# Machine Learning and Imaging

## BME 548L
Roarke Horstmeyer

Lecture 7: Gradient descent and going beyond linear classification

# Summary of machine learning pipeline:
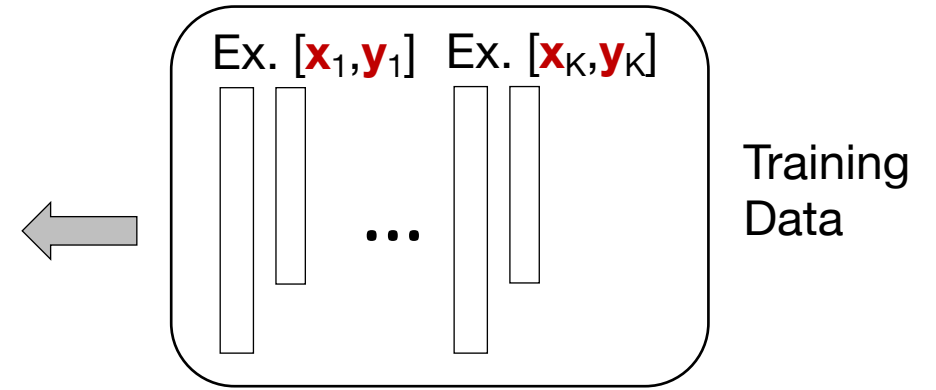
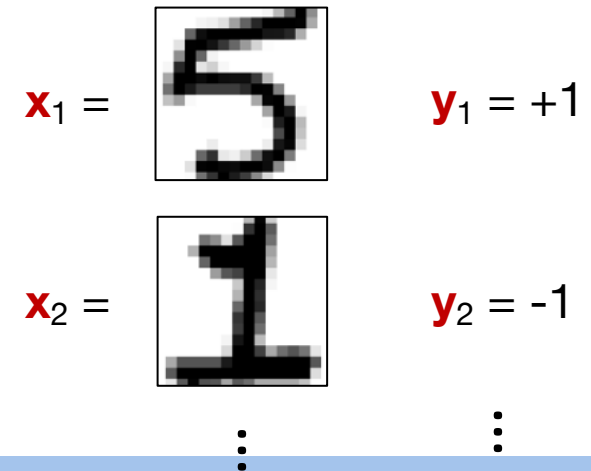## 1. Network Training

What we need for network training:

**1. Labeled examples**

# Summary of machine learning pipeline:

## 1. Network Training

Ex. [$\mathbf{x}_1$,$\mathbf{y}_1$]  Ex. [$\mathbf{x}_K$,$\mathbf{y}_K$]

...

Training
Data

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

What we need for network training:

**1. Labeled examples**

$\mathbf{x}_1 =$      $\mathbf{y}_1 = +1$

$\mathbf{x}_2 =$      $\mathbf{y}_2 = -1$

# Summary of machine learning pipeline:

## 1. Network Training



Output
$\mathbf{y^*}$

**Model**

$f(\mathbf{W,x})$

Ex. $[\mathbf{x}_1,\mathbf{y}_1]$   Ex. $[\mathbf{x}_K,\mathbf{y}_K]$
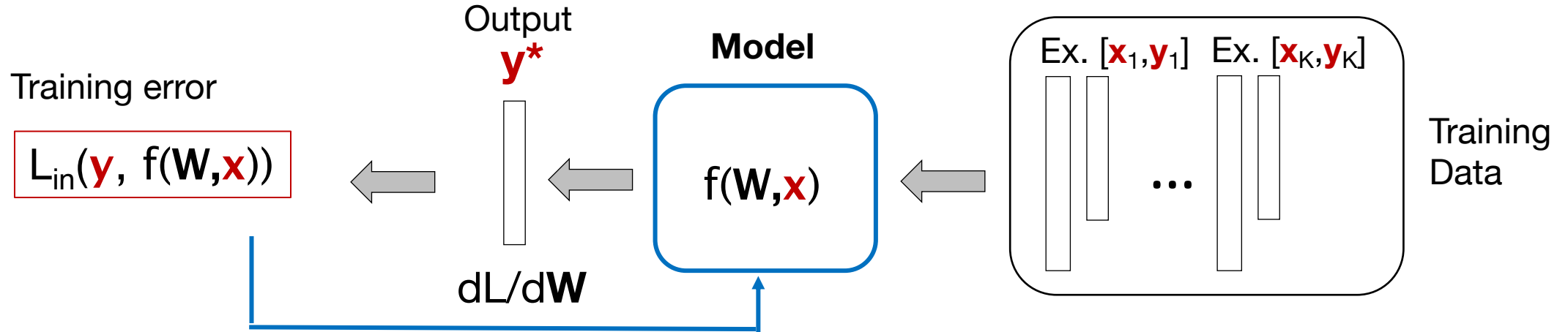
...

Training
Data

What we need for network training:

**1. Labeled examples**

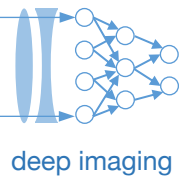**2. A model and loss function**

# Summary of machine learning pipeline:

## 1. Network Training



Training error

$$L_{in}(\mathbf{y}, f(\mathbf{W},\mathbf{x}))$$

Output $\mathbf{y^*}$

**Model**

$f(\mathbf{W},\mathbf{x})$

dL/d**W**

Ex. $[\mathbf{x}_1,\mathbf{y}_1]$   Ex. $[\mathbf{x}_K,\mathbf{y}_K]$

...

Training Data

What we need for network training:

**1. Labeled examples**

**2. A model and loss function**

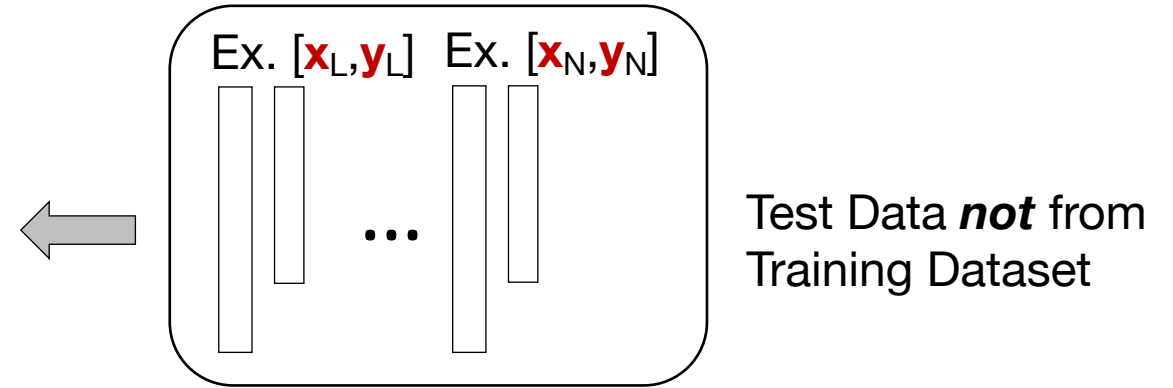**3. A way to minimize the loss function L**

**Summary of machine learning pipeline:**

2. Network Testing

What we need for network testing:

# Summary of machine learning pipeline:
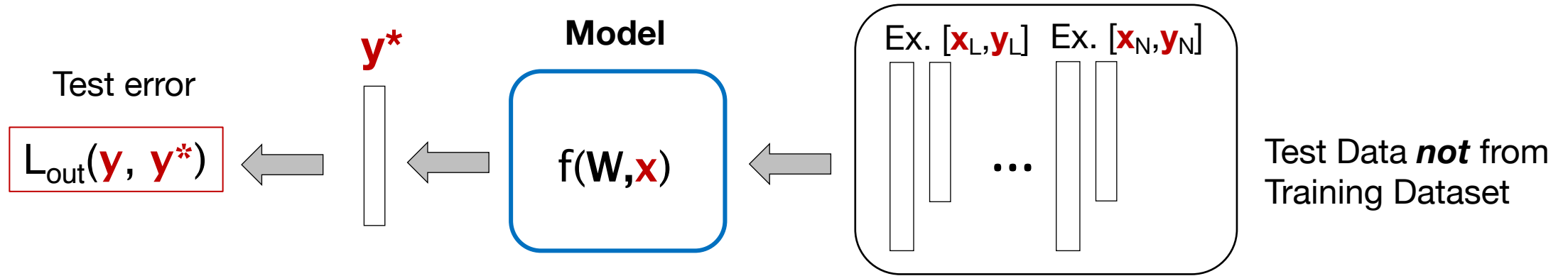
## 2. Network Testing

Ex. [$\mathbf{x}_L$,$\mathbf{y}_L$]   Ex. [$\mathbf{x}_N$,$\mathbf{y}_N$]

...

Test Data **not** from Training Dataset

What we need for network testing:

**4. *Unique* labeled test data**

# Summary of machine learning pipeline:

## 2. Network Testing

**Model**

$y^*$

Test error

$$L_{out}(y, y^*)$$

$f(W,x)$

Ex. $[x_L, y_L]$    Ex. $[x_N, y_N]$

...

