathcal{H} should we use

Let's use hyperplane $\underline{w}^T \underline{x} = 0$, and make prediction based on the distance b/w the point and the hyperplane

How do you want to model your predicted likelihood based on d ?

what is the distance?
Assuming $\|\underline{w}\|=1$,
 $d = \underline{w}^T \underline{x}$

Bingo! \mathcal{H} is

■ So, assuming $\|\underline{w}\|=1$, the distance b/w \underline{x}_n and the hyperplane is $d = \underline{w}^T \underline{x}_n$

■ Nice! we now know how to formulate the distance.

■ Now, let's see how to formulate our decision rule $\hat{P}_{\underline{w}}(+1 | \underline{x}_n)$ based on this distance.

■ We should formulate $\hat{P}_{\underline{w}}(+1 | \underline{x}_n)$ in a way that

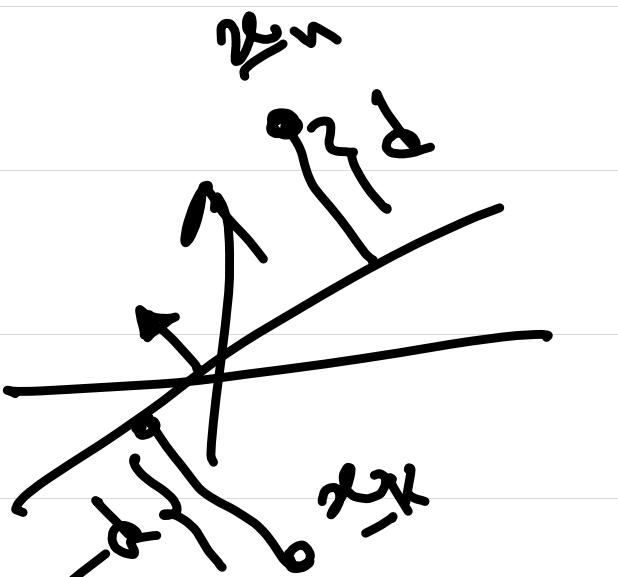
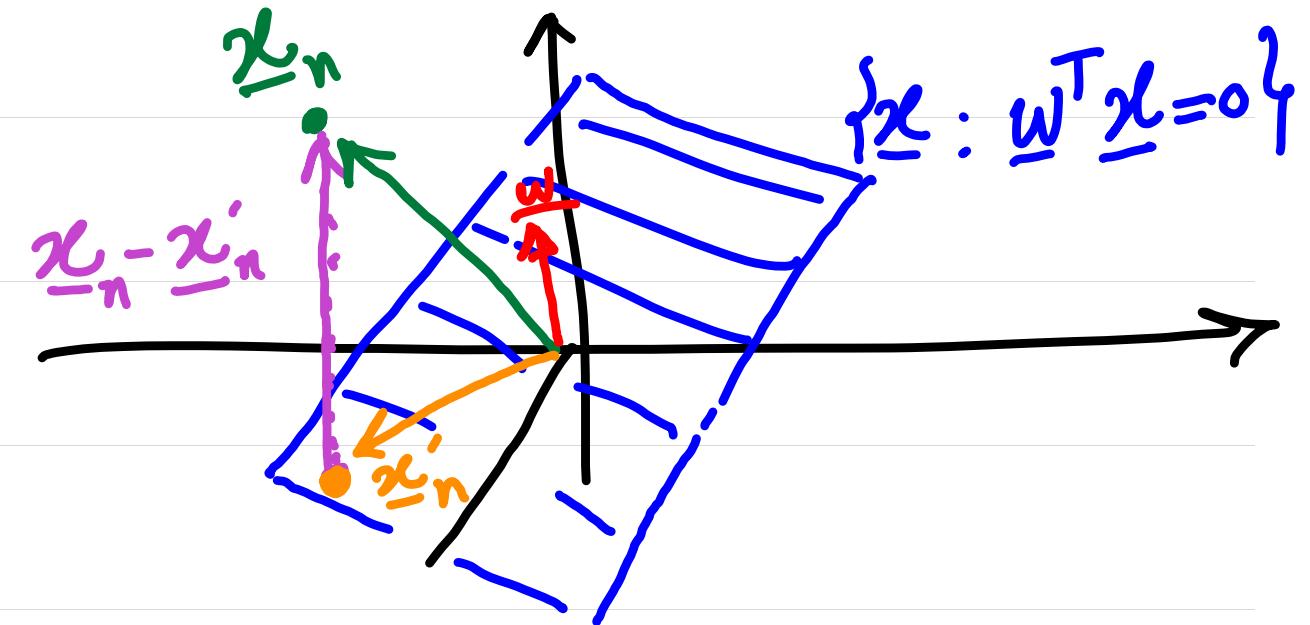
□ When $d = \underline{w}^T \underline{x}_n \rightarrow +\infty$, $\hat{P}_{\underline{w}}(+1 | \underline{x}_n) \rightarrow 1$

