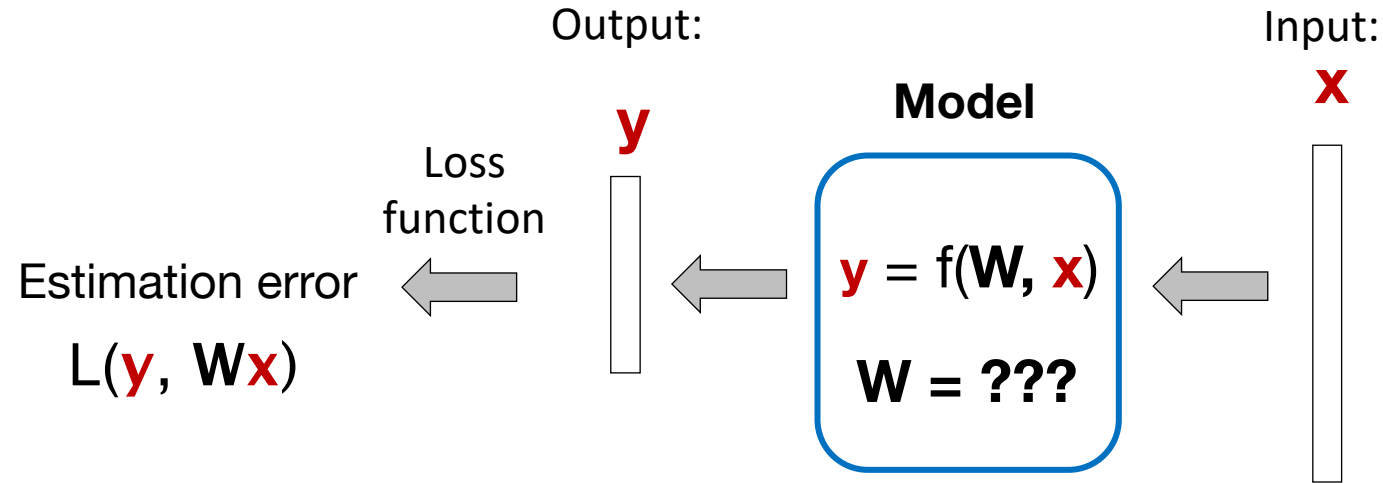


Lecture 6, Part 2: Ingredients for Machine Learning

Machine Learning and Imaging

BME 548L
Roarke Horstmeyer

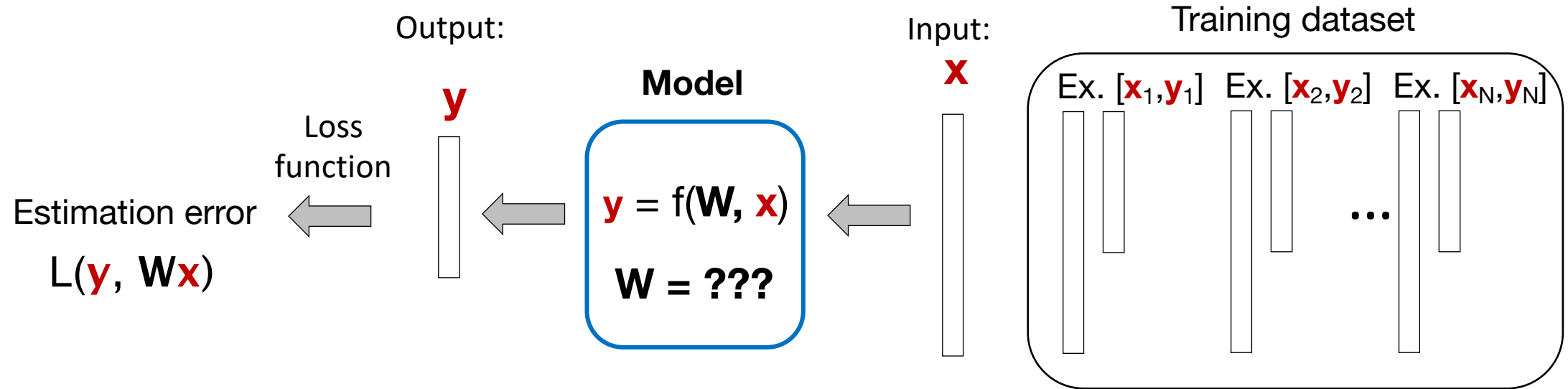
Pipeline for machine learning



Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix W)

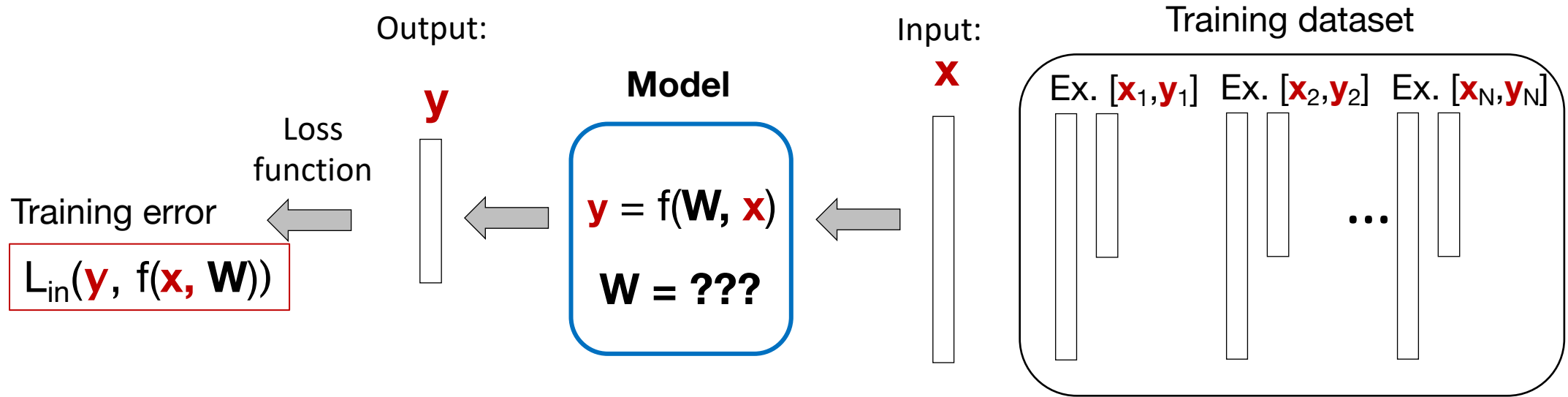
Pipeline for machine learning



Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix W)
2. Using large set of "training" data to first determine mapping $f(x, W)$

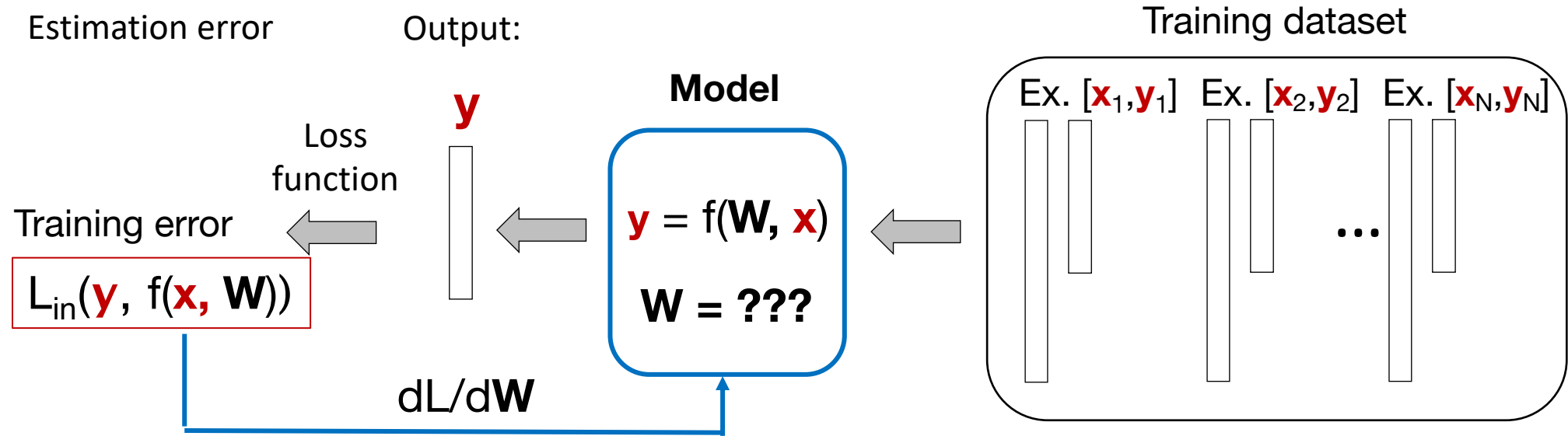
Pipeline for machine learning



Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix W)
2. Using large set of “training” data to first determine mapping $f(x, W)$
3. To do so, use a *loss function* L that depends upon the training inputs (x, y) and the model (W)

Pipeline for machine learning



Changes for machine learning framework:

1. Now must establish the mapping from inputs to outputs (here, matrix W)
2. Using large set of “training” data to first determine mapping $f(x, W)$
3. To do so, use a *loss function* L that depends upon the training inputs (x, y) and the model (W)
4. Find optimal mapping (W) using the training data, guided by gradient descent on L

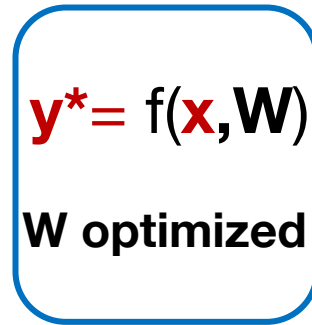
Pipeline for machine learning

Output:

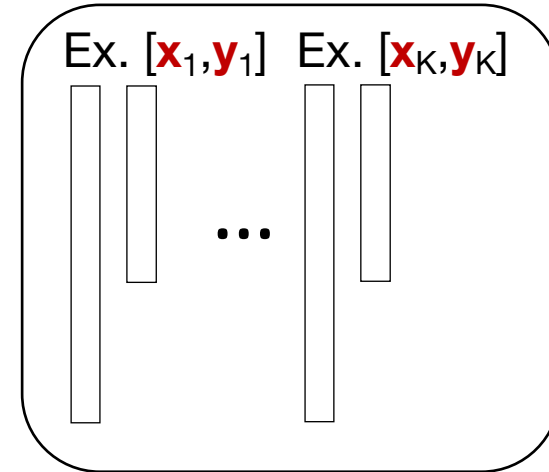
y^*



Model



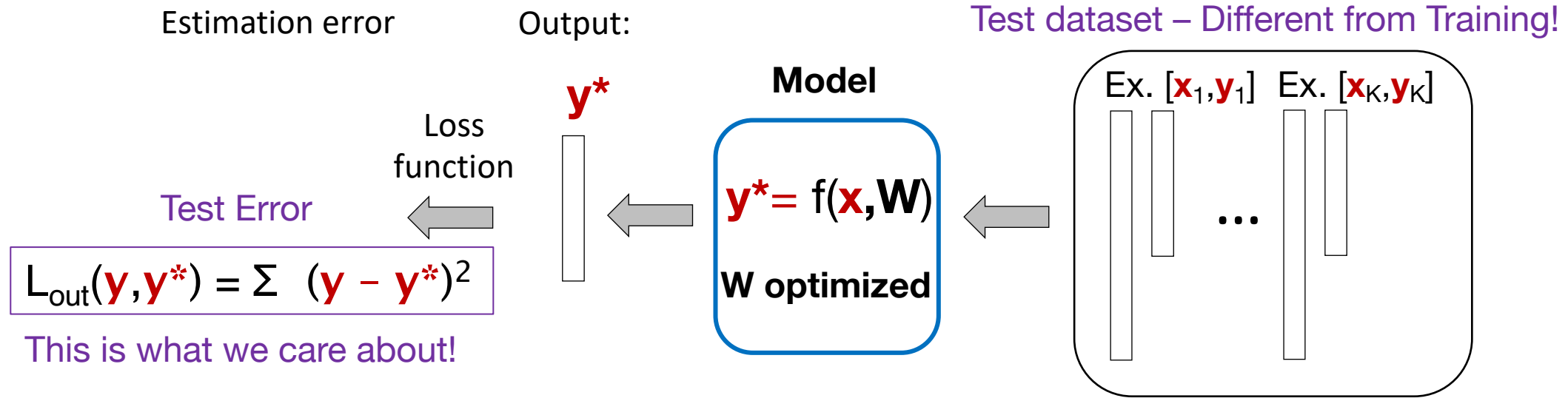
Test dataset – Different from Training!



