

# Machine Learning from Data

Lecture 8: Spring 2021

# Today's Lecture

- Linear Models
  - Classification and
  - Regression

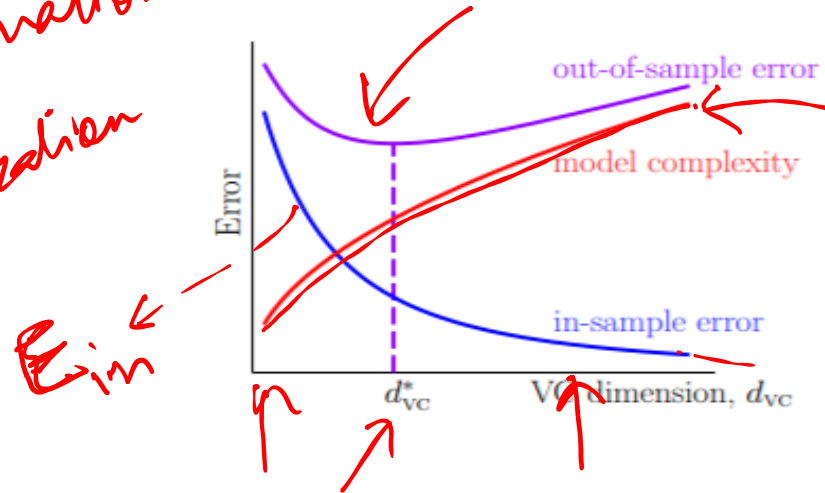
# TRADE-OFF

## VC Analysis

$$E_{\text{out}} \leq E_{\text{in}} + \Omega(d_{\text{VC}})$$

1. Did you fit your data well enough ( $E_{\text{in}}$ )?
2. Are you confident your  $E_{\text{in}}$  will generalize to  $E_{\text{out}}$

Approximation  
vs  
generalization



## The VC Insurance Co.

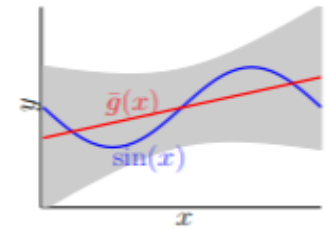
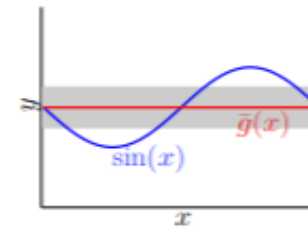
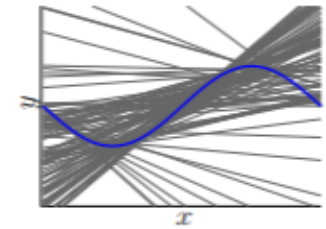
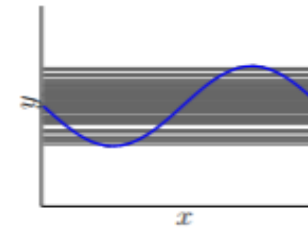
The VC warranty had conditions for becoming void:

- You can't look at your data *before* choosing  $\mathcal{H}$ .
- Data must be generated i.i.d from  $P(\mathbf{x})$ .
- Data and test case from *same*  $P(\mathbf{x})$  (same bin).

## Bias-Variance Analysis

$$E_{\text{out}} = \text{bias} + \text{var}$$

1. How well can you fit your data (**bias**)?
2. How close to that best fit can you get (**var**)?



$\mathcal{H}_0$

bias = 0.50;  
var = 0.25.

$$E_{\text{out}} = 0.75 \quad \checkmark$$

$\mathcal{H}_1$

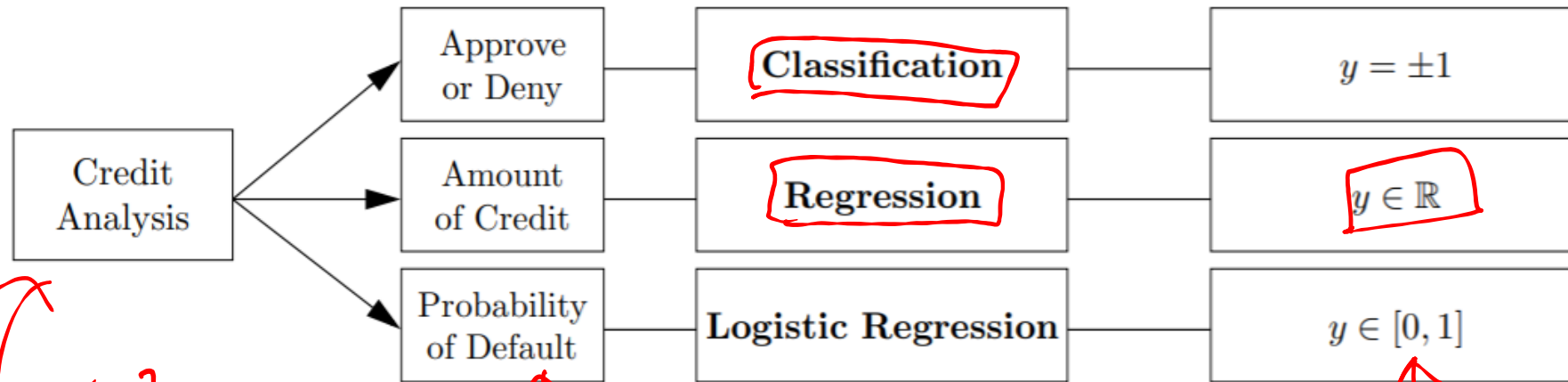
bias = 0.21;  
var = 1.69.