Test Data *not* from Training Dataset

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
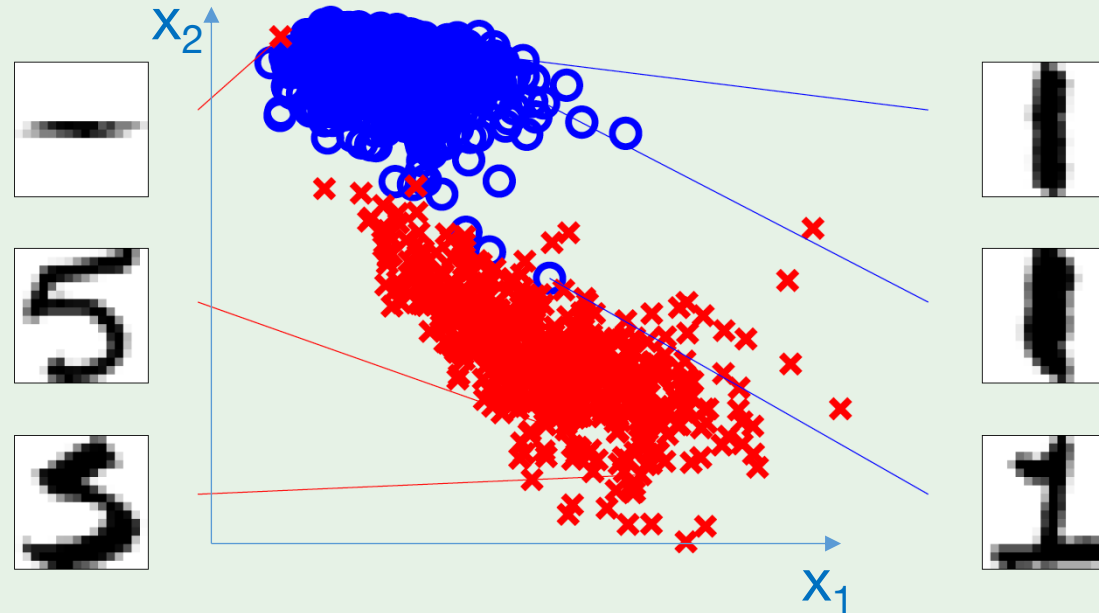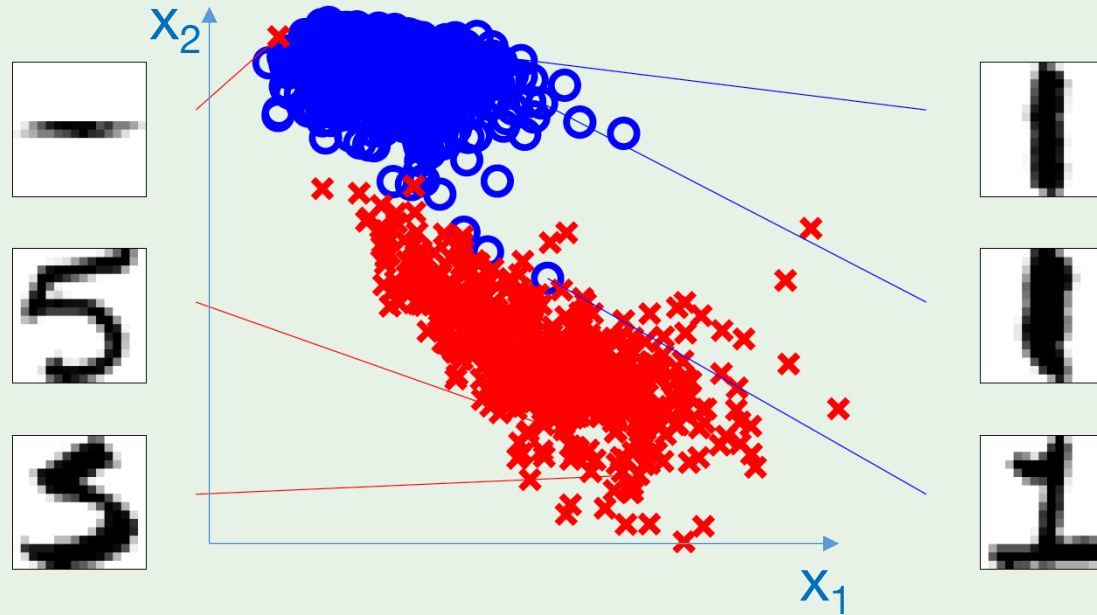
Illustration of features

$\mathbf{x} = (x_0, x_1, x_2)$    $x_1$: intensity    $x_2$: symmetry
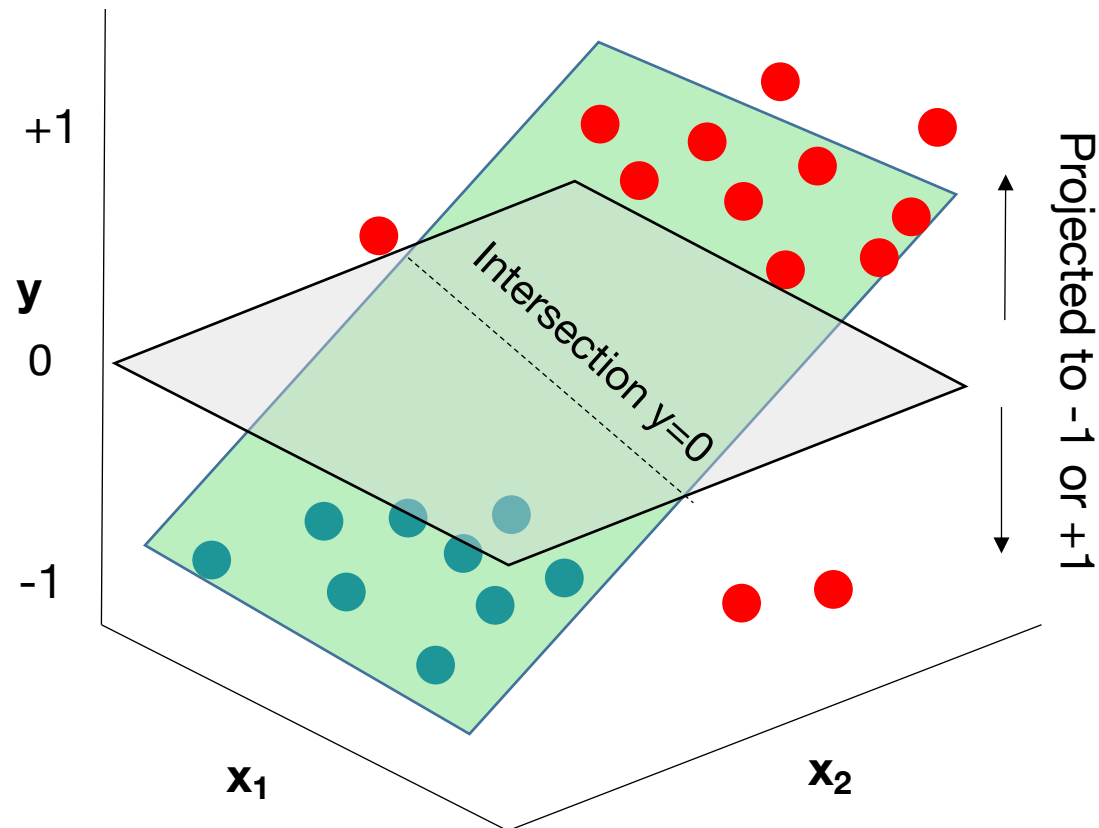
**Linear classification:**

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 sgn() achieve classification?



With sgn() operation:

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

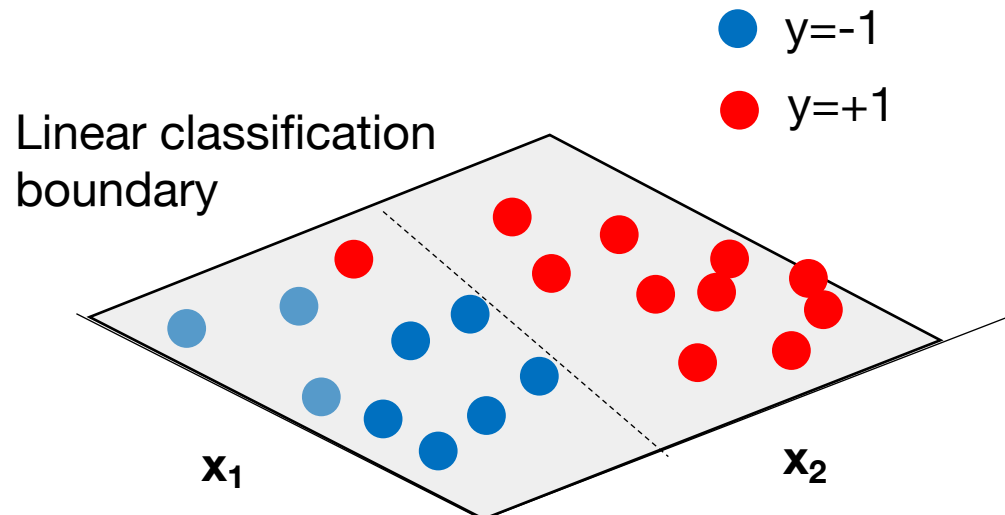$$L = \frac{1}{N}\sum_{i=1}^{N}(w^Tx_i - y_i)^2$$

- 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 sgn() achieve classification?