In a *separate step*, we then need to do the following to test the network:

1. Evaluate model accuracy by sending *new x* through – need *new, unique data with label*

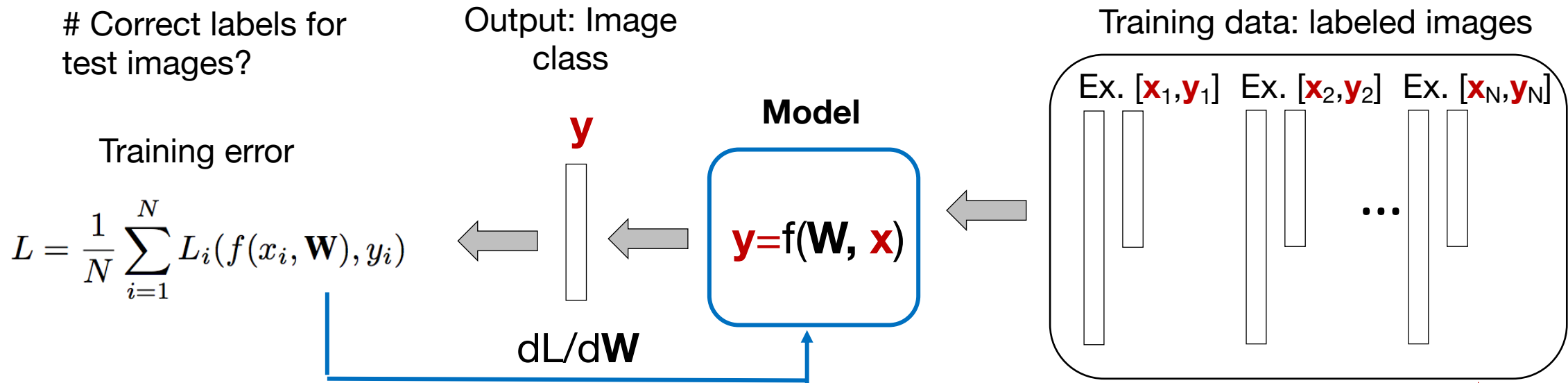
Pipeline for machine learning



In a *separate step*, we then need to do the following to test the network:

1. evaluate model accuracy by sending *new* x through – need *new, unique* data with label
2. Compare output y^* to known “test data” label y
3. Evaluate performance with an error equation L_{out}

Example: machine learning for image classification



Let's consider a simple example – image classification. What do we need for training?

1. Labeled examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Example: machine learning for image classification



https://en.wikipedia.org/wiki/MNIST_database

MNIST image set: <http://yann.lecun.com/exdb/mnist/>

Example: MNIST image dataset



\mathbf{X} = 28x28 pixel matrix

$\mathbf{x} = \text{vec}[\mathbf{X}] = 784\text{-long vector}$

Linear model would require $W = 784$ element matrix

Start simple: use $\mathbf{x} = (x_0, x_1, x_2)$ to describe **intensity** and **symmetry** of image X

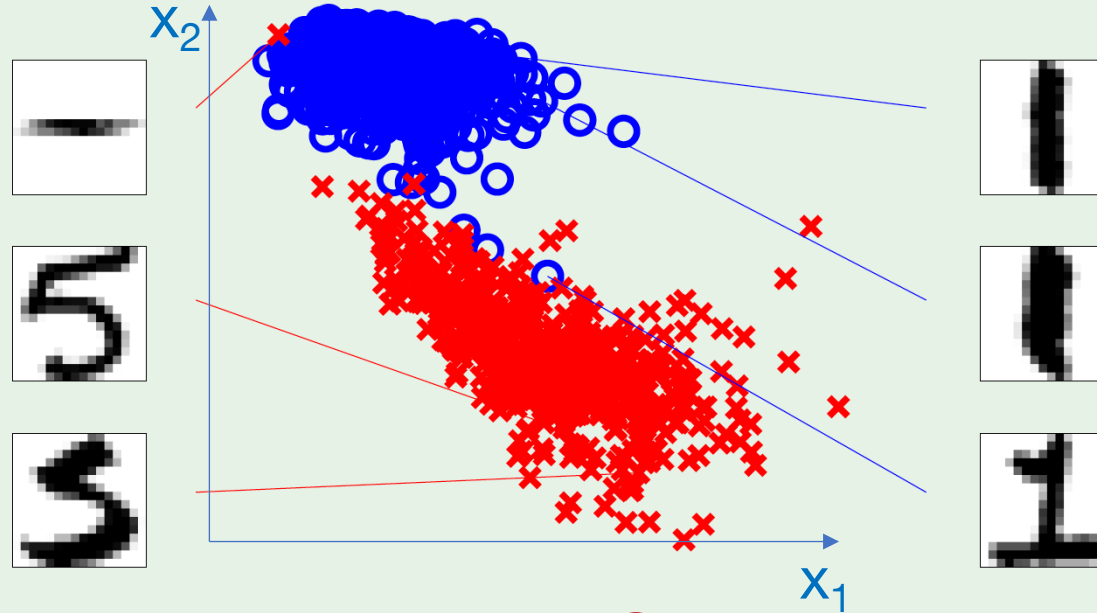
Linear model can now use smaller $w = (w_0, w_1, w_2)$

Illustration of features

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

x_1 : intensity

x_2 : symmetry

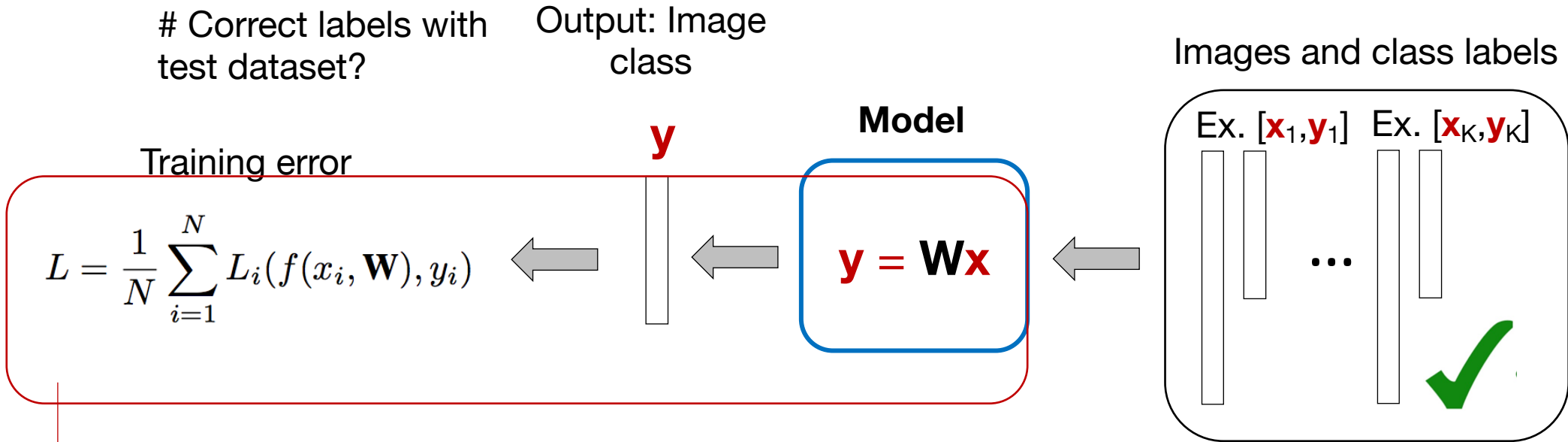


Dataset: 1000 examples of 1's and 5's mapped to $\mathbf{x}_j = (1, x_1, x_2)$, with associated label $y_j = 1$ or -1

This transforms $w\mathbf{x} + b$ into $\mathbf{w}\mathbf{x}$

$$\begin{bmatrix} w_0 & w_1 & w_2 \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} = w_0 + w_1x_1 + w_2x_2$$

Example: machine learning for image classification



Let's consider a simple example – image classification. What do we need for training?

✓ **1. Labeled examples**

$$\{(x_i, y_i)\}_{i=1}^N$$

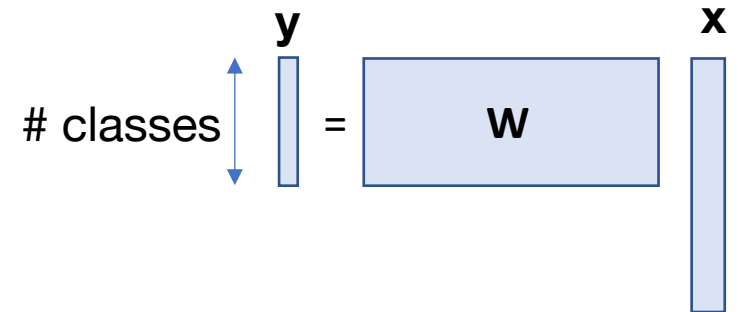
➔ **2. A model and loss function**

Let's start with a simpler approach: linear regression

General linear model:

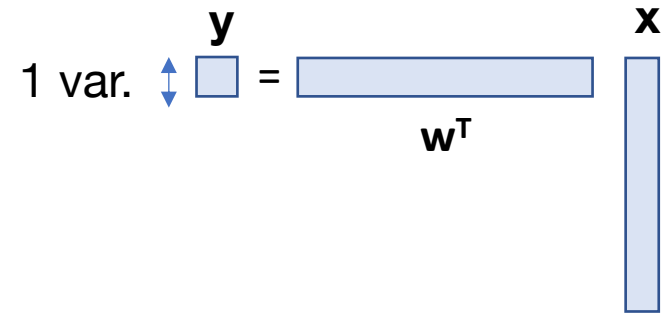
$$L = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, \mathbf{W}), y_i)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i(\mathbf{W}\mathbf{x}_i, y_i)$$



