

Perceptron Learning Algorithm (1957, Rosenblatt)

Input: Training set  $D$  that is linearly separable

Output:  $w \in \mathbb{R}^{d+1}$  that achieves  $E_n(w) = 0$

Initialization

Initialize  $w$

in any arbitrary manner

$$w = (0, 0, 0, \dots, 0) \in \mathbb{R}^{d+1}$$

Step 1,

Check if  $E_n(w) = 0$

if yes, stop and output  $w$

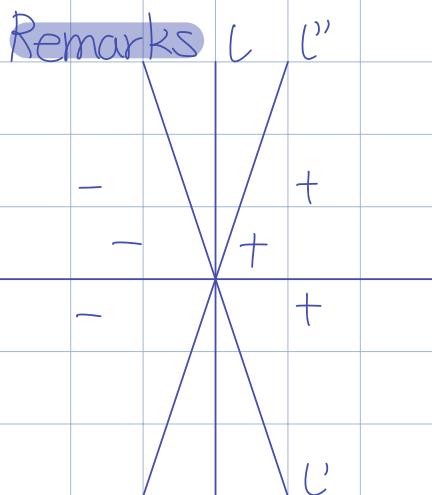
Step 2,

( $x_n, y_n$ ) be any misclassified point

$$w \leftarrow w + x_n \text{ if } y_n = +1$$

$$w \leftarrow w - x_n \text{ if } y_n = -1$$

Go to step 1

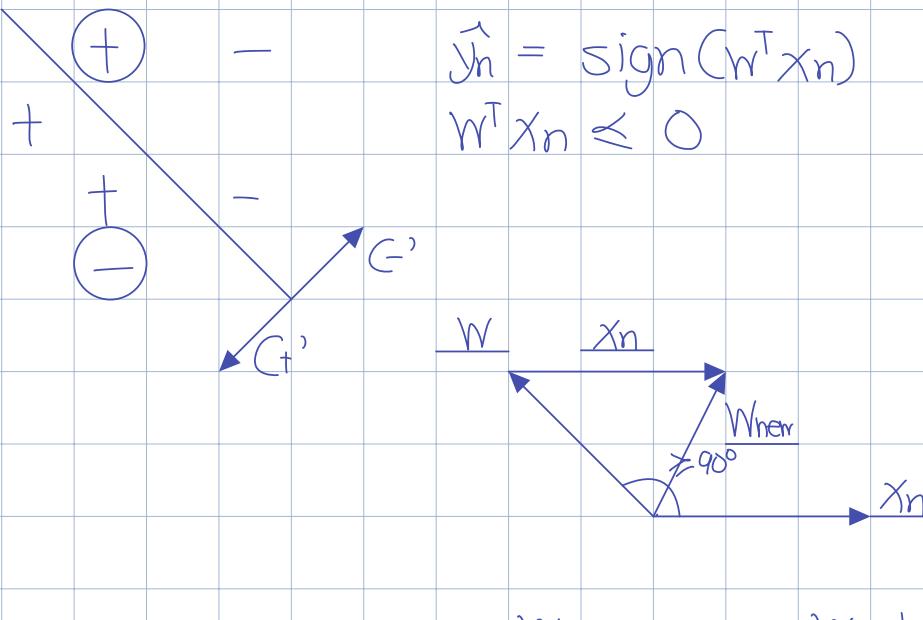


- 1) output is determined by the sequence of  $(x_n, y_n)$  in step 2
- 2) update rule is related gradient descent

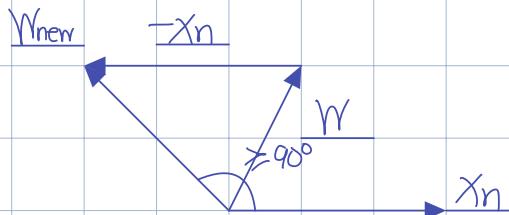
### Geometric Intuition for Update Rule

1)  $(x_n, y_n)$  is a misclassified point

$$y_n = +1 \quad \hat{y}_n = -1 \text{ or } x_n \text{ is on boundary}$$



$$y_n = -1 \quad \hat{y}_n = +1 \quad w^T x_n \geq 0$$



Problem 1.3 in textbook shows that PLA will converge in finitely many steps  $E_n(w) = 0$

### Pocket Algorithm

Idea: At each iteration, say  $t$  keep aside the "best weight vector" witnessed so far in iteration 1, 2, ...,  $t-1$

Do the update in PLA and see if it yields a better weight vector (if so update accordingly)

Stop after sufficiently many iterations

### Initialization

Initialize the "best weight vector" in some arbitrary manner

$$\hat{w} = w(0) = (0, 0, 0, \dots, 0) \in \mathbb{R}^{d+1}$$

for  $t=0, 1, 2, \dots, T-1$  do :

run PLA for one update to compute  $w(t+1)$

evaluate  $E_n(w(t+1))$

if  $E_n(w(t+1)) < E_n(\hat{w})$ :

$\hat{w} \leftarrow w(t+1)$

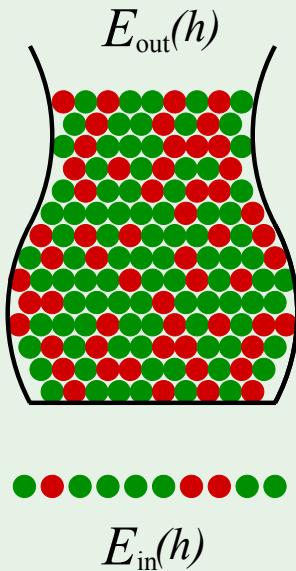
endif

output  $\hat{w}$

## Review of Lecture 2

Is Learning feasible?