With 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=-1

● y=+1

Linear classification boundary

$\mathbf{x_1}$

$\mathbf{x_2}$

Closed-form solution available for this boundary line **w** via pseudo-inverse (see last lecture's notes)

Let's consider some other strategies to solve for **w**...

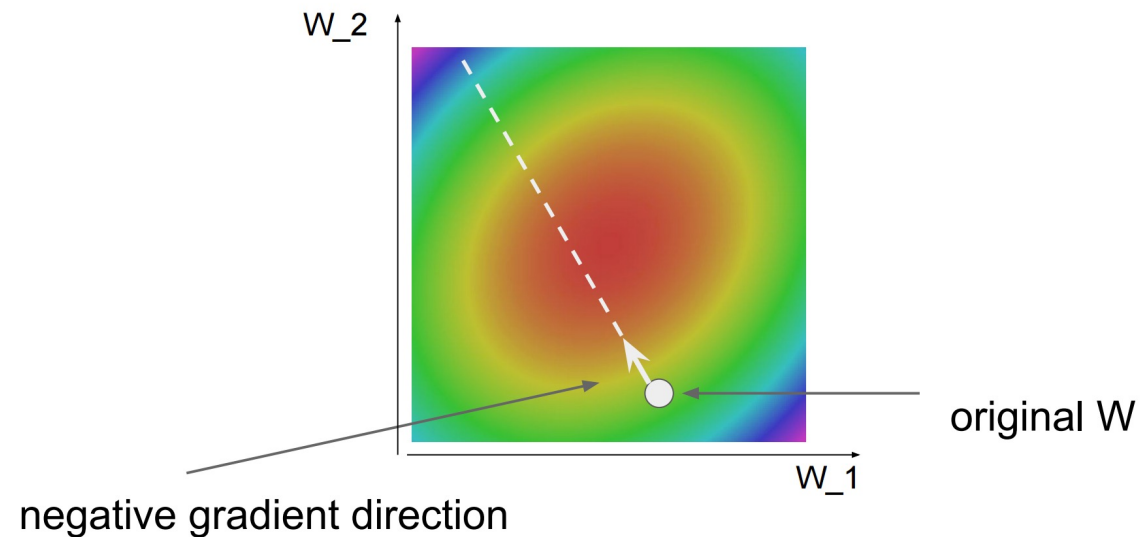**3 methods to solve for w$^T$ in the case of linear regression:**

(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 **W**

Until the gradient does not change very much:
      dL/d**W** = evaluate_gradient(**W**, **x** ,y ,L)
      **W** = **W** – step_size * dL/d**W**

evaluate_gradient can
be achieved numerically
or algebraically



W_2

W_1

original W

negative gradient direction

# 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 $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \varepsilon^{(0)}\mathbf{g}^{(0)}$

4. Repeat: $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \varepsilon^{(n)}\mathbf{g}^{(n)}$, until $|\mathbf{x}^{(n+1)}-\mathbf{x}^{(n)}| <$ 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

deep imaging

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 $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \varepsilon^{(0)}\mathbf{g}^{(0)}$

4. Repeat: $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \varepsilon^{(n)}\mathbf{g}^{(n)}$, until $|\mathbf{x}^{(n+1)}-\mathbf{x}^{(n)}| <$ 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 $\boldsymbol{\varepsilon}^{(n)}$?

# Steepest descent and the best step size ε

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

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

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

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

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

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

Then, evaluate at the next step:

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

# Steepest descent and the best step size ε

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

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

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

Then, evaluate at the next step:

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

Solve for optimal step (when Hessian is positive): $\boxed{\epsilon^* = \dfrac{\boldsymbol{g}^\top \boldsymbol{g}}{\boldsymbol{g}^\top \boldsymbol{H} \boldsymbol{g}}}$

J. R. Shewchuck, "An Introduction to the Conjugate Gradient Method Without the Agonizing Pain"
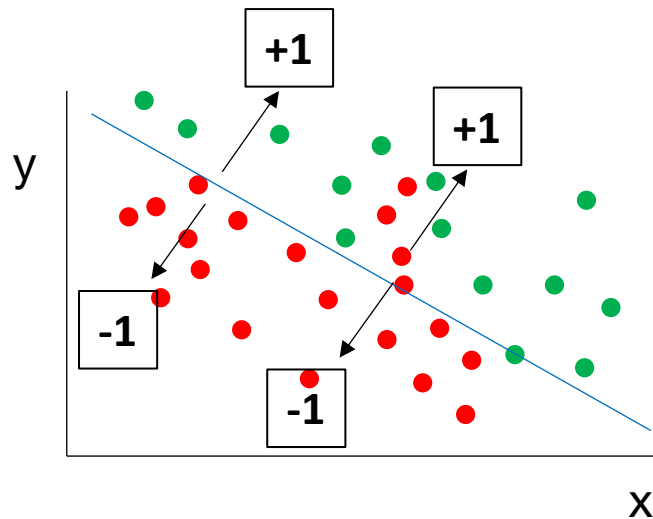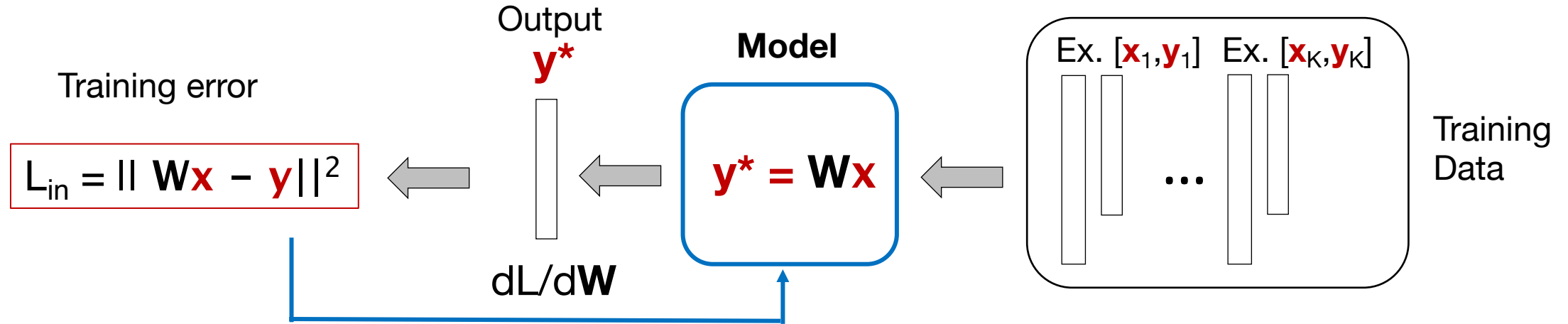
# The linear classification model – what's not to like?



Training error

$$L_{in} = \| \mathbf{W}\mathbf{x} - \mathbf{y} \|^2$$

Output
$\mathbf{y}^*$

dL/d$\mathbf{W}$

**Model**

$$\mathbf{y}^* = \mathbf{W}\mathbf{x}$$

Ex. $[\mathbf{x}_1, \mathbf{y}_1]$  Ex. $[\mathbf{x}_K, \mathbf{y}_K]$

...

Training Data

deep imaging

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



Training error

$$L_{in} = || \mathbf{W}\mathbf{x} - \mathbf{y} ||^2$$

Output $\mathbf{y}^*$

dL/d$\mathbf{W}$

**Model**

$$\mathbf{y}^* = \mathbf{W}\mathbf{x}$$

Ex. [$\mathbf{x}_1$,$\mathbf{y}_1$]   Ex. [$\mathbf{x}_K$,$\mathbf{y}_K$]

...

Training Data

Solution hypothesis

y

x

1. Can only separate data with lines (hyper-planes)…
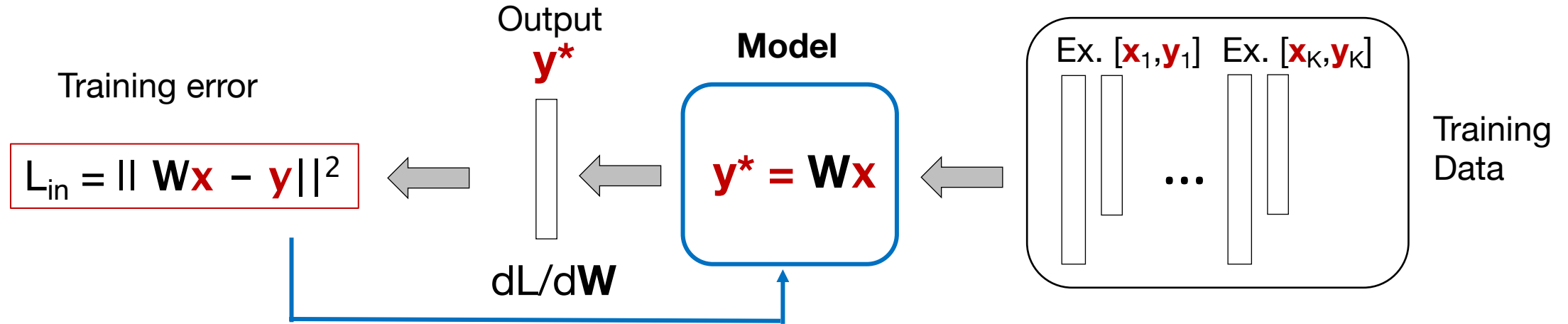
# The linear classification model – what's not to like?