Assume 1 class = 1 linear fit

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



Use MSE error model

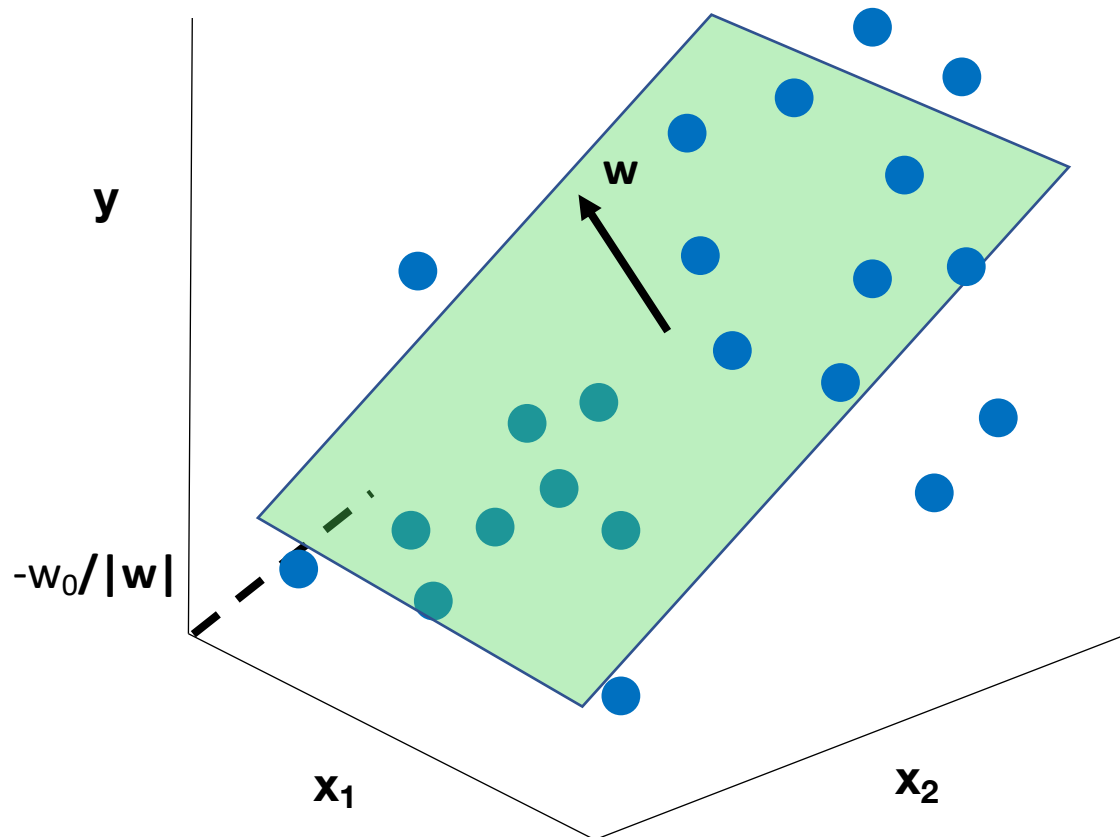
$$L = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2$$

Where labels determined by thresholding

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

$$\text{sgn}(x) := \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases}$$

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



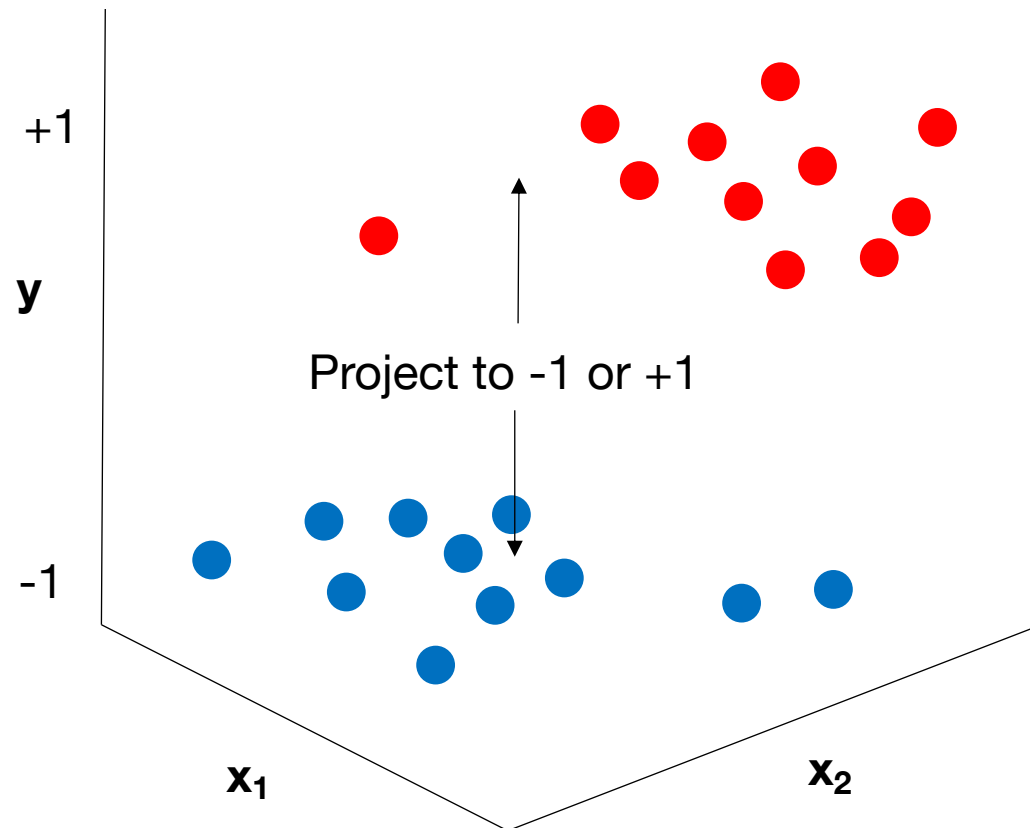
Without $\text{sgn}()$: regression for best fit

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

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

- If y_i can be anything, minimizing L makes \mathbf{w} the plane of best fit

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



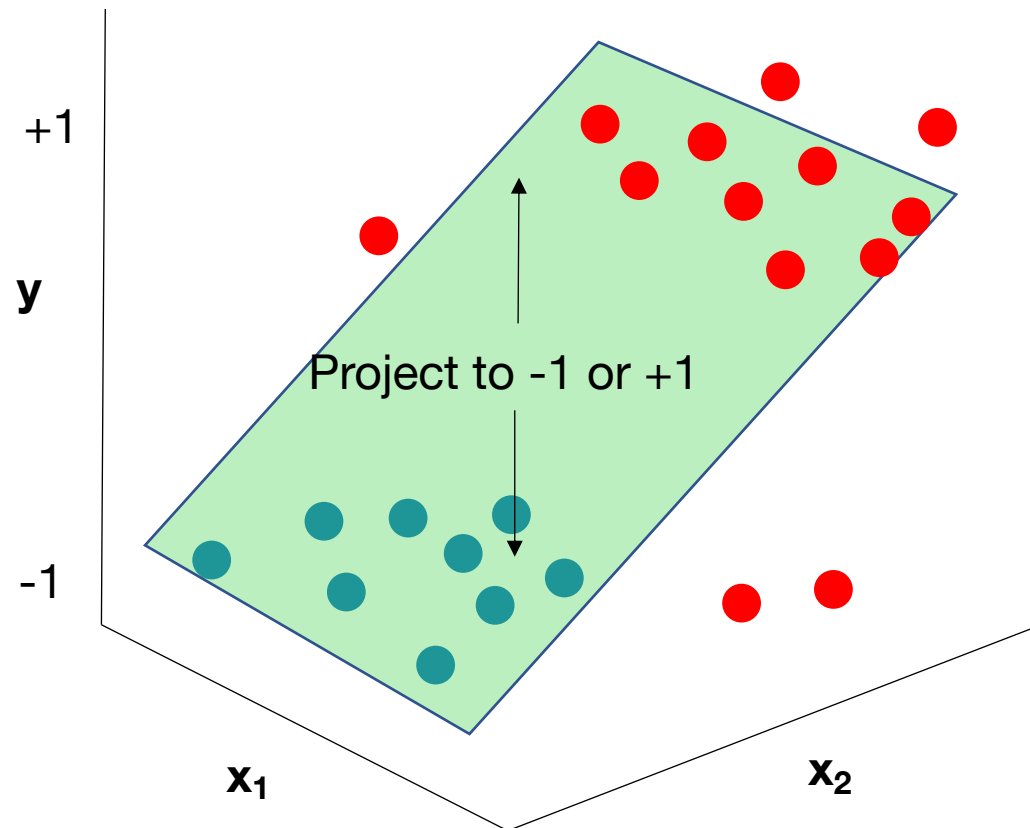
Without $\text{sgn}()$: regression for best fit

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

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

- y_i can only be -1 or +1, which defines its class

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



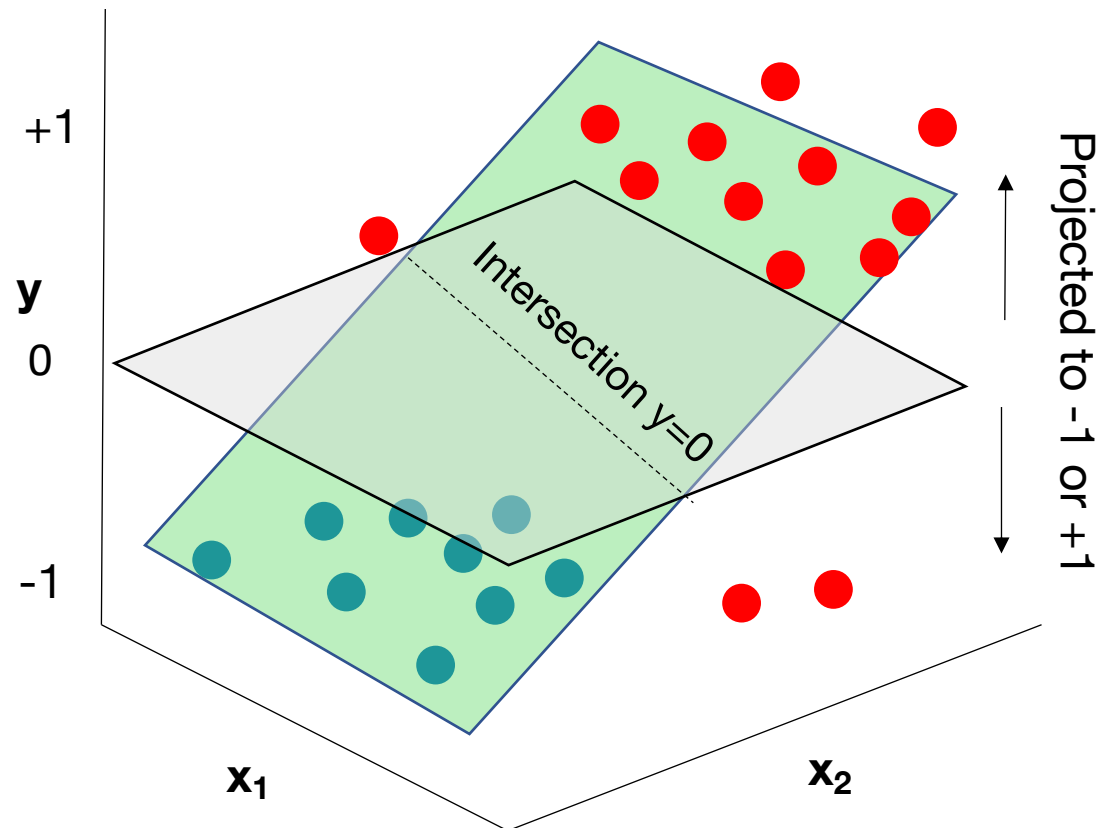
Without $\text{sgn}()$: regression for best fit

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$$

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

- y_i can only be -1 or +1, which defines its class
- Can still find plane of best fit