Yes, in a **probabilistic** sense.



$$\mathbb{P} [ |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon ] \leq 2e^{-2\epsilon^2 N}$$

Since  $g$  has to be one of  $h_1, h_2, \dots, h_M$ , we conclude that

If:

$$|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon$$

Then:

$$|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| > \epsilon \quad \text{or}$$

$$|E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| > \epsilon \quad \text{or}$$

...

$$|E_{\text{in}}(h_M) - E_{\text{out}}(h_M)| > \epsilon$$

This gives us an added  $M$  factor.

# Learning From Data

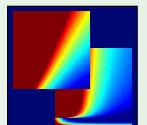
Yaser S. Abu-Mostafa  
*California Institute of Technology*

## Lecture 3: Linear Models I



Sponsored by Caltech's Provost Office, E&AS Division, and IST

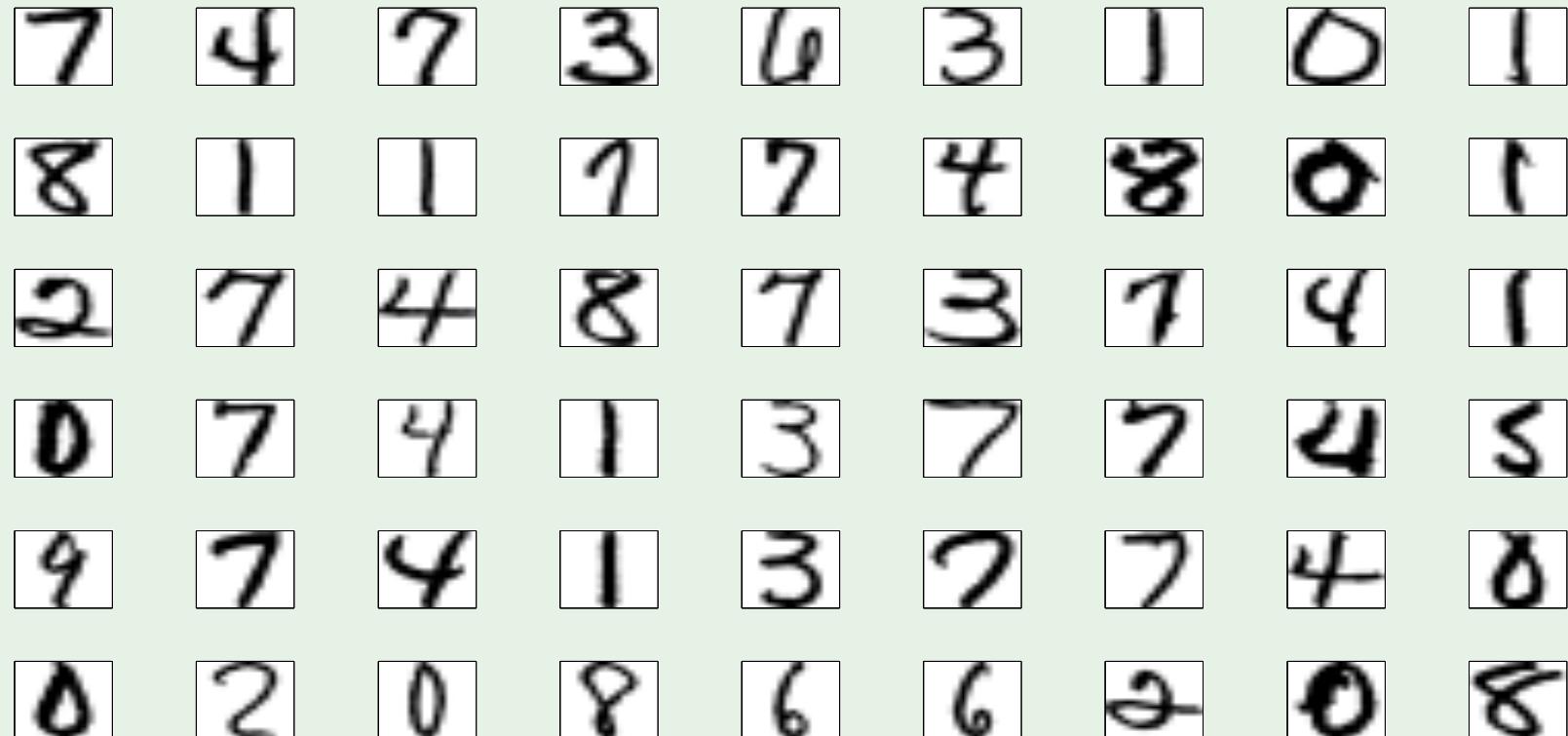
• Tuesday, April 10, 2012



# Outline

- Input representation
- Linear Classification
- Linear Regression
- Nonlinear Transformation

