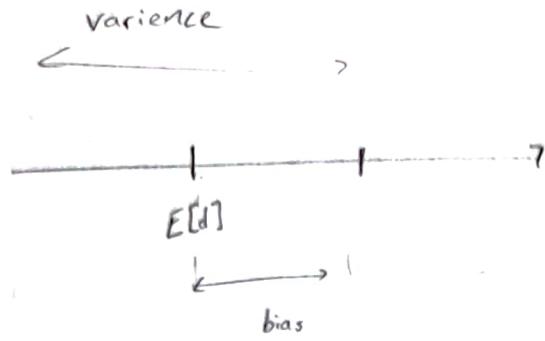


# Bias and Variance



$$\text{Bias} = E[d] - \theta_1$$

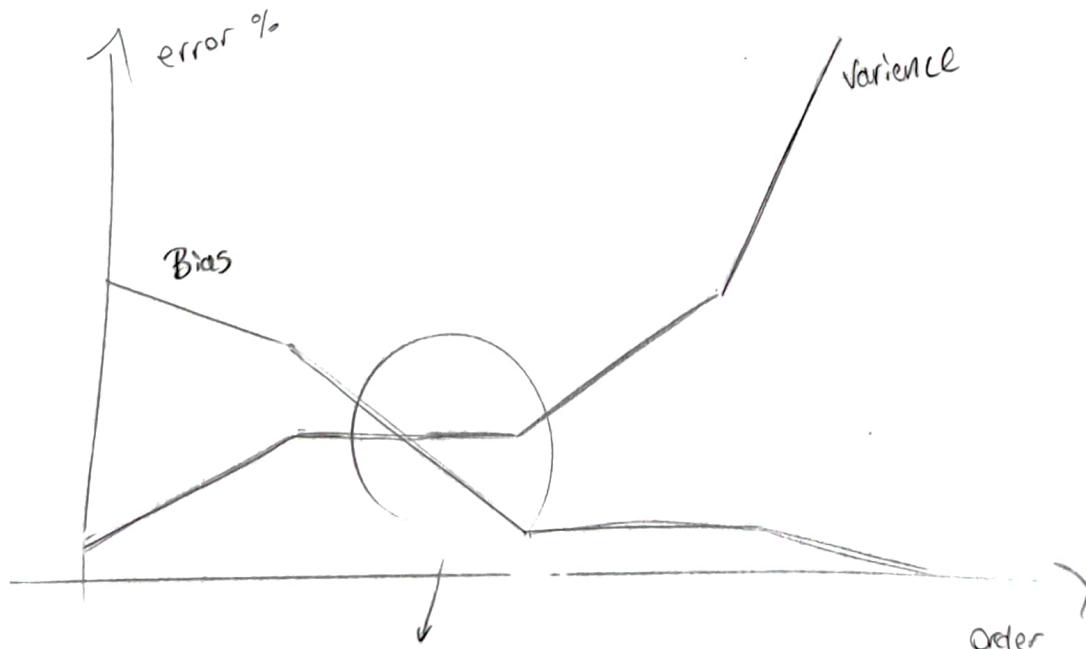
$$\text{Variance} = E[(d - E[d])^2]$$

mean square error

$$E[\theta - \theta_1]^2 = \underbrace{(E[d] - \theta_1)^2}_{\text{bias}^2} + \underbrace{E[(d - E[d])^2]}_{\text{Variance}}$$

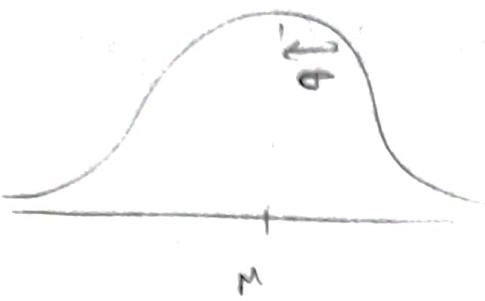
bias / variance dilemma

few parameters	high variance	/	lots of parameters
low bias	high bias		low variance



We use cross-validation to find this sweet spot.