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



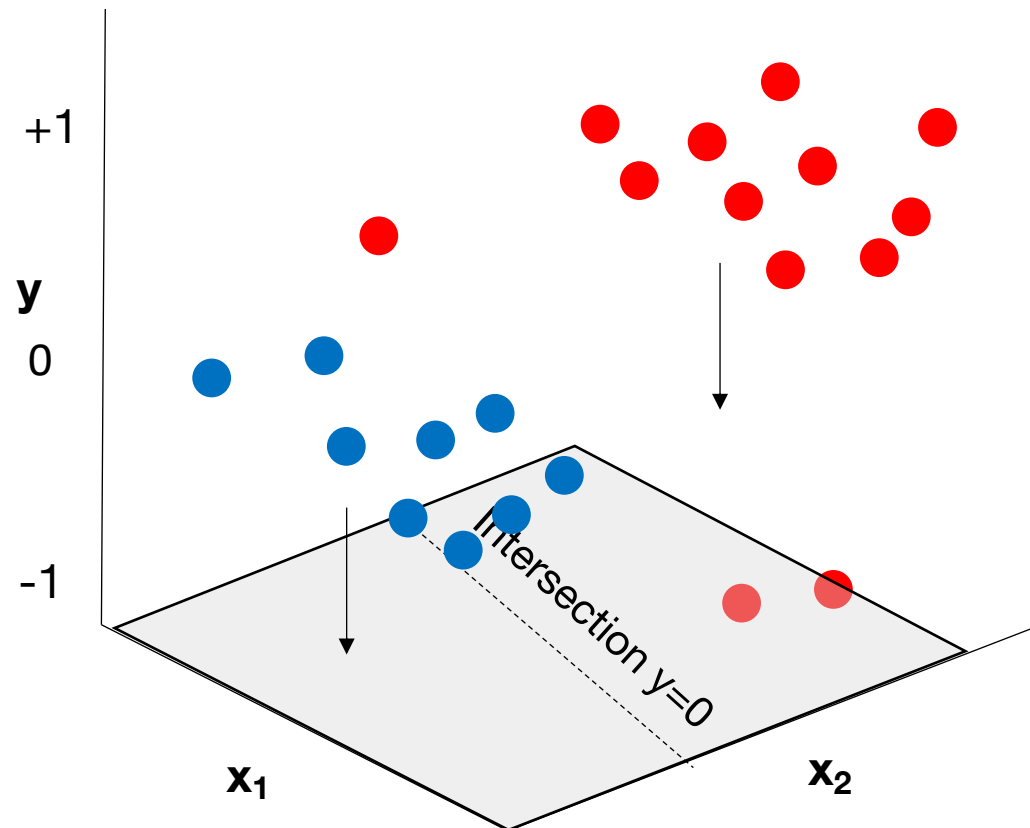
With $\text{sgn}()$ operation:

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

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

- 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:

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

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

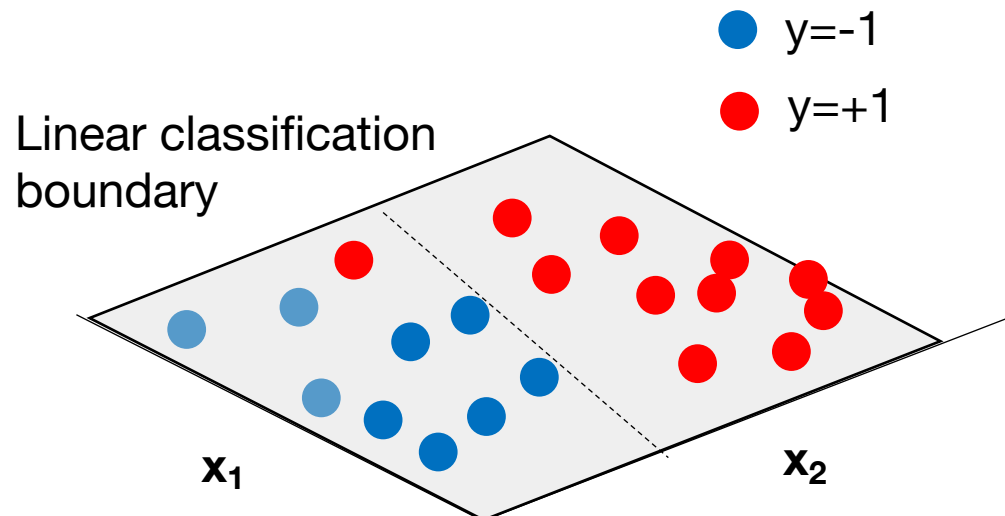
- y axis isn't really needed now & can view this decision boundary in 2D

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

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

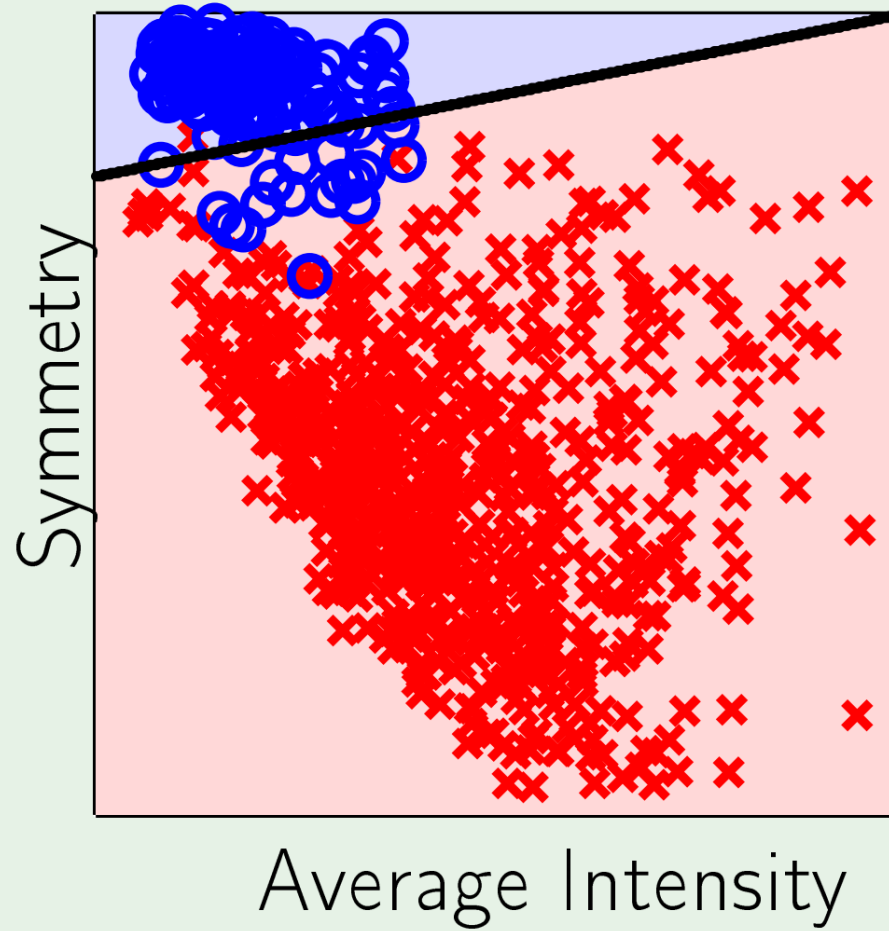
$$f(\mathbf{x}_i) = y_i^* = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$$

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

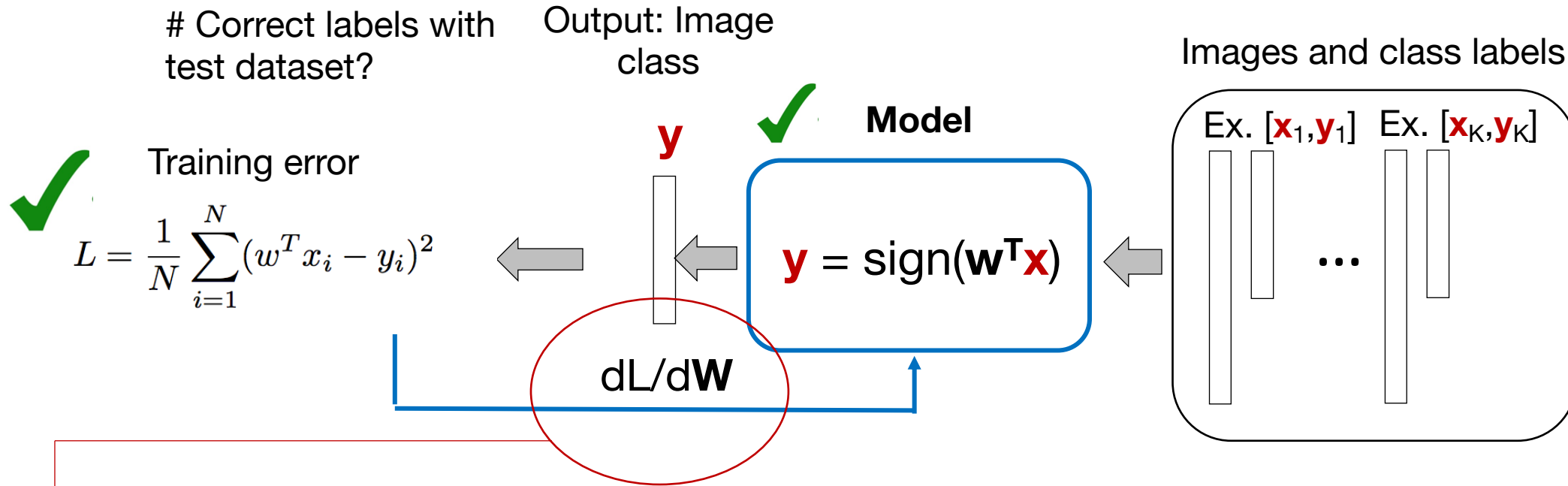


Sign operation takes linear regression and makes it a classification operation!

Linear regression boundary



Example: machine learning for image classification




Let's consider a simple example – image classification. What do we need for training?

1. Labeled examples

✓ 2. A model and loss function

3. A way to minimize the loss function L

3 methods to solve for w^T in the case of linear regression:

- (easier)
1. Pseudo-inverse (this is one of the few cases with a closed-form solution)
 2. Numerical gradient descent
 3. Gradient descent on the cost function with respect to W
- (harder)
- 
- A blue arrow pointing downwards, indicating the increasing difficulty of the methods listed.

3 methods to solve for w^T in the case of linear regression:

(easier)

1. Pseudo-inverse (this is one of the few cases with a closed-form solution)

2. Numerical gradient descent

3. Gradient descent on the cost function with respect to W

(harder)

1. Turning linear regression *for unknown weights W* into a pseudo-inverse:

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

We are multiplying many x_i 's with the same w and are adding them up - let's make a matrix!

1. Turning linear regression for unknown weights W into a pseudo-inverse:

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

We are multiplying many x_i 's with the same w and are adding them up - let's make a matrix!

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

Each training image is 1 row of X

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Each training label is 1 entry of y

1. Turning linear regression for unknown weights W into a pseudo-inverse:

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

We are multiplying many x_i 's with the same w and are adding them up - let's make a matrix!

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

Each training image is 1 row of X

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Each training label is 1 entry of y

$$L = \frac{1}{N} \|Xw - y\|^2$$

This is the same form as the pseudo-inverse we were working with before, but now we want to solve for w

Write this out as a matrix equation:

$$L = \frac{1}{N} \|Xw - y\|^2$$

Note: Training data goes into “dictionary” matrix

$$L = \frac{1}{N} (w^T X^T X w - 2w^T X^T y + yy^T)$$

Write this out as a matrix equation:

$$L = \frac{1}{N} \|Xw - y\|^2$$

Note: Training data goes into “dictionary” matrix

$$L = \frac{1}{N} (w^T X^T X w - 2w^T X^T y + yy^T)$$

Take derivative wrt w and set to 0:

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

Solution is pseudo-inverse:

$$w_o = (X^T X)^{-1} X^T y$$

Write this out as a matrix equation: $L = \frac{1}{N} \|Xw - y\|^2$ **Note: Training data goes into “dictionary” matrix**

$$L = \frac{1}{N} (w^T X^T X w - 2w^T X^T y + yy^T)$$

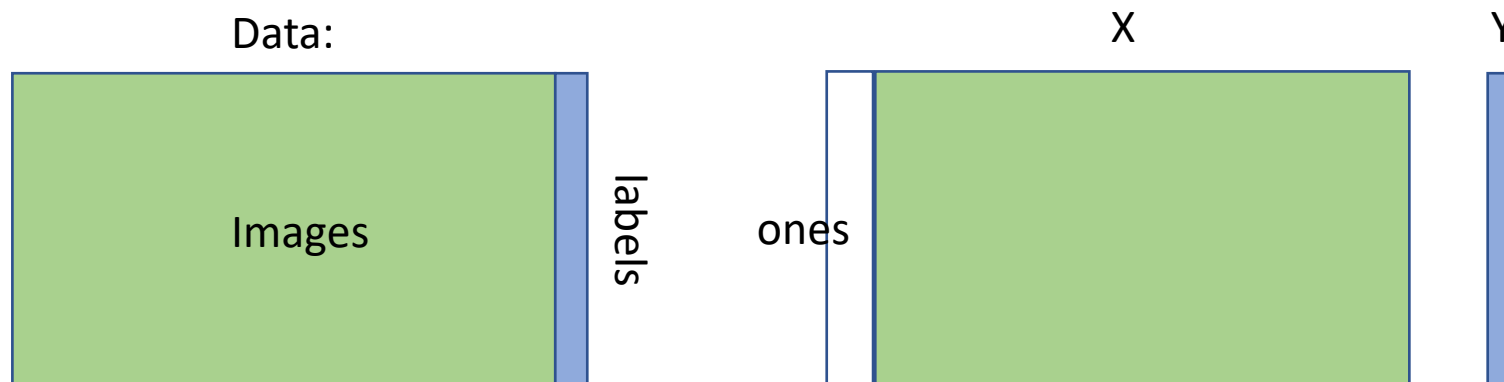
Take derivative wrt w and set to 0: $\nabla L(w) = \frac{2}{N} X^T (Xw - y) = 0$

Solution is pseudo-inverse: $w_o = (X^T X)^{-1} X^T y$

- Steps for Pseudo-inverse:
1. Construct matrix X and vector y from data
 - Each training image is 1 row of X
 - Each training label is 1 entry of y
 2. Compute solution for w_o via above equation

Example pseudo-code

```
data = np.loadtxt('train_data.txt', dtype=int)
X = numpy.zeros((data.shape[0], data.shape[1]-1))
X[:, 0]=1
Y = numpy.zeros((data.shape[0], 1))
for row in m:
    X[row, 1:X.shape[1]-1] = data[row, 0:data.shape[1]:1]
    Y[row] = data[row, data.shape[1]-1]
X_dagger = np.linalg.pinv(X)
w = np.matmul(X_dagger, Y)
```



3 methods to solve for w^T in the case of linear regression:

(easier) 1. Pseudo-inverse (this is one of the few cases with a closed-form solution)

2. Numerical gradient descent

3. Gradient descent on the cost function with respect to W

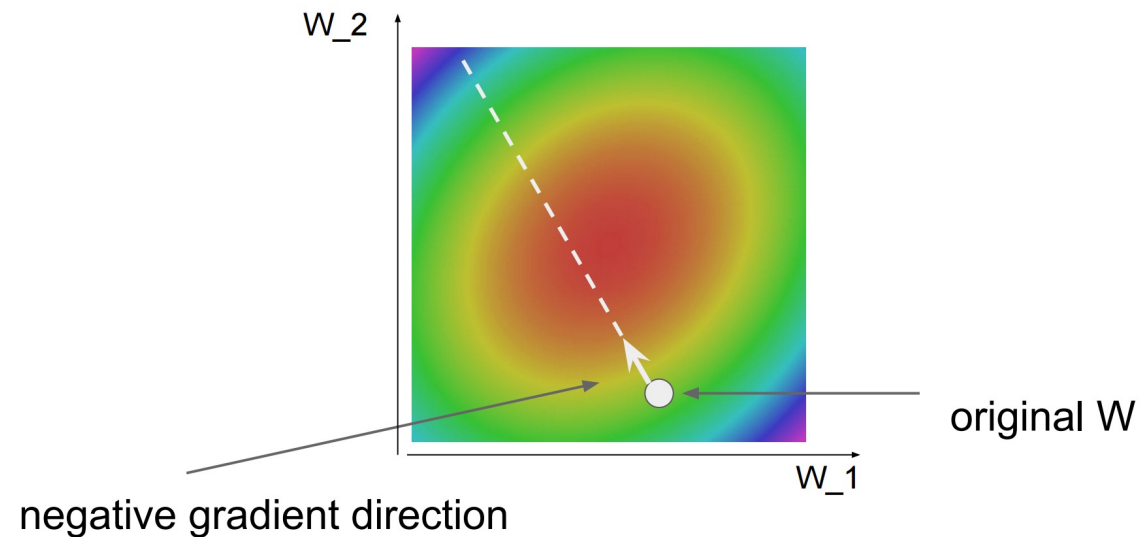
(harder)

Gradient descent: The iterative recipe

Initialize: Start with a guess of \mathbf{W}

Until the gradient does not change very much:
 $dL/d\mathbf{W} = \text{evaluate_gradient}(\mathbf{W}, \mathbf{x}, y, L)$
 $\mathbf{W} = \mathbf{W} - \text{step_size} * dL/d\mathbf{W}$

evaluate_gradient can be achieved numerically or algebraically



Gradient descent: Numerical evaluation example



With a matrix, compute this for each entry:

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

Gradient descent: Numerical evaluation example



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:

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

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

Gradient descent: Numerical evaluation example



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:

$$W = [1,2;3,4]$$
$$L(W, x, y) = 12.79$$

$$W_{1+h} = [1.001,2;3,4]$$
$$L(W_{1+h}, x, y) = 12.8$$

Gradient descent: Numerical evaluation example



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:

$$W = [1,2;3,4]$$
$$L(W, x, y) = 12.79$$

$$W_{1+h} = [1.001,2;3,4]$$
$$L(W_{1+h}, x, y) = 12.8$$

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

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

Gradient descent: Numerical evaluation example



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:

$$W = [1,2;3,4]$$
$$L(W, x, y) = 12.79$$

$$W_{1+h} = [1.001,2;3,4]$$
$$L(W_{1+h}, x, y) = 12.8$$

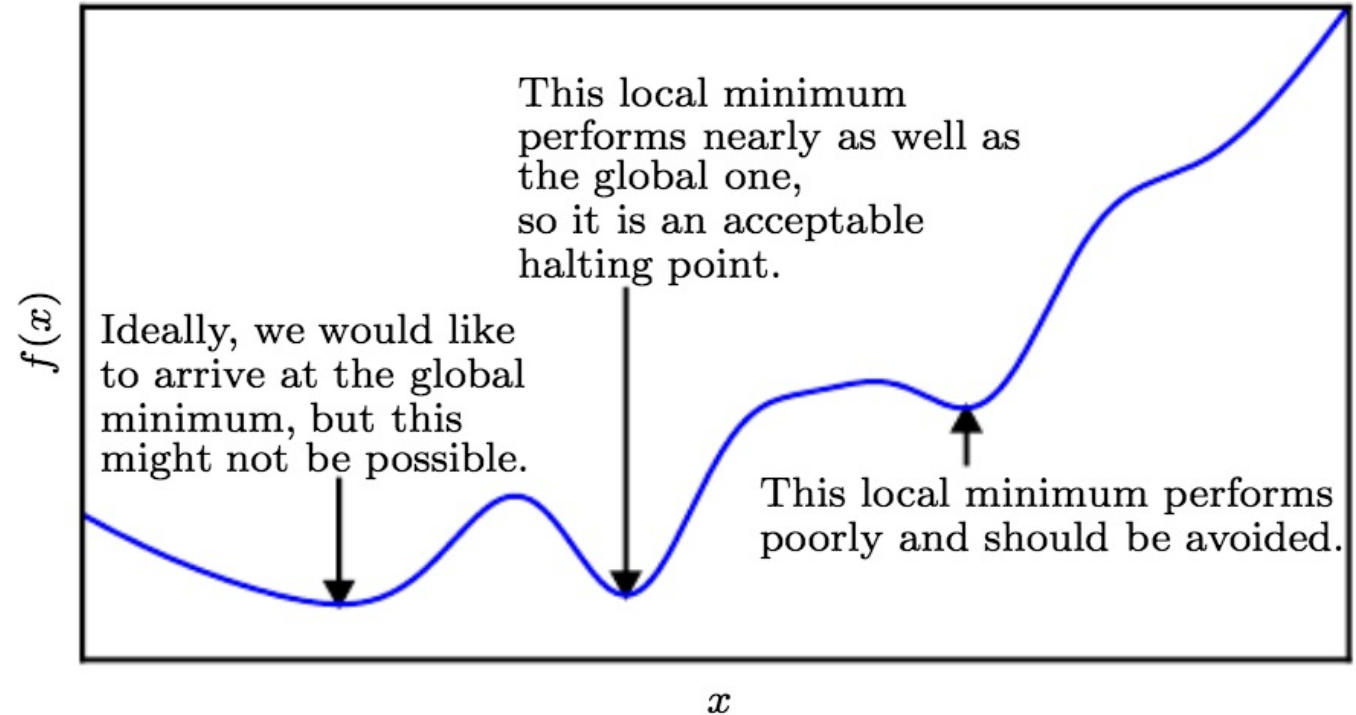
$$dL(W_1)/dW_1 = 12.8 - 12.79 / .001$$

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

- Repeat for all entries of \mathbf{W} , $dL/d\mathbf{W}$ will have $N \times M$ entries for $N \times M$ matrix
- **This is a “brute force” approach – not ideal, but sometimes helpful**

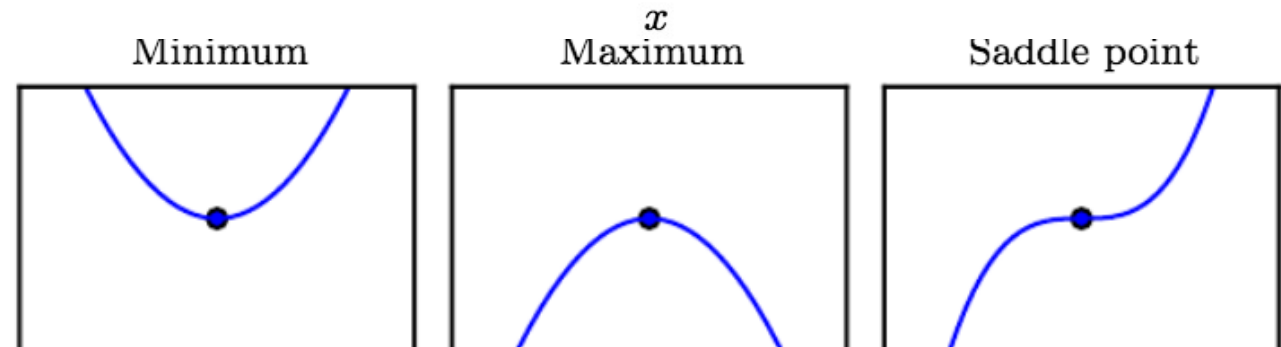
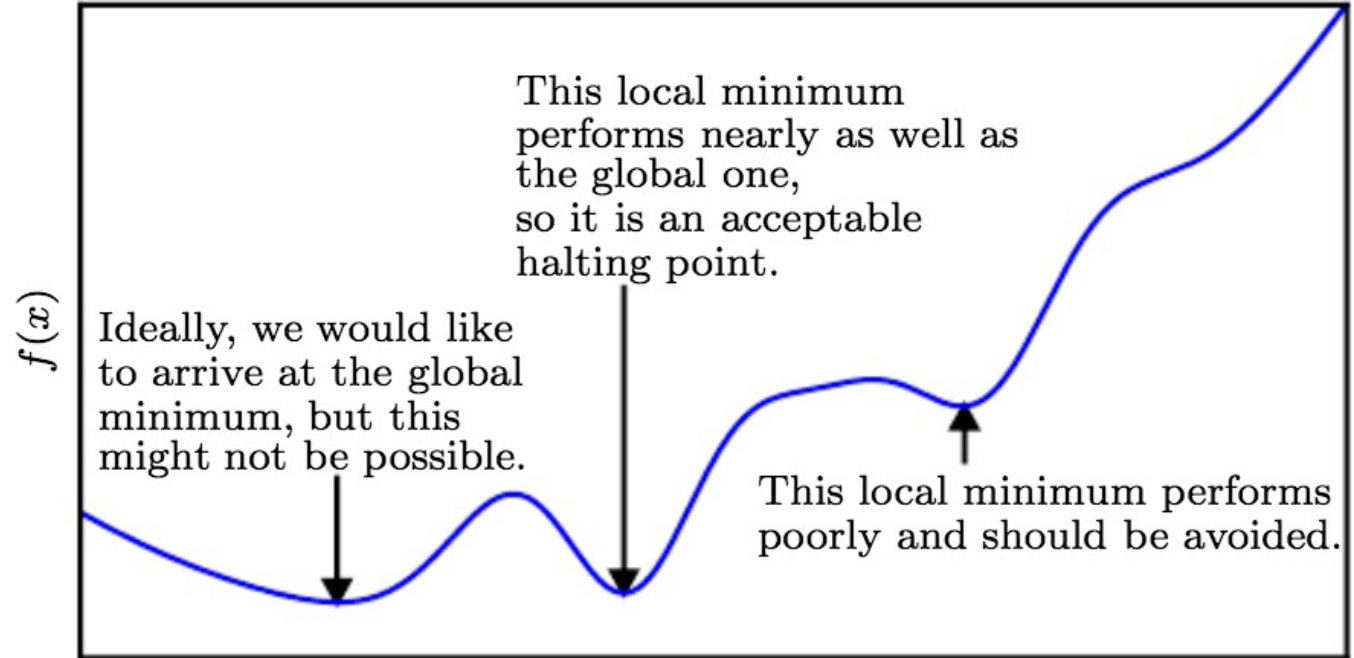
Some quick details about gradient descent

- For non-convex functions, local minima can obscure the search for global minima
- Analyzing critical points (plateaus) of function of interest is important



Some quick details about gradient descent

- For non-convex functions, local minima can obscure the search for global minima
- Analyzing critical points (plateaus) of function of interest is important
- Critical points at $df/dx = 0$
- 2nd derivative d^2f/dx^2 tells us the type of critical point:
 - Minima at $d^2f/dx^2 > 0$
 - Maxima at $d^2f/dx^2 < 0$



Some quick details about gradient descent

Often we'll have functions of m variables

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{e.g., } f(\mathbf{x}) = \sum (\mathbf{Ax}-\mathbf{y})^2)$$

We take partial derivatives $\frac{\partial}{\partial x_i} f(\mathbf{x})$ and put them in **gradient vector** $\mathbf{g} = \nabla_{\mathbf{x}} f(\mathbf{x})$

Some quick details about gradient descent

Often we'll have functions of m variables

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \quad (\text{e.g., } f(\mathbf{x}) = \sum (\mathbf{Ax}-\mathbf{y})^2)$$

We take partial derivatives $\frac{\partial}{\partial x_i} f(\mathbf{x})$ and put them in **gradient vector** $\mathbf{g} = \nabla_x f(\mathbf{x})$

We have many second derivatives: $\mathbf{H}(f)(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$ **Hessian Matrix**

Some quick details about gradient descent

Often we'll have functions of m variables

$$f : \mathbb{R}^m \rightarrow \mathbb{R} \quad (\text{e.g., } f(\mathbf{x}) = \sum (\mathbf{Ax}-\mathbf{y})^2)$$

We take partial derivatives $\frac{\partial}{\partial x_i} f(\mathbf{x})$ and put them in **gradient vector** $\mathbf{g} = \nabla_{\mathbf{x}} f(\mathbf{x})$

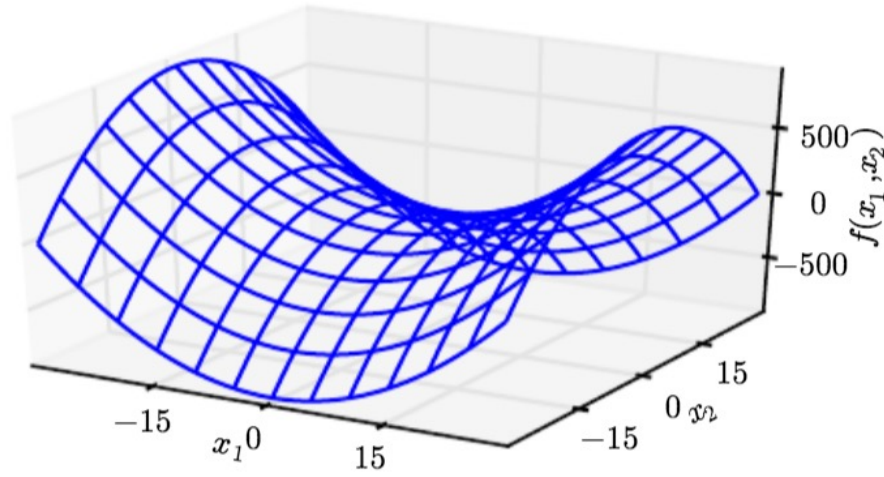
We have many second derivatives: $\mathbf{H}(f)(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$ **Hessian Matrix**

In general, we'll have functions that map m variables to n variables

$$\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n \quad (\text{e.g., } \mathbf{f}(\mathbf{x}) = \mathbf{Wx}, \mathbf{W} \text{ is } n \times m)$$

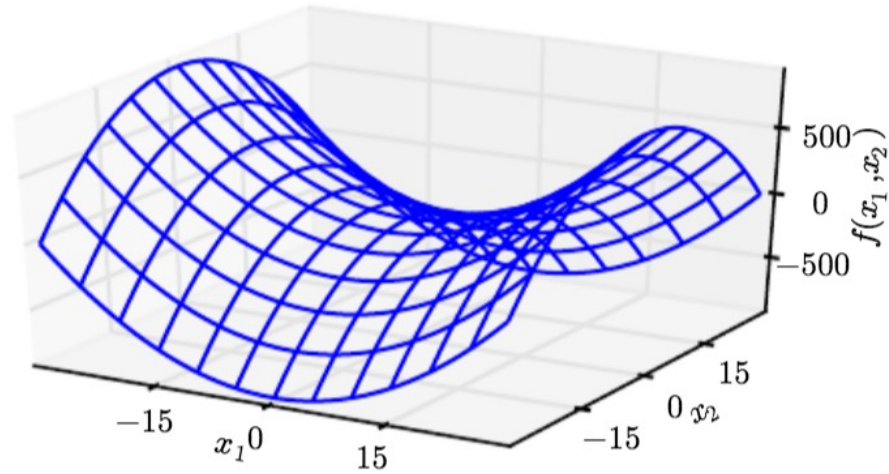
$\mathbf{J} \in \mathbb{R}^{n \times m}$ of \mathbf{f} : $J_{i,j} = \frac{\partial}{\partial x_j} f(\mathbf{x})_i$ **Jacobian Matrix**

Quick example



$$f(\mathbf{x}) = x_1^2 - x_2^2$$

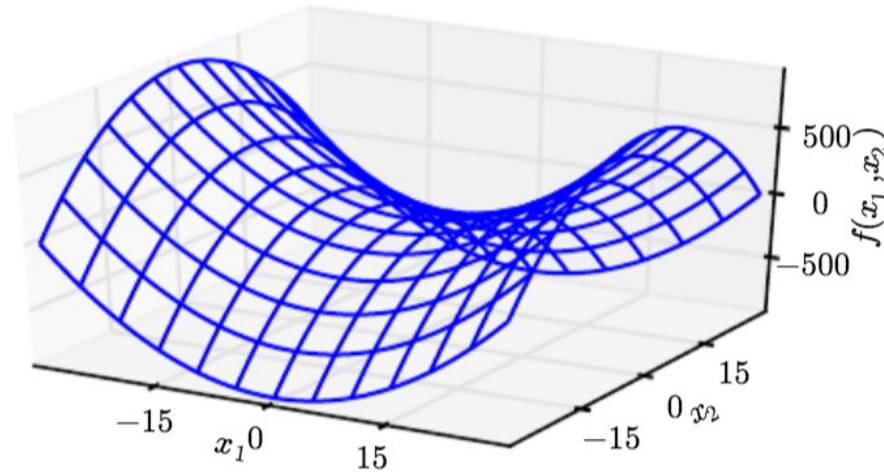
Quick example



$$f(\mathbf{x}) = x_1^2 - x_2^2$$

$$\mathbf{g} = \begin{bmatrix} 2x_1 \\ -2x_2 \end{bmatrix}$$

Quick example



$$f(\mathbf{x}) = x_1^2 - x_2^2$$

$$\mathbf{g} = \begin{bmatrix} 2x_1 \\ -2x_2 \end{bmatrix}$$

$$\mathbf{H} = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

- Convex functions have positive semi-definite Hessians (Trace ≥ 0)
- Trace/eigenvalues of Hessian are useful evaluate critical points & guide optimization

Steepest descent and the best step size ϵ

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 $|\mathbf{x}^{(n+1)} - \mathbf{x}^{(n)}| < \text{threshold } t$

```
while previous_step_size > precision and iters < max_iters:  
    prev_x = cur_x  
    cur_x -= epsilon * df(prev_x)  
    previous_step_size = abs(cur_x - prev_x)  
    **Update epsilon - see next slide  
    iters+=1
```

Steepest descent and the best step size ϵ

We computed this – computers can too in interesting ways

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 -= epsilon * df(prev_x)  
    previous_step_size = abs(cur_x - prev_x)  
    **Update epsilon - see next slide  
    iters+=1
```

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

To find out, take 2nd order Taylor expansion of f (a good approx. for nearby points):

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

To find out, take 2nd order Taylor expansion of f (a good approx. for nearby points):

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

Then, evaluate at the next step:

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

Steepest descent and the best step size ϵ

What is a good step size $\epsilon^{(n)}$?

To find out, take 2nd order Taylor expansion of f (a good approx. for nearby points):

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

Then, evaluate at the next step:

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

Solve for optimal step (when Hessian is positive):

$$\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}$$

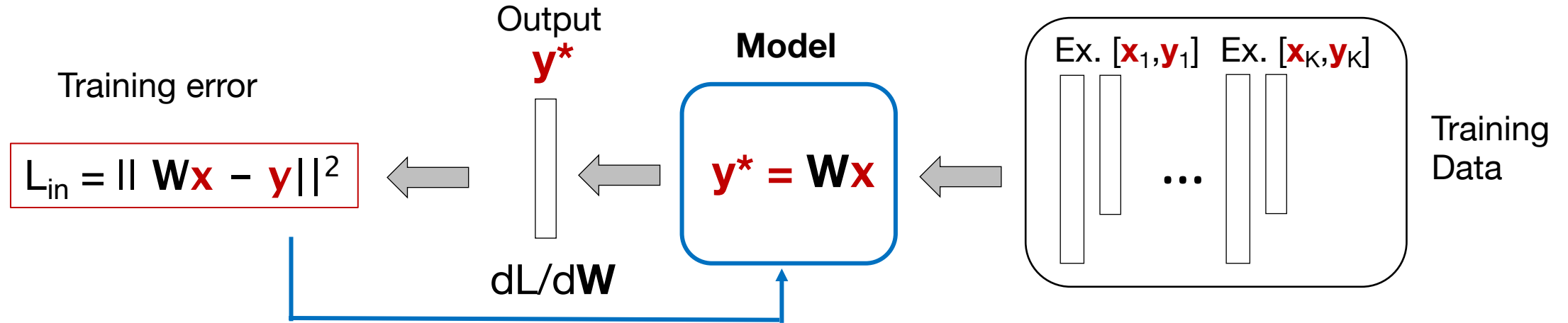
J. R. Shewchuck, [“An Introduction to the Conjugate Gradient Method Without the Agonizing Pain”](#)

3 methods to solve for w^T in the case of linear regression:

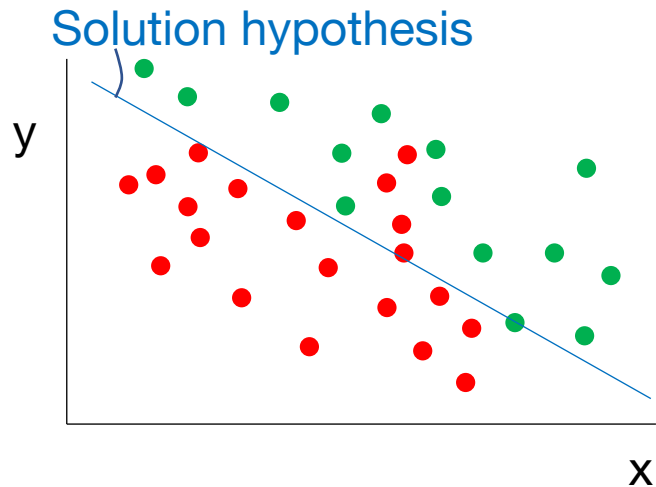
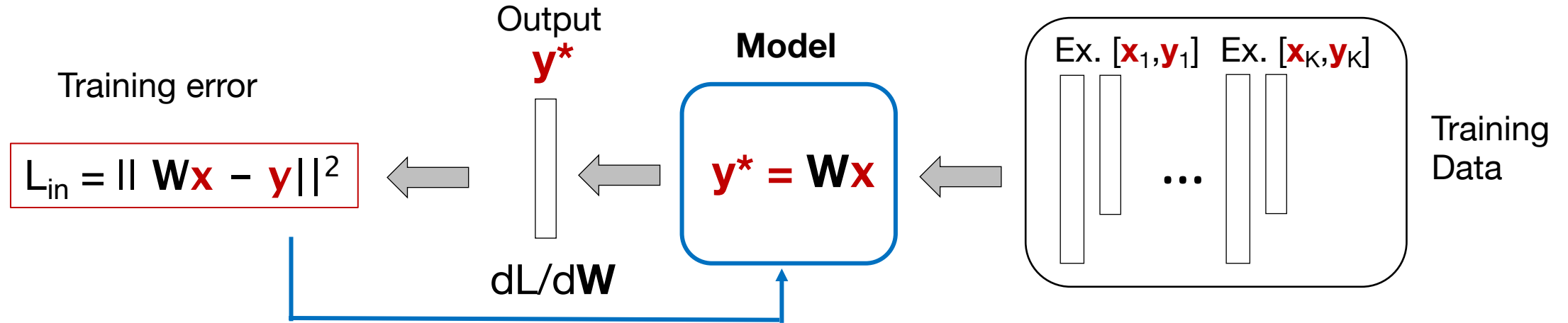
- (easier)
1. Pseudo-inverse (this is one of the few cases with a closed-form solution)
 2. Numerical gradient descent
 3. Gradient descent on the cost function with respect to W
- (harder)
- A blue arrow points downwards from the word '(easier)' to the word '(harder)', indicating the relative difficulty of the methods listed.

Next : We'll understand why linear regression doesn't work so well, and extend things beyond this simple starting point

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

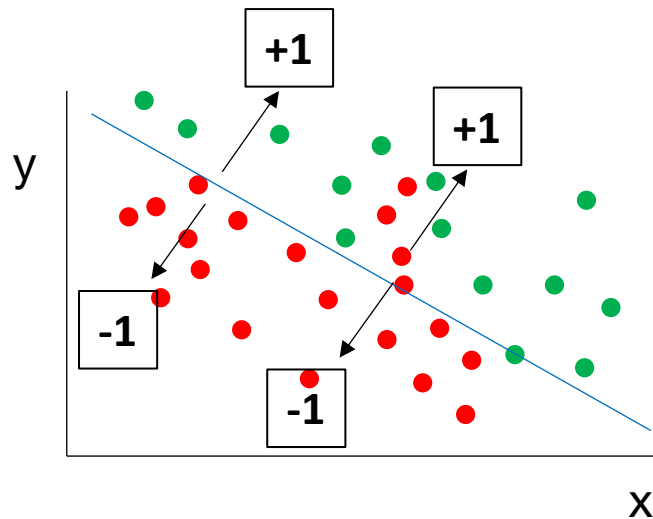
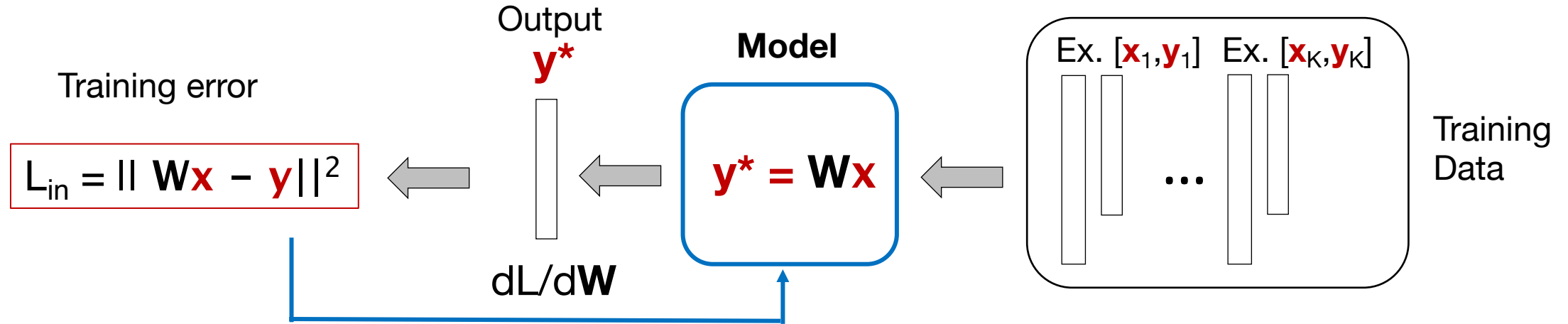


The linear classification model – what’s not to like?



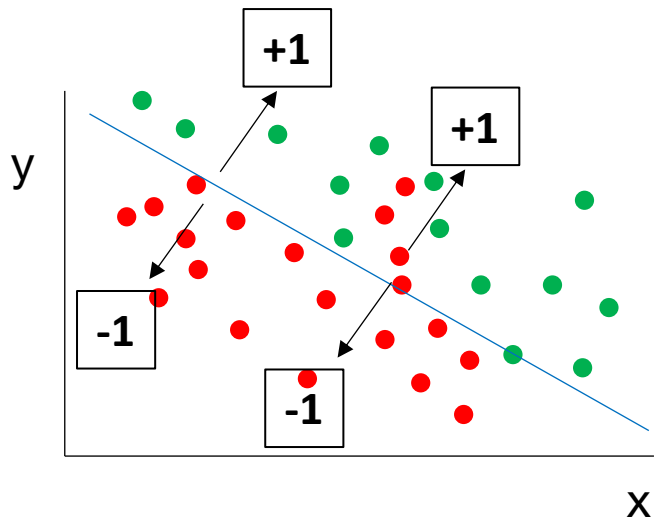
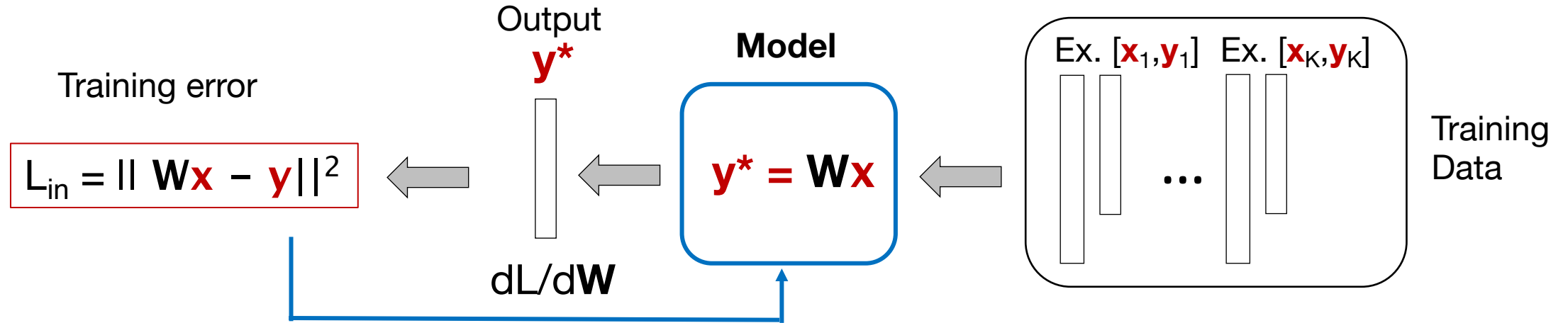
1. Can only separate data with lines (hyper-planes)...

The linear classification model – what’s not to like?



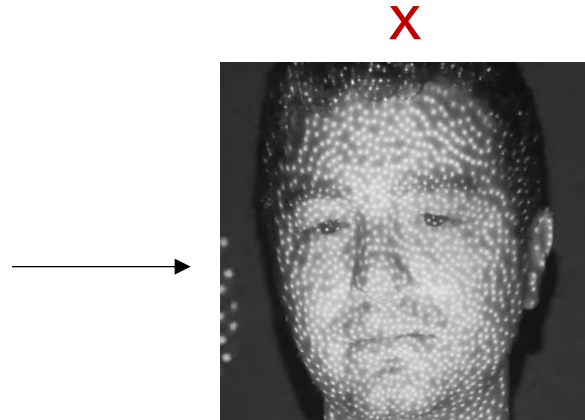
1. Can only separate data with lines (hyper-planes)...