$$E_{\text{out}} = 1.90$$

# Linear Model

---

$y \rightarrow$  outcome variable



Chapter 2

A  $\rightarrow$  0.4 or 0.9

# Linear Model: FUNDAMENTAL

i) On its own it's able to solve most ML problems

ii) It's a building block for other models.  
NN, SVM (robust linear model)

$$S = \underbrace{w^T}_{\text{linear in weights}} \underbrace{x}_{\text{linear in } x}$$

where

$$w \in \mathbb{R}^{d+1}$$

$$x \in 1 \times \mathbb{R}^{d+1}$$

linear in weights.

linear in  $x$

hyperplane  $\xrightarrow{2d}$  line

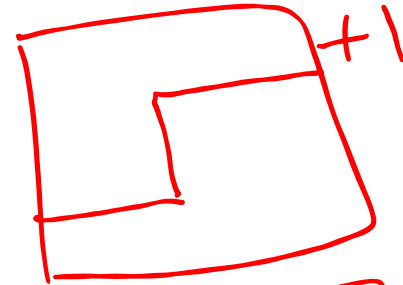
Is it possible to solve any kind of problem using linear model?

TRANSFORM

$$S = \underbrace{W^T x}_{\text{linear signal}}$$

linear signal

sign function  
(classification)



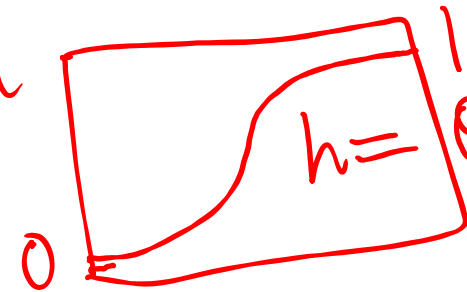
$$h(x) = \text{sign}(w^T x)$$

Identity Transform  
(Regression)

$$h = s$$

$$h(x) = w^T x$$

Sigmoid function  
(Logistic Regression)



$$h = \theta(s)$$
$$h(x) = \theta(w^T x)$$

# Classification (Perceptron)

$$\mathcal{H}_{\text{linear}} = \left\{ h(x) = \text{sign}(w^T x) \mid w \in \mathbb{R}^{d+1} \right\}$$

•  $\checkmark$  i)  $d_{VC} = d+1$ ,  $\therefore E_{\text{out}} \leq E_{\text{in}} + O\left(\sqrt{d \frac{\ln N}{N}}\right)$

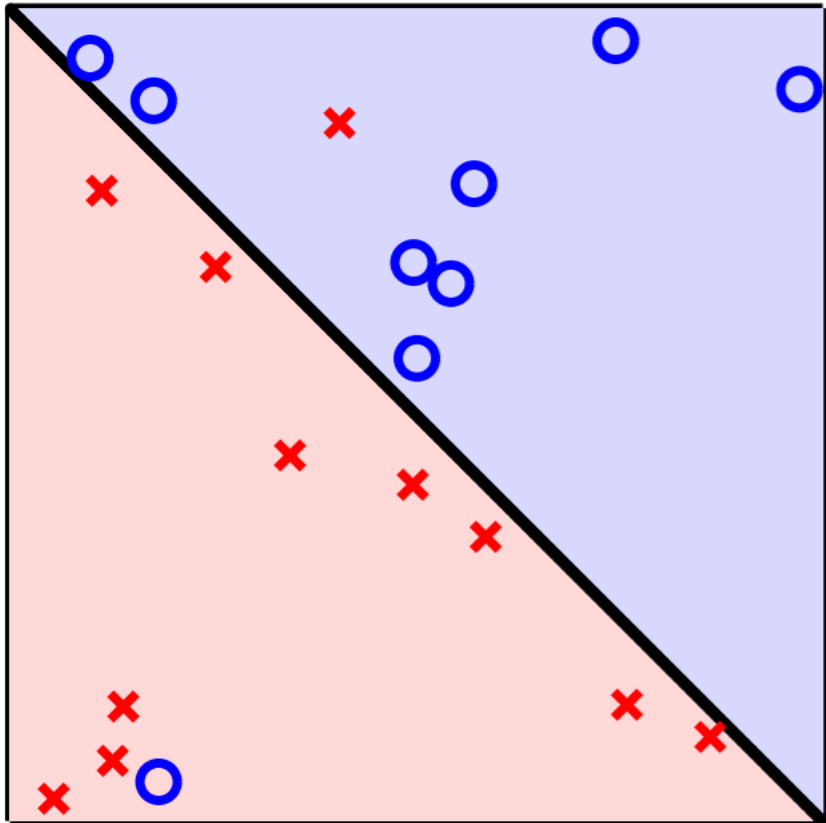
$\checkmark$  ii)  $E_{\text{in}} \approx 0$ , PLA  $\longrightarrow$  accomplish that.  $E_{\text{in}}(g) \approx 0$

Q1. What happens if data is not linearly separable?