## Parameters Estimation



$$x^* \sim N(\mu, \sigma^2)$$

Gaussian (Normal) Distribution

$$p(x) = N(\mu, \sigma^2)$$

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

maximum likelihood estimate

we will use it for classification

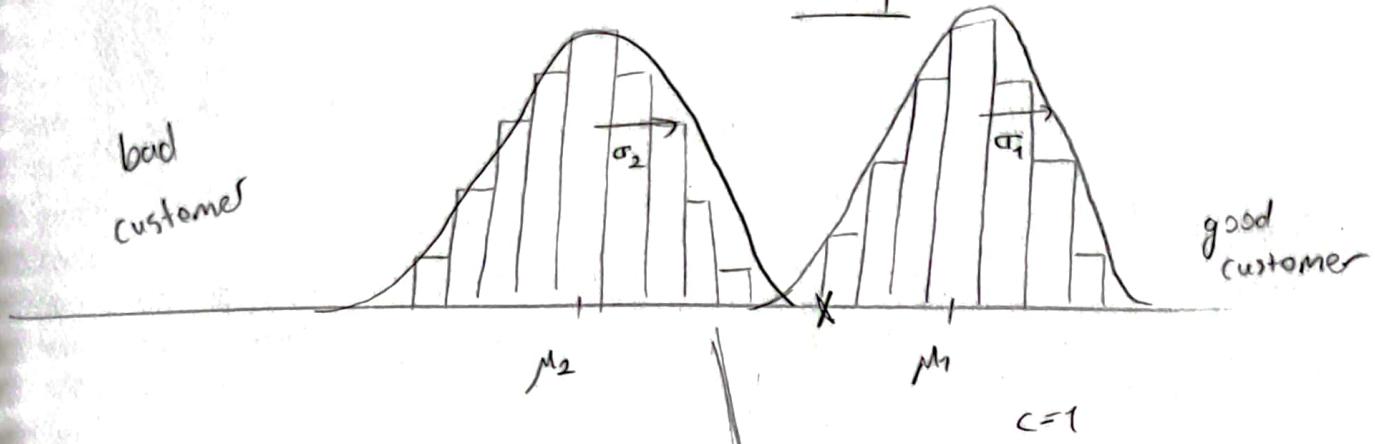
$$\underbrace{P(c/x)}_{\text{posterior prob.}} = \frac{\underbrace{P(x/c) \cdot P(c)}_{\text{likelihood}}} {\underbrace{P(x)}_{\text{evidence}}} \underbrace{P(c/x)}_{\text{prior}}$$

$$P(c/x) \propto P(x/c) P(c)$$

proportional

Bayes agent

Classification  
example



$c=0$

$$P(x | c=0)$$

$$\frac{\sum x^t}{N} = \mu_2$$

$$P(x | c=0) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left(-\frac{(x^t - \mu_2)^2}{2\sigma^2}\right)$$

$c=1$

$$P(x | c=1)$$

$$\frac{\sum x^t}{N} = \mu_1$$

$$P(x | c=1) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x^t - \mu_1)^2}{2\sigma^2}\right)$$

$$P(c_i | x) =$$

$$g_i(x) = \log P(x | c_i) \quad P(c_i) = \log P(x | c_i) + \log P(c_i)$$

$$g_i(x) = \log \left[ \frac{1}{\sqrt{2\pi\sigma_i^2}} \cdot \exp\left(-\frac{(x - \mu_i)^2}{2\sigma^2}\right) \right] + \log P(c_i)$$

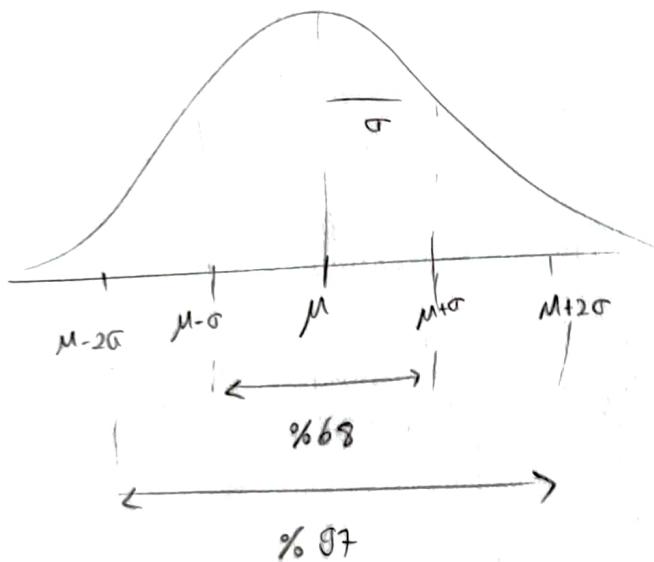
Error function = bias<sup>2</sup> + variance

	Bias	Variance
Simple model	large	low
2nd		↓
3rd		↑

To find the minimum error; we use cross-validation.

a brief explanation ↑ of last lesson

Gaussian Dist.  
(Normal)



$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

it gives a density  
not probability.