□ When $d = \underline{w}^T \underline{x}_n \rightarrow -\infty$, $\hat{P}_{\underline{w}}(+1 | \underline{x}_n) \rightarrow 0$

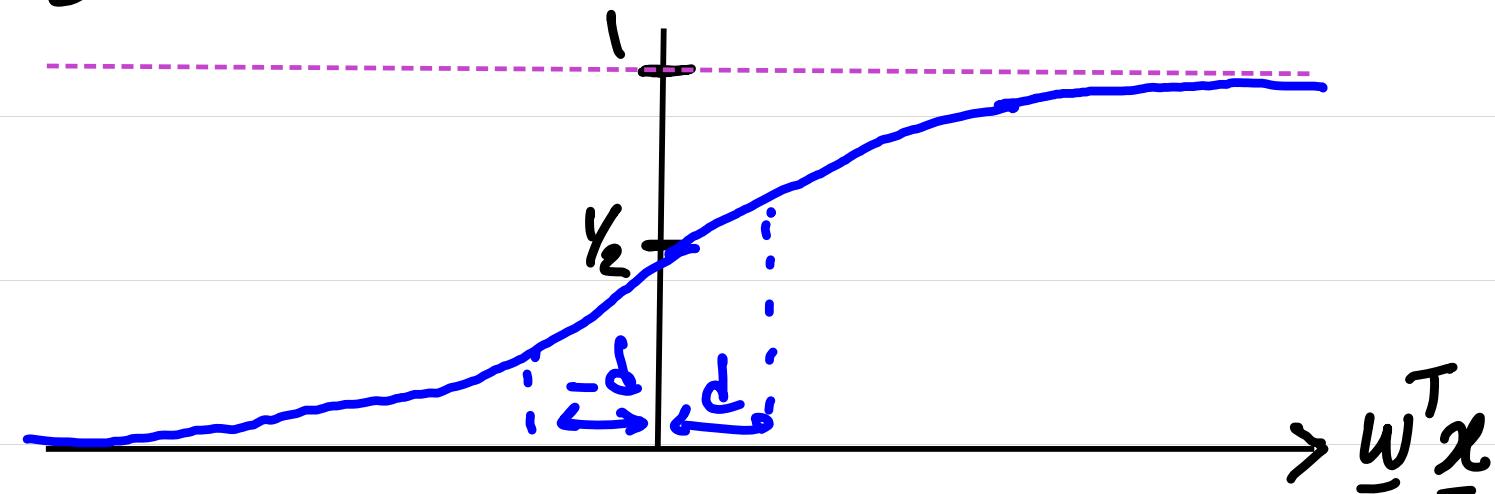
□ When $d = \underline{w}^T \underline{x}_n \rightarrow 0$, $\hat{P}_{\underline{w}}(+1 | \underline{x}_n) \rightarrow 0.5$