Training error

$$L_{in} = || \mathbf{W}\mathbf{x} - \mathbf{y}||^2$$

Output
$\mathbf{y^*}$

**Model**

$$\mathbf{y^* = Wx}$$

dL/d**W**

Ex. [$\mathbf{x}_1$,$\mathbf{y}_1$]  Ex. [$\mathbf{x}_K$,$\mathbf{y}_K$]
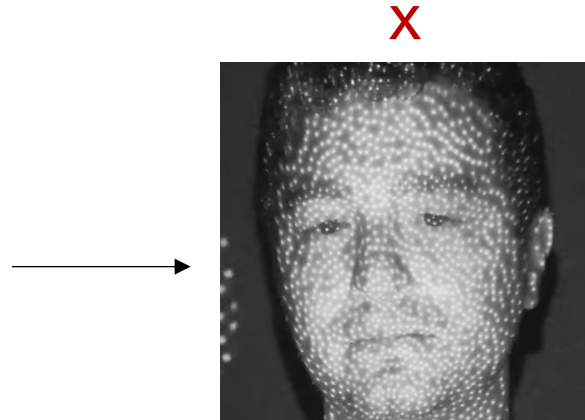
...

Training
Data



+1

+1

-1

-1

y

x

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

2. We only allowed for binary labels (y = +/- 1)

deep imaging

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

Training error

$$L_{in} = || \mathbf{W}\mathbf{x} - \mathbf{y}||^2$$

Output
$\mathbf{y}^*$

dL/d**W**

**Model**

$$\mathbf{y}^* = \mathbf{W}\mathbf{x}$$

Ex. [$\mathbf{x}_1$,$\mathbf{y}_1$]   Ex. [$\mathbf{x}_K$,$\mathbf{y}_K$]

...

Training Data



+1

+1

-1

-1

y

x

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

deep imaging

# Cost functions matter: a simple example


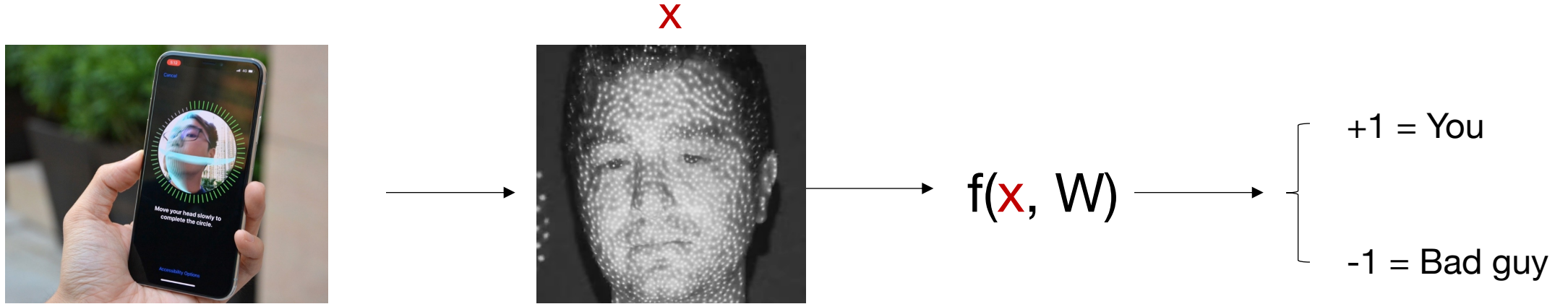
X

$f(x, W)$

+1 = You

-1 = Bad guy

# Cost functions matter: a simple example

x

f(x, W)

+1 = You

-1 = Bad guy

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

x

$f(x, W)$

+1 = You

-1 = Bad guy

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 |
| **-1** | False accept | No Error |

y

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

Letting an intruder in

# Cost functions matter: a simple example

deep imaging

x

$f(x, W)$

+1 = You

-1 = Bad guy

<u>Two types of error</u>: false accept and false reject

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

ReLU(x)

ReLU(x)  = 0, x < 0
         = x, x >= 0

x

$f(x, W)$

|  | +1 | -1 |
|---|---|---|
| **+1** | No Error | False reject |
| **-1** | False accept | No Error |

y

It's you, but you can't get in…

Letting an intruder in

# Cost functions matter: a simple example

x

f(x, W)

+1 = You

-1 = Bad guy

Two types of error: false accept and false reject

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

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

Penalty for intruder

Large penalty for annoyance…
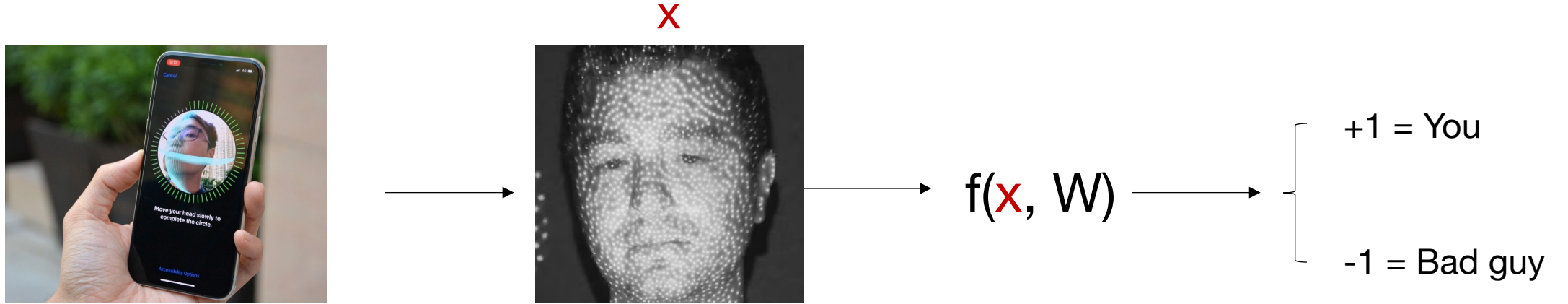
f(x, W)

|  | +1 | -1 |
|---|---|---|
| **+1** | No Error | False reject |
| **-1** | False accept | No Error |

y

It's you, but you can't get in…

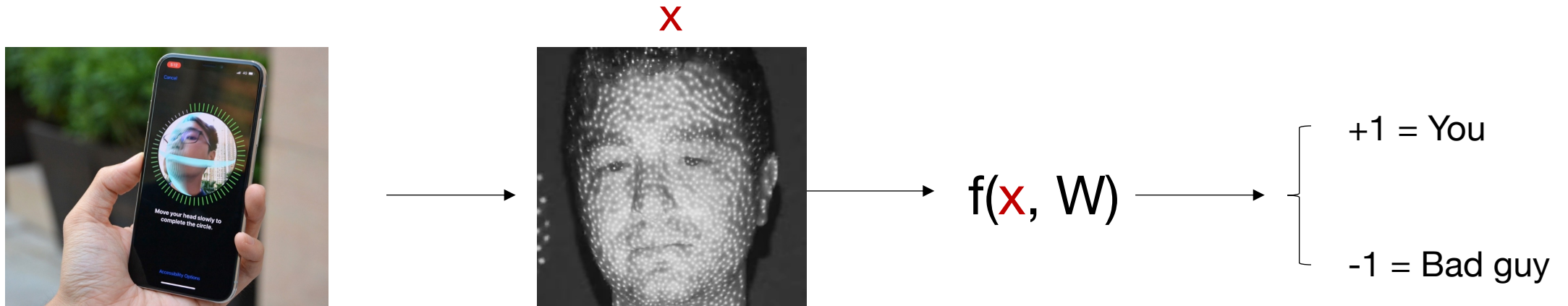Letting an intruder in

# Cost functions matter: a simple example

x

f(x, W)

+1 = You

-1 = Bad guy

What if you're a CIA agent?

f(x, W)

|  | +1 | -1 |
|---|---|---|
| **+1** | No Error | False reject |
| **-1** | False accept | No Error |

y

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

Letting an intruder in

# Cost functions matter: a simple example

x

f(x, W)

+1 = You

-1 = Bad guy

What if you're a CIA agent?