To achieve probability we  
integrate for that point.

$x_1$	$x_2$	$c$
$x_1^{(1)}$	$x_2^{(1)}$	
$x_1^{(2)}$	$x_2^{(2)}$	
$\vdots$	$\vdots$	
$x_1^{(n)}$	$x_2^{(n)}$	

$$P(c=1 | x_1, x_2) = \frac{P(x_1, x_2 | c) \cdot P(c)}{P(x_1, x_2)}$$

$$P(x_1, x_2) = P(x_1) \cdot P(x_2) \quad (\text{if they are independent})$$

$$x_1 \sim N(\mu_1, \sigma_1^2)$$

$$x_2 \sim N(\mu_2, \sigma_2^2)$$

$$P(x_1, x_2) = \left( \frac{1}{\sqrt{2\pi}\sigma_1} \right) \cdot \exp \left( -\frac{1}{2} \cdot \frac{(x_1 - \mu_1)^2}{\sigma_1^2} \right) \cdot \left( \frac{1}{\sqrt{2\pi}\sigma_2} \right) \cdot \exp \left( -\frac{1}{2} \cdot \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right)$$

$$\frac{1}{2\pi\sigma_1\sigma_2} \cdot \exp \left[ \frac{1}{4} \cdot \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right]$$

$$\begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$m = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$(x-m)^T \cdot (x-m)$$

$$\begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 \end{bmatrix} \underbrace{\begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix}}_{\Sigma^{-1}} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}$$

covariance matrix

$$\Sigma^{-1}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$\Sigma = \Sigma^T \quad (\text{symmetric})$$

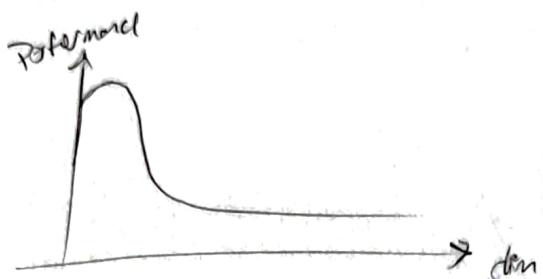
\* Dimensionality Reduction → reduces

- time complexity
- space complexity
- cost of observing feature

↓  
data visualization

## curse of dimensionality

there is a point on number of features that the bigger we go above that point, performance starts to degrading instead of improving.



## feature selection vs extraction

### class labels

Selection → correlation with class labels

more  
less redundancy

\* minimum redundancy  
\* maximum relevant

mRMR algorithm

extraction → project the original  $d$  to new  $k$  dim

2d

\* Principal Component

\* Analysis

PCA

(unsupervised model)

NO class labels

eigen  
natur latwors

$A_{d \times d}$

$$Ax = \lambda x$$

$$(A - \lambda I)x = 0$$

$$\det(\lambda I - A) = 0$$

$$A \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 & x_2 & \dots & x_n \\ f_1 & f_2 & \dots & f_m \end{bmatrix}$$

$$A = X \Sigma X^{-1} \quad (\text{if } X^{-1} \text{ exists})$$

$$\begin{array}{c} \downarrow \text{columns} \\ \xrightarrow{\text{rows}} \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \\ \text{IMDB} \quad \uparrow \quad \uparrow \quad \uparrow \quad \text{movies} \\ \text{users} \quad \text{titles} \quad \text{titles} \quad \text{genres} \end{array}$$

matrix factorization

+ singular value decomposition

$$A_{n \times d} = U \Sigma V^T$$

$\sum_{\text{diagonal}}$   
 $U \perp V$

$$U \cdot U^T = I$$
$$V \cdot V^T = I$$

$$A^T \cdot A = (U \Sigma V^T)^T \cdot U \Sigma V^T = U \Sigma V^T \cdot U \Sigma V^T = \Sigma^2 V^T$$

$$A^T \cdot A = B = \Sigma^2 V^T$$

$$A \cdot A^T = C = U \Sigma^2 U^T$$

## Subset selection

- forward selection

add the best feature at each step until it stops getting better

$$f_1 \rightarrow f_1, f_4 \rightarrow f_1, f_4, f_7 \rightarrow \dots$$

$O(d^2)$

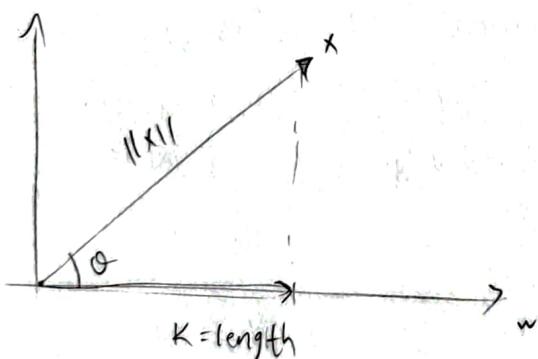
- backward selection

same as forward but reverse

$$f_1, f_2, \dots, f_d \rightarrow f_1, f_2, f_3, \dots, f_d \rightarrow \dots$$

## Principal Component Analysis

Project the data on to  $z$  where the variance is maximized.



$$\cos \theta = \frac{K}{\|x\|}$$

$$x^T w = \|x\| \cdot \|w\| \cdot \cos \theta$$

$$w^T x = \|x\| \cdot \|w\| \cdot \frac{K}{\|x\|}$$

$$K = \frac{w^T x}{\|w\|} \approx 1$$

$$K = w^T x$$

after projection:  $z = w^T x$

$$\|w\| = 1$$

find the  $w$  value that maximizes  $z$ .

$$\arg \max E \left[ (w^T x - w^T \mu)^2 \right]$$

$$= (w^T x - w^T \mu) \cdot (w^T x - w^T \mu) = (w^T x - w^T \mu) (x^T w - \mu^T w)$$

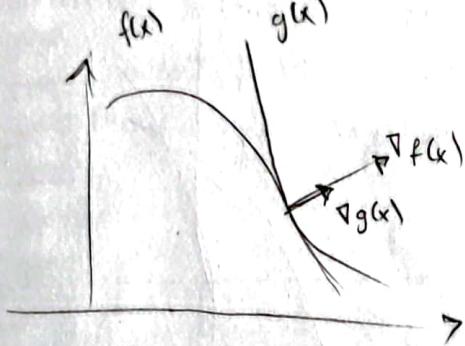
$$= w^T (x - \mu) (x^T - \mu^T) w = E \left[ w^T (x - \mu) (x - \mu)^T w \right]$$

$$w^T \cdot E \left[ (x - \mu) (x - \mu)^T \right] w$$

subject to  $\|w\|=1$

$\Sigma \rightarrow$  covariance matrix

constraint optimization



Both functions have the gradients that have same direction.

Lagrange Multipliers

$$\nabla f(x) = \lambda \nabla g(x)$$

$$\nabla f(x) - \lambda \nabla g(x) = 0$$

$$\frac{\partial L(w, \lambda)}{\partial w} = 0$$

partial  
check out: derivative of vector (least square)

$$\Sigma w = \lambda w$$

$\downarrow$   
covariance matrix  
largest eigenvalue  
corresponding eigenvector

$$\begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_{N \times d}$$

$$\Sigma_{d \times d}$$

d eigen values

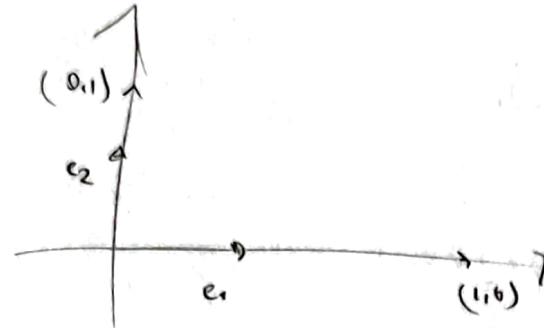
$$d_1, d_2, d_3, \dots, d_d$$

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}_{d \times 1}$$

$$z_i = x \cdot w_i$$

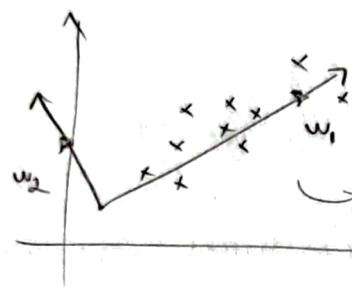
$$z_k = x \cdot w_k$$

\* PCA



$$[u] = 5e_1 + 4e_2$$

$e_1$  and  $e_2$  are bases of dimensions



these eigenvectors can be our new bases so that we can visualize data better and that would reduce dimensions to whatever we want.

### -Covariance Matrix-

x	y
1	1
1	1
1	1

# samples = N

mean of X =  $\mu_x$

mean of Y =  $\mu_y$

Step 1: find  $\mu_x$  and  $\mu_y$

Step 2: var of X  $\rightarrow$   $\frac{[(\mu_x - x_1)^2 + (\mu_x - x_2)^2 + \dots + (\mu_x - x_n)^2]}{N}$

Step 3: var of Y  $\rightarrow$  - - -

$$\text{cov}(x,y) = \frac{[(\mu_x - x_1) \cdot (\mu_y - y_1) + (\mu_x - x_2) \cdot (\mu_y - y_2) + \dots + (\mu_x - x_n) \cdot (\mu_y - y_n)]}{(N-1)}$$

$$\text{cov} \sum = \begin{bmatrix} \text{var}(x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(x,y) & \text{var}(y) & \text{cov}(y,z) \\ \text{cov}(x,z) & \text{cov}(y,z) & \text{var}(z) \end{bmatrix}$$

PCA → covariance matrix ( $\leq$ )

→ eigenvalue and eigenvectors of  $\leq$

$$A \cdot x = \lambda x$$

→ normalize eigenvector

$$\rightarrow e_1 = \begin{bmatrix} a \\ b \end{bmatrix} \rightarrow e_1 = \begin{bmatrix} \frac{a}{\sqrt{a^2+b^2}} \\ \frac{b}{\sqrt{a^2+b^2}} \end{bmatrix}$$

→ derive dataset

$$\begin{array}{|c|c|} \hline x & y \\ \hline 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ \hline \end{array}$$

$$P_1 \quad P_2$$



$$P_{11} = e_1^\top \begin{bmatrix} x_i - \bar{x} \\ y_i - \bar{y} \end{bmatrix}$$

$$P_{12} = \dots$$

SVD

$$\text{let } A = \begin{bmatrix} 3 & 3 & 2 \\ 2 & 3 & -2 \end{bmatrix}$$

$$\textcircled{1} \quad A \cdot A^\top = \begin{bmatrix} 3 & 3 & 2 \\ 2 & 3 & -2 \end{bmatrix} \cdot \begin{bmatrix} 3 & 2 \\ 3 & 3 \\ 2 & -2 \end{bmatrix} = \begin{bmatrix} 17 & 8 \\ 8 & 17 \end{bmatrix}$$

$$A \cdot A^\top - \lambda I = 0$$

$$w - \lambda I = 0$$

$$\hookrightarrow \lambda^2 - 34 + 225 = 0$$

$$(\lambda - 25)(\lambda - 9) = 0$$

$$\lambda_1 = 25 \quad \lambda_2 = 9$$

$$\textcircled{2} \quad A \cdot A^\top - \lambda_1 I = \begin{bmatrix} 1 & -10 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$



$$y = w^\top \times$$



$$w = [e_1 \ e_2]$$

regression  $\rightarrow$  gradient descent

$$g_i(x) \rightarrow$$

hypothesis  $\rightarrow$

bias = variance

normal dist

$$\frac{1}{2}$$

green 0,1

$$S=1$$

terorist = 0,1

citizen = 0,1

$$P(S=1 \mid T=1) = 3\%$$

$$P(S=0 \mid T=0) = 97\%$$

$$P(T=1 \mid S=1) = \frac{P(S=1 \mid T=1) \cdot P(T)}{P(S=1)}$$

$$P(T=1 \mid S=1) = \frac{3\% \cdot 1/100}{1} =$$

$$P(T=1) = P(T=1 \mid S=0) + P(T=1 \mid S=1)$$

$$89\%$$

$$P(S=1) = P(S=1 \mid T=1) \cdot P(T=1) + P(S=1 \mid T=0) \cdot P(T=0)$$

$$P(A=0 \mid B=0) = \%x \rightarrow P(A=1 \mid B=0) = \% (100-x)$$

$$P(A=1) = P(A=1 \mid B=0) \cdot P(B=0) + P(A=1 \mid B=1) \cdot P(B=1)$$

## Gradient Des

$$w^{k+1} = w^k - \alpha \cdot \nabla g(w^k)$$

linear  $D_m = \frac{-2}{n} \sum_{i=0}^n x_i(y_i - \bar{y}_i)$

$$D_C = \frac{-2}{n} \sum_{i=0}^n (y_i - \bar{y}_i)$$

bias - variance

$$\text{error } f = \text{bias}^2 + \text{variance}$$

$$E[(d - \hat{\phi}_i)^2] = \underbrace{(E[d] - \phi_i)^2}_{\text{bias}} + E[(d - E[d])^2]$$

Gaussian

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

$$g_i(x) = \log \left[ \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left[-\frac{(x-\mu_i)^2}{2\sigma_i^2}\right] \cdot P(c_i) \right]$$

$$\frac{d}{dx}(a^x) = a^x \cdot \ln a$$

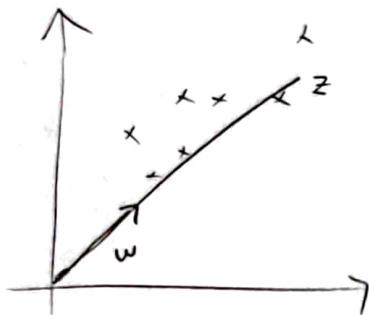
$$\frac{d}{dx}(\log_a g(x)) = \frac{g'(x)}{g(x) \cdot \ln a}$$

$$+\frac{d}{dx}(e^{g(x)}) = e^{g(x)} \cdot g'(x)$$

$$+\frac{d}{dx}(a^{g(x)}) = \ln(a) \cdot a^{g(x)} \cdot g'(x)$$

PCA

After Midterm



$$z = w^T \cdot x \quad , \quad \|w\|=1$$

$\arg \max$  variance

$$\arg \max \underbrace{E[(z - E(z))^2]}_{}$$

$$E \left[ (w^T \cdot x - \underbrace{E[w^T \cdot x]}_{\mu})^2 \right]$$

$$\begin{aligned} L(w, d) &= E \left[ (w^T \cdot x - E[w^T \cdot x])^2 \right] - \lambda \cdot w^T \cdot w \\ &= E \left[ (w^T \cdot x - \underbrace{w^T \cdot \mu}_{\mu})^2 \right] - \lambda \cdot w^T \cdot w \end{aligned}$$

$w^T E[x] \rightarrow x$  is variable  
 $\mu$

$$\begin{bmatrix} w_1, w_2, \dots, w_n \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \frac{w_1 \cdot x_1 + w_2 \cdot x_2 -}{\dots + w_n \cdot x_n} = [x_1, x_2, \dots, x_n] \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} w$$

Hatrlatma  $\rightarrow w^T \cdot x = x^T \cdot w$

$$E \left[ (w^T \cdot x - w^T \cdot \mu) \cdot (w^T \cdot x - w^T \cdot \mu) \right]$$

$$E \left[ w^T (x - \mu) \cdot (x^T \cdot w - \mu^T \cdot w) \right]$$

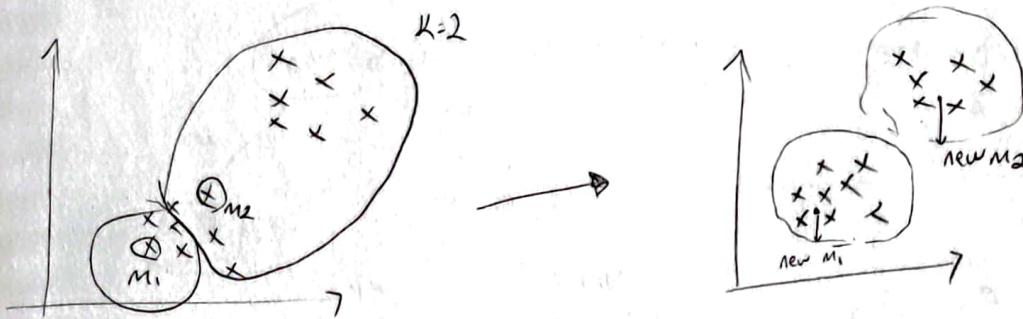
$$E \left[ w^T (x - \mu) \cdot (x^T - \mu^T) \cdot w \right]$$

$$\underbrace{w^T \cdot E \left[ (x - \mu) \cdot (x^T - \mu^T) \right] \cdot w}_{\sum} \quad \boxed{L(w, d) = w^T \cdot \Sigma \cdot w - \lambda \cdot w^T \cdot w}$$

## Unsupervised Learning

- unlabelled data
  - discovers groups/clusters or patterns within the data
  - dimensionality reduction, clustering → two most popular.
- classes vs clusters

### K-means clustering



\*\* It is an iterative algorithm not a closed form formula.

Initialization of  $m_i$ 's:

- Random labelling (not recommended)
  - Mean + small random vectors
  - Random data points as centers
  - PCA and take the means of these groups.
- \* Multiple restart for more accurate results.

drawbacks of K-means

it may stuck in local minima



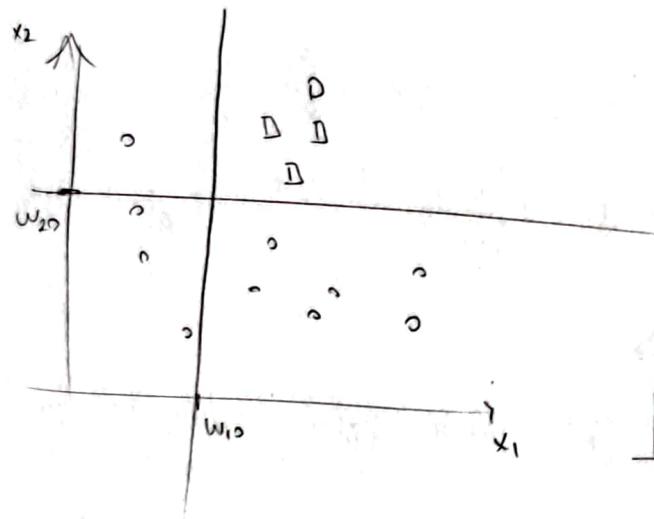
$$\frac{d}{dx} (a^{g(x)}) = a^{g(x)} \cdot \ln(a) \cdot g'(x)$$

Anil K. Jain  
youtube talk

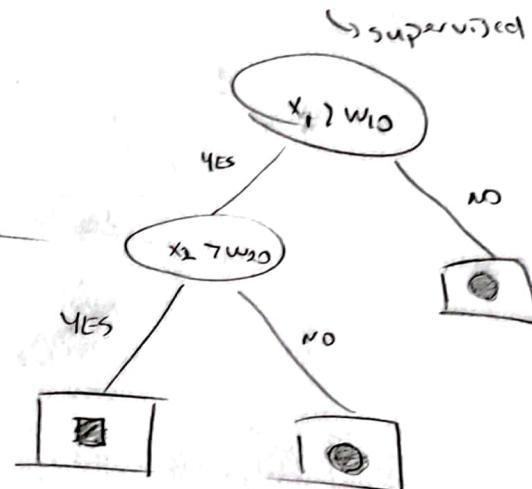
$$\frac{d}{dx} (\log_a g(x)) = \frac{g'(x)}{g(x) \cdot \ln(a)}$$

HATIRJIMA

Trees / Decision Trees



(classification trees)



splitting may vary

→ YES NO  
→  $x=10, x=12, x=16$

N cases

$$\text{Gini Index} = P \cdot (-P) \rightarrow$$

$x_1$

3	1+	5
1+	2+	1+
0-	1-	3-

	$x_1$	$x_2$	$x_3$
+	3	1	m
+	4	y	f
+	4	y	m
+	5	y	f
-	5	1	f
-	5	y	f
-	5	1	m
-	5	y	m

$$Gini = \frac{1}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} + \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{3}{4}$$

Impurity

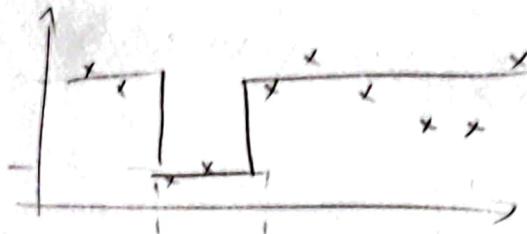
$$(\text{expected value}) \quad \frac{1}{4} \cdot \frac{1}{8} + \frac{2}{4} \cdot \frac{3}{8} + \frac{3}{4} \cdot \frac{4}{8}$$

$x_2$

$$\begin{matrix} 1+ \\ 3- \\ 1+ \\ 1- \\ \frac{1}{4} \cdot \frac{3}{4} \\ \frac{3}{4} \cdot \frac{1}{4} \end{matrix}$$

$$Imp = \frac{3}{16} \cdot \frac{4}{8} + \frac{3}{16} \cdot \frac{4}{8}$$

## Regression Trees



we choose a cluster pair  
that minimizes the Error f.

