

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 assignment

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

Week 03 - Part 01

Review: Linear classification & linear Regression

■ Given: $\mathcal{D} = \{(\underline{x}_n, y_n)\}_{n=1}^N$

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

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

Linear classification

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

■ Illustrating as Neuron

Linear Regression

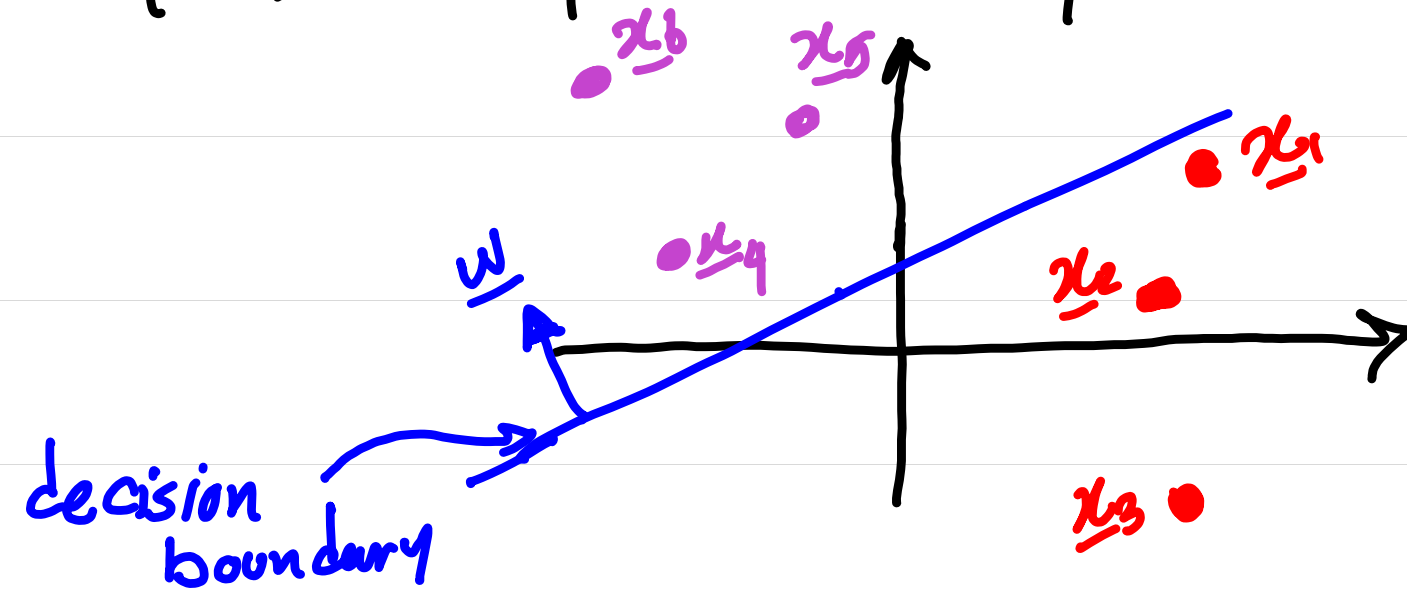
$$\underline{x} \in \mathbb{R}^{d+1}, y \in \mathbb{R}, \hat{y} = \underline{w}^T \underline{x}$$

■ Illustrating as Neuron

■ So far, We have studied deterministic hypothesis.

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

□ E.g.: $d=2$

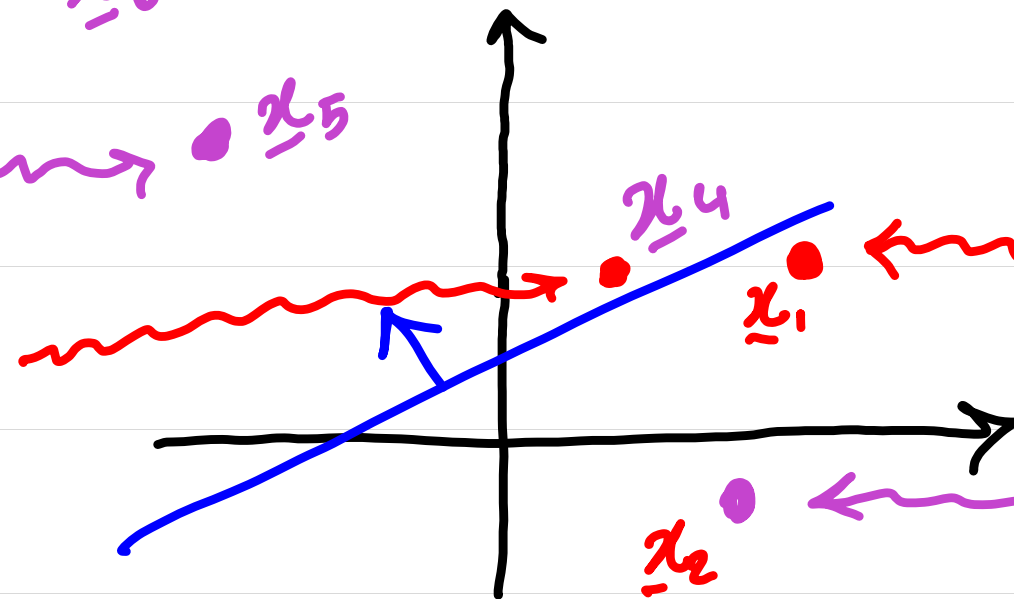


■ What if we want randomness in our prediction?