## A real data set



16x16 Images

## Input representation

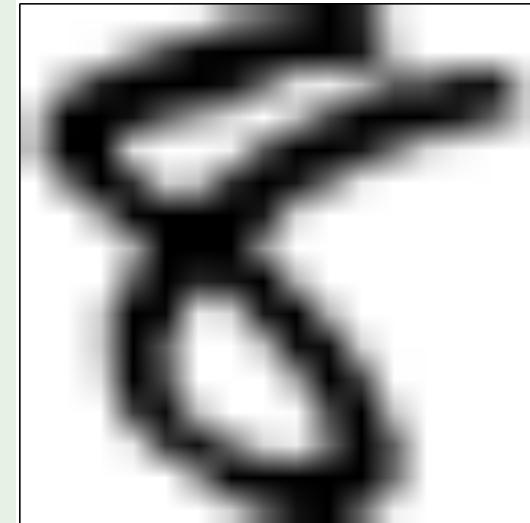
$$\text{'raw' input } \mathbf{x} = (x_0, x_1, x_2, \dots, x_{256}) \in \mathbb{R}^{256+1}$$

$$\text{linear model: } (w_0, w_1, w_2, \dots, w_{256})$$

**Features:** Extract useful information, e.g.,

$$\text{intensity and symmetry } \mathbf{x} = (x_0, x_1, x_2) \in \mathbb{R}^{2+1}$$

$$\text{linear model: } (w_0, w_1, w_2)$$

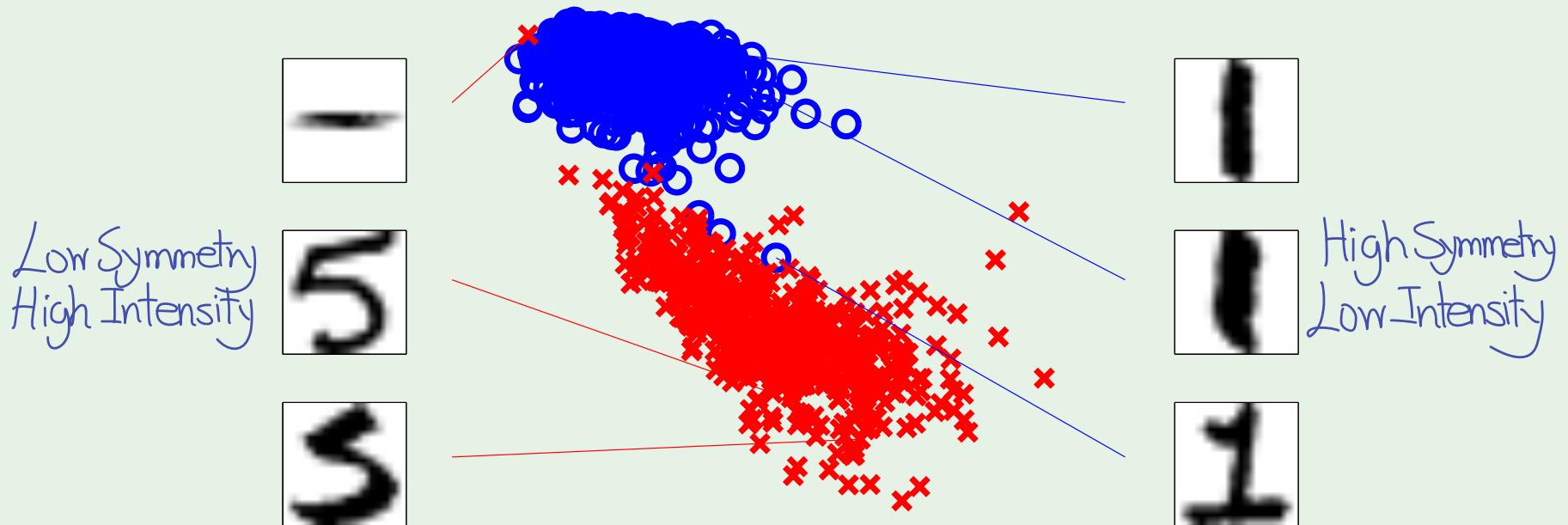


## Illustration of features

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

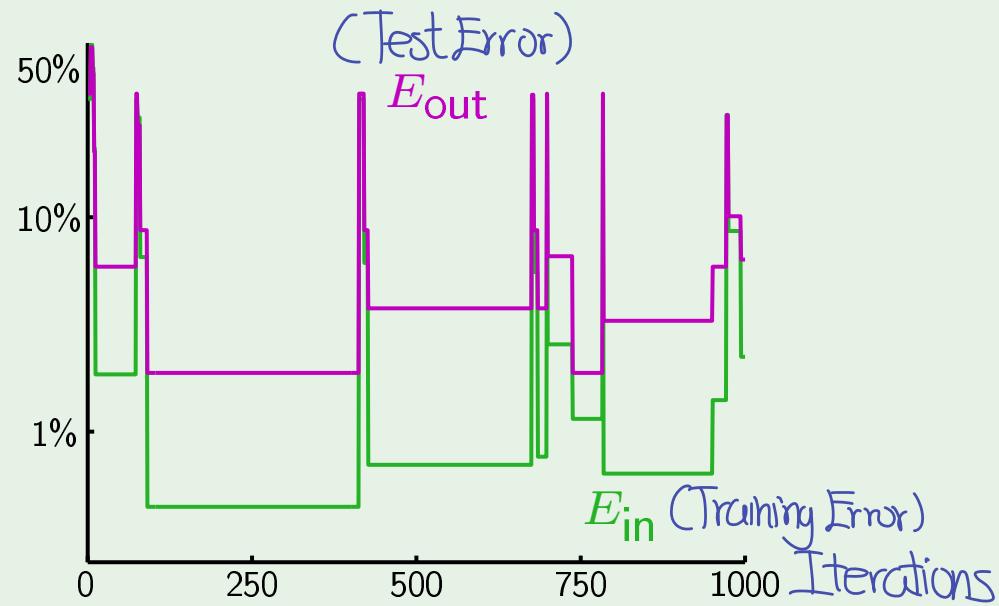
$x_1$ : intensity

$x_2$ : symmetry

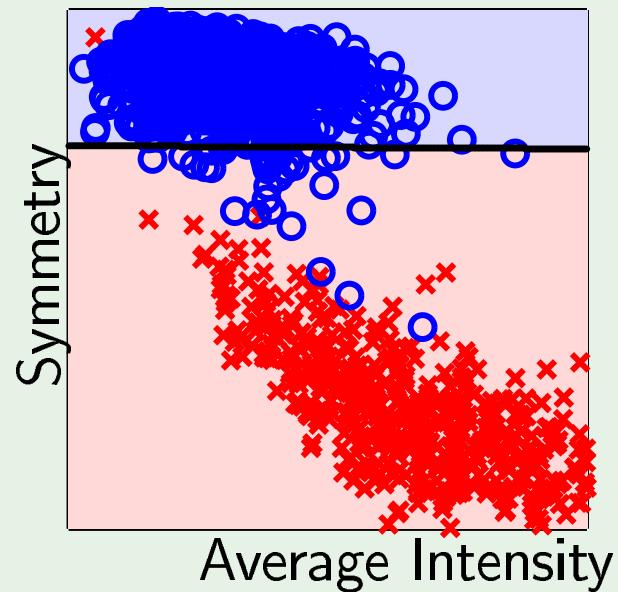


# What PLA does

Evolution of  $E_{\text{in}}$  and  $E_{\text{out}}$

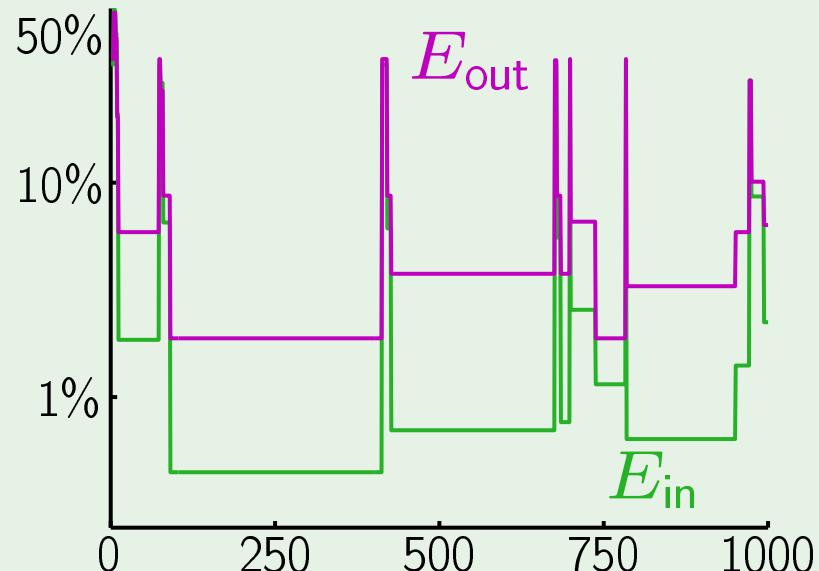


Final perceptron boundary

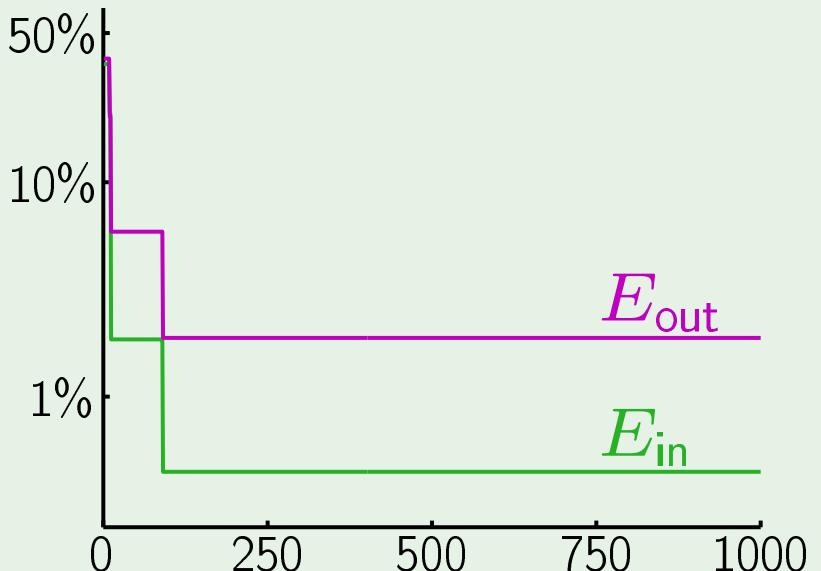


# The ‘pocket’ algorithm

PLA:

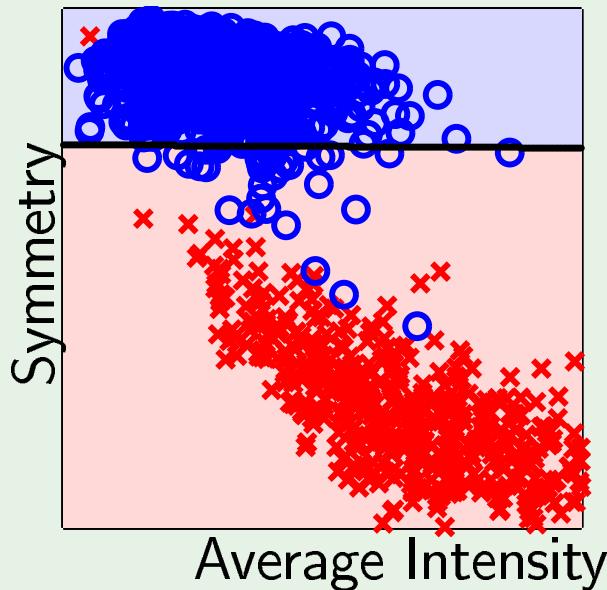


Pocket:

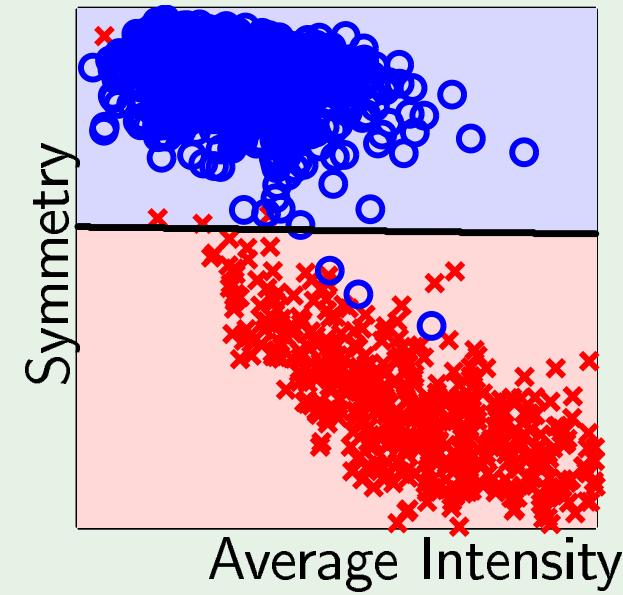


## Classification boundary - PLA versus Pocket

PLA:



Pocket:



# Outline

- Input representation
- Linear Classification
- Linear Regression      **regression  $\equiv$  real-valued output**
- Nonlinear Transformation

## Linear Regression

Given  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_d, y_d)\}$

$$x_i = (x_{i0=1}, x_{i1}, x_{i2}, \dots, x_{id})$$

$$y_i \in \mathbb{R}$$

Output a linear prediction function

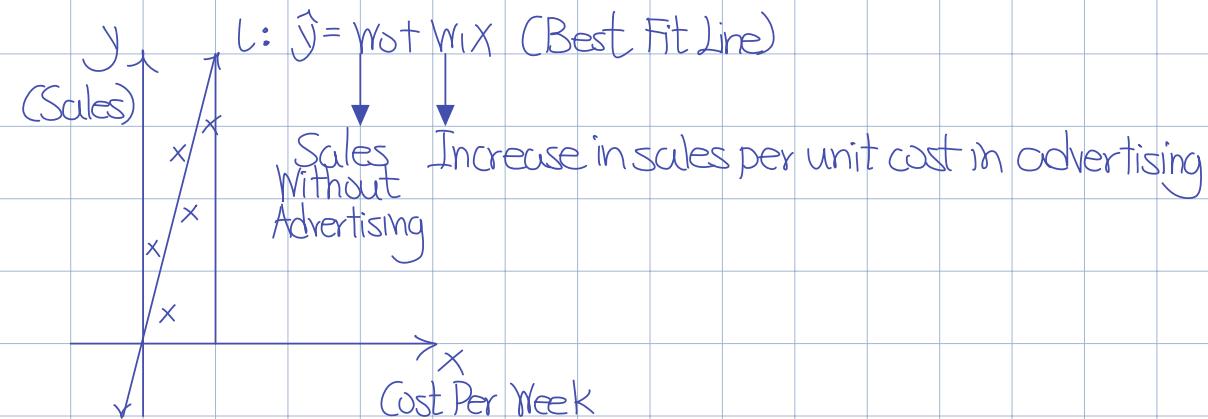
Given  $x = (x_0=1, x_1, x_2, \dots, x_d)$

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

Example |  $d=1$ )

$x$  = advertising cost in /week

$y$  = sales in that week



$$d = 3$$

$$x = (x_0, x_1, x_2, x_3)$$

|

$x_1$  = cost per week in TV advertising

$x_2$  = cost per week in Radio advertising

$x_3$  = cost per week in Internet advertising

$$\hat{y} = x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

### Notation

#### Data Vector

$$x_i = \begin{bmatrix} x_{i0} = 1 \\ x_{i1} \\ \vdots \\ x_{id} \end{bmatrix} \in \mathbb{R}^{d+1}$$

Arbitrary Scaling

#### Data Matrix

$$X = \begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_N^T \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$$

#### Observation Vector

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^{N \times 1}$$

## Weight Vector

$$w = w_0 \in \mathbb{R}^{d+1}$$
$$\vdots$$
$$w_d$$

## Prediction Vector

$$\hat{y} = \begin{matrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{matrix} \in \mathbb{R}^{N \times 1}$$

$$= \begin{matrix} x_1^T w \\ x_2^T w \\ \vdots \\ x_N^T w \end{matrix} = \begin{matrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{matrix} w = \cancel{x} w$$

Given  $x_i$

$$\hat{y}_i = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_d x_{id}$$
$$= w^T x_i = x_i^T w$$

Given  $\mathcal{f}$ ,  $Y$  compute  $w$  that outputs  $\hat{y} = \mathcal{f} w$

## Loss Function

Ideal World

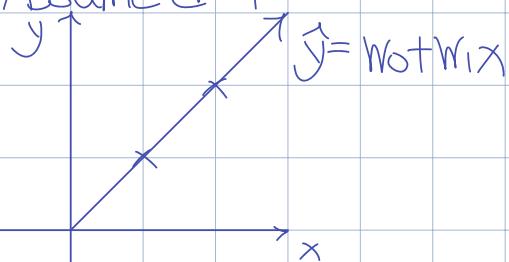
Find  $w$  such that all labels in training set are perfectly predicted

$$y_1 = \hat{y}_1 = x_1^T w$$

$$y_N = \hat{y}_N = x_N^T w$$

$$y_2 = \hat{y}_2 = x_2^T w$$

Assume  $d=1$



$$N \leq d+1$$

$$\begin{aligned} y_1 &= w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_d x_{1d} \\ y_2 &= w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_d x_{2d} \\ &\vdots \\ y_N &= w_0 + w_1 x_{N1} + w_2 x_{N2} + \dots + w_d x_{Nd} \end{aligned} \quad \left. \begin{array}{l} \text{Equation 1} \\ \text{Equation 2} \\ \vdots \\ \text{Equation } N \end{array} \right\}$$

Solve for  $w_0, w_1, w_2, \dots, w_d$

# of variable  $\geq$  # of constraints

$$d+1 \geq N$$

In practice,  $N \gg d$

Typically,  $N \approx 10d$

### Loss Function

$$\begin{aligned} \text{Ein}(w) &= \|y - \hat{y}\|^2 = \sum_{i=1}^N (y_i - \hat{y}_i)^2 \\ &\quad \begin{matrix} \text{observation vector} & \text{prediction vector} \end{matrix} \\ &= \|y - \cancel{\text{or}} w\|^2 \end{aligned}$$

(Squared Error Loss)

$$w^* = \underset{w \in \mathbb{R}^{d+1}}{\operatorname{argmin}} E_n(w)$$

$$w^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$\mathbb{R}^{(d+1) \times (d+1)}$        $\mathbb{R}^{(d+1) \times N}$        $\mathbb{R}^{N \times 1}$

Optimal Weight Vector  
Solved Analytically

Special Case ( $d=1$ )

$$w = (w_0, w_1)$$

$$\hat{y} = w_0 + w_1 x \quad \hat{y}_i = w_0 + w_1 x_i$$

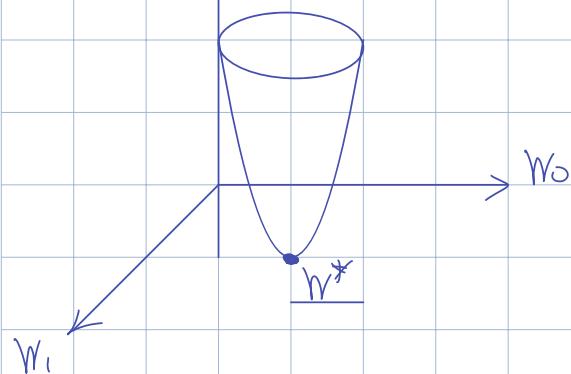
$$E_n(w) = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$= f(w)$$

$$= \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

minimize  $E_n(w)$

$$f(w_0, w_1)$$



Convex Function  
Paraboloid

$$f(w) = \text{constant}$$



Contour Plot

Ellipse

To minimize,  $f(\underline{w}) \rightarrow (w_0, w_1, w_2, \dots, w_d)$

$$\nabla_{\underline{w}} f(\underline{w}) = 0$$

Gradient of  $f(\cdot)$  wrt  $w$

$$\frac{\partial f}{\partial w_0}$$

$$\frac{\partial f}{\partial w_1} = 0$$

$\vdots$

$$\frac{\partial f}{\partial w_d}$$

Gradient

D Linear function

$$g(\underline{w}) = w_0 v_0 + w_1 v_1 + w_2 v_2 + \dots + w_d v_d$$

$$= \underline{w}^T \underline{v} = \underline{v}^T \underline{w}$$

$v$  : fixed vector

$$\nabla_{\underline{w}} g(\underline{w}) = \begin{matrix} \frac{\partial g}{\partial w_0} \\ \vdots \\ \frac{\partial g}{\partial w_d} \end{matrix} = \begin{matrix} v_0 \\ \vdots \\ v_d \end{matrix} = \underline{v}$$

$$g(\underline{w}) = \underline{w}^\top A \underline{w}$$

$A \in \mathbb{R}^{(d+1) \times (d+1)}$  matrix

$$\nabla g(\underline{w}) = (A + A^\top) \underline{w}$$

$$= 2A \underline{w} \quad \text{if } A = A^\top$$

$$\begin{aligned} f(\underline{w}) &= \|Y - \underline{Y}\|^2 \\ &= \|Y - \cancel{\mathcal{X}} \underline{w}\|^2 \\ &= (Y - \cancel{\mathcal{X}} \underline{w})^\top (Y - \cancel{\mathcal{X}} \underline{w}) \\ &= Y^\top Y - Y^\top \cancel{\mathcal{X}} \underline{w} - (\cancel{\mathcal{X}} \underline{w})^\top Y + \\ &\quad (\cancel{\mathcal{X}} \underline{w})^\top (\cancel{\mathcal{X}} \underline{w}) \\ &= Y^\top Y - 2 \underline{w}^\top \cancel{\mathcal{X}}^\top Y + \underline{w}^\top \cancel{\mathcal{X}}^\top \cancel{\mathcal{X}} \underline{w} \end{aligned}$$

$$\nabla_{\underline{w}} f(\underline{w}) = \cancel{\nabla}_{\underline{w}} (Y^\top Y) - 2 \cancel{\nabla}_{\underline{w}} (\underline{w}^\top \cancel{\mathcal{X}}^\top Y) +$$

$$\cancel{\nabla}_{\underline{w}} (\underline{w}^\top \cancel{\mathcal{X}}^\top \cancel{\mathcal{X}} \underline{w}) = 0$$

$$= -2 \cancel{\mathcal{X}}^\top Y + 2 \cancel{\mathcal{X}}^\top \cancel{\mathcal{X}} \underline{w} = 0$$

$$\cancel{\mathcal{X}}^\top \cancel{\mathcal{X}} \underline{w} = \cancel{\mathcal{X}}^\top Y$$

$$\underline{w} = (\cancel{\mathcal{X}}^\top \cancel{\mathcal{X}})^\top \cancel{\mathcal{X}}^\top Y$$

# Credit again

**Classification:** Credit approval (yes/no)

**Regression:** Credit line (dollar amount)

Input:  $\mathbf{x} =$

age	23 years
annual salary	\$30,000
years in residence	1 year
years in job	1 year
current debt	\$15,000
...	...

Linear regression output:  $h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^\top \mathbf{x}$

# The data set

Credit officers decide on credit lines:

$$(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_N, \mathbf{y}_N)$$

$y_n \in \mathbb{R}$  is the credit line for customer  $\mathbf{x}_n$ .

Linear regression tries to replicate that.

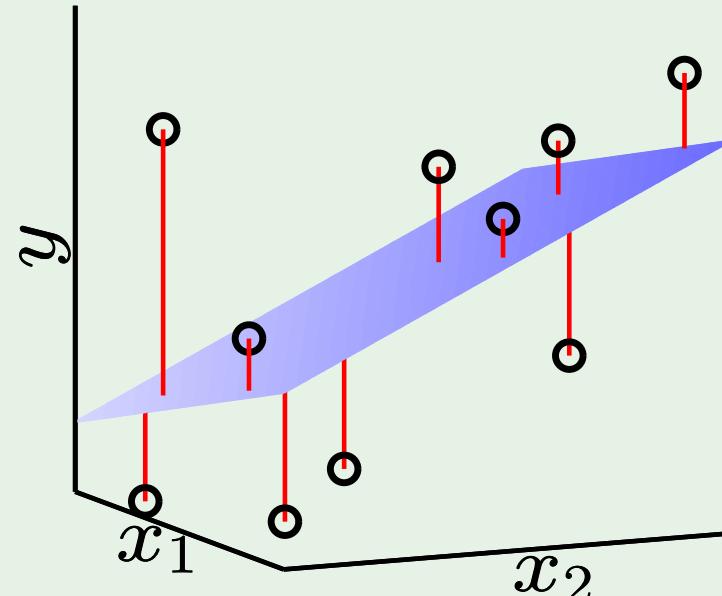
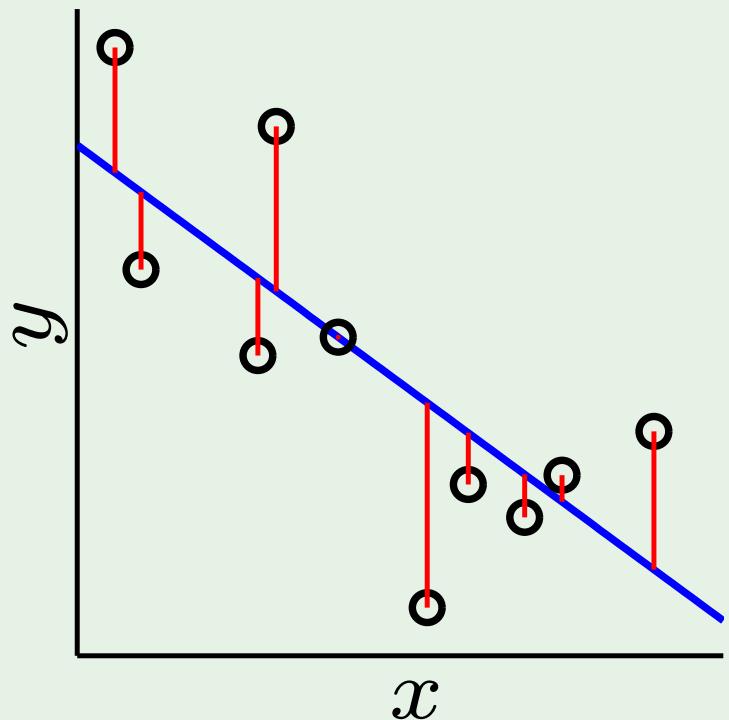
## How to measure the error

How well does  $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  approximate  $f(\mathbf{x})$ ?

In linear regression, we use squared error  $(h(\mathbf{x}) - f(\mathbf{x}))^2$

**in-sample error:**  $E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$

## Illustration of linear regression



## The expression for $E_{\text{in}}$

$$\begin{aligned} E_{\text{in}}(\mathbf{w}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - y_n)^2 \\ &= \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \end{aligned}$$

where  $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

## Minimizing $E_{\text{in}}$

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla E_{\text{in}}(\mathbf{w}) = \frac{2}{N} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0}$$

$$\mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y}$$

$$\mathbf{w} = \mathbf{X}^\dagger \mathbf{y} \quad \text{where} \quad \mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$\mathbf{X}^\dagger$  is the '**pseudo-inverse**' of  $\mathbf{X}$

# The pseudo-inverse

$$\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$$\left( \begin{array}{c|cc} & & \\ \hline & \begin{bmatrix} & \end{bmatrix} & \\ & \underbrace{\hspace{1cm}}_{d+1 \times d+1} & \\ & \begin{bmatrix} & \end{bmatrix} & \\ \hline & \begin{bmatrix} & \end{bmatrix} & \\ & \underbrace{\hspace{1cm}}_{d+1 \times N} & \end{array} \right)^{-1}$$

# 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  $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$ .

## Linear regression for classification

Linear regression learns a real-valued function  $y = f(\mathbf{x}) \in \mathbb{R}$

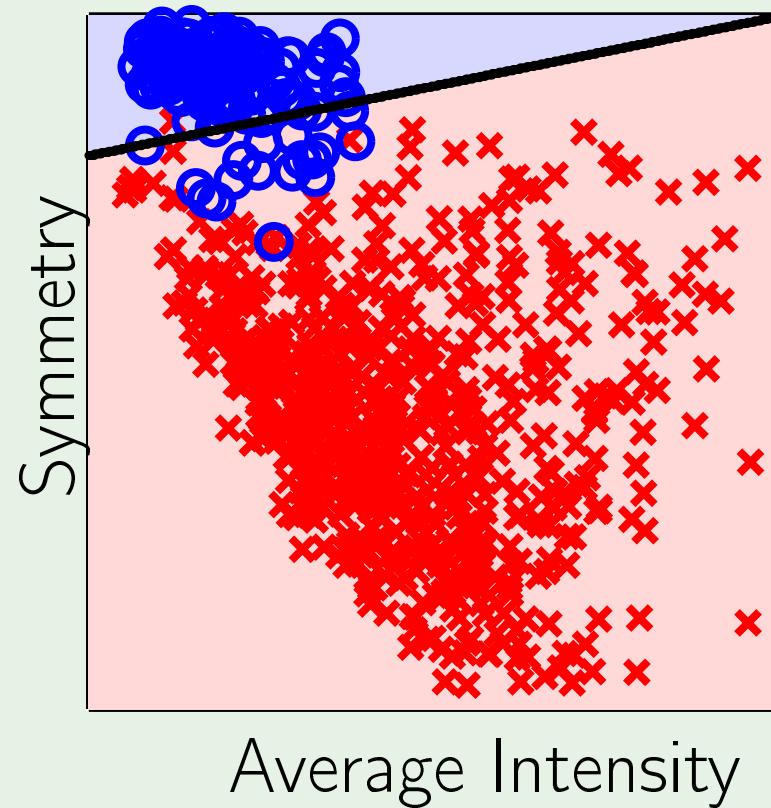
Binary-valued functions are also real-valued!  $\pm 1 \in \mathbb{R}$

Use linear regression to get  $\mathbf{w}$  where  $\mathbf{w}^\top \mathbf{x}_n \approx y_n = \pm 1$

In this case,  $\text{sign}(\mathbf{w}^\top \mathbf{x}_n)$  is likely to agree with  $y_n = \pm 1$

Good initial weights for classification

## Linear regression boundary

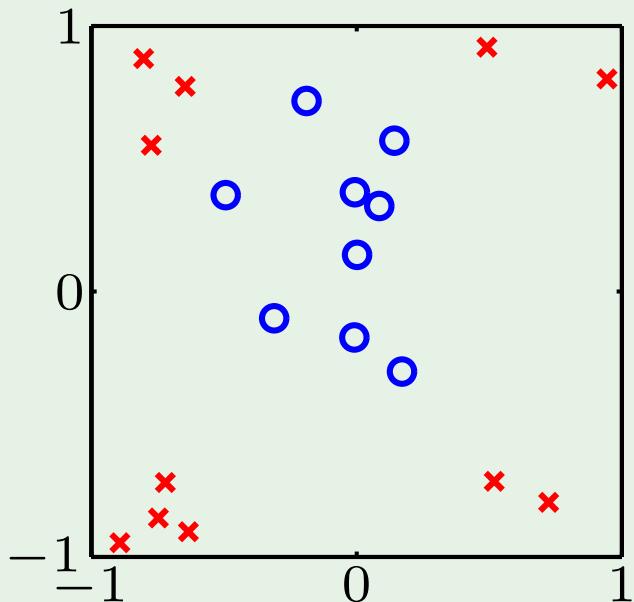


# Outline

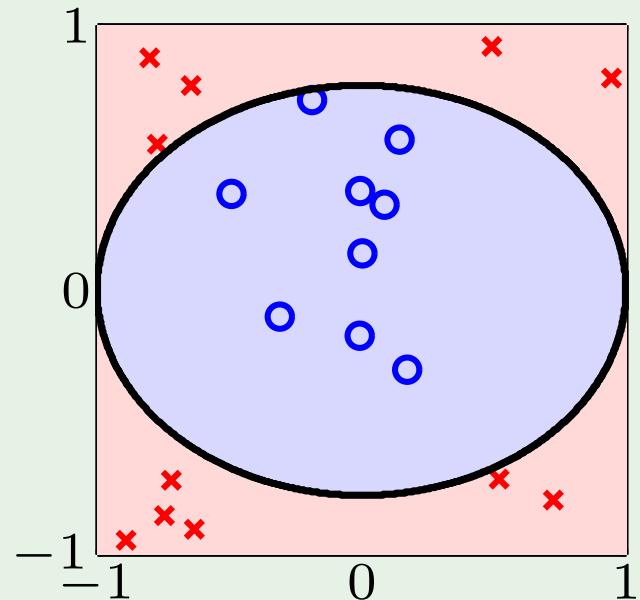
- Input representation
- Linear Classification
- Linear Regression
- Nonlinear Transformation

## Linear is limited

Data:



Hypothesis:



## Another example

Credit line is affected by ‘years in residence’

but **not** in a linear way!

Nonlinear  $[[x_i < 1]]$  and  $[[x_i > 5]]$  are better.

Can we do that with linear models?

## Linear in what?

Linear regression implements

$$\sum_{i=0}^d \textcolor{red}{w}_i x_i$$

Linear classification implements

$$\text{sign} \left( \sum_{i=0}^d \textcolor{red}{w}_i x_i \right)$$

Algorithms work because of **linearity in the weights**

## Transform the data nonlinearly

$$(x_1, x_2) \xrightarrow{\Phi} (x_1^2, x_2^2)$$