Q2. How can we ensure  $E_{\text{in}}(g) \approx 0$

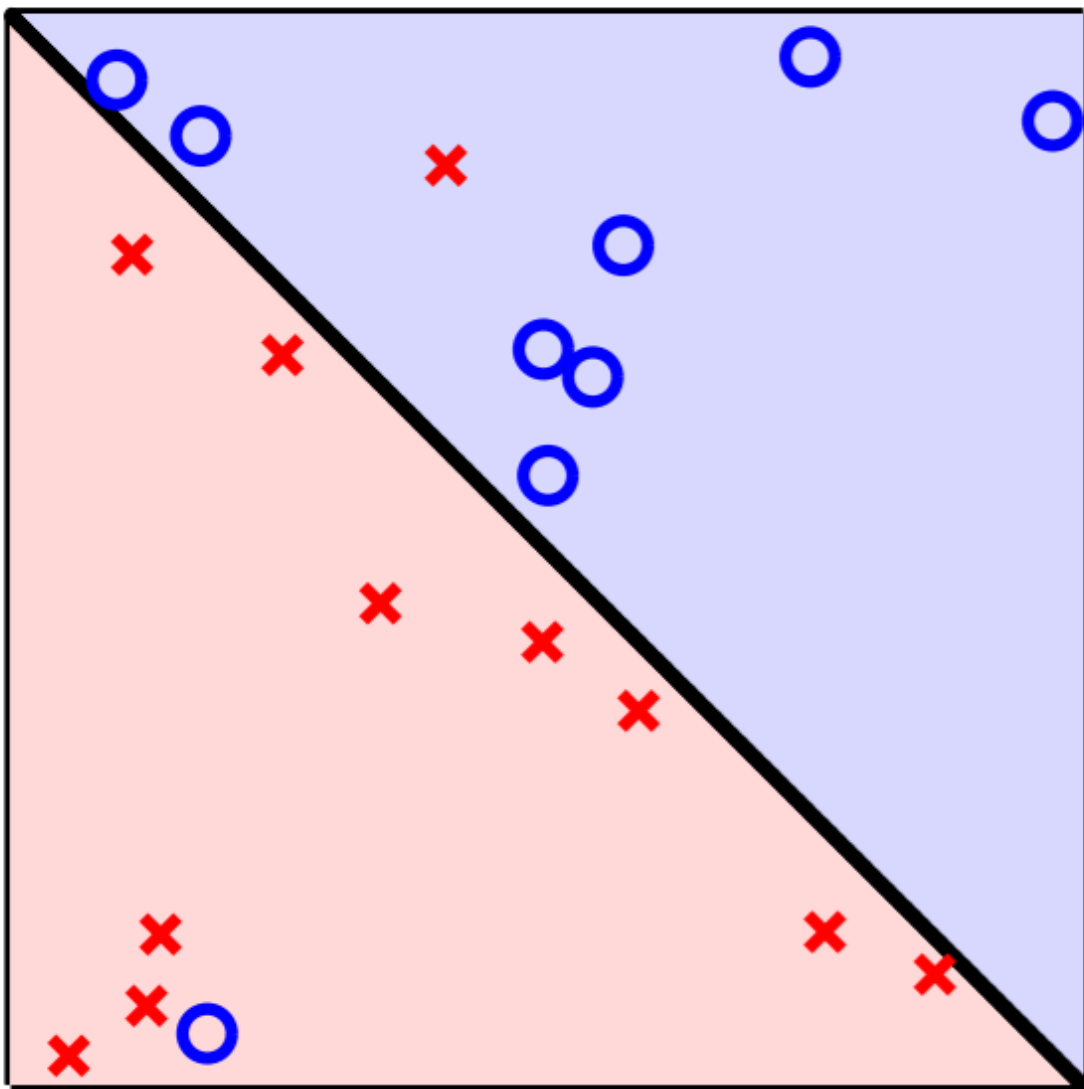
Feature Set

$\downarrow$   
To learn error  
(Pocket Algorithm)



PLA : Perceptron Algorithm  
 $w(t+1) \leftarrow \underline{w(t)} + y_* x_*$   
 $\downarrow$   
 $w_i \rightarrow E_{in}(best)$   
 $\downarrow$   
Pocket  
 $\downarrow$   
Update