95% sure +1 \rightsquigarrow x_6

80% sure it is +1 \rightsquigarrow x_5

Very uncertain,
52% sure it
is +1



Very uncertain, 52% sure it is -1

80% sure it is -1

x_3 \rightsquigarrow 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)

- \underline{y} = heart attack

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

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

Lecture Outline

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

What hypothesis Set
 \mathcal{H} should we use

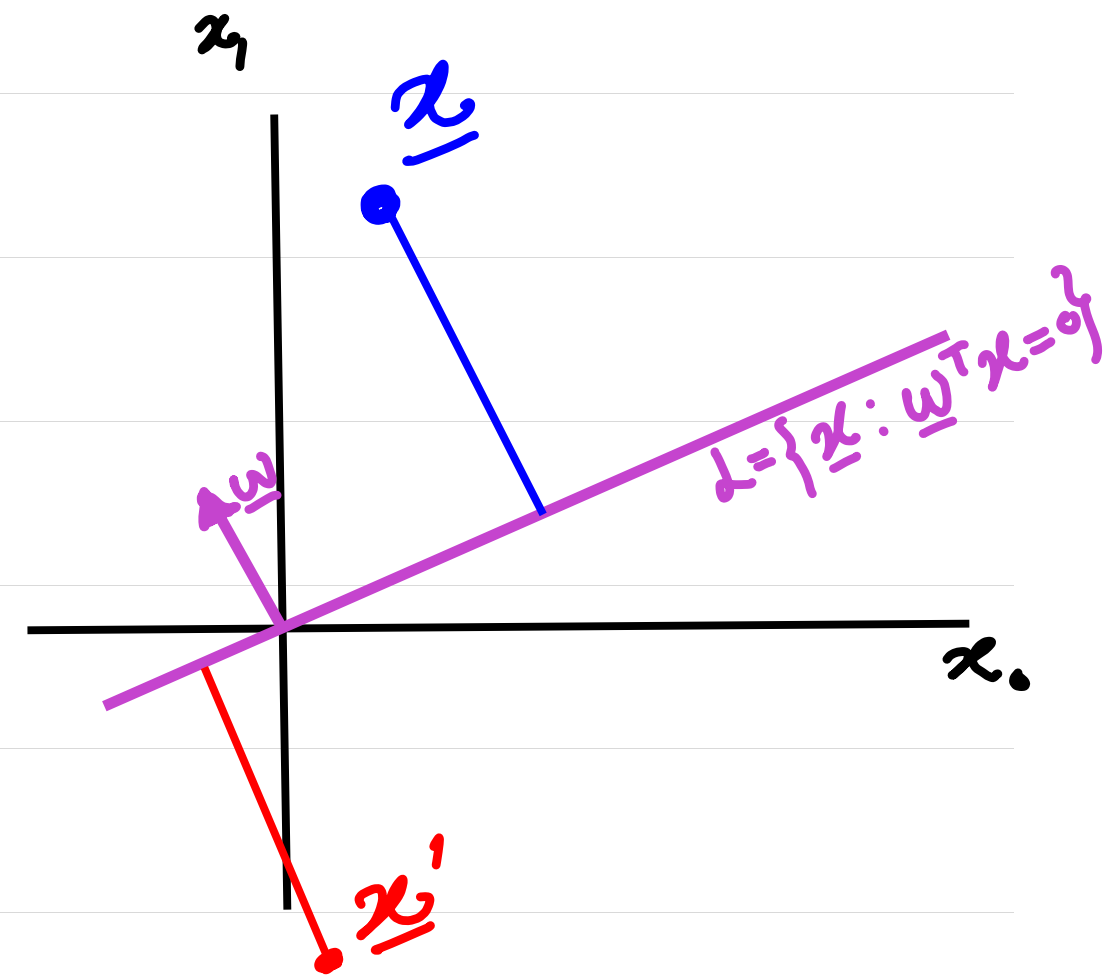
Bingo! \mathcal{H} is

Today: Logistic Regression

■ Predicts the probability of label: $\hat{P}_{\underline{w}}(+1|x)$

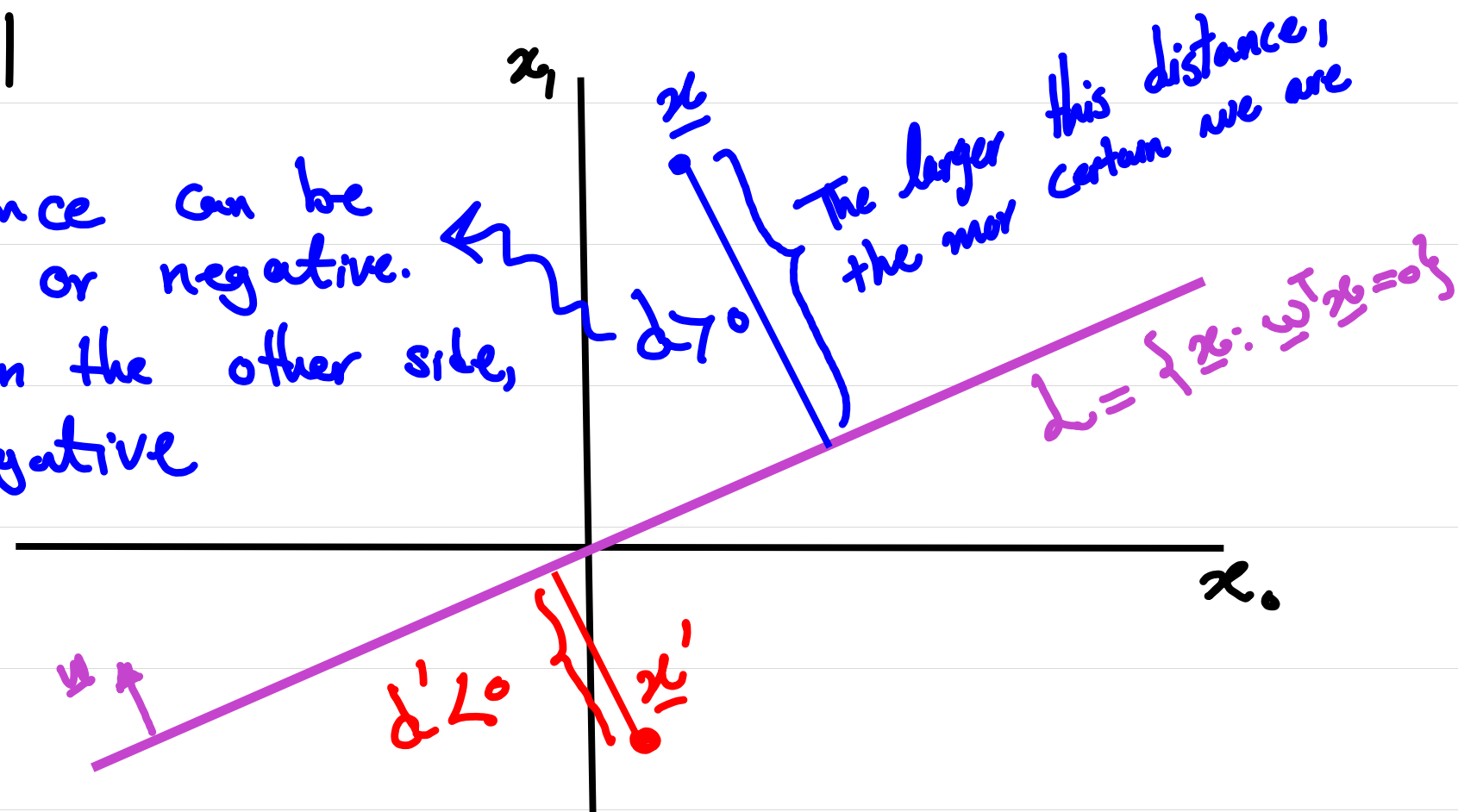
■ Let's see what could be a reasonable formulation for logistic regression model to output probability of label, $\hat{P}_{\underline{w}}(+1|x)$

□ We are looking for a decision rule that its prediction depends on



□ 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} and the decision boundary indicates the probability of our Prediction.

● As $d \rightarrow +\infty$, $\hat{P}_{\underline{w}}(+1 | \underline{x}) \rightarrow$

● As $d \rightarrow -\infty$, $\hat{P}_{\underline{w}}(+1 | \underline{x}) \rightarrow$

● When $d \approx 0$, $\hat{P}_{\underline{w}}(+1 | \underline{x}) \approx$

Lecture Outline

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

What hypothesis Set
 \mathcal{H} should we use

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

What is the distance?

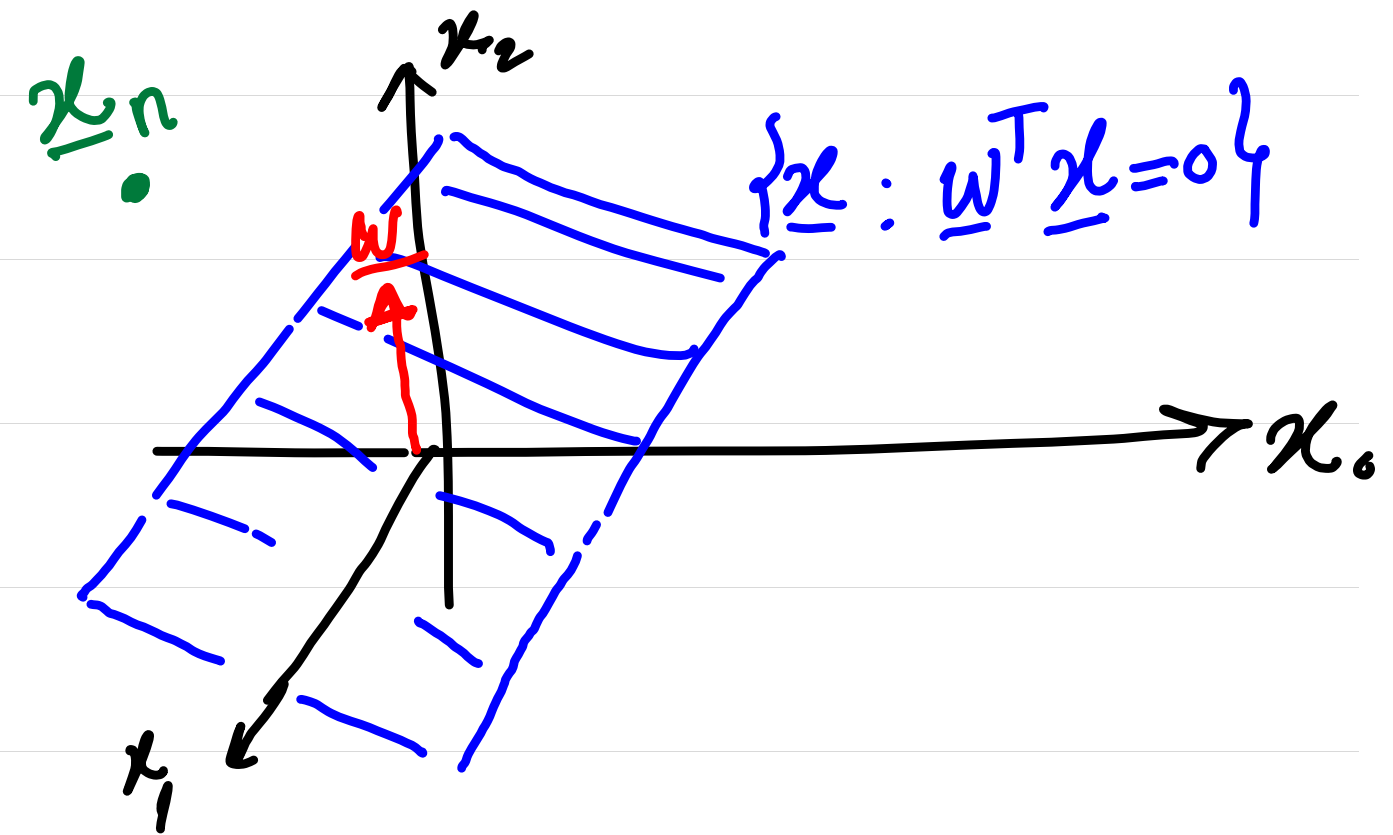
Bingo! \mathcal{H} is

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 $\|\underline{w}\|_2 = 1$.

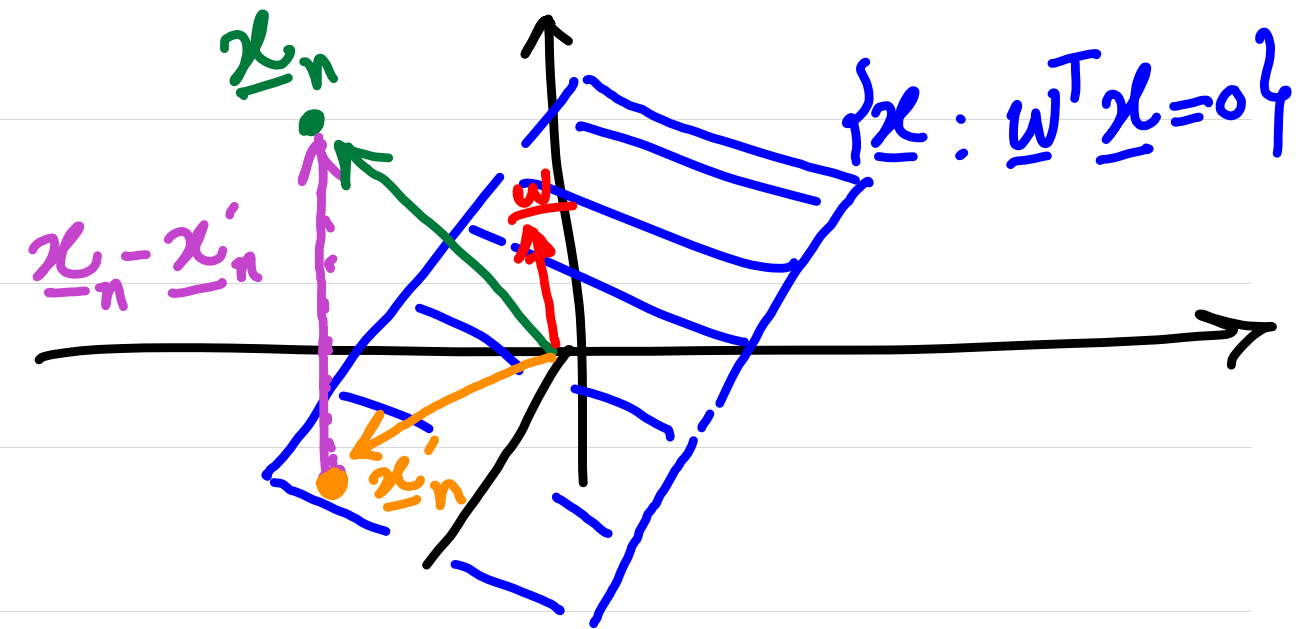
□ Recall: \underline{w} is orthogonal to the hyperplane $\{\underline{x} : \underline{w}^T \underline{x} = 0\}$



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

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

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



Lecture Outline

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

What hypothesis Set
 \mathcal{H} should we use

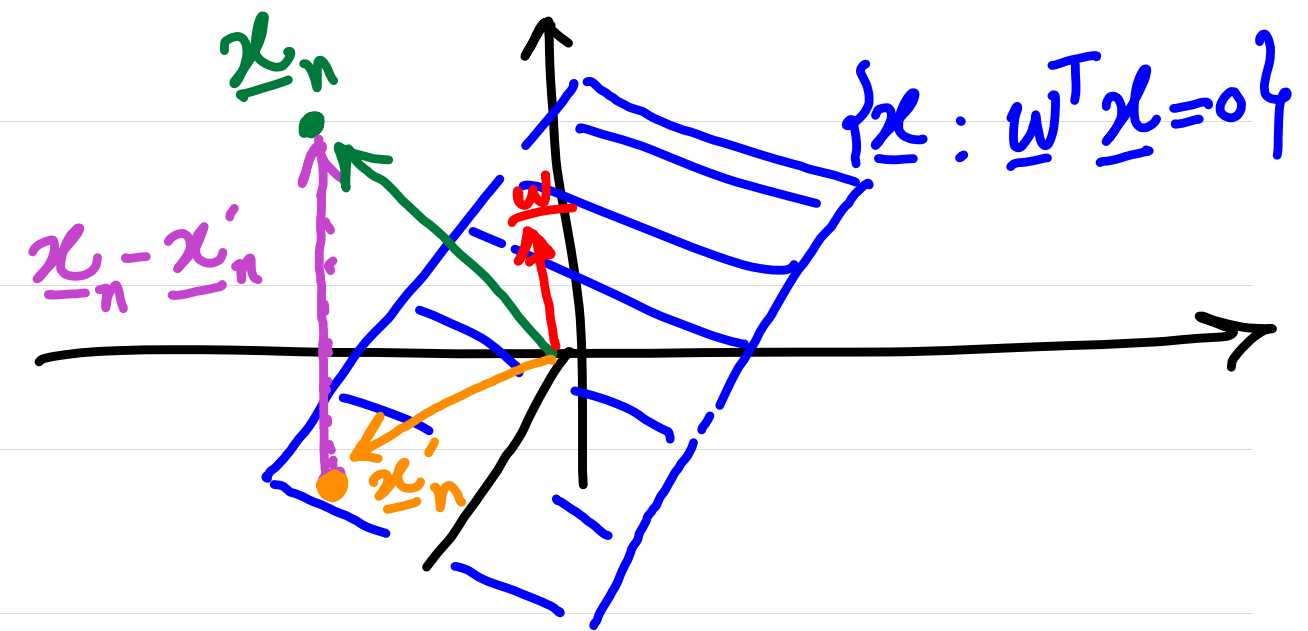
Let's use hyperplane
 $\underline{w}^T \underline{x} = 0$, and
make prediction
based on the
distance b/w
the point and
the hyperplane

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

Bingo! \mathcal{H} is

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

■ 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

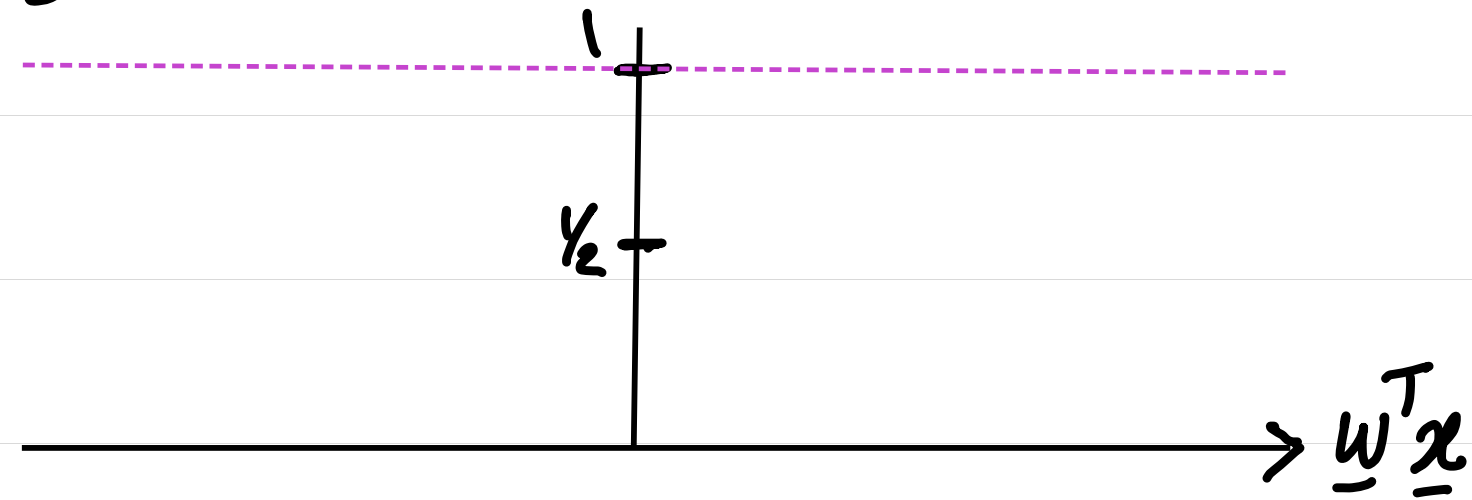
□ When $d = \underline{w}^T \underline{x}_n \rightarrow$, $\hat{P}_{\underline{w}}(+1|\underline{x}_n) \rightarrow$

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

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

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

■ So, we are looking for a function like this



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



Lecture Outline

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

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 ?

- $\underline{w}^T \underline{x} \rightarrow +\infty, \hat{P}_{\underline{w}}(1|\underline{x}) \rightarrow 1$
- $\underline{w}^T \underline{x} \rightarrow 0, \hat{P}_{\underline{w}}(1|\underline{x}) \rightarrow 1/2$
- $\underline{w}^T \underline{x} \rightarrow -\infty, \hat{P}_{\underline{w}}(1|\underline{x}) \rightarrow 0$
- Symmetry

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

Bingo! \mathcal{H} is

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

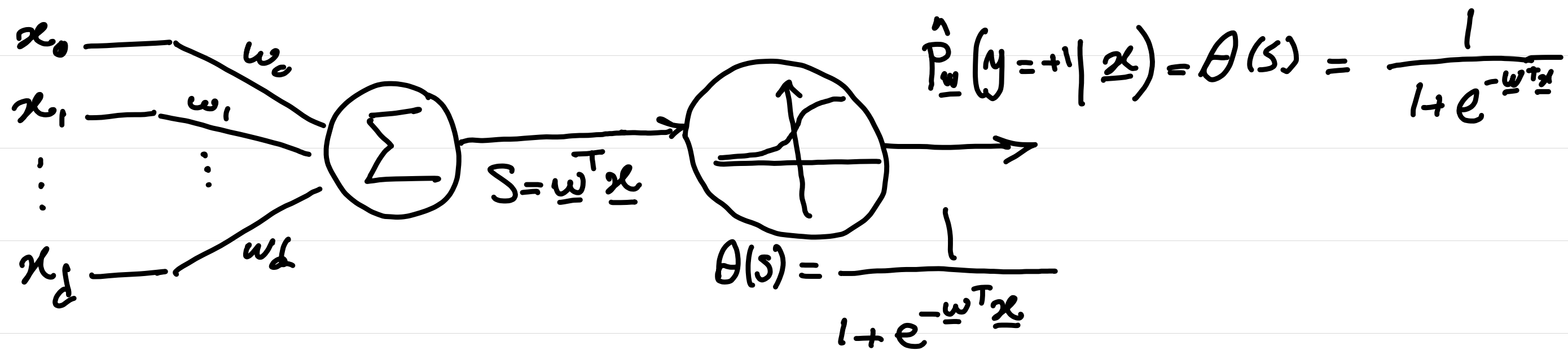
What is the distance?
Assuming $\|\underline{w}\|=1$,
 $d = \underline{w}^T \underline{x}$

Illustrating Logistic Regression Model as Neuron

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

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

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



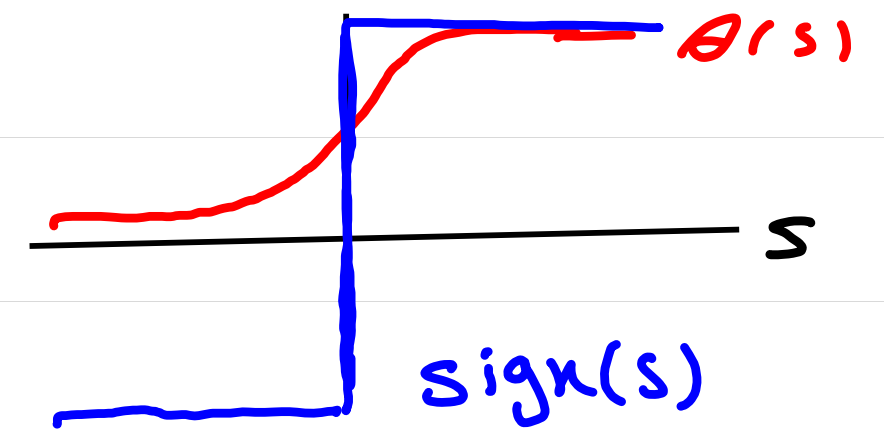
■ $\theta(s)$ is called the "logistic function."

■ Logistic function is an example of "Sigmoid functions"

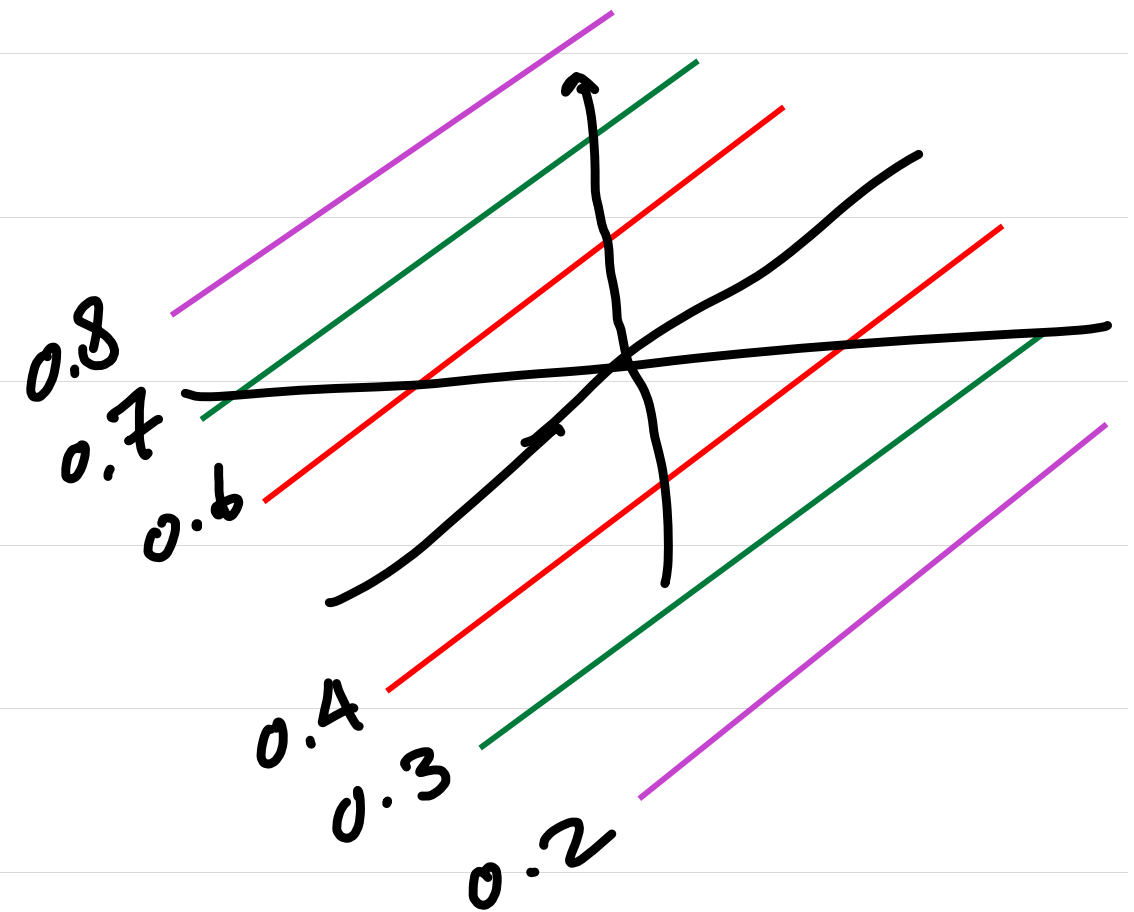
□ "Sigmoid" functions is the class of functions that are "S" shaped.