2. We only allowed for binary labels ($y = +/- 1$)

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

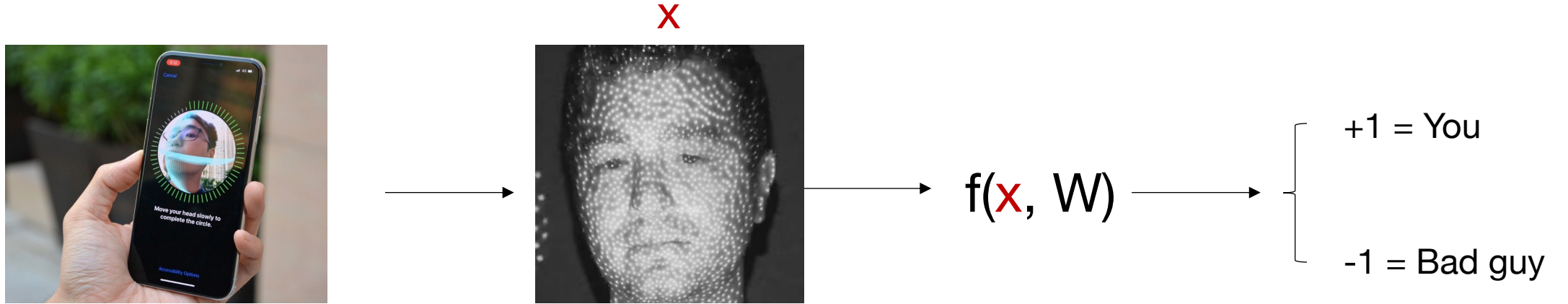
Cost functions matter: a simple example



$$f(x, W)$$

- +1 = You
- 1 = Bad guy

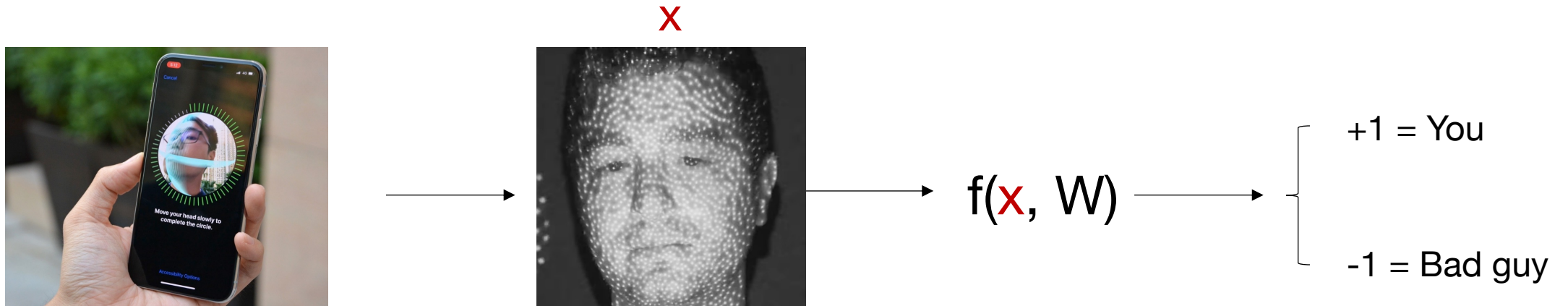
Cost functions matter: a simple example



Two types of error: false accept and false reject

		$f(x, W)$	
		+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



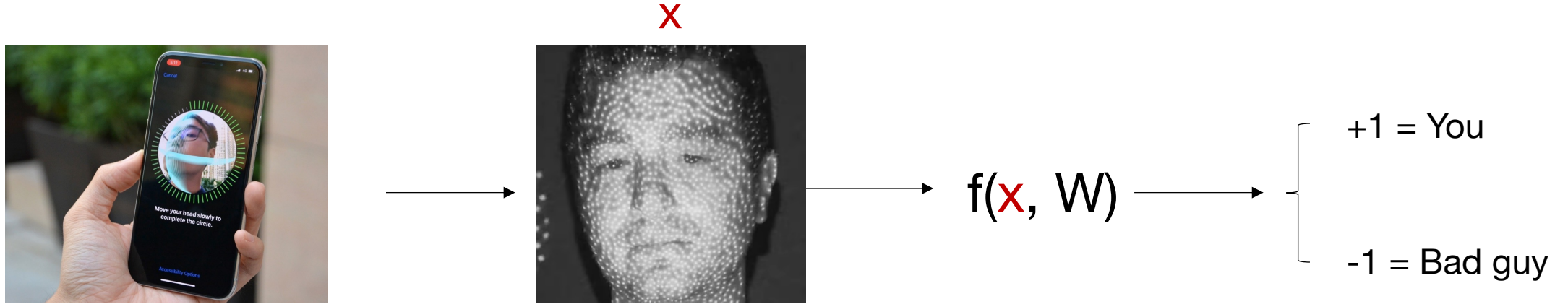
Two types of error: false accept and false reject

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

		$f(x, W)$		
		+1	-1	
y	+1	No Error	False reject	It's you, but you can't get in...
	-1	False accept	No Error	

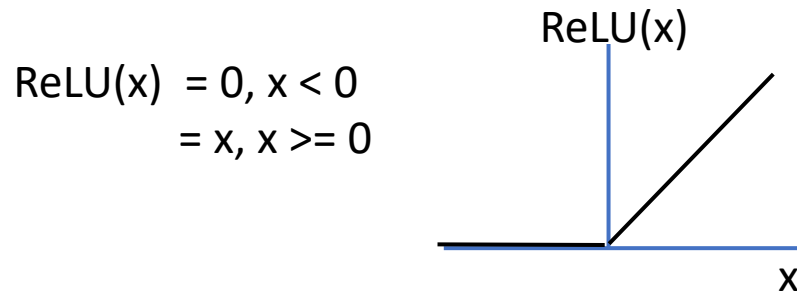
Letting an intruder in

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?



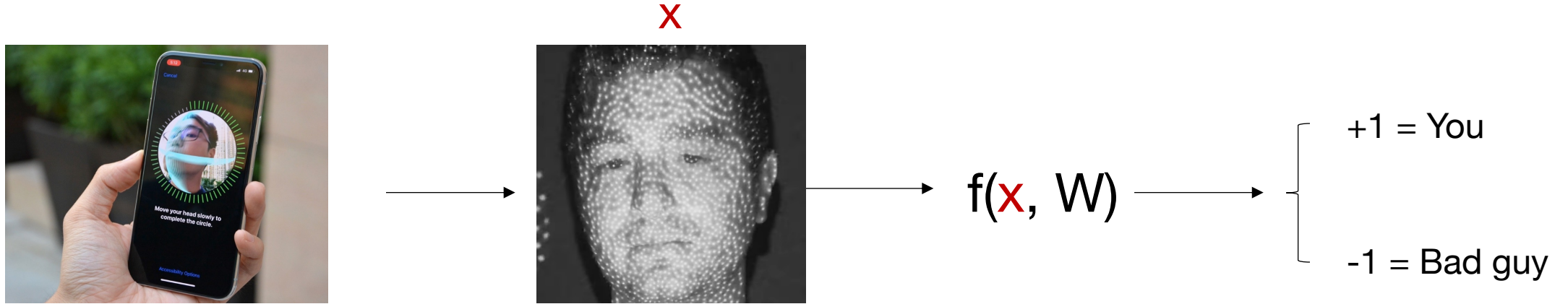
$f(x, W)$

	+1	-1	
+1	No Error	False reject	It's you, but you can't get in...
-1	False accept	No Error	

y

Letting an intruder in

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?

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

Penalty for intruder

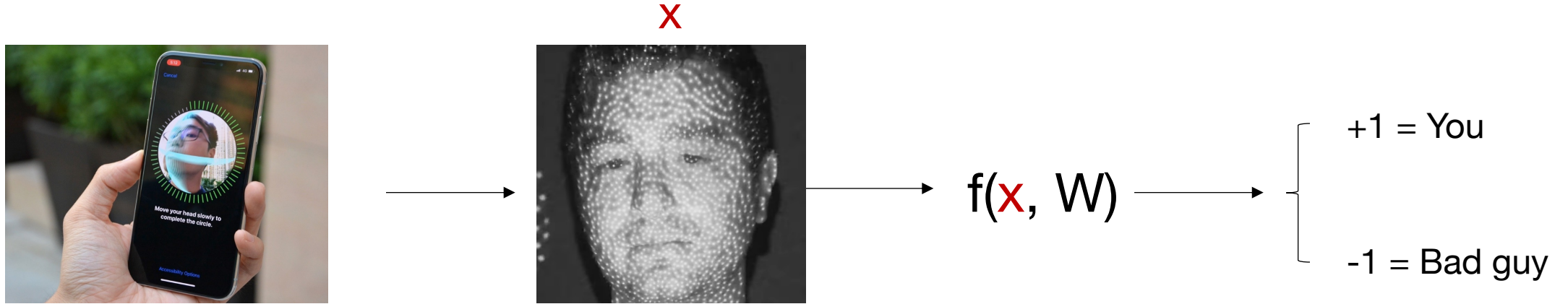
Large penalty for annoyance...

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

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

Letting an intruder in

Cost functions matter: a simple example

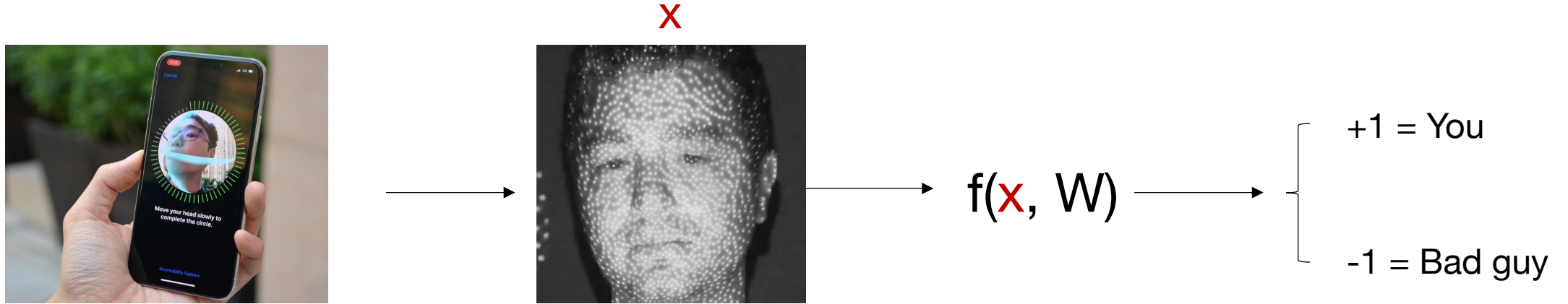


What if you're a CIA agent?

		$f(x, W)$		
		+1	-1	
y	+1	No Error	False reject	It's you, but you can't get in...
	-1	False accept	No Error	

Letting an intruder in

Cost functions matter: a simple example



What if you're a CIA agent?

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

BIG penalty for intruder

Don't mind about annoyance...

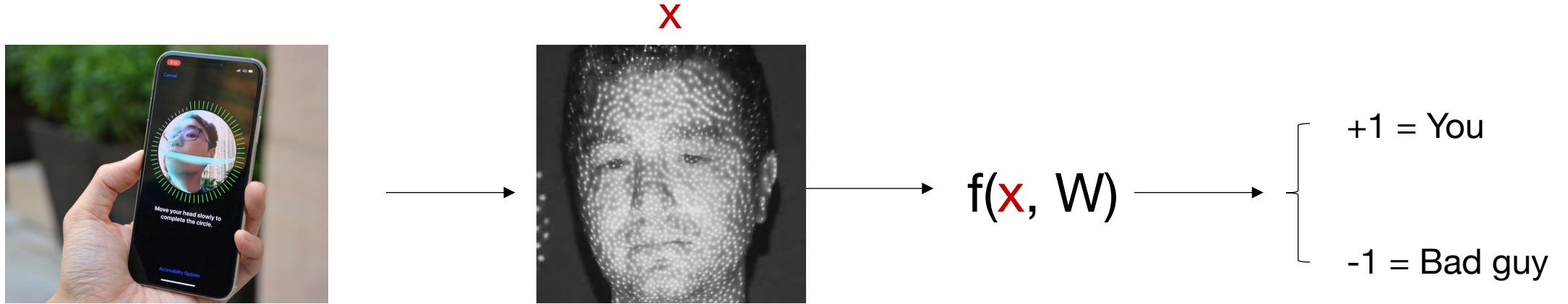
$f(x, W)$

	+1	-1	
+1	No Error	False reject	It's you, but you can't get in...
-1	False accept	No Error	

y

Letting an intruder in

Cost functions matter: a simple example



Establishing cost function tied to conditional probabilities:

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

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

Establish L, W to balance and minimize these probabilities

		$f(x, W)$		
		+1	-1	
y	+1	No Error	False reject	It's you, but you can't get in...
	-1	False accept	No Error	

Letting an intruder in