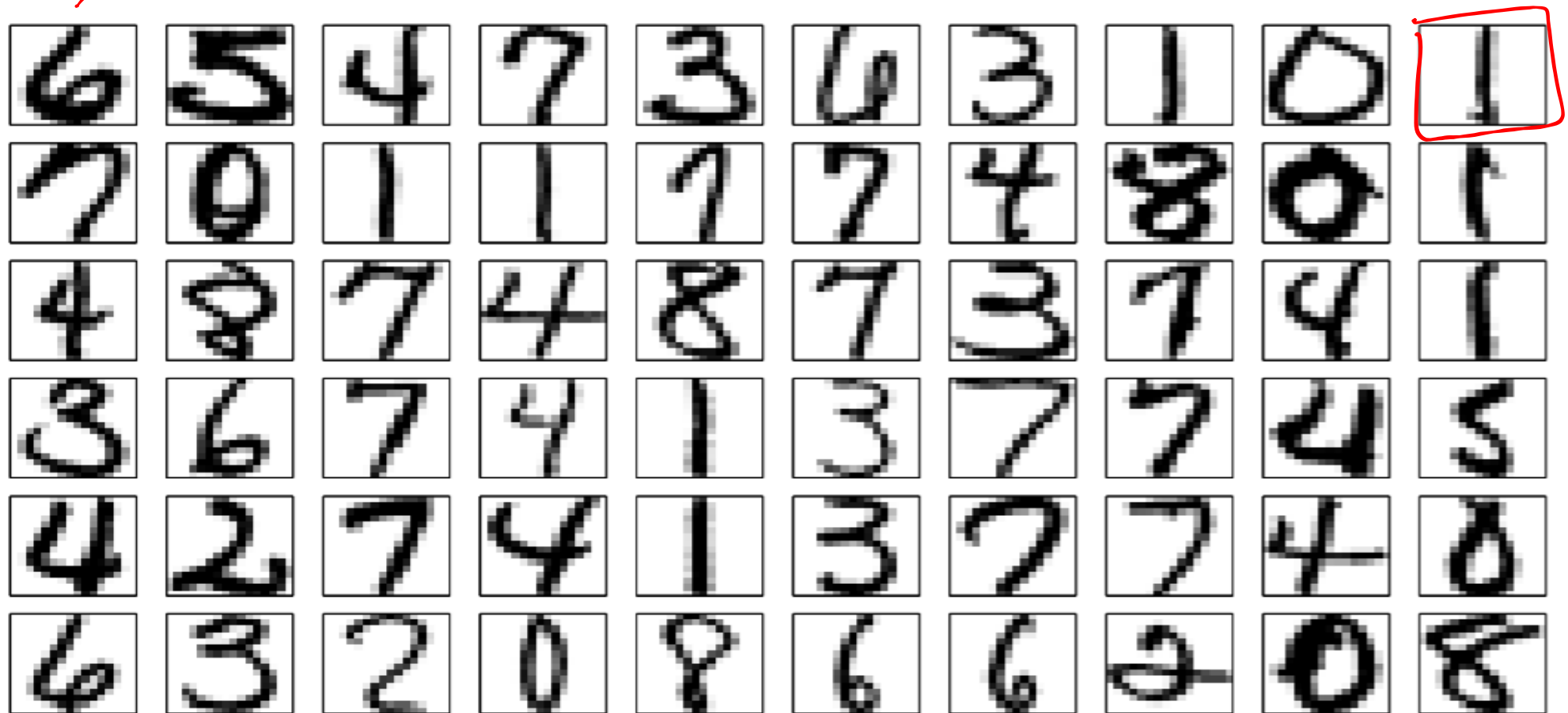
→ Get a better starting point.  
 → Algorithms.  
 → Minimizing  $E_{in}$  is a hard combinatorial problem.

## The Pocket Algorithm

- Run PLA ✓
- At each step keep the best  $E_{in}$  (and  $\mathbf{w}$ ) so far.

(It's not rocket science, but it works.)

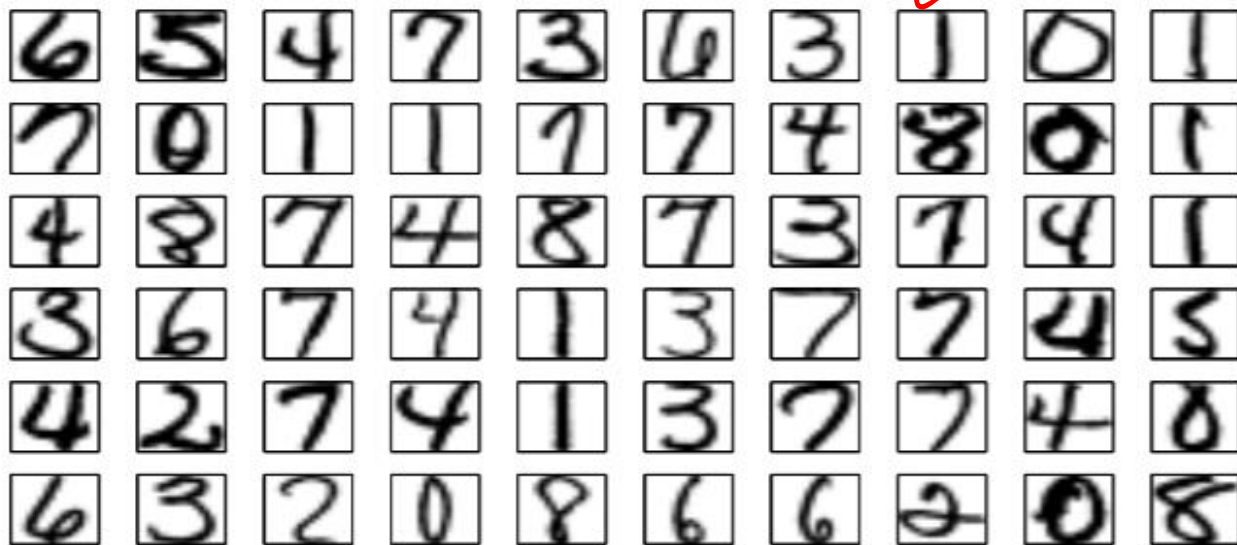
Stopping Criteria  
 10000 iterations



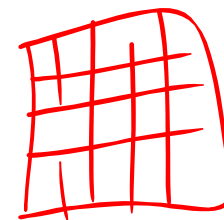
Each digit is a  $16 \times 16$  image.

*Pre-processing.*

$(E_{in} \approx E_{out}) \rightarrow E_{in} \approx 0 \checkmark$



$d = 256$   
 $d_{vc} = 257$



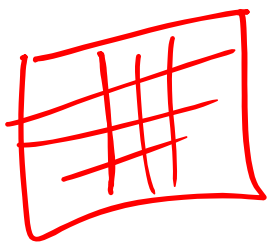
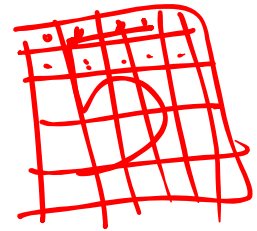
Each digit is a  $16 \times 16$  image.

$-1$  to  $+1$



$\begin{bmatrix} -1 & -1 & -1 & -1 & -1 & -1 & -1 & -0.63 & \underline{0.86} & \underline{-0.17} & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -0.99 & 0.3 & 1 & 0.31 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ -1 & -1 & -0.41 & 1 & 0.99 & -0.57 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -0.68 & 0.83 & 1 & 0.56 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -0.94 & 0.54 \\ 1 & 0.78 & -0.72 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 0.1 & 1 & 0.92 & -0.44 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -0.26 & 0.95 & 1 & -0.16 & -1 & -1 \\ -1 & -0.99 & -0.71 & -0.83 & -1 & -1 & -1 & -1 & -1 & -0.8 & 0.91 & 1 & 0.3 & -0.96 & -1 & -1 & -0.55 & 0.49 & 1 & 0.88 & 0.09 & -1 & -1 & -1 & -1 & 0.28 & 1 & 0.88 & -0.8 & -1 \\ -0.9 & 0.14 & 0.97 & 1 & 1 & 1 & 0.99 & -0.74 & -1 & -1 & -0.95 & 0.84 & 1 & 0.32 & -1 & -1 & 0.35 & 1 & 0.65 & -0.10 & -0.18 & 1 & 0.98 & -0.72 & -1 & -1 & -0.63 & 1 & 1 \\ 0.07 & -0.92 & 0.11 & 0.96 & 0.30 & -0.88 & -1 & -0.07 & 1 & 0.64 & -0.99 & -1 & -1 & -0.67 & 1 & 1 & 0.75 & 0.34 & 1 & 0.70 & -0.94 & -1 & -1 & 0.54 & 1 & 0.02 & -1 & -1 \\ -1 & -0.90 & 0.79 & 1 & 1 & 1 & 1 & 0.53 & 0.18 & 0.81 & 0.83 & 0.97 & 0.86 & -0.63 & -1 & -1 & -1 & -1 & -0.45 & 0.82 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0.13 & -1 & -1 & -1 & -1 & -1 \\ -1 & -0.48 & 0.81 & 1 & 1 & 1 & 1 & 1 & 1 & 0.21 & -0.94 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -0.97 & -0.42 & 0.30 & 0.82 & 1 & 0.48 & -0.47 & -0.99 & -1 & -1 & -1 & -1 \end{bmatrix}$

# Feature Construction

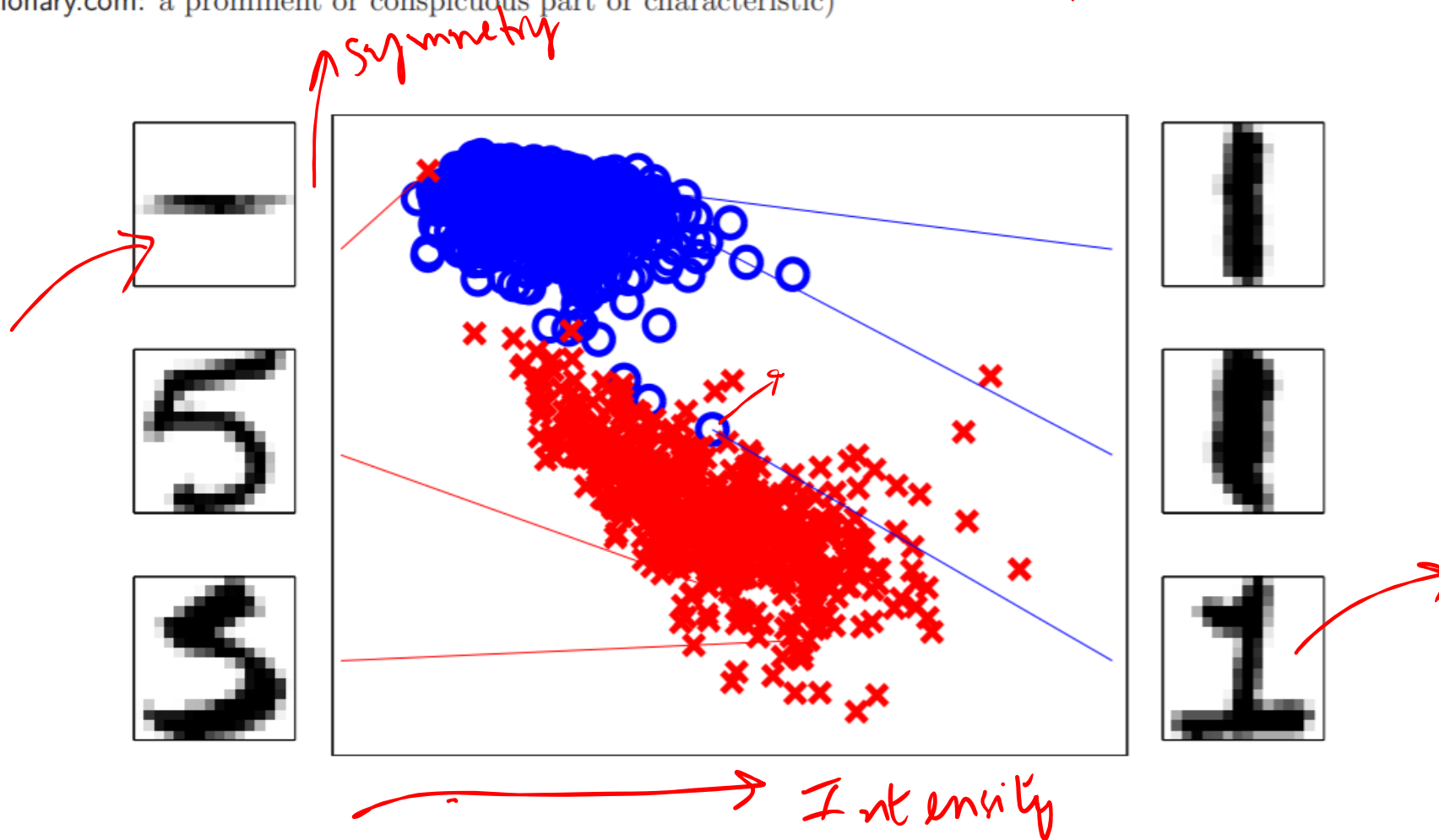
- Summarize the input into features that are useful  
↓  
dvc
  - i) Reduce the dimension (dvc) significantly
  - ii) What is really important?  
 $E_{in} \approx E_{out}$
- 1      5      ✓ Intensity 
- ↗      ↖      ✓ Symmetry 
- Mathematical form.      ↘      Algorithm
- $dvc = 3!$

# Linear Model

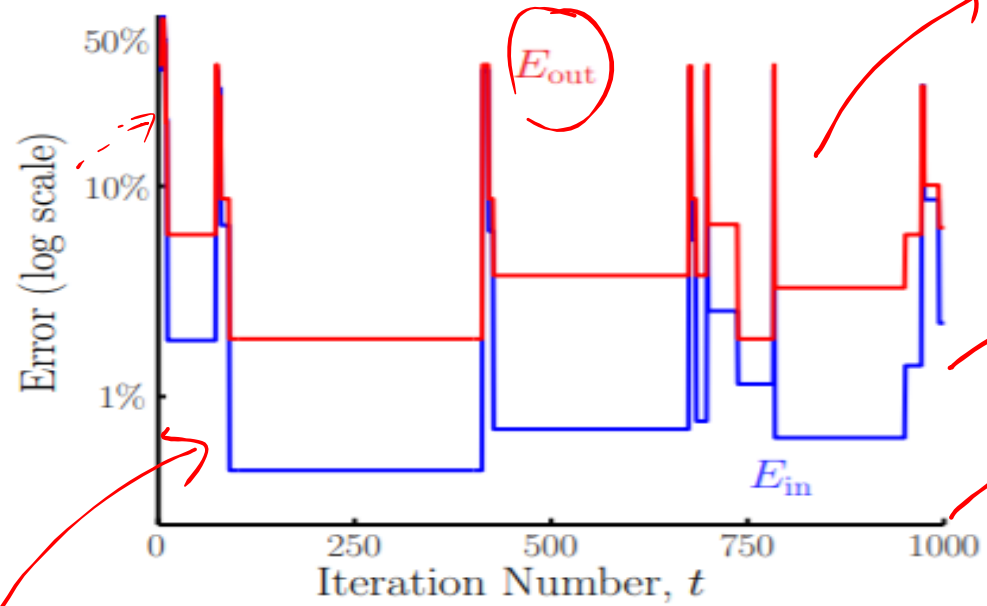
- Powerful .  $E_{in} \approx E_{out}$   
Feature Construction.

## Intensity and Symmetry Features

feature: an important property of the input that you think is useful for classification.  
(dictionary.com: a prominent or conspicuous part or characteristic)



# PLA

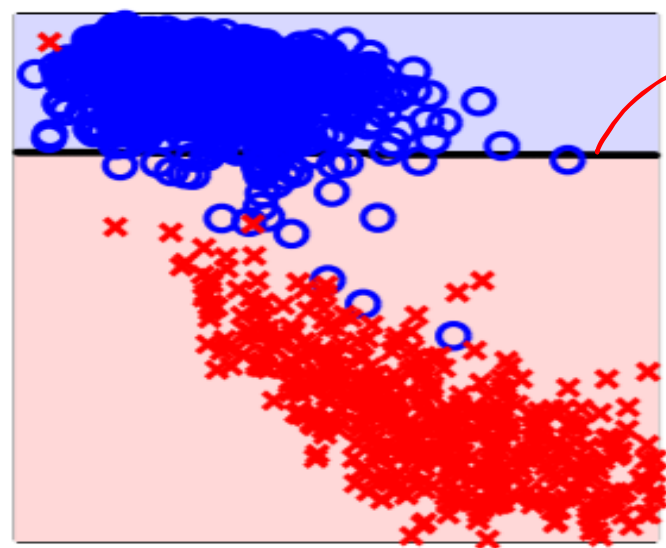


sync

500

situations

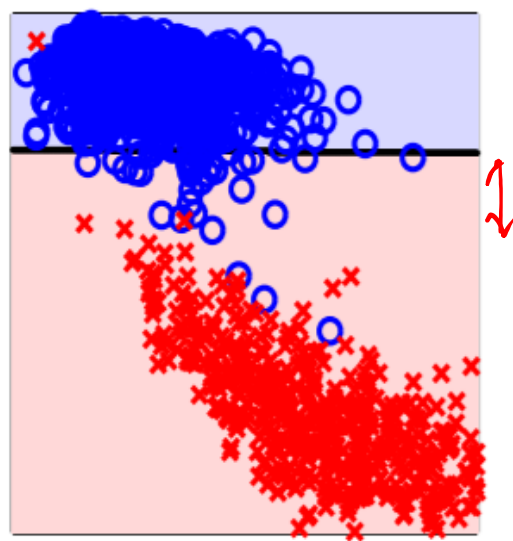
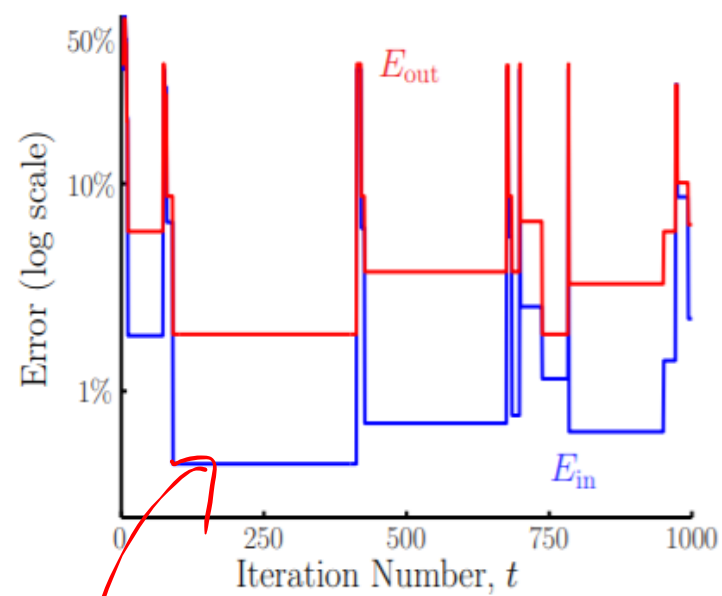
classifier



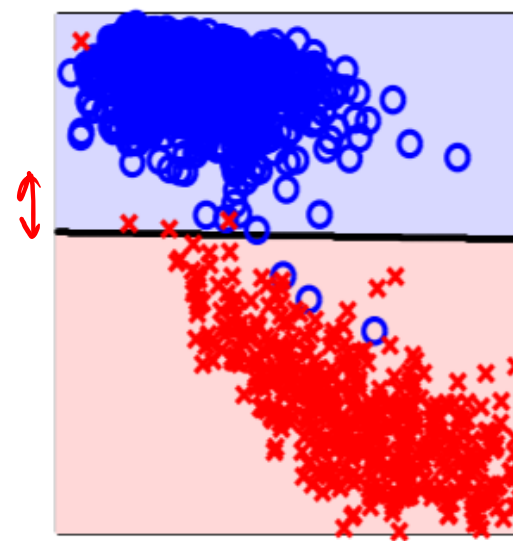
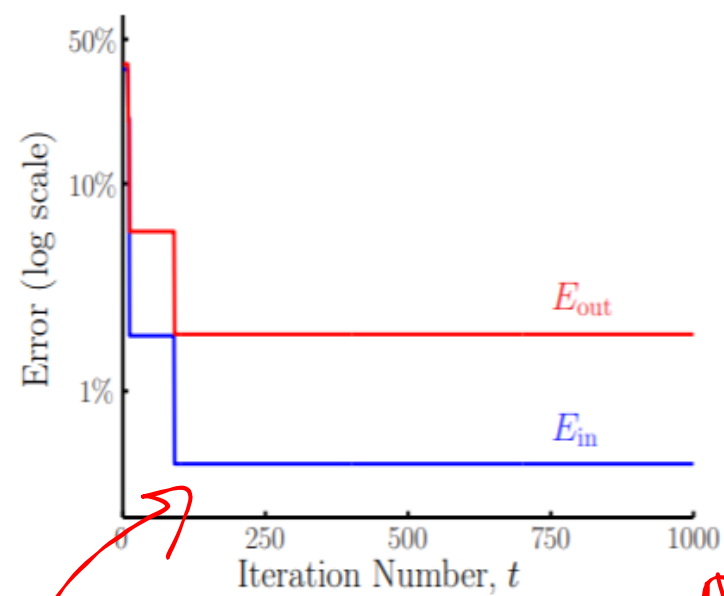
$d_{vc} = 3$

$E_{in} \approx E_{out}$   
tight

# PLA



# Pocket





# Linear Regression

(Credit Card)

•

age	32 years
gender	male
salary	40,000
debt	26,000
years in job	1 year
years at home	3 years
...	...

**Classification:** Approve/Deny

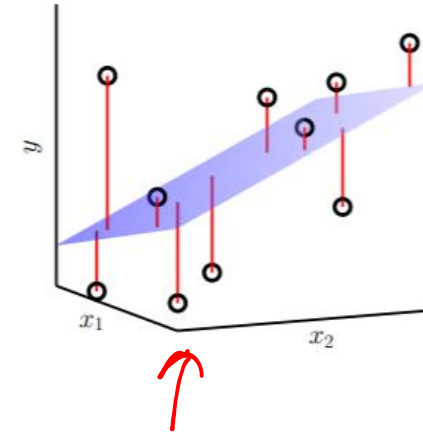
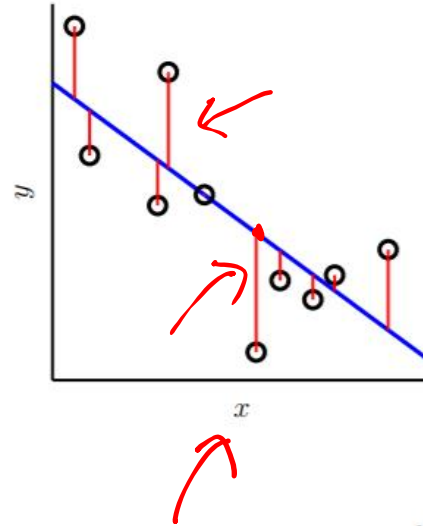
**Regression:** Credit Line (dollar amount)

regression  $\equiv y \in \mathbb{R}$

$$h(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$$

# Least Squares Linear Regression

$h(\mathbf{x})$



$$y = f(\mathbf{x}) + \epsilon$$

← noisy target  $P(y|\mathbf{x})$

in-sample error

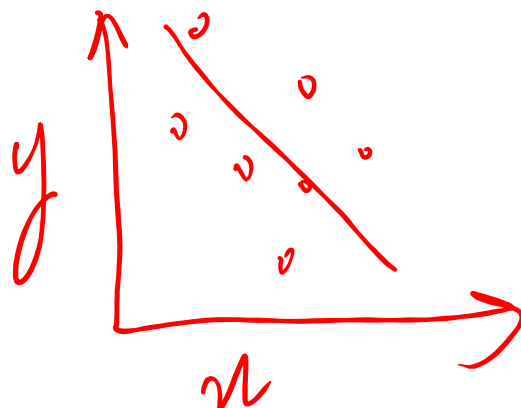
$$E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$$

→ MSE

out-of-sample error

$$E_{\text{out}}(h) = \mathbb{E}_{\mathbf{x}}[(h(\mathbf{x}) - y)^2]$$

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



$$h(x) = \boxed{\omega^T x}$$

$$E_{\text{in}}(\omega) = \frac{1}{N} \sum_{i=1}^N (h(x_i) - y_i)^2$$

(std. regression error)

$$= \frac{1}{N} \sum_{i=1}^N (\underbrace{\omega^T x_i}_{\hat{y}_i} - y_i)^2$$

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

$\left\{ \begin{array}{l} \hat{y}_i = \omega^T x_i \\ \text{error} \end{array} \right.$

MATRIX NOTATION

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \quad \underline{\hat{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}$$

$$X = \begin{bmatrix} \text{---} u_1^T \text{---} \\ \text{---} u_2^T \text{---} \\ \vdots \\ \text{---} u_n^T \text{---} \end{bmatrix}$$

$$X \in \mathbb{R}^{N \times (d+1)} \quad (\hat{y}_i - y_i)^2$$

$$E_{in}(\omega) = \frac{1}{N} \|\hat{y} - y\|^2$$

$$\begin{aligned}
 \hat{y}_1 &= \omega^T x_1 = x_1^T \omega \longrightarrow \left[ \text{---} x_1^T \text{---} \right] \begin{bmatrix} \omega \end{bmatrix} \\
 \hat{y}_2 &= \omega^T x_2 = x_2^T \omega \longrightarrow \left[ \text{---} x_2^T \text{---} \right] \begin{bmatrix} \omega \end{bmatrix} \\
 &\vdots \\
 \hat{y}_n &= \left[ \text{---} x_n^T \text{---} \right] \begin{bmatrix} \omega \end{bmatrix} \checkmark
 \end{aligned}$$

$\underbrace{\hspace{10em}}_{X \text{ (data)}}$

$$\hat{y} = X \omega \quad (\text{Verify})$$

$$\therefore E_{in}(\omega) = \frac{1}{N} \| X \omega - y \|^2 \longrightarrow \text{LA def.}$$

$$\|z\|^2 = z^T z$$

$$\begin{aligned} \therefore \underline{\underline{E_{in}(w)}} &= \frac{1}{N} (Xw - y)^T (Xw - y) \\ &= \frac{1}{N} (w^T X^T X w - \underbrace{w^T X^T y - y^T X w}_{-2w^T X^T y} + y^T y) \\ &= \frac{1}{N} (w^T X^T X w - 2w^T X^T y + y^T y) \end{aligned}$$

$$\therefore \text{Set } \nabla E_{in}(w) \rightarrow 0$$

$$E_{in}(w) = \frac{1}{N} \left( \underbrace{w^T X^T X w}_{\text{quadratic}} - \underbrace{2 w^T X^T y}_{\text{linear}} + \underbrace{y^T y}_{\text{constant}} \right)$$

$$\underline{ax^2 + bx + c} \rightarrow 2ax + b$$

$$\underline{\nabla E_{in}(w)} = \frac{1}{N} \left( 2X^T X w - 2X^T y \right) \rightarrow 0$$

Normal eqns.

$$\underbrace{X^T X w}_{\text{}} = X^T y$$

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

$X^T X$  is not invertible  
 $\rightarrow$  pseudo inverse

g

Pseudo Inverse Algo.

# Summarize Pseudo Inverse Algorithm

① Put your data in matrix form

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} \text{---} & x_1^T & \text{---} \\ \text{---} & x_2^T & \text{---} \\ \text{---} & x_n^T & \text{---} \end{bmatrix}$$

② Apply pseudo inverse Algorithm.

$$o(Nd^2 + d^3) \xrightarrow{\quad} \begin{aligned} w_{\text{lin}} &= (X^T X)^{-1} X^T y \\ w_{\text{lin}} &= X^+ y \end{aligned} \quad ; \quad X^+ = \underbrace{(X^T X)^{-1}}_{\text{learn}} X^T$$

③ 
$$g(x) = w_{\text{lin}}^T x$$
$$g(x_i) \approx y_i$$

$$E_{\text{out}}(g) = E_{\text{in}}(g) + O\left(\frac{d}{N}\right)$$



# Pocket Algorithm

$E_{in} \approx 0$  (Hard)

Pseudo Inverse Algorithm  $\rightarrow w_{lin}$

$y$ 's  $\rightarrow +1 -1 \rightarrow w_{lin}$

$$\boxed{\text{sign}(w_{lin}^T x_i) \approx y_i}$$

- $\rightarrow$  Get  $w_{lin}$  for classification problem.
  - $\rightarrow$  Make this your starting point
- Regression for Classification.

Thanks!