Remarks:

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



② $\hat{P}_w(y=-1 | \underline{x}) =$



③ Notation:

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

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

■ We can combine them:

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

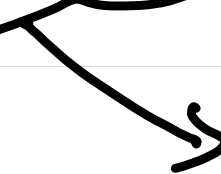
Lecture Outline

We want to predict
with randomness



Let's use \mathcal{H} as

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

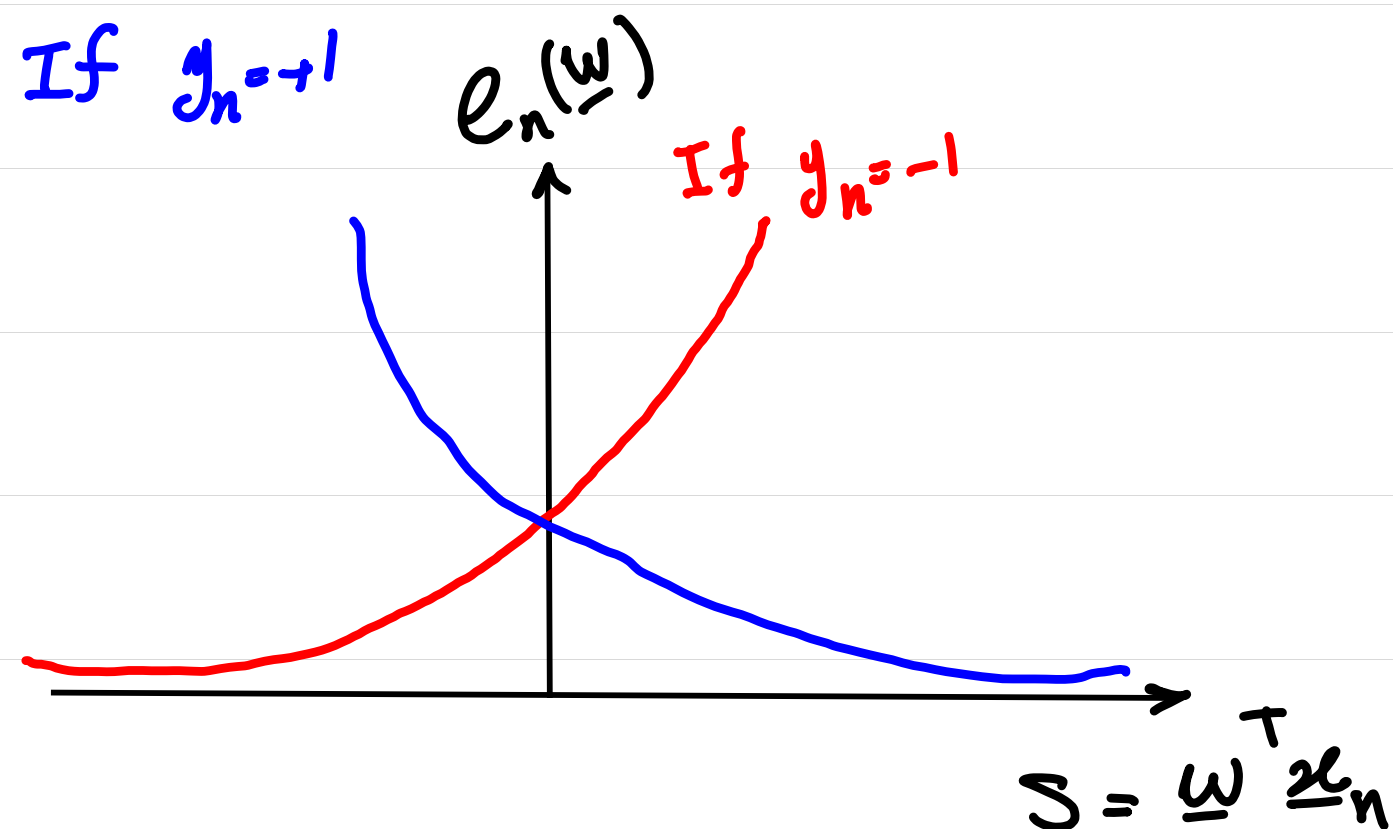


Now, we need to come up with
a meaningful $E_{in}(\underline{w})$ to
minimize

What is the Error Criterion?

For the n_{th} example in train data,

$$e_n(\underline{w}) =$$



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

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

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

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

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

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

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

□ Large $e_n(\underline{w})$

□ Small $e_n(\underline{w})$

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

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

□ If $y_n = +1$, $e_n(\underline{w}) = -\log \hat{P}_{\underline{w}}(y_n/\underline{x}_n) =$

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

□ larger $e_n(\underline{w})$ □ smaller $e_n(\underline{w})$

$$e_n(\underline{w}) =$$

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

Summary: Logistic Regression

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

■ find $\underline{w} \in \mathbb{R}^{d+1}$,

■ to minimize $E_{in}(\underline{w}) = \frac{1}{N} \sum_{n=1}^N e_n(\underline{w}) = \frac{1}{N} \sum_{n=1}^N \log(1 + e^{-y_n \underline{w}^T \underline{x}_n})$.

■ For a new input \underline{x} , predict the probability of being in class y as $\hat{P}_{\underline{w}}(y | \underline{x}) = \frac{1}{1 + e^{-y \underline{w}^T \underline{x}}}$