

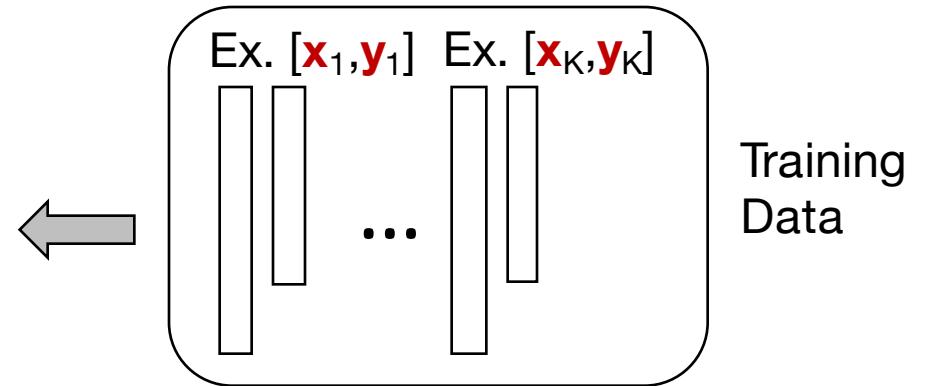
# Machine Learning in Imaging

BME 590L  
Roarke Horstmeyer

Lecture 6: Beyond linear classification and model complexity

# Summary of machine learning pipeline:

## 1. Network Training



E.g., images of 1's and 5's with labels:

What we need for network training:

### 1. Labeled examples

$$x_1 = \begin{matrix} \text{5} \end{matrix} \quad y_1 = +1$$

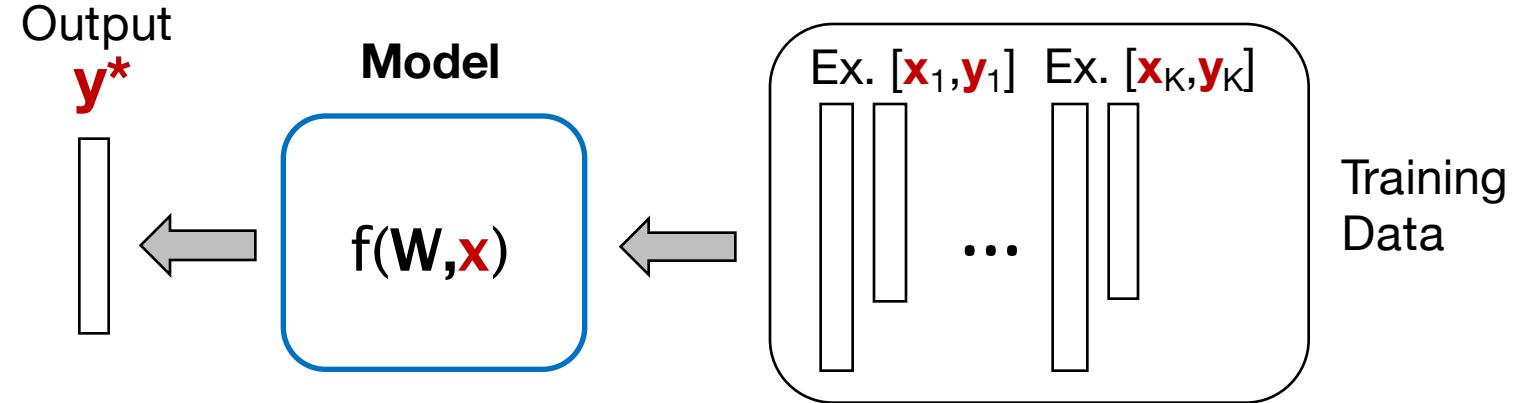
$$x_2 = \begin{matrix} \text{1} \end{matrix} \quad y_2 = +1$$

⋮

⋮

# Summary of machine learning pipeline:

## 1. Network Training

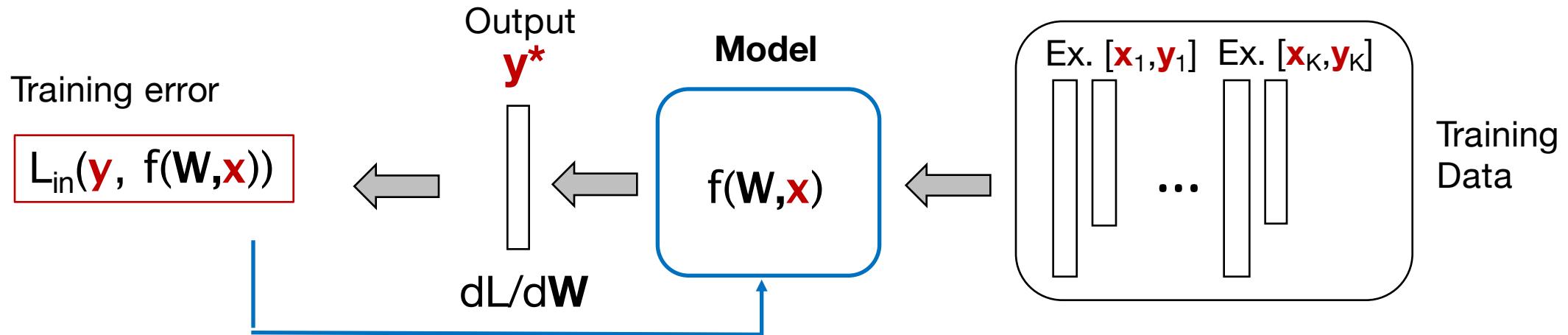


What we need for network training:

- 1. Labeled examples**
- 2. A model and loss function**

# Summary of machine learning pipeline:

## 1. Network Training

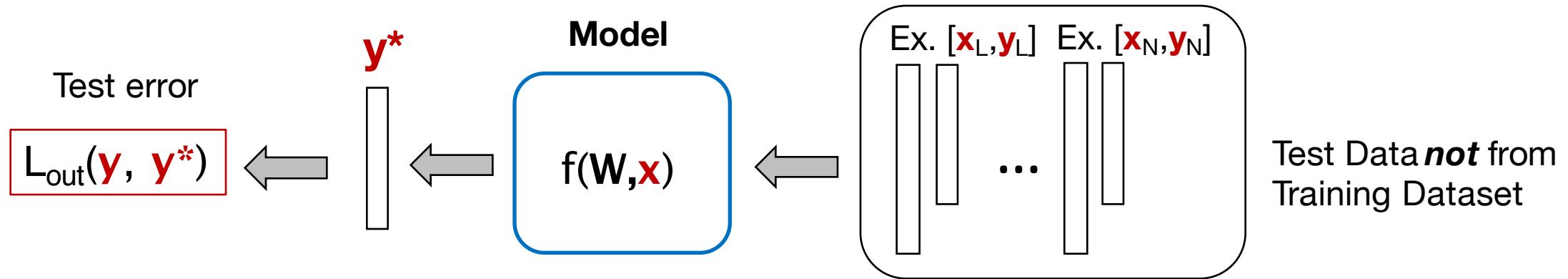


What we need for network training:

- 1. Labeled examples**
- 2. A model and loss function**
- 3. A way to minimize the loss function  $L$**

# Summary of machine learning pipeline:

## 2. Network Testing



What we need for network testing:

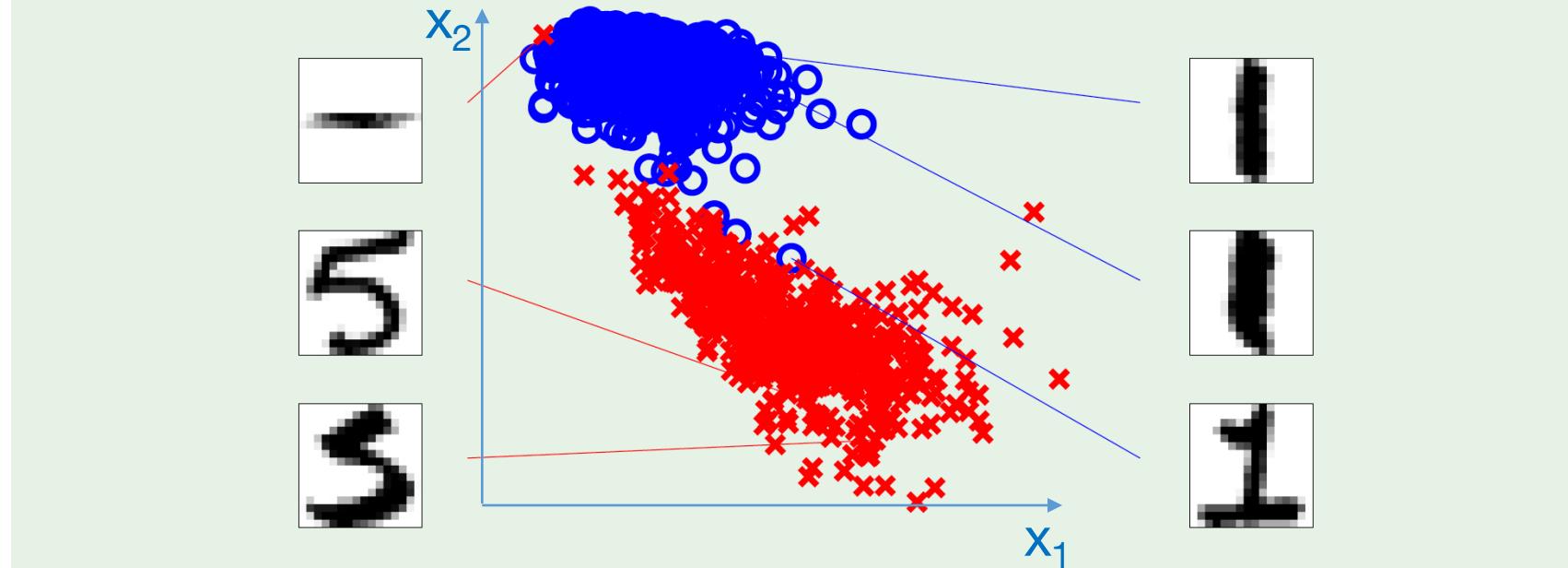
4. ***Unique* labeled test data**
5. **Evaluation of model error**

## Illustration of features

$$\mathbf{x} = (x_0, x_1, x_2)$$

$x_1$ : intensity

$x_2$ : symmetry



**Linear classification:**

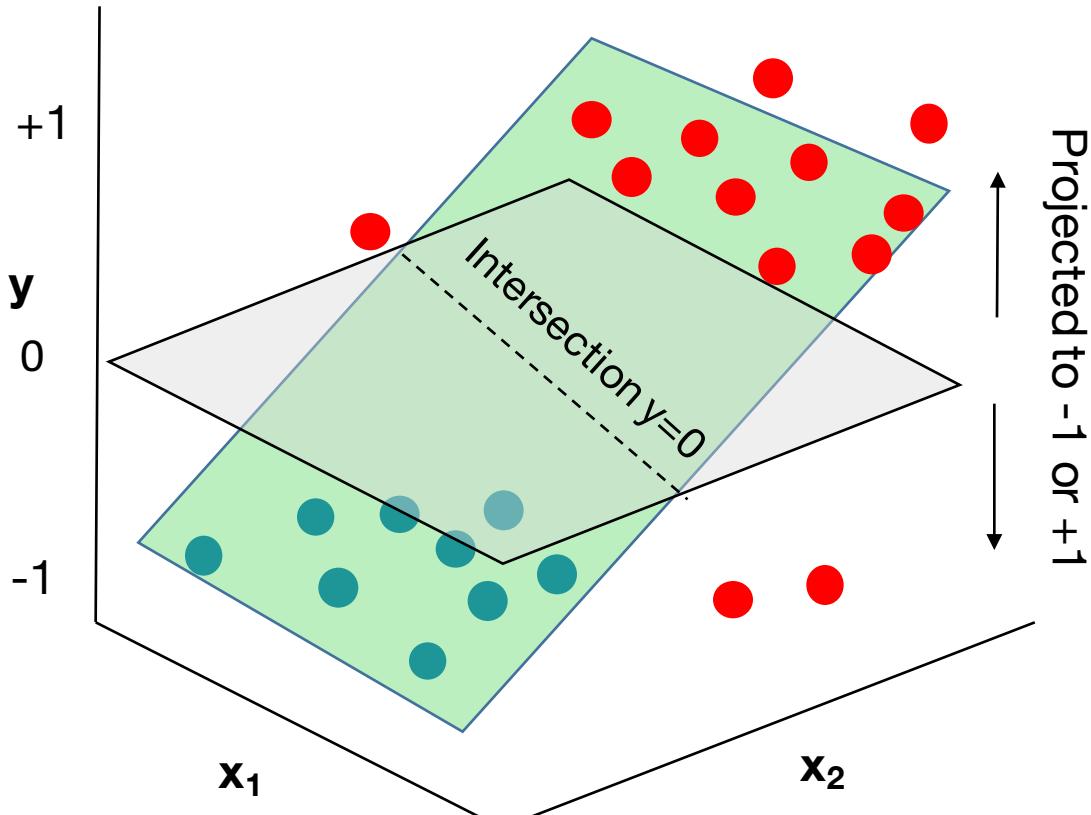
Use MSE error model

Where labels  
determined by  
thresholding

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

$$f(\mathbf{x}_i) = y_i^* = \operatorname{sgn}(\mathbf{w}^T \mathbf{x}_i) \quad \operatorname{sgn}(x) := \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases}$$

## Why does linear regression with $\text{sgn}()$ achieve classification?



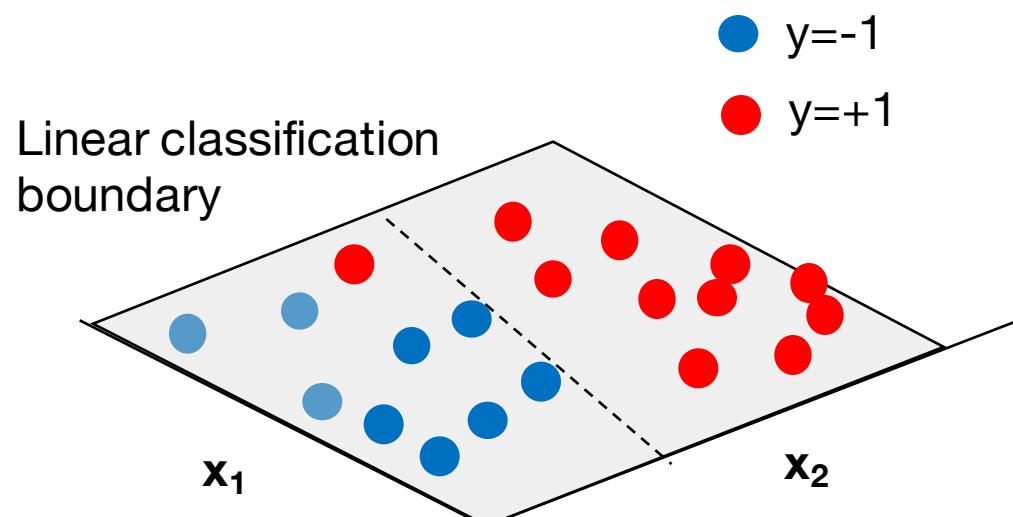
With  $\text{sgn}()$  operation:

$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

- Anything point to one side of  $y=0$  intersection is class  $+1$ , anything on the other side of intersection is class  $-1$

# Why does linear regression with $\text{sgn}()$ achieve classification?



With  $\text{sgn}()$  operation:

$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

Closed-form solution available!

# The linear regression algorithm

- 1: Construct the matrix  $\mathbf{X}$  and the vector  $\mathbf{y}$  from the data set  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  as follows

$$\mathbf{X} = \underbrace{\begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}}_{\text{input data matrix}}, \quad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

- 2: Compute the pseudo-inverse  $\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ .
- 3: Return  $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$ .

# In the rest of this class: solve via gradient descent



With a matrix, compute this for each entry:

$$\frac{dL(W_i)}{dW_i} = \lim_{h \rightarrow 0} \frac{L(W_i + h) - L(W_i)}{h}$$

Example:

$$W = [1,2;3,4]$$

$$L(W, x, y) = 12.79$$

$$W_1 + h = [1.001, 2; 3, 4]$$

$$L(W_1 + h, x, y) = 12.8$$

$$dL(W_1)/dW_1 = 12.8 - 12.79 / .001$$

$$dL(W_1)/dW_1 = 10$$

- Repeat for all entries of  $W$ ,  $dL/dW$  will have  $N \times M$  entries for  $N \times M$  matrix

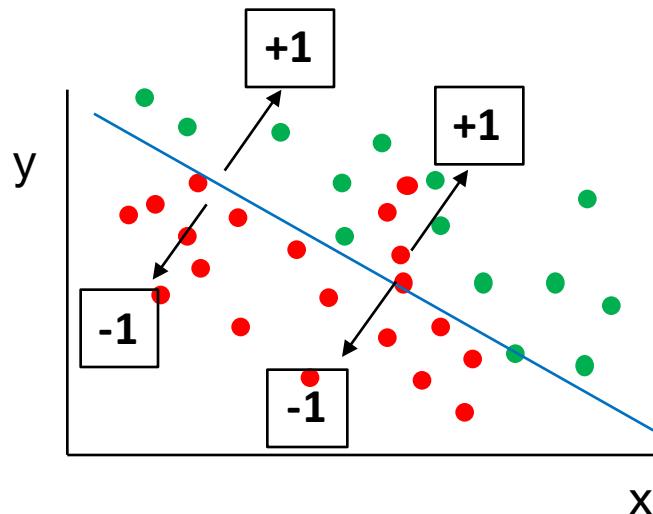
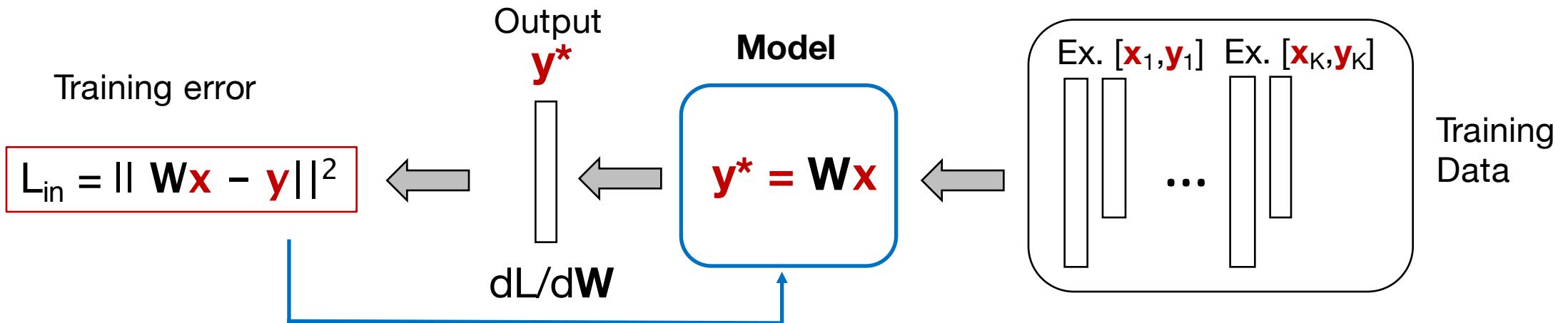
## Steepest descent and the best step size $\epsilon$

1. Evaluate function  $f(\mathbf{x}^{(0)})$  at an initial guess point,  $\mathbf{x}^{(0)}$
2. Compute gradient  $\mathbf{g}^{(0)} = \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$
3. Next point  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \epsilon^{(0)} \mathbf{g}^{(0)}$
4. Repeat –  $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \epsilon^{(n)} \mathbf{g}^{(n)}$ , until  $|\mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}| < \text{threshold } t$

$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$
$$\nabla L(w) = \frac{2}{N} X^T (Xw - y) = 0$$

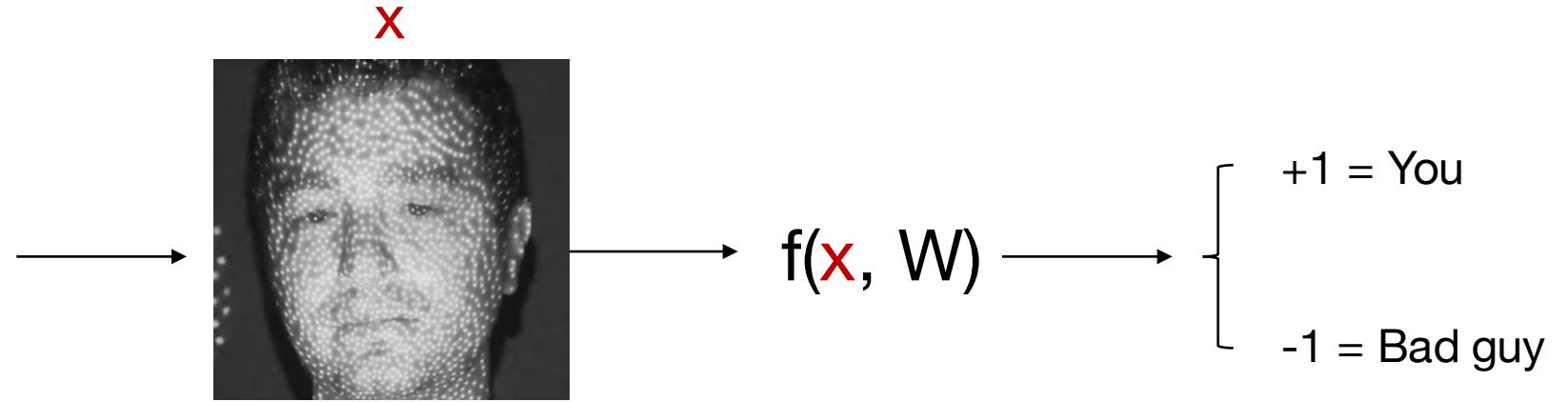
```
while previous_step_size > precision and iters < max_iters:  
    prev_x = cur_x  
    cur_x -= gamma * df(prev_x)  
    previous_step_size = abs(cur_x - prev_x)  
    iters+=1
```

# The linear classification model – what's not to like?



1. Can only separate data with lines (hyper-planes)...
2. We only allowed for binary labels ( $y = +/- 1$ )
3. Error function  $L_{in}$  inherently makes assumptions about statistical distribution of data

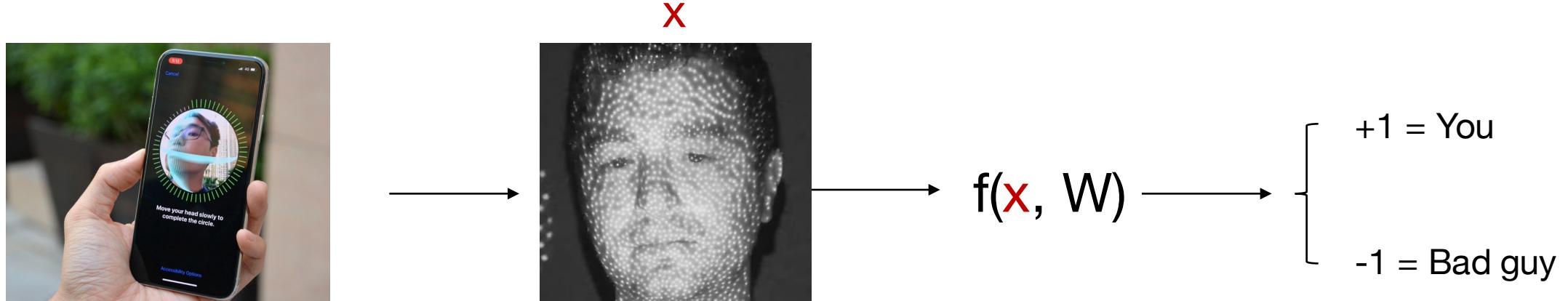
## Cost functions matter: a simple example



Two types of error: false accept and false reject

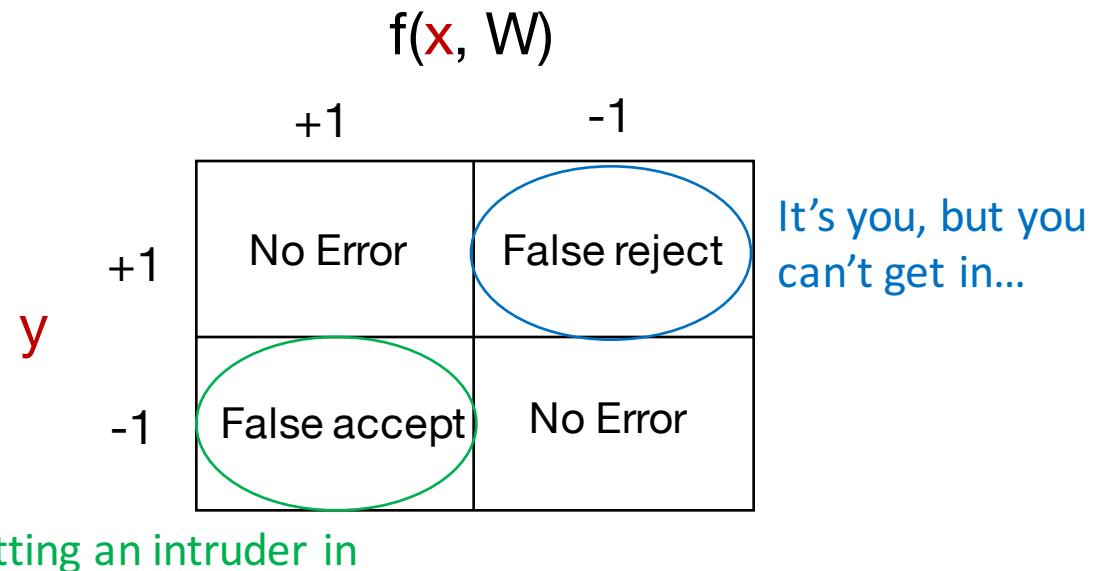
		$f(\mathbf{x}, \mathbf{W})$	
		+1	-1
$y$	+1	No Error (you/you)	False reject
	-1	False accept	No Error (bad guy/bad guy)

# Cost functions matter: a simple example

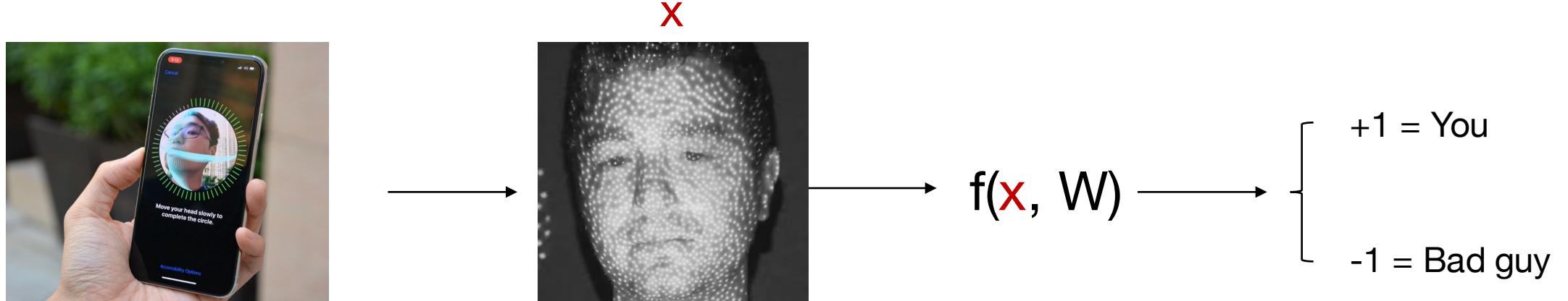


Two types of error: false accept and false reject

On a standard phone, what's a good cost function?



# Cost functions matter: a simple example



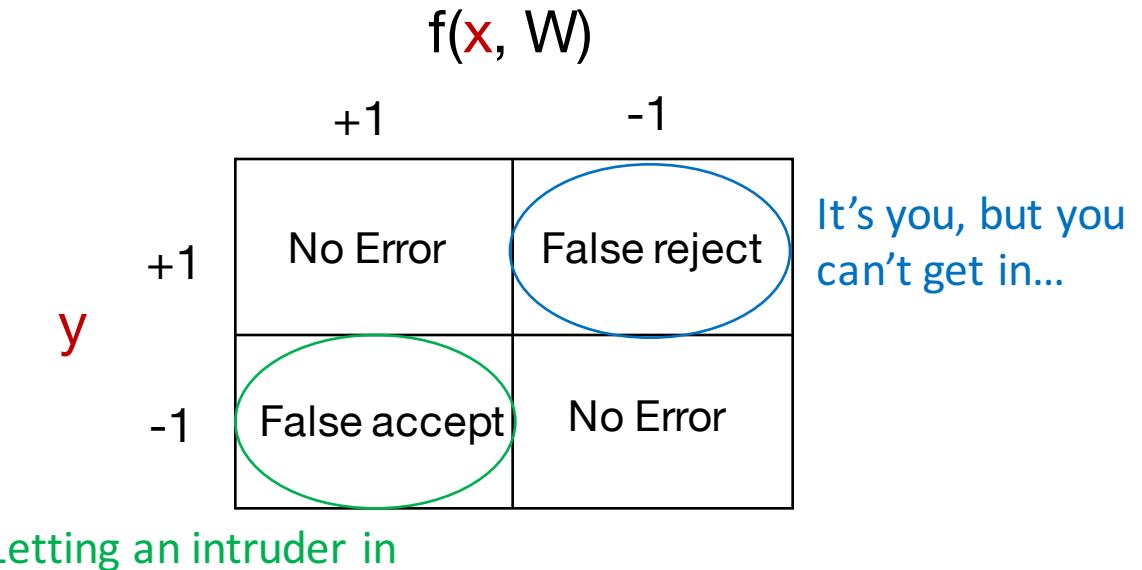
Two types of error: false accept and false reject

On a standard phone, what's a good cost function?

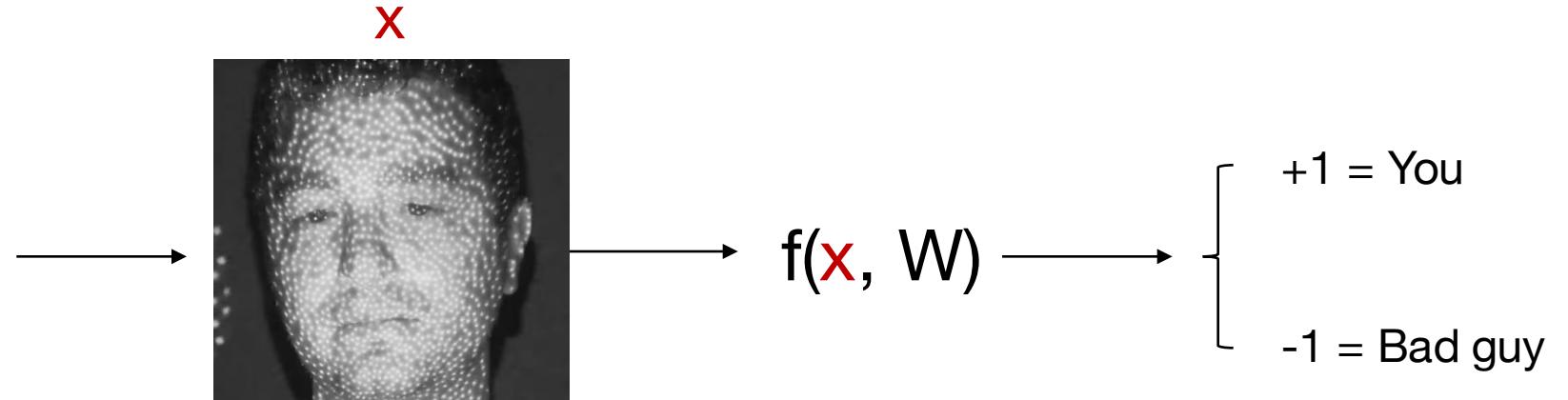
$$L_{in} = \text{ReLU}[f(x, W) - y] + 10 \text{ReLU}[y - f(x, W)]$$

Penalty for  
intruder

Large penalty for  
annoyance...



# Cost functions matter: a simple example



What if you're a CIA agent?

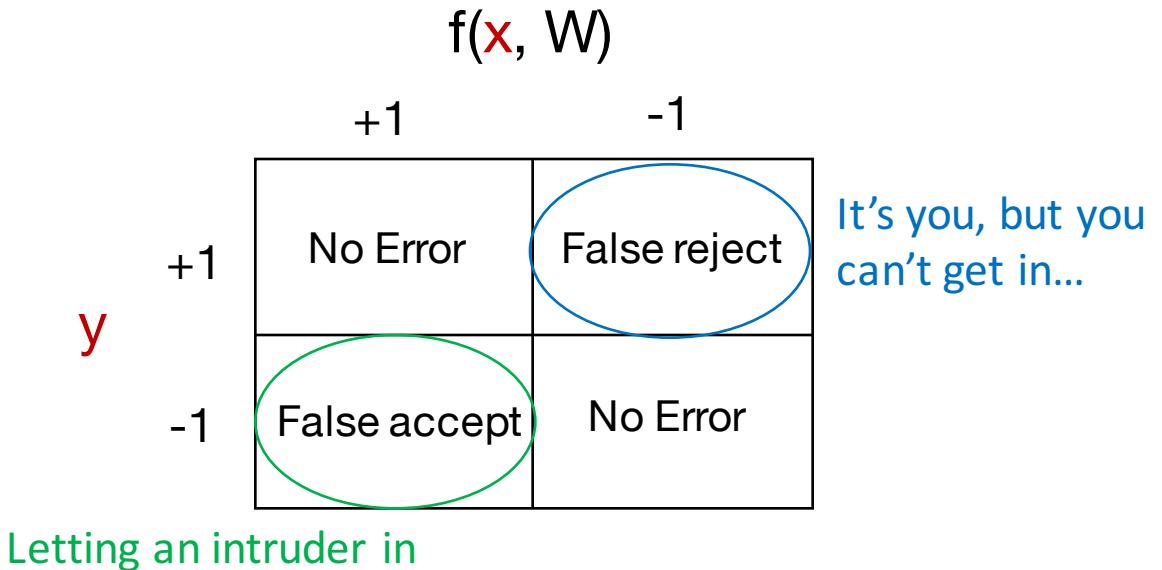
$$L_{in} = 100,000 \operatorname{ReLU}[f(\mathbf{x}, \mathbf{W}) - y] + \operatorname{ReLU}[y - f(\mathbf{x}, \mathbf{W})]$$

---

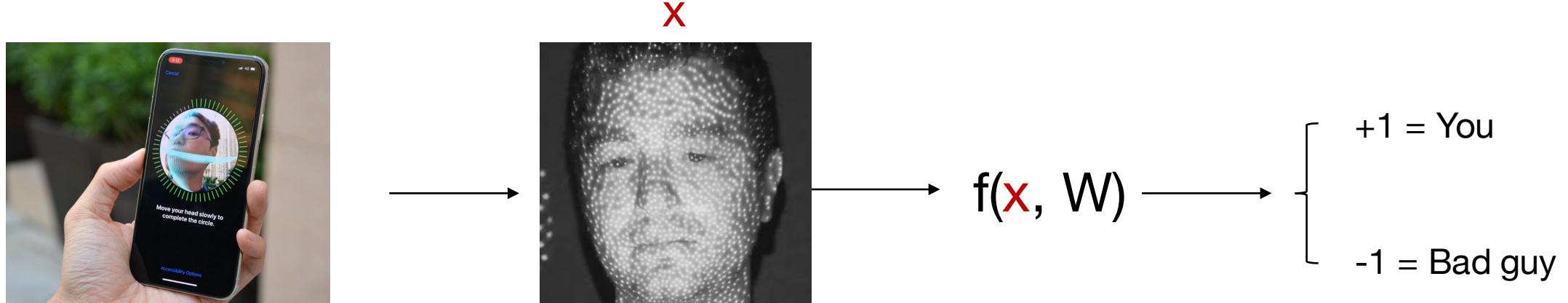
BIG penalty  
for intruder

---

Don't mind about  
annoyance...



# Cost functions matter: a simple example



Establishing cost function tied to conditional probabilities:

$$P(y = -1 \mid f(\mathbf{x}, \mathbf{W}) = +1)$$

$$P(y = +1 \mid f(\mathbf{x}, \mathbf{W}) = -1)$$

Establish  $L, W$  to balance these probabilities

		$f(\mathbf{x}, \mathbf{W})$	
		+1	-1
$y$	+1	No Error	False reject
	-1	False accept	No Error

Letting an intruder in

It's you, but you can't get in...

# Machine learning and probability

- Probability measures help determine when to use certain cost functions
- In previous case, we can measure probability of seeing a particular label, given model & data:

$$p(\mathbf{y}|\mathbf{w}, \mathbf{x})$$

- In general, to find the best model, we'd like to infer  $\mathbf{w}$ , having at hand our labeled data:

Maximum Likelihood Estimation

$$p(\mathbf{w}|\mathbf{x}_1, \dots, \mathbf{x}_N; \mathbf{y}_1, \dots, \mathbf{y}_N)$$

- These two quantities are connected via Bayes' Theorem

$$p(\mathbf{w}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{x})}$$

With 2  
conditioned  
variables:

$$p(\mathbf{w}|\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{x}, \mathbf{y})}$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{y}) \propto p(\mathbf{x}, \mathbf{y}|\mathbf{w}) \propto p(\mathbf{y}|\mathbf{x}, \mathbf{w})$$

# Machine learning and probability

- Probability measures help determine when to use certain cost functions
- In previous case, we can measure probability of seeing a particular label, given model & data:

$$p(\mathbf{y}|\mathbf{w}, \mathbf{x})$$

- In general, to find the best model, we'd like to infer  $\mathbf{w}$ , having at hand our labeled data:

Maximum Likelihood Estimation

$$p(\mathbf{w}|\mathbf{x}_1, \dots, \mathbf{x}_N; \mathbf{y}_1, \dots, \mathbf{y}_N)$$

- These two quantities are connected via Bayes' Theorem

$$p(\mathbf{w}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{x})}$$

With 2  
conditioned  
variables:

$$p(\mathbf{w}|\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y}|\mathbf{w})p(\mathbf{w})}{p(\mathbf{x}, \mathbf{y})}$$

What you want: but hard to  
vary data to find model...

$$p(\mathbf{w}|\mathbf{x}, \mathbf{y}) \propto p(\mathbf{x}, \mathbf{y}|\mathbf{w}) \propto p(\mathbf{y}|\mathbf{x}, \mathbf{w})$$

What you can do: test the model,  
check the result

## Linear classification is the maximum likelihood for Gaussian data

- Given a close relationship between  $p(\mathbf{w}|\mathbf{x}, \mathbf{y}) \longleftrightarrow p(\mathbf{y}|\mathbf{x}, \mathbf{w})$  :

Maximum likelihood estimation asks the question,

For what  $\mathbf{w}$  is  $p(\mathbf{y}_1, \dots, \mathbf{y}_N | \mathbf{x}_1, \dots, \mathbf{x}_N; \mathbf{w})$  maximized?

For what  $\mathbf{w}$  is  $\prod_{i=1}^N p(\mathbf{y}_i, | \mathbf{x}_i, \mathbf{w})$  maximized?

- Let's assume our labels are a "noisy" Gaussian process that surround the correct label:

$$y_i = \mathbf{w}^T \mathbf{x}_i + n \quad (n \text{ is zero-mean Gaussian noise})$$

- Then, the above cond. prob. for labels can be expressed as a multivariate Gaussian

## Linear classification is the maximum likelihood for Gaussian data

For what  $\mathbf{w}$  is  $\prod_{i=1}^N p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{w})$  maximized?



For what  $\mathbf{w}$  is  $\prod_{i=1}^N \exp\left(\frac{-(\mathbf{y}_i - \mathbf{w}^T \mathbf{x})^2}{2\sigma^2}\right)$  maximized?



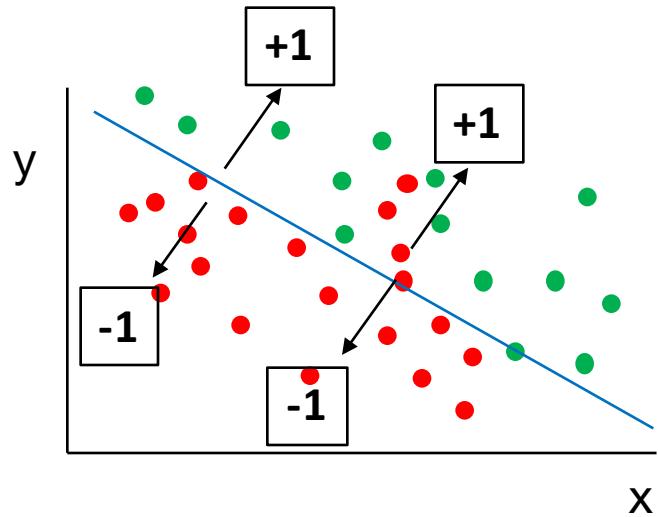
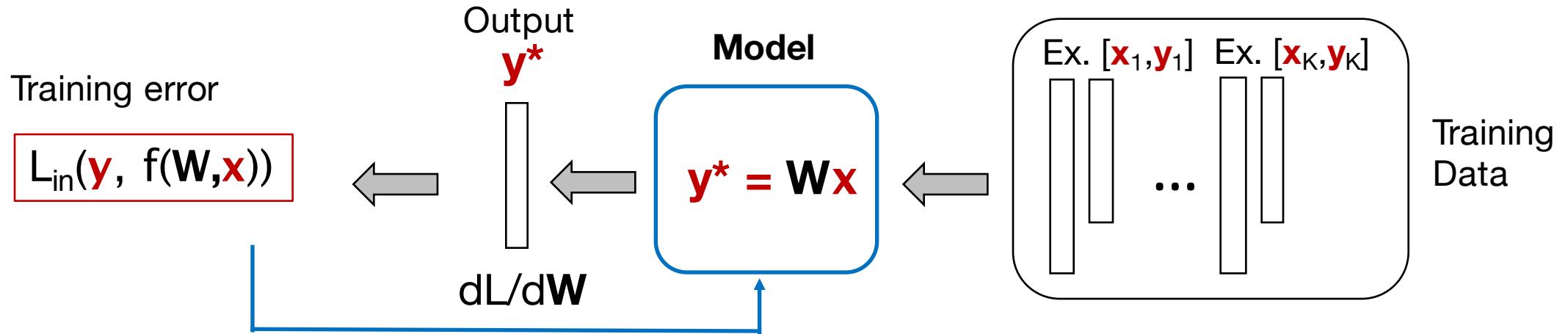
For what  $\mathbf{w}$  is  $\sum_{i=1}^N \frac{-(\mathbf{y}_i - \mathbf{w}^T \mathbf{x})^2}{2\sigma^2}$  maximized?



For what  $\mathbf{w}$  is  $\sum_{i=1}^N (\mathbf{y}_i - \mathbf{w}^T \mathbf{x})^2$  minimized?

Summary: Linear classification assumes labels are a Gaussian random process, and then thresholds them to either -1 or +1

# The linear classification model – what's not to like?



1. Can only separate data with lines (hyper-planes)...
2. We only allowed for binary labels ( $y = +/- 1$ )
3. Error function  $L_{in}$  inherently makes assumptions about statistical distribution of data

## Let's think about the labels as a probabilistic measure:

Linear regression:

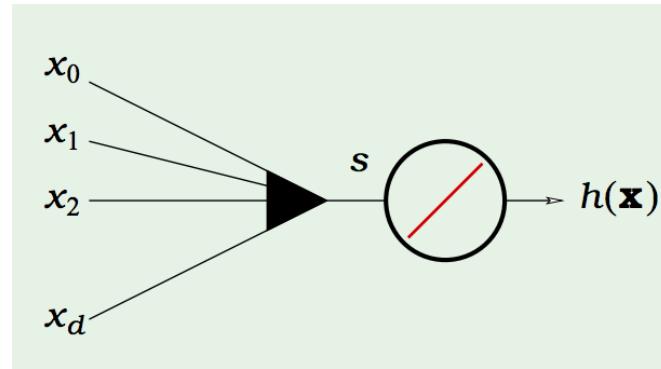


$$h(x) = s \in (-\infty, \infty)$$

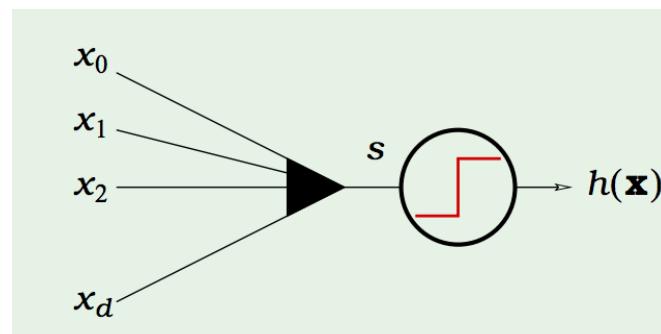
Linear classifier:

$$h(x) = \text{sign}(w_o^T x_j)$$

$$h(x) = \text{sign}(s) \in \{0, 1\}$$



Not a  
probabilistic  
mapping



Probabilistic, but  
all-or-nothing

## Let's think about the labels as a probabilistic measure:

Linear regression:

$$h(x) = s \in (-\infty, \infty)$$

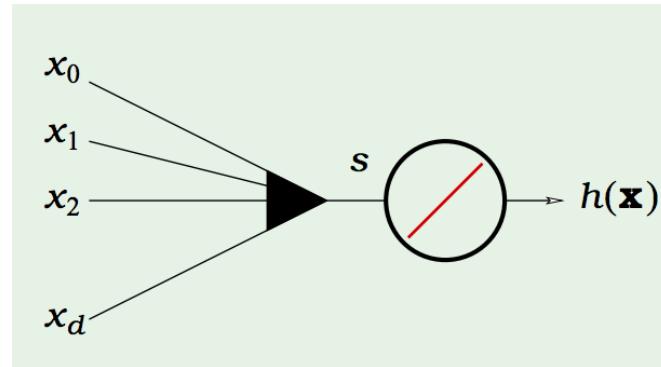
Linear classifier:

$$h(x) = \text{sign}(w_o^T x_j)$$

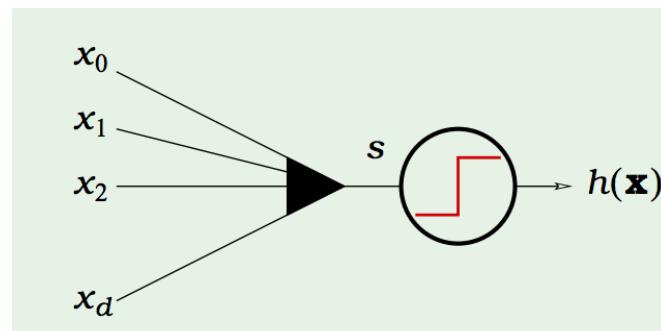
$$h(x) = \text{sign}(s) \in \{0, 1\}$$

Logistic classifier:

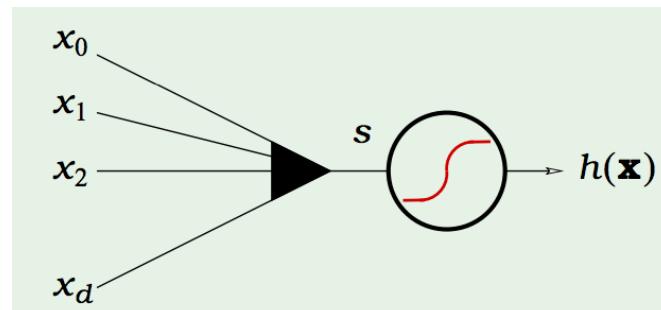
$$h(x) = \theta(s) \in [-1, 1]$$



Not a  
probabilistic  
mapping



Probabilistic, but  
all-or-nothing



Probabilistic

## Probability interpretation

$h(\mathbf{x}) = \theta(s)$  is interpreted as a probability

**Example.** Prediction of heart attacks

Input  $\mathbf{x}$ : cholesterol level, age, weight, etc.

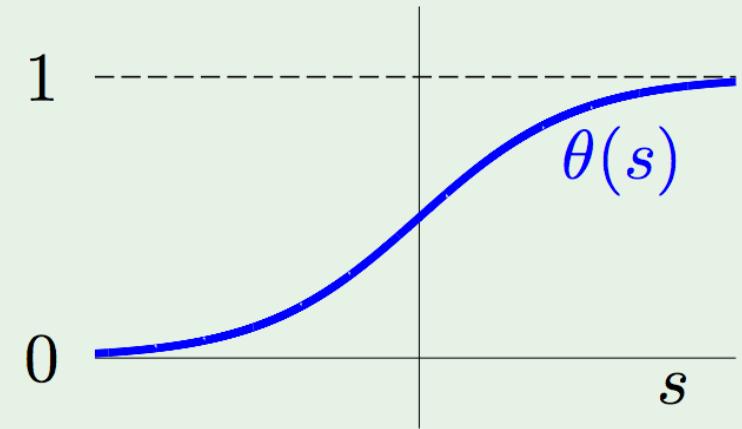
$\theta(s)$ : probability of a heart attack

The signal  $s = \mathbf{w}^\top \mathbf{x}$       “risk score”

## The logistic function $\theta$

The formula:

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



soft threshold: uncertainty

sigmoid: flattened out 's'

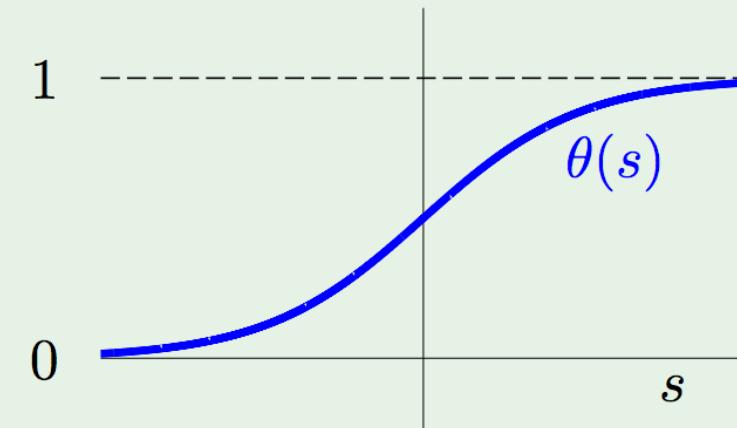
## From linear classification to logistic classification

- Let's re-derive a cost function for the case where labels are treated as probabilities
  - You'll use derivatives of this more often than not in Tensorflow
- During learning, we will again have two classes (in this simple example),  $y = +/- 1$
- map these binary values onto a  $[0,1]$  probability distribution,  $f(x)$ :

### Formula for likelihood

$$P(y | \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

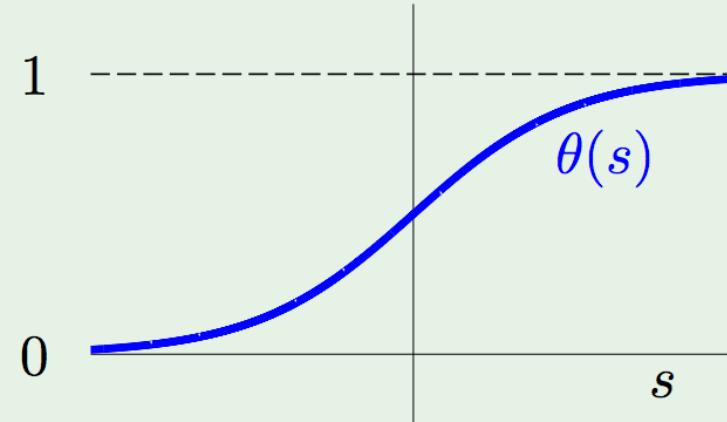
$h(\mathbf{x}) = \theta(s)$  is interpreted as a probability



# Formula for likelihood

$$P(y | \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

$h(\mathbf{x}) = \theta(s)$  is interpreted as a probability



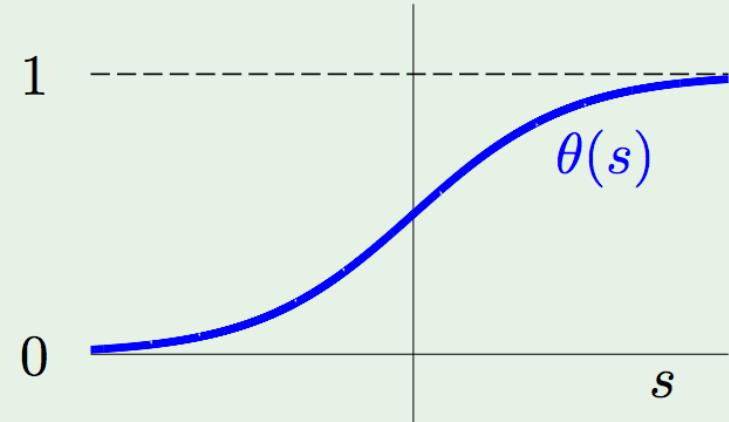
Rough Interpretation:

- If score  $x$  is high, then should map to  $y=+1$  with high probability
- $h(x) = \theta(x)$  is large for large value of  $x$
- If score  $x$  is small, then should map to  $y=-1$  with high probability
- $h(x) \sim 0$  for small values of  $x$ , so  $1 - h(x) \sim 1$  is high probability to  $y=-1$  mapping

## Formula for likelihood

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute  $h(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x})$ , noting  $\theta(-s) = 1 - \theta(s)$



$$P(y \mid \mathbf{x}) = \theta(y \mathbf{w}^\top \mathbf{x})$$

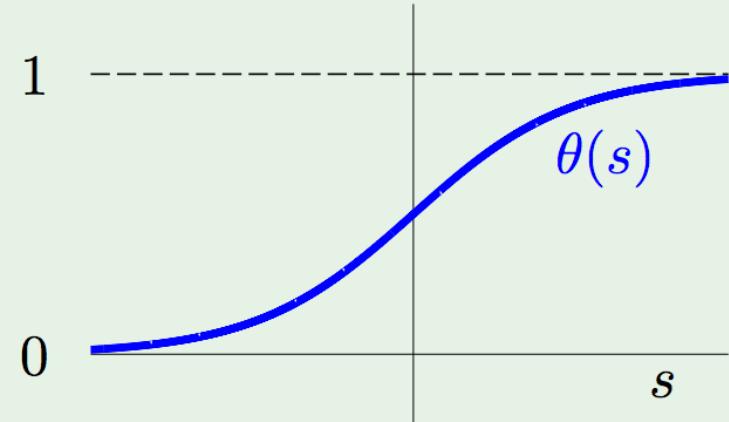
Likelihood of  $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n)$$

## Formula for likelihood

$$P(y \mid \mathbf{x}) = \begin{cases} h(\mathbf{x}) & \text{for } y = +1; \\ 1 - h(\mathbf{x}) & \text{for } y = -1. \end{cases}$$

Substitute  $h(\mathbf{x}) = \theta(\mathbf{w}^\top \mathbf{x})$ , noting  $\theta(-s) = 1 - \theta(s)$



This is the probability of the labels, given the data. We'd like to maximize this probability!

\*Like the linear regression case, but now the probability of classes given the data is not Gaussian distributed, but instead follows the sigmoid curve (is bound to  $[0,1]$ , which is more realistic)

$$P(y \mid \mathbf{x}) = \theta(y \mathbf{w}^\top \mathbf{x})$$

Likelihood of  $\mathcal{D} = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$  is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n) = \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n)$$

## Maximizing the likelihood

Minimize

$$-\frac{1}{N} \ln \left( \prod_{n=1}^N \theta(y_n \mathbf{w}^\top \mathbf{x}_n) \right)$$

$$= \frac{1}{N} \sum_{n=1}^N \ln \left( \frac{1}{\theta(y_n \mathbf{w}^\top \mathbf{x}_n)} \right)$$

$$\left[ \theta(s) = \frac{1}{1 + e^{-s}} \right]$$

$$L_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \underbrace{\ln \left( 1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right)}_{e(h(\mathbf{x}_n), y_n)}$$

“cross-entropy” error

## How to minimize $L_{\text{in}}$

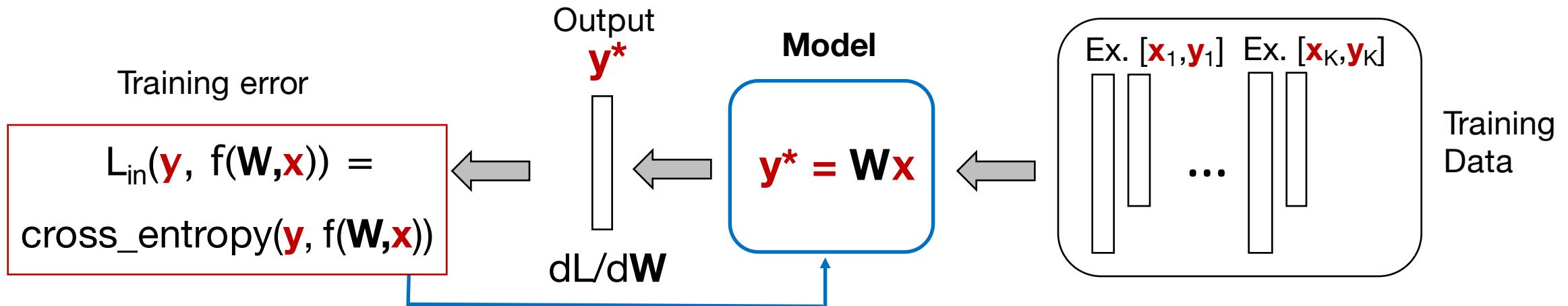
For logistic regression,

$$L_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln \left( 1 + e^{-y_n \mathbf{w}^\top \mathbf{x}_n} \right) \quad \longleftarrow \text{iterative solution}$$

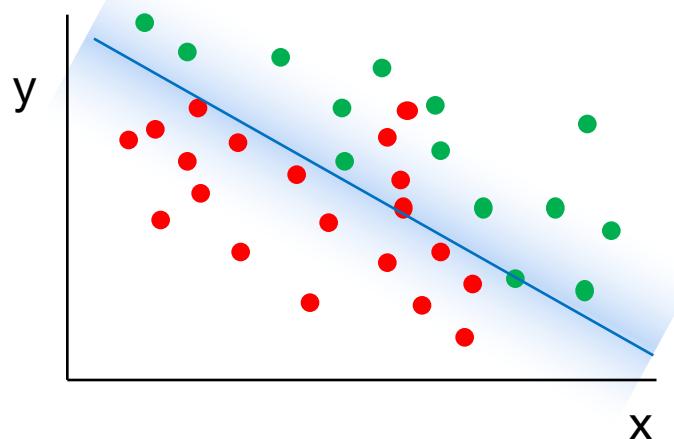
Compare to linear regression:

$$L_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 \quad \longleftarrow \text{closed-form solution}$$

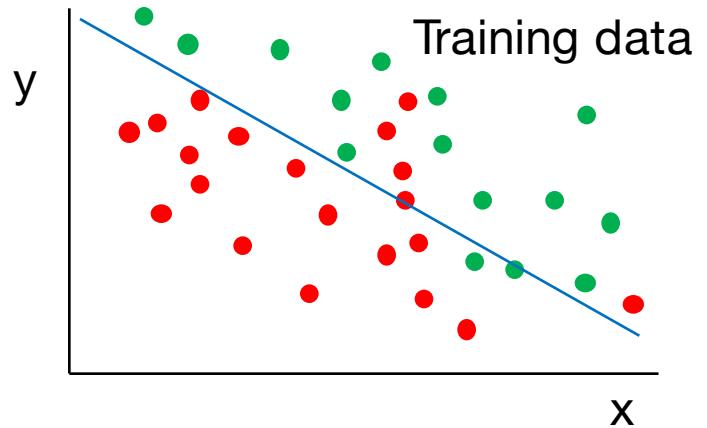
# The linear classification model – what's not to like?



Probabilistic mapping to  $y$

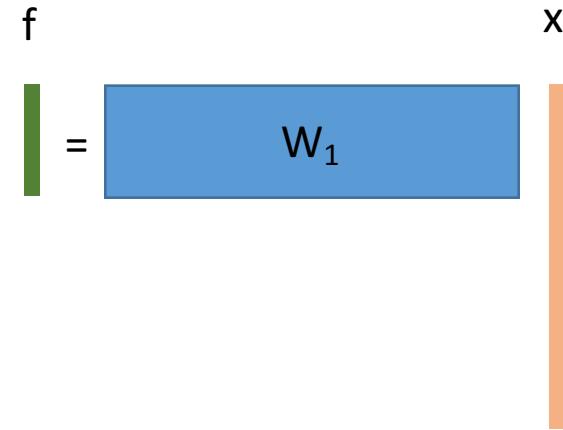


1. Can only separate data with lines (hyper-planes)...
2. We only allowed for binary labels ( $y = +/- 1$ )
3. Error function  $L_{in}$  inherently makes assumptions about statistical distribution of data



$$f = W_1x$$

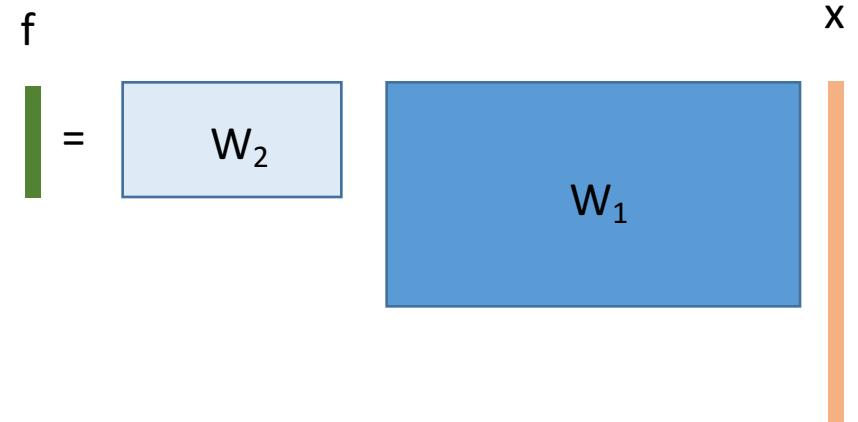
Learned  $f$ : not flexible



Can we add flexibility by multiplying with another weight matrix?

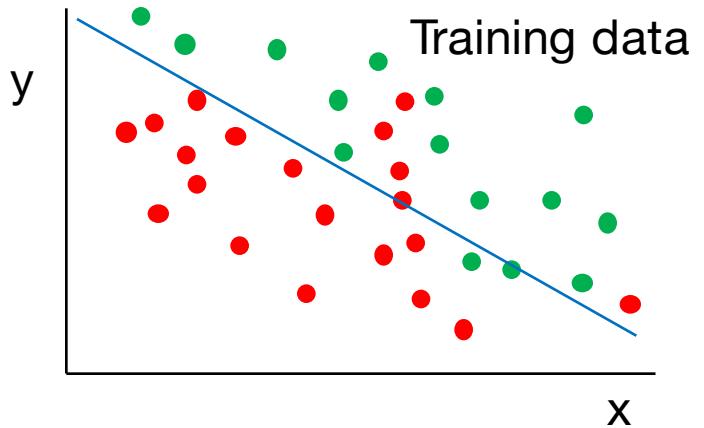
$$\begin{cases} f_1 = W_1x + b_1 \\ f_2 = W_2f_1 + b_2 \end{cases}$$

$$f_2 = W_2(W_1x + b_1) + b_2$$



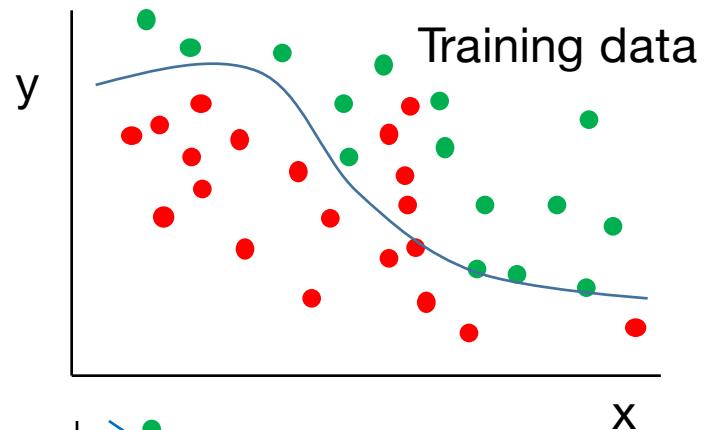
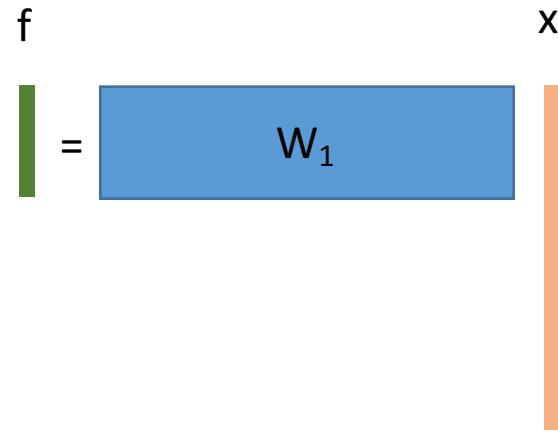
$$f_2 = W'x + b'$$

Unfortunately not...



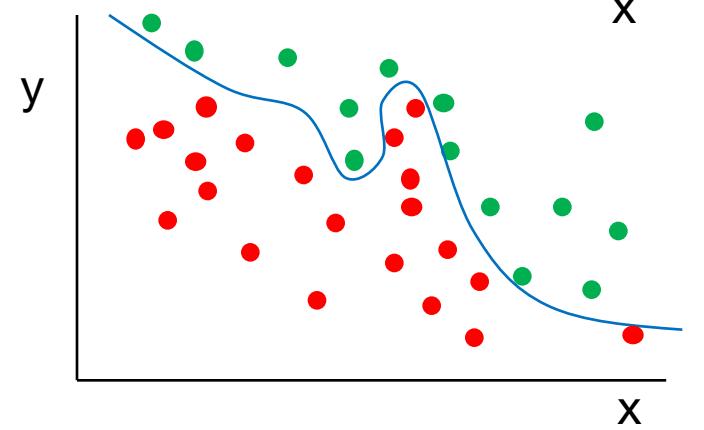
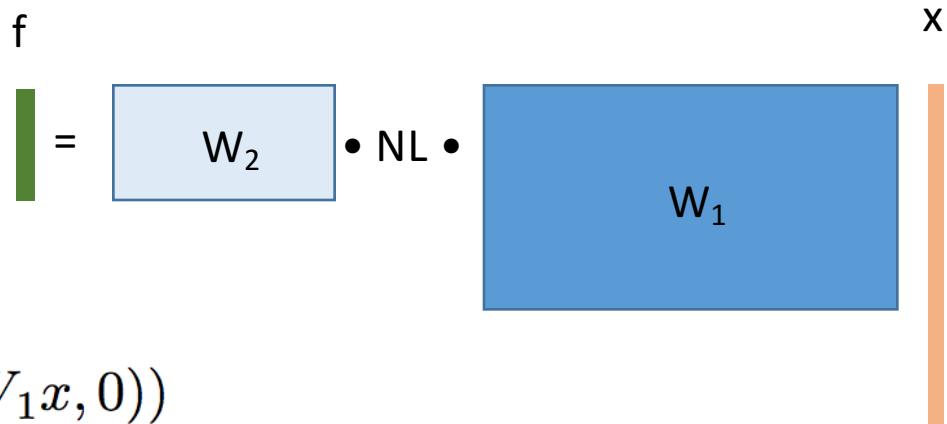
$$f = W_1x$$

Learned  $f$ : not flexible



$$f = W_2 \max(W_1x, 0)$$

Learned  $f$ : a bit flexible



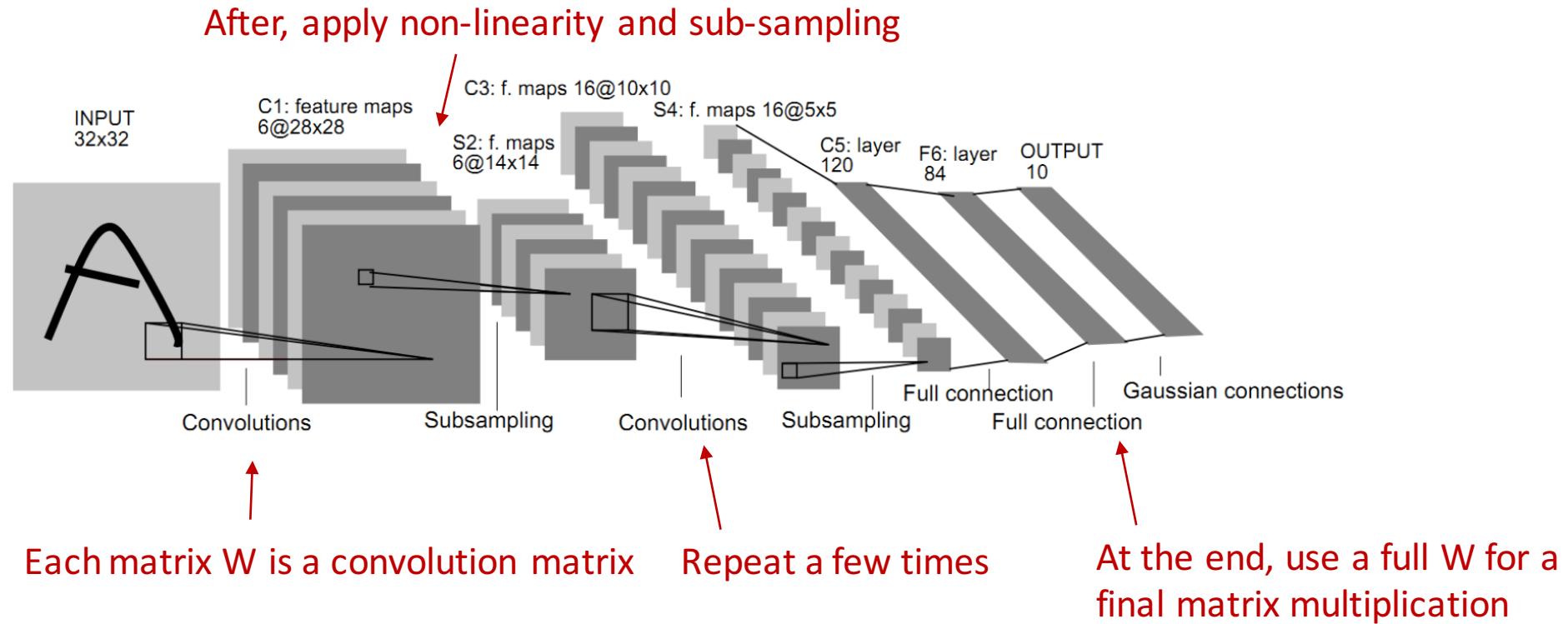
$$f = W_3 \max(0, W_2 \max(W_1x, 0))$$

Learned  $f$ : more flexible

Does it generalize???

We can keep adding  
these “layers”...

# Getting us to Convolutional Neural Networks



# Getting us to Convolutional Neural Networks