□ We want symmetry, i.e., if $d = \underline{w}^T \underline{x}_n$ and $-d = \underline{w}^T \underline{x}_K$, then

$$\hat{P}_{\underline{w}}(+1 | \underline{x}_n) = \hat{P}_{\underline{w}}(-1 | \underline{x}_K) = 1 - \hat{P}_{\underline{w}}(+1 | \underline{x}_K), \text{ if } \underline{w}^T \underline{x}_n = -\underline{w}^T \underline{x}_K$$



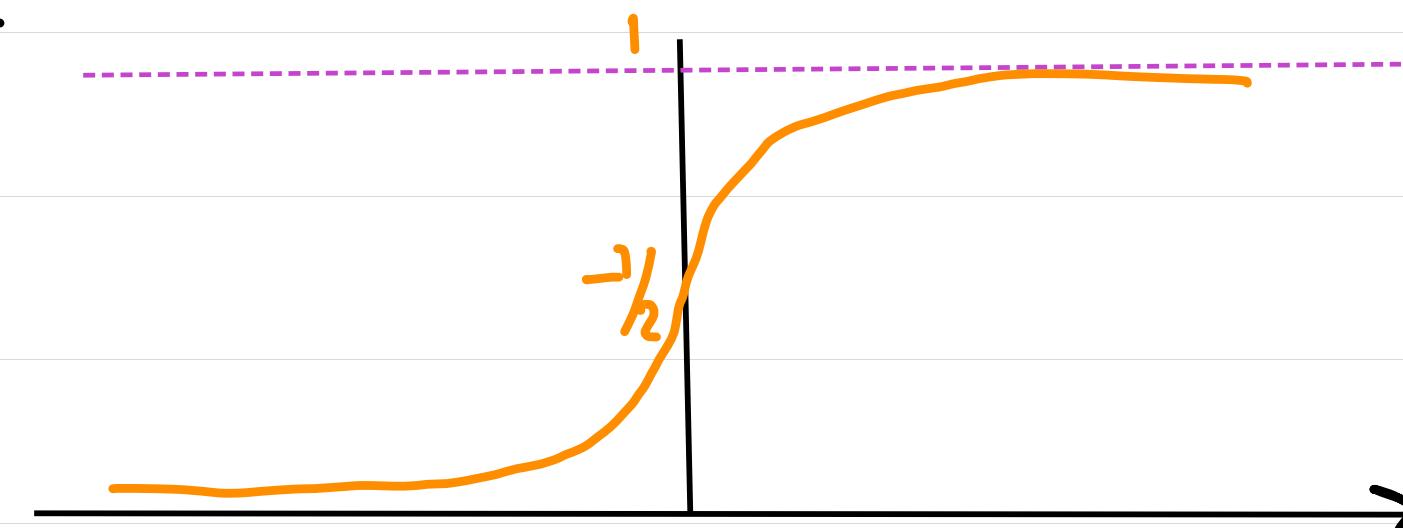
■ So, we are looking for a function like this



■ There is a nice function that we can use to model such behaviour.

□ The "logistic" function

$$\theta(s) = \frac{1}{1 + e^{-s}}$$



$$\hat{P}_{\underline{w}}(1|\underline{x}) = \theta(s) = \frac{1}{1 + e^{-s}} = \frac{1}{1 + e^{-w^T \underline{x}}}$$

$$s = w^T \underline{x}$$

Lecture Outline

We want to predict with randomness
 $\hat{P}_w(y=+1|\underline{x})$

- * $\underline{w}^T \underline{x} \rightarrow +\infty, \hat{P}_w(1|\underline{x}) \rightarrow 1$
- * $\underline{w}^T \underline{x} \rightarrow -\infty, \hat{P}_w(1|\underline{x}) \rightarrow 0$
- * Symmetry

What hypothesis set \mathcal{H} should we use

How do you want to model your predicted likelihood based on d ?

logistic function is a suitable choice
 $\delta(s) = \frac{1}{1+e^{-s}}$

Bingo! \mathcal{H} is

$$\hat{P}_w(1|\underline{x}) = \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$$

What is the distance?
Assuming $\|\underline{w}\|=1$,
 $d = \underline{w}^T \underline{x}$

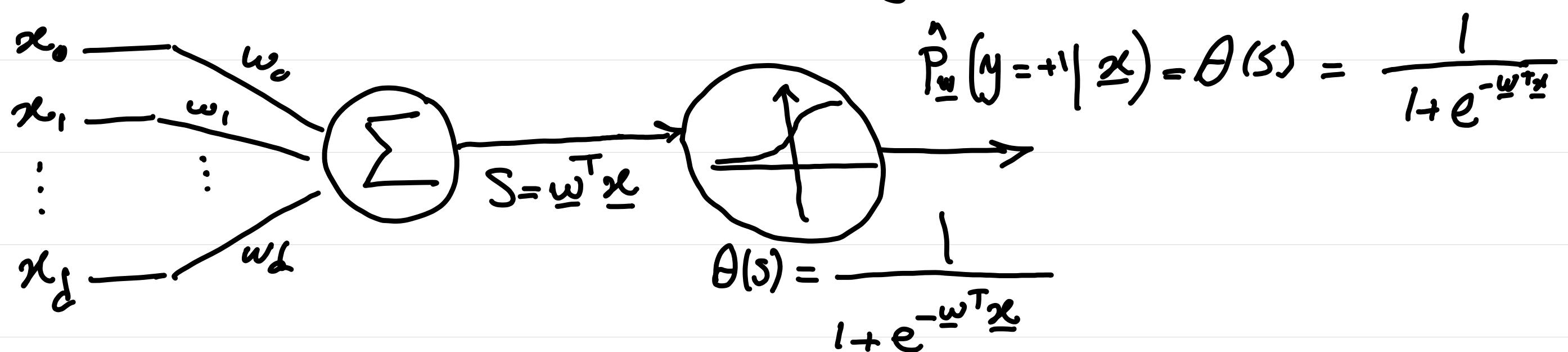
Let's use hyperplane $\underline{w}^T \underline{x} = 0$, and make prediction based on the distance b/w the point and the hyperplane

Illustrating Logistic Regression Model as Neuron

Given : $\mathcal{D} = \{(\underline{x}_n, y_n)\}_{n=1}^N$, $\underline{x}_n \in \mathbb{R}^{d+1}$, $y_n \in \{+1, -1\}$

Unknown Target Function : $y = f(\underline{x})$

Hypothesis Set : $\hat{y} = \theta(\underline{w}^T \underline{x}) = \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$



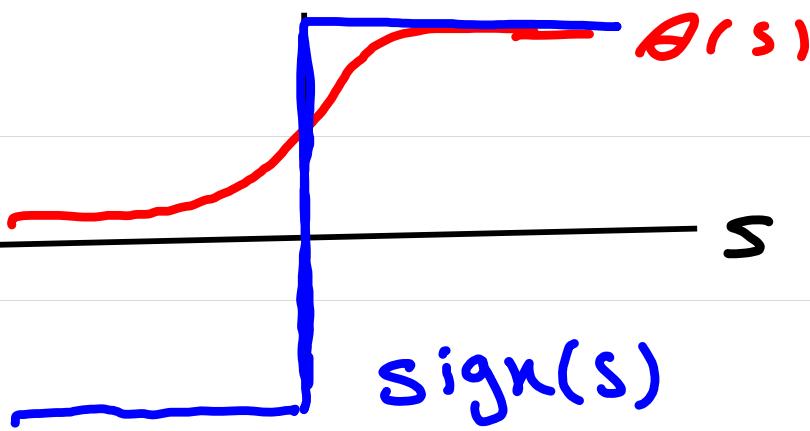
$\theta(s)$ is called the "logistic function."

Logistic function is an example of "Sigmoid functions"

↳ "Sigmoid" functions is the class of functions that are "S" shaped.

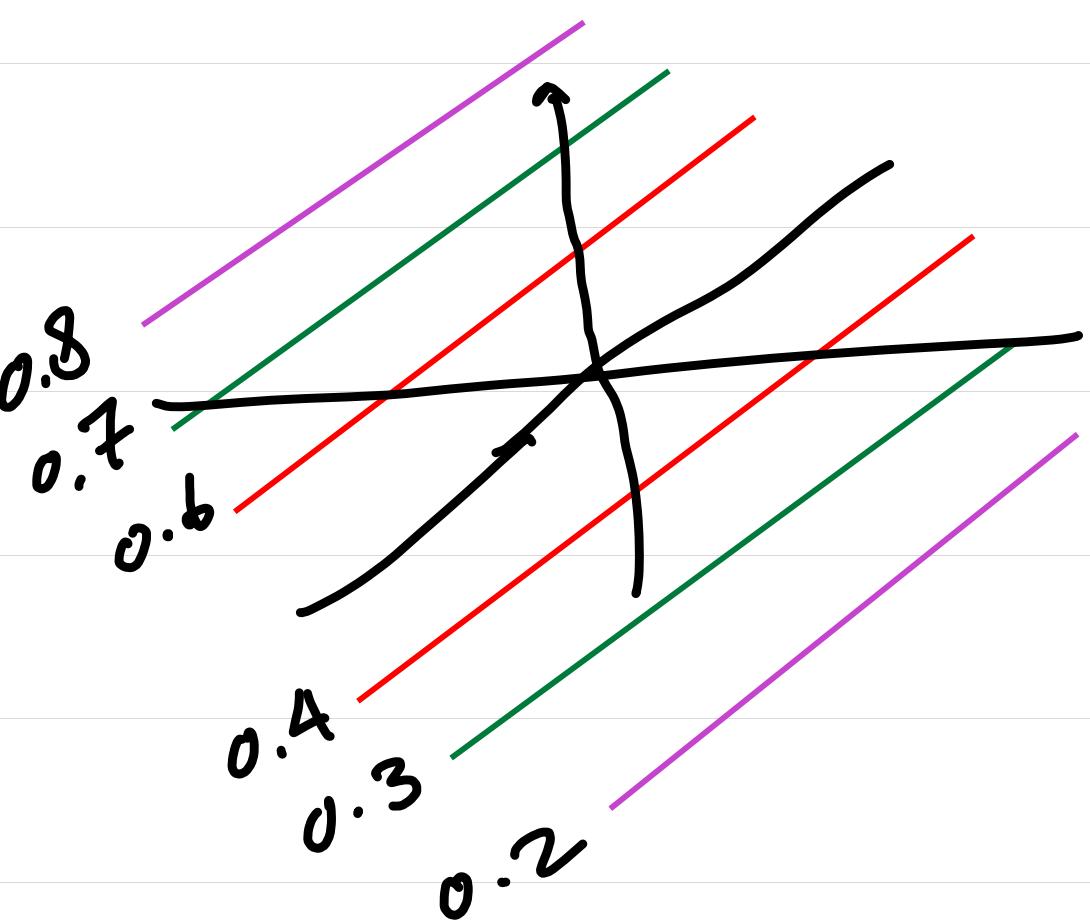
Remarks:

- ① This is a soft threshold,
in contrast to the hard threshold used by
linear classification



② $\hat{P}_{\underline{w}}(-1 | \underline{x}) = 1 - \hat{P}_{\underline{w}}(+1 | \underline{x}) = 1 - \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$

$$= \frac{1}{1 + e^{\underline{w}^T \underline{x}}} = \sigma(-\underline{w}^T \underline{x})$$



③ Notation:

$$\hat{P}_{\underline{w}}(1 | \underline{x}) = \theta(\underline{w}^\top \underline{x}) \quad \text{estimate} \quad P[y=+1 | \underline{x}]$$

$$\hat{P}_{\underline{w}}(-1 | \underline{x}) = \theta(-\underline{w}^\top \underline{x})$$

■ We can combine them:

$$\hat{P}_{\underline{w}}(y | \underline{x}) = \theta(y \underline{w}^\top \underline{x}) = \frac{1}{1 + e^{-y \underline{w}^\top \underline{x}}} , \text{ where } y \in \{+1, -1\}$$

Lecture Outline

We want to predict
With randomness



Let's use \mathcal{H} as

$$\hat{P}_{\underline{w}}(y|\underline{x}) = \frac{1}{1 + e^{-y\underline{w}^T \underline{x}}}$$



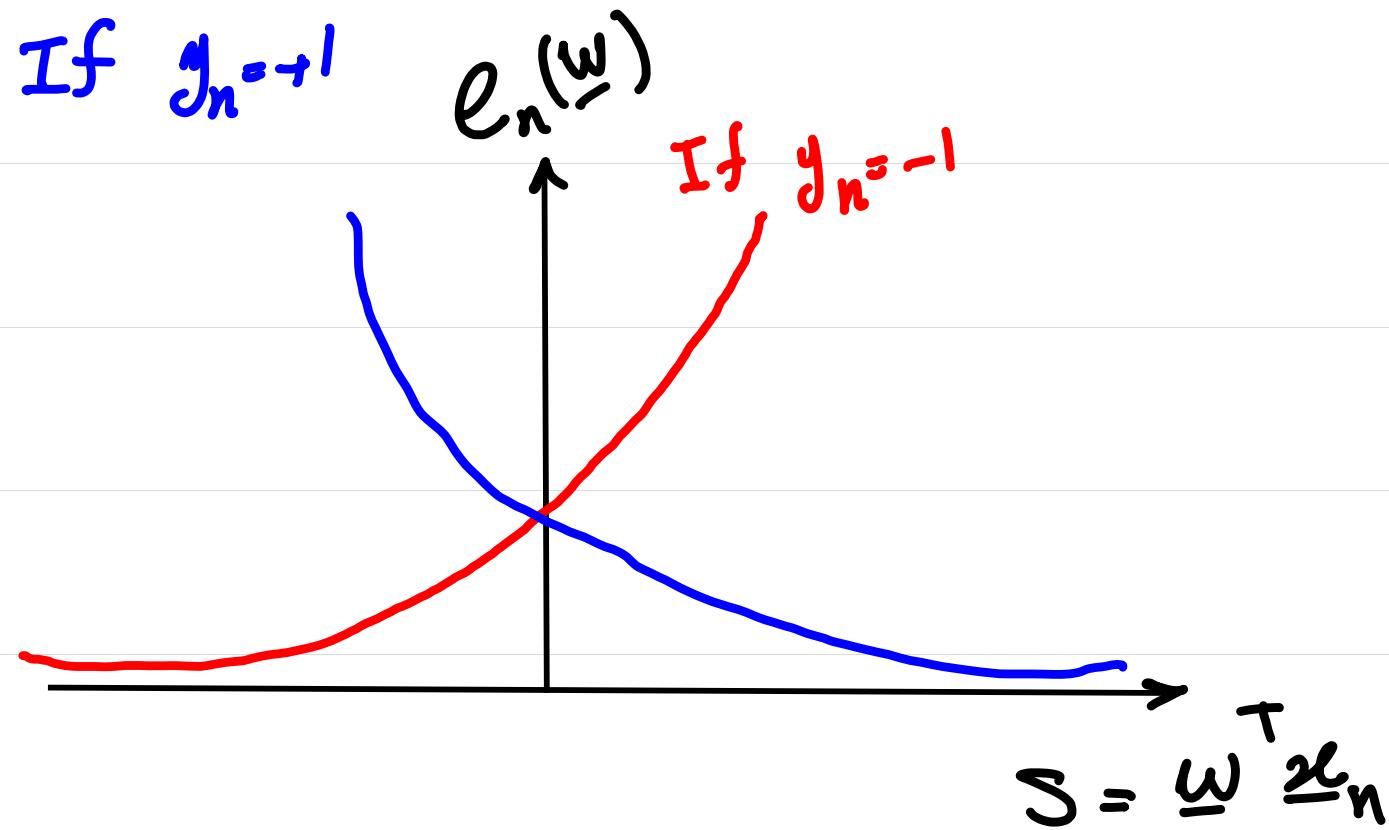
Now, we need to come up with
a meaningful $E_{in}(\underline{w})$ to
minimize

What is the Error Criterion?

For the n_{th} example in train data,
natural log

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = \log \left(1 + e^{-y_n \underline{w}^T \underline{x}_n} \right)$$

"log-loss" function



Note:

- If $\underline{w}^T \underline{x}_n \gg 0$ and $y_n = +1 \rightarrow \text{loss} = 0$
- If $\underline{w}^T \underline{x}_n \ll 0$ and $y_n = -1 \rightarrow \text{loss} = 0$

Lecture Outline

We want to predict
With randomness

Let's use \mathcal{H} as
 $\hat{P}_{\underline{w}}(y|\underline{x}) = \frac{1}{1 + e^{-y\underline{w}^T \underline{x}}}$

We use log-loss for
our error criterion

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n)$$

Now, we need to come up with
a meaningful $E_{in}(\underline{w})$ to
minimize

Why is
log-loss
a reasonable
choice?

E.g.: Suppose logistic regression outputed a \underline{w} such that

$$\hat{P}_{\underline{w}}(+1 | \underline{x}_n) = 0.999$$

$$(\Rightarrow \hat{P}_{\underline{w}}(-1 | \underline{x}_n) = 0.001)$$

The "machine" is very confident that $\hat{y}_n = +1$.

◻ If $y_n = +1$, i.e., the true label was $+1$:

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = -\log(0.999) \approx 10^{-4}$$

◻ what if $y_n = -1$. What do you expect to see for $e_n(\underline{w})$?

• The "machine" is confident in its incorrect prediction.

Large $e_n(\underline{w})$

Small $e_n(\underline{w})$

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = -\log(0.001) \approx 3$$

E.g.: Suppose $\hat{P}_{\underline{w}}(+1 | \underline{x}_n) = 0.2$ ($\hat{P}_{\underline{w}}(-1 | \underline{x}_n) = 0.8$)

◻ If $y_n = +1$, $\ell_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = -\log(0.2) \approx 1.39$

◻ If $y_n = -1$, what do we expect to see?

◻ larger $\ell_n(\underline{w})$ smaller $\ell_n(\underline{w})$

$$\ell_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n) = -\log(0.8) \approx 0.22$$

Lecture Outline

We want to predict with randomness

Let's use \mathcal{H} as

$$\hat{P}_{\underline{w}}(y|\underline{x}) = \frac{1}{1 + e^{-y \underline{w}^T \underline{x}}}$$

We use log-loss for our error criterion

$$e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n | \underline{x}_n)$$

Now, we need to come up with a meaningful $E_{in}(\underline{w})$ to minimize

Why is log-loss a reasonable choice?

Numerical E.g. shows that it makes sense

Seriously? Justifying based on Numerical example? We are better than this

Summary: Logistic Regression

- Given $\mathcal{D} = \{(\underline{x}_n, y_n)\}_{n=1}^N$, $\underline{x}_n \in \mathbb{R}^{d+1}$ $y_n \in \{+1, -1\}$,
- find $\underline{w} \in \mathbb{R}^{d+1}$,
- to minimize $E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w}) = \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \underline{w}^\top \underline{x}_n})$.
- For a new input \underline{x} , predict the probability of being in class y as $\hat{P}_{\underline{w}}(y | \underline{x}) = \frac{1}{1 + e^{-y \underline{w}^\top \underline{x}}}$