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

BIG penalty
for intruder

Don't mind about
annoyance…

f(x, W)

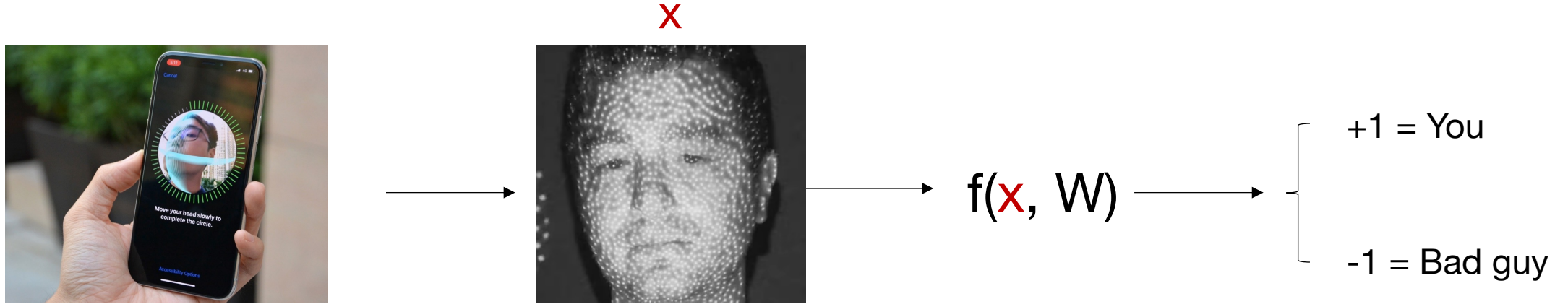|  | +1 | -1 |
|---|---|---|
| **+1** | No Error | False reject |
| **-1** | False accept | No Error |

y

It's you, but you
can't get in…

Letting an intruder in

# Cost functions matter: a simple example

deep imaging

x

$f(x, W)$

+1 = You

-1 = Bad guy

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 |
|---|---|---|
| +1 | No Error | False reject |
| -1 | False accept | No Error |

y

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

Letting an intruder 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:

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

# 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 **w**, having at hand our labeled data:

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

# 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 **w**, having at hand our labeled data:

Maximum Likelihood Estimation $\qquad p(\mathbf{w}|\mathbf{x}_1, ...\mathbf{x}_N; \mathbf{y}_1, ...\mathbf{y}_N)$

- 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 **w**, having at hand our labeled data:

Maximum Likelihood Estimation   $p(\mathbf{w}|\mathbf{x}_1, ...\mathbf{x}_N; \mathbf{y}_1, ...\mathbf{y}_N)$

- 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,

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

# 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,

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

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

# 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,

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

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

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

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

- 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?
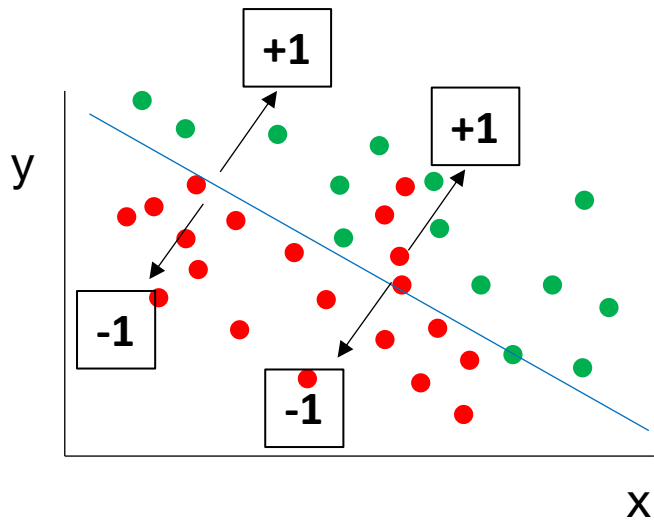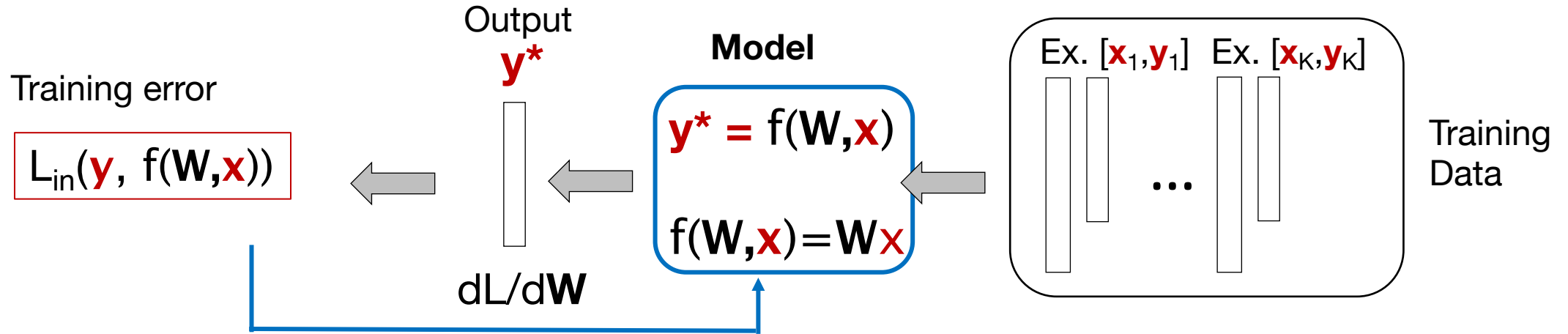
# Linear classification is the maximum likelihood for Gaussian data

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?

# 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?

$\downarrow$

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?

$\downarrow$

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

$\downarrow$

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

Summary: Linear classification with MSE assumes model output deviates from true labels via a Gaussian random process. Is this fair, given that labels are ether -1 or +1?

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

Training error

$L_{in}(\textcolor{red}{y}, f(\textbf{W},\textcolor{red}{x}))$

Output
$\textcolor{red}{y^*}$

**Model**

$\textcolor{red}{y^*} = f(\textbf{W},\textcolor{red}{x})$

$f(\textbf{W},\textcolor{red}{x}) = \textbf{W}\textcolor{red}{x}$

dL/d**W**

Ex. [$\textbf{x}_1,\textbf{y}_1$]  Ex. [$\textbf{x}_K,\textbf{y}_K$]

...

Training Data

+1

+1

y

-1

-1

x

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

deep imaging

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

Linear regression: predict some output h(x) from x



$x_0$
$x_1$
$x_2$
$x_N$

h(x)

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

Linear regression: predict some output h(x) from x

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

Linear regression

$x_0$
$x_1$
$x_2$
$x_N$

h(x)

Not a probabilistic mapping

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

Linear regression: predict some output h(x) from x

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

Linear classifier: predict binary output h(x) from x

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

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

Linear regression

$x_0$
$x_1$　　　　h(x)
$x_2$
$x_N$

Not a probabilistic mapping

Sign(x)

$x_0$
$x_1$　　　　h(x)
$x_2$
$x_N$

Probabilistic, but all-or-nothing: either 0, or 1

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

Linear regression:

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

Linear regression

$x_0$
$x_1$
$x_2$
$x_N$

h(x)

Not a probabilistic mapping

Linear classifier:

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

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

Sign(x)

$x_0$
$x_1$
$x_2$
$x_N$

h(x)

Probabilistic, but all-or-nothing: either 0, or 1

Logistic classifier:

$$h(x) = \theta(x) \in [0, 1]$$

Any value between 0 and 1

$\theta(x)$

$x_0$
$x_1$
$x_2$
$x_N$

h(x)

Probabilistic: continuous value between 0 and 1

Probabilistic interpretation of function that maps outputs to labels, h(x) = θ(x)

**Example**: You are trying to predict the probability that a patient may have a certain form a disease, θ(**x**), given a number of observations and measurements, **x**

**Example**: You are trying to predict the probability of rain tomorrow, θ(**x**), given a set of satellite image data, **x**

Probabilistic interpretation of function that maps outputs to labels, h(x) = θ(x)

**Example**: You are trying to predict the probability that a patient may have a certain form a disease, θ(**x**), given a number of observations and measurements, **x**

**Example**: You are trying to predict the probability of rain tomorrow, θ(**x**), given a set of satellite image data, **x**

**The Logistic Function $\theta$**

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

Also called Sigmoid function



θ(x)

- Use soft threshold to map any number to [0,1] range
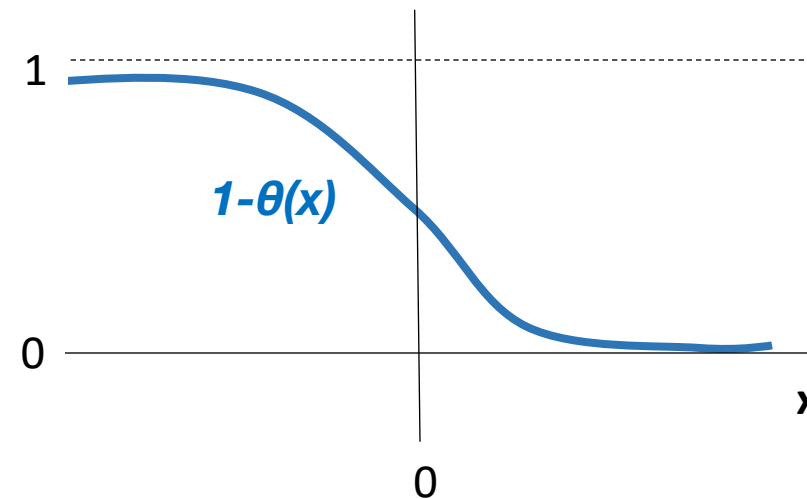- Sigmoid "flattens out" x

# From linear classification to logistic classification

- Let's re-derive a cost function for the case where labels are treated *as probabilities*

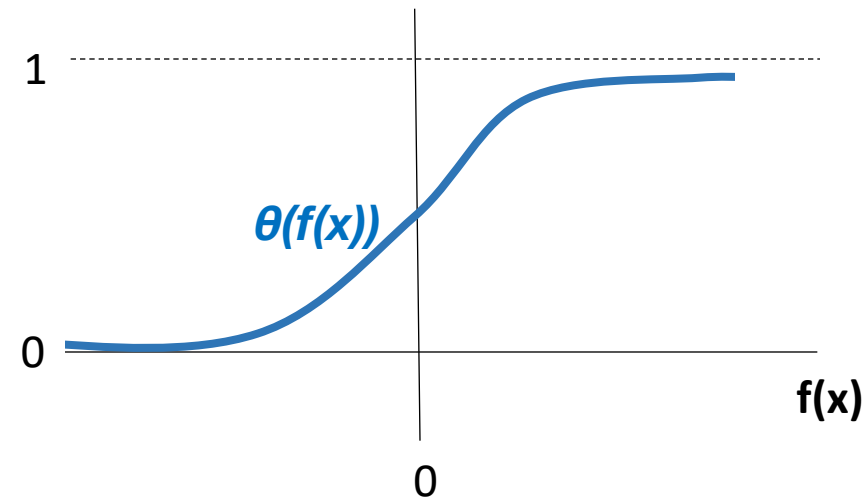  - You'll use this approach more often than not in Tensorflow…

# 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 this approach 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

deep imaging

- Let's re-derive a cost function for the case where labels are treated *as probabilities*

  - You'll use this approach 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

Formula for likelihood using the logistic function θ, given binary labels

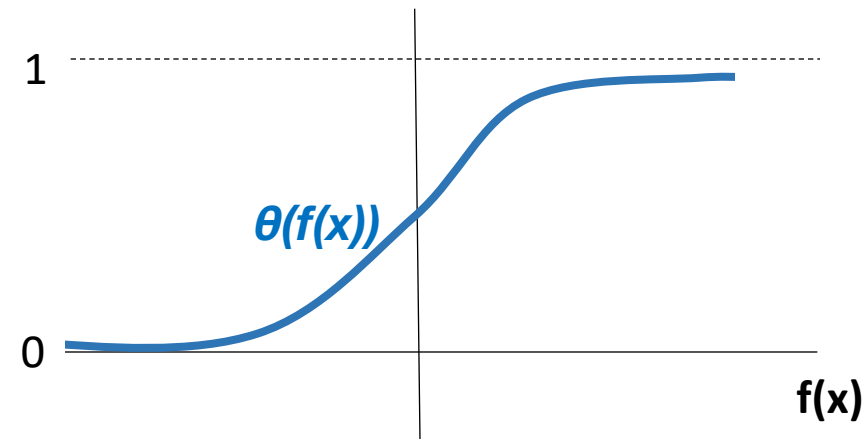$$P(y \mid \mathbf{x}) = \begin{cases} \theta(\mathbf{x}) & \text{For y = +1} \end{cases}$$
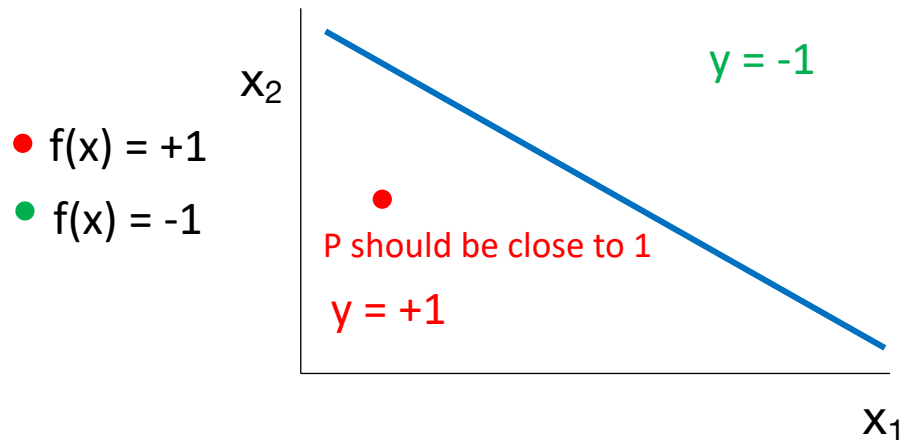
# From linear classification to logistic classification

- Let's re-derive a cost function for the case where labels are treated *as probabilities*

  - You'll use this approach 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

Formula for likelihood using the logistic function θ, given binary labels

$$P(y \mid \mathbf{x}) = \begin{cases} \theta(\mathbf{x}) & \text{For } y = +1 \\ 1 - \theta(\mathbf{x}) & \text{For } y = -1 \end{cases}$$



*1-θ(x)*

# 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 this approach 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

Formula for likelihood using the logistic function θ, given binary labels

$$P(y \mid f(\mathbf{x})) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$
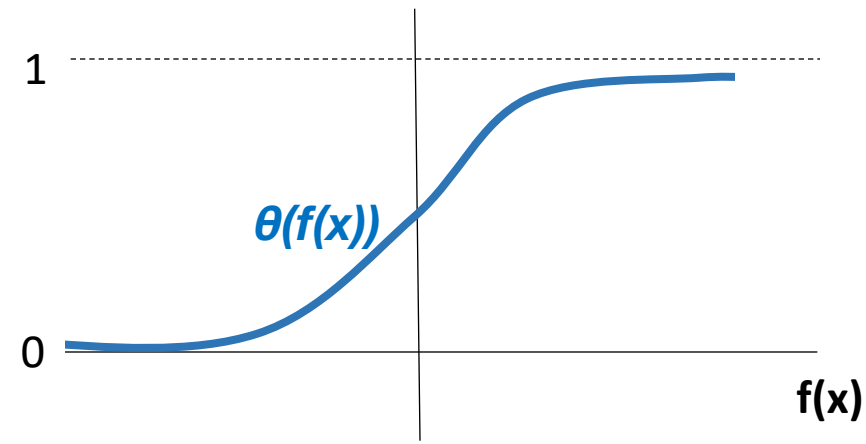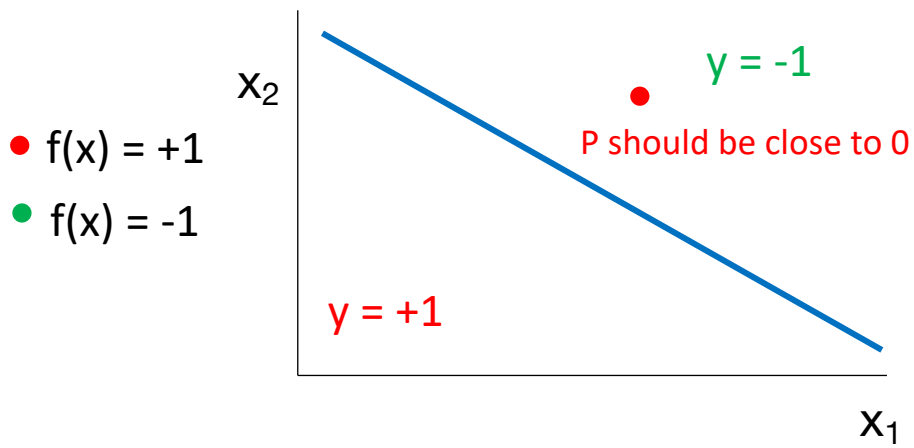
# From linear classification to logistic classification

- Let's re-derive a cost function for the case where labels are treated as probabilities

  - You'll use this approach 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

Formula for likelihood using the logistic function θ, given binary labels

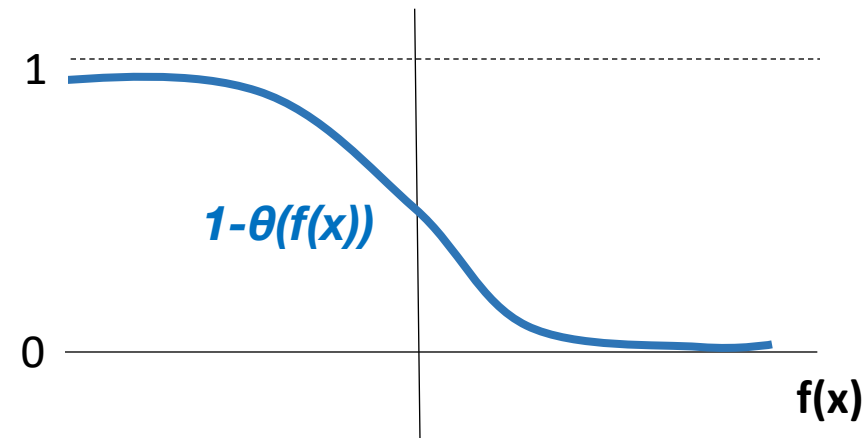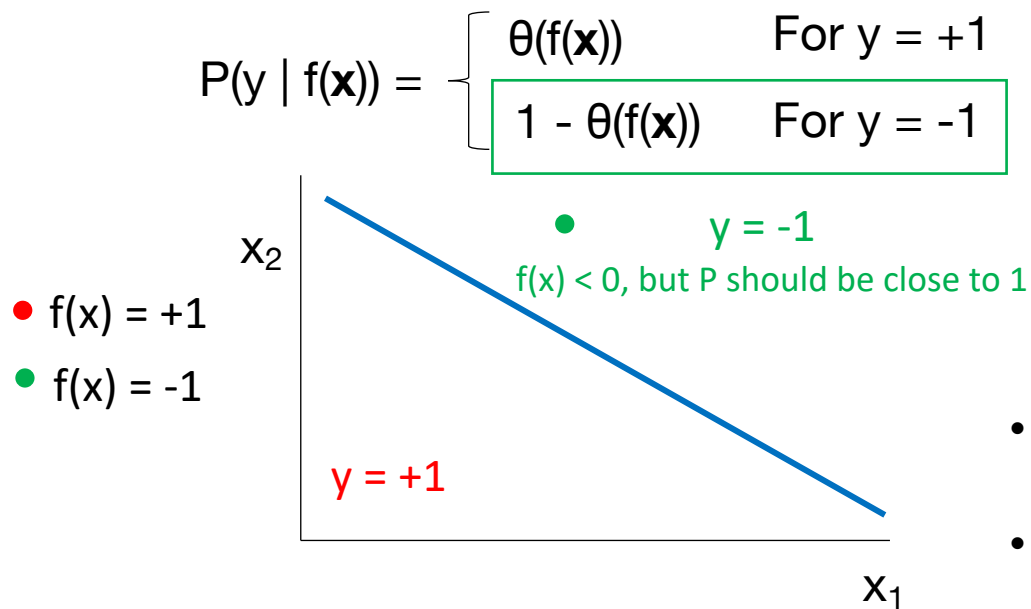$$P(y \mid f(\mathbf{x})) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$



$x_2$

● f(x) = +1

● f(x) = -1

y = -1

P should be close to 1

y = +1

$x_1$

1

$\theta(f(x))$

0

f(x)

- If network output f(x) is large, then should map to y=+1 with high probability
- $\theta(f(x))$ is large for large value of x

deep imaging

# 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 this approach 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

Formula for likelihood using the logistic function θ, given binary labels

$$P(y \mid f(\mathbf{x})) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$



$x_2$

• f(x) = +1

• f(x) = -1

y = -1

P should be close to 0

y = +1

$x_1$

1

$\theta(f(x))$

0

f(**x**)

# 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 this approach 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

Formula for likelihood using the logistic function θ, given binary labels

$$P(y \mid f(\mathbf{x})) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$

● y = -1
f(x) < 0, but P should be close to 1

$x_2$

● f(x) = +1
● f(x) = -1

y = +1

$x_1$

*1-θ(f(x))*

f(x)

- If network output f(x) is small, then should map to y=-1 with high probability
- θ(f(x)) ~0 for small values of f(x), so 1- θ(f(x)) ~1 is high probability to y=-1 mapping

**Deriving cost function for logistic classification for probabilistic outputs**

Instead of mapping f(**x**) to either +1 or -1 with the sign operator, let's use θ to map it to lie between 0 and 1:

$$P(y \mid \mathbf{x}) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$

**Deriving cost function for logistic classification for probabilistic outputs**

Instead of mapping f(**x**) to either +1 or -1 with the sign operator, let's use θ to map it to lie between 0 and 1:

$$P(y \mid \mathbf{x}) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$

We'll stick with the case of linear classification, where $f(\mathbf{x}) = \mathbf{w}^\mathsf{T}\mathbf{x}$

Also, please note that for the logistic function, $\theta(-a) = 1 - \theta(a)$.

**Deriving cost function for logistic classification for probabilistic outputs**

Instead of mapping f(**x**) to either +1 or -1 with the sign operator, let's use θ to map it to lie between 0 and 1:

$$P(y \mid \mathbf{x}) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$

We'll stick with the case of linear classification, where $f(\mathbf{x}) = \mathbf{w}^\mathsf{T}\mathbf{x}$

Also, please note that for the logistic function, θ(-a) = 1 - θ(a).

So, we can summarize the case-based definition above with a single function,

$$P(y \mid \mathbf{x}) = \theta(y\, f(\mathbf{x}))$$

**Deriving cost function for logistic classification for probabilistic outputs**

Instead of mapping f(**x**) to either +1 or -1 with the sign operator, let's use θ to map it to lie between 0 and 1:

$$P(y \mid \mathbf{x}) = \begin{cases} \theta(f(\mathbf{x})) & \text{For } y = +1 \\ 1 - \theta(f(\mathbf{x})) & \text{For } y = -1 \end{cases}$$

We'll stick with the case of linear classification, where f(**x**) = **w**$^T$**x**

Also, please note that for the logistic function, θ(-a) = 1 - θ(a).

So, we can summarize the case-based definition above with a single function,

$$P(y \mid \mathbf{x}) = \theta(y \, f(\mathbf{x}))$$
$$P(y \mid \mathbf{x}) = \theta(y \, \mathbf{w}^T\mathbf{x})$$

Here, y = +/-1 flips θ(**w**$^T$**x**) to be either θ(**w**$^T$**x**) or θ(-**w**$^T$**x**) = 1- θ(**w**$^T$**x**)

**Deriving cost function for logistic classification for probabilistic outputs**

Similar to the linear classification case, the likelihood of observing *N* independent outputs is given by,

$$P(y_1, y_2 \dots y_N \mid \mathbf{x}_1, \mathbf{x}_2, \dots \mathbf{x}_N) = \prod_{n=1}^{N} P(y_n \mid \mathbf{x}_n)$$

$$= \prod_{n=1}^{N} \theta(y_n \, \mathbf{w}^T \mathbf{x}_n)$$

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)

$$\text{Maximize} \quad P(y_1, y_2 \dots y_N \mid \mathbf{x}_1, \mathbf{x}_2, \dots \mathbf{x}_N) = \prod_{n=1}^{N} \theta(y_n \, \mathbf{w}^T \mathbf{x}_n)$$

**Deriving cost function for logistic classification for probabilistic outputs**

Maximize $P(y_1, y_2 \ldots y_N \mid \mathbf{x}_1, \mathbf{x}_2, \ldots \mathbf{x}_N) = \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}_n)$

Minimize $-\dfrac{1}{N} \ln \left( \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}) \right)$

**Deriving cost function for logistic classification for probabilistic outputs**

Maximize $P(y_1, y_2 ... y_N \mid \mathbf{x}_1, \mathbf{x}_2, ... \mathbf{x}_N) = \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}_n)$

Minimize $\quad -\dfrac{1}{N} \ln \left( \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}) \right)$

Minimize $\quad \dfrac{1}{N} \sum_{n=1}^{N} \ln \left( \dfrac{1}{\theta(y_n \mathbf{w}^T \mathbf{x})} \right)$

Use relationship $\quad \theta(a) = \dfrac{1}{1 + e^{-a}}$

**Deriving cost function for logistic classification for probabilistic outputs**

Maximize $\quad P(y_1, y_2 \dots y_N \mid \mathbf{x}_1, \mathbf{x}_2, \dots \mathbf{x}_N) = \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}_n)$

Minimize $\quad -\dfrac{1}{N} \ln \left( \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}) \right)$

Minimize $\quad \dfrac{1}{N} \sum_{n=1}^{N} \ln \left( \dfrac{1}{\theta(y_n \mathbf{w}^T \mathbf{x})} \right)$ $\qquad$ Use relationship $\qquad \theta(a) = \dfrac{1}{1 + e^{-a}}$

Minimize $\quad L_{in}(\mathbf{w}) = \dfrac{1}{N} \sum_{n=1}^{N} \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}})$

Cross entropy error for logistic classification

Typically requires iterative solution to minimize

**Deriving cost function for logistic classification for probabilistic outputs**

Maximize   $P(y_1, y_2 ... y_N \mid \mathbf{x}_1, \mathbf{x}_2, ... \mathbf{x}_N) = \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}_n)$

Minimize   $-\dfrac{1}{N} \ln \left( \prod_{n=1}^{N} \theta(y_n \mathbf{w}^T \mathbf{x}) \right)$

Minimize   $\dfrac{1}{N} \sum_{n=1}^{N} \ln \left( \dfrac{1}{\theta(y_n \mathbf{w}^T \mathbf{x})} \right)$    Use relationship    $\theta(a) = \dfrac{1}{1 + e^{-a}}$

Minimize   $L_{in}(\mathbf{w}) = \dfrac{1}{N} \sum_{n=1}^{N} \ln(1 + e^{-y_n \mathbf{w}^T \mathbf{x}})$        $L_{in}(\mathbf{w}) = \dfrac{1}{N} \sum_{n=1}^{N} (y_n - \mathbf{w}^T \mathbf{x})^2$

Cross entropy error for logistic classification            Mean-square error for linear classification

Typically requires iterative solution to minimize          Closed form solution available

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



Training error

$$L_{in}(\textbf{y}, f(\textbf{W},\textbf{x})) = \text{cross\_entropy}(\textbf{y}, f(\textbf{W},\textbf{x}))$$

Output
$\textbf{y}^*$

Model

$$\textbf{y}^* = \textbf{W}\textbf{x}$$

dL/d**W**

Training Data

Ex. $[\textbf{x}_1,\textbf{y}_1]$   Ex. $[\textbf{x}_K,\textbf{y}_K]$

...

Probabilistic mapping to y

y

x

deep imaging

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



**Training error**

$$L_{in}(\mathbf{y}, f(\mathbf{W},\mathbf{x})) =$$

$$cross\_entropy(\mathbf{y}, f(\mathbf{W},\mathbf{x}))$$

Output
$\mathbf{y}^*$

dL/d**W**

**Model**

$$\mathbf{y}^* = \mathbf{Wx}$$

Ex. [$\mathbf{x}_1$,$\mathbf{y}_1$]  Ex. [$\mathbf{x}_K$,$\mathbf{y}_K$]

...

Training
Data

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

Training data

$$f = W_1 x$$

Learned *f:* not flexible

f $\quad$ x

$$| = | \quad W_1 \quad |$$

Training data

$$f = W_1 x$$

Learned *f:* not flexible

Can we add flexibility by multiplying with another weight matrix?

$$\begin{cases} f_1 = W_1 x + b_1 \\ f_2 = W_2 f_1 + b_2 \end{cases}$$

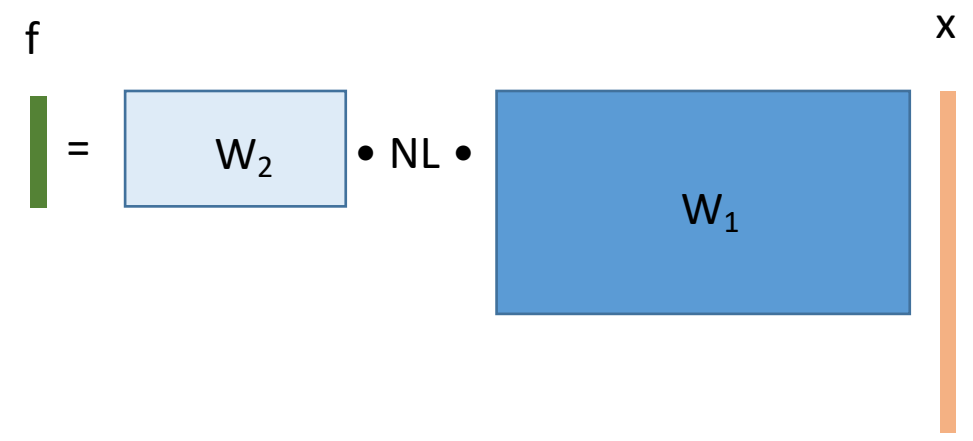$$f_2 = W_2(W_1 x + b_1) + b_2$$

$$f_2 = W' x + b'$$   Unfortunately not…

Training data

Learned *f:* not flexible

$$f = W_1 x$$

Training data

Learned *f:* a bit flexible

$$f = W_2 \max(W_1 x, 0)$$

Training data

$y$

$x$

$f = W_1 x$

Learned *f:* not flexible

f     x

$= \quad W_1$

Training data

$y$

$x$

$f = W_2 \mathrm{max}(W_1 x, 0)$

Learned *f:* a bit flexible

f     x

$= \quad W_2 \quad \bullet \mathrm{NL} \bullet \quad W_1$

$y$

$x$

$f = W_3 \mathrm{max}(0, W_2 \mathrm{max}(W_1 x, 0))$
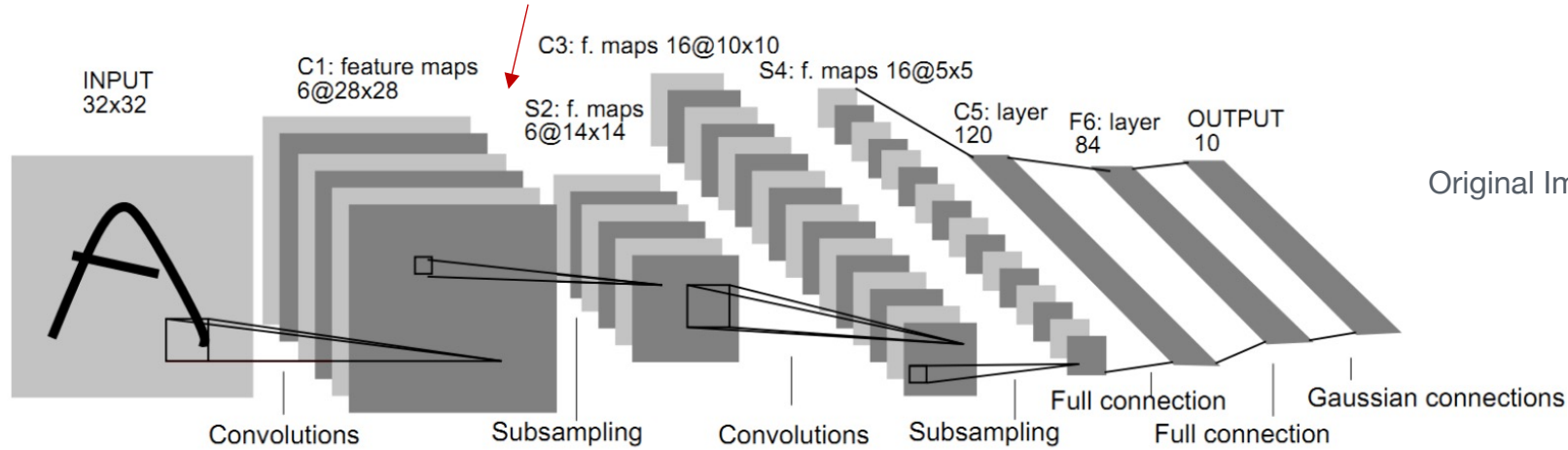
Learned *f:* more flexible

Does it generalize???

We can keep adding these "layers"…

deep imaging

# Getting us to Convolutional Neural Networks

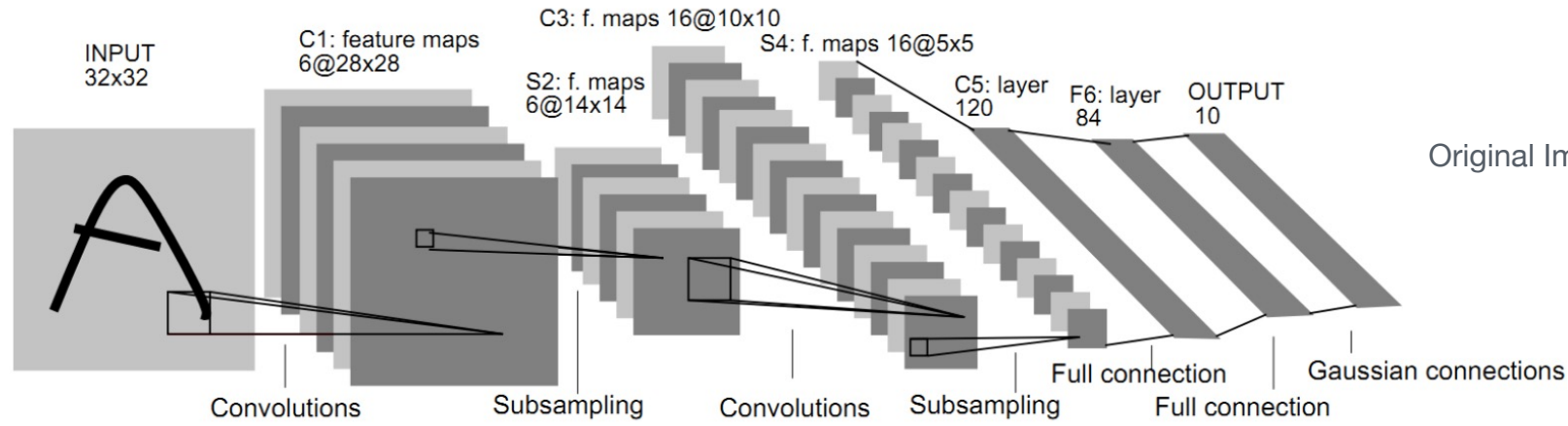After, apply non-linearity and sub-sampling



Original Image published in [LeCun et al., 1998]

Each matrix W is a convolution matrix

Repeat a few times

At the end, use a full W for a final matrix multiplication

# Getting us to Convolutional Neural Networks

Original Image published in [LeCun et al., 1998]

In practice, this process is repeated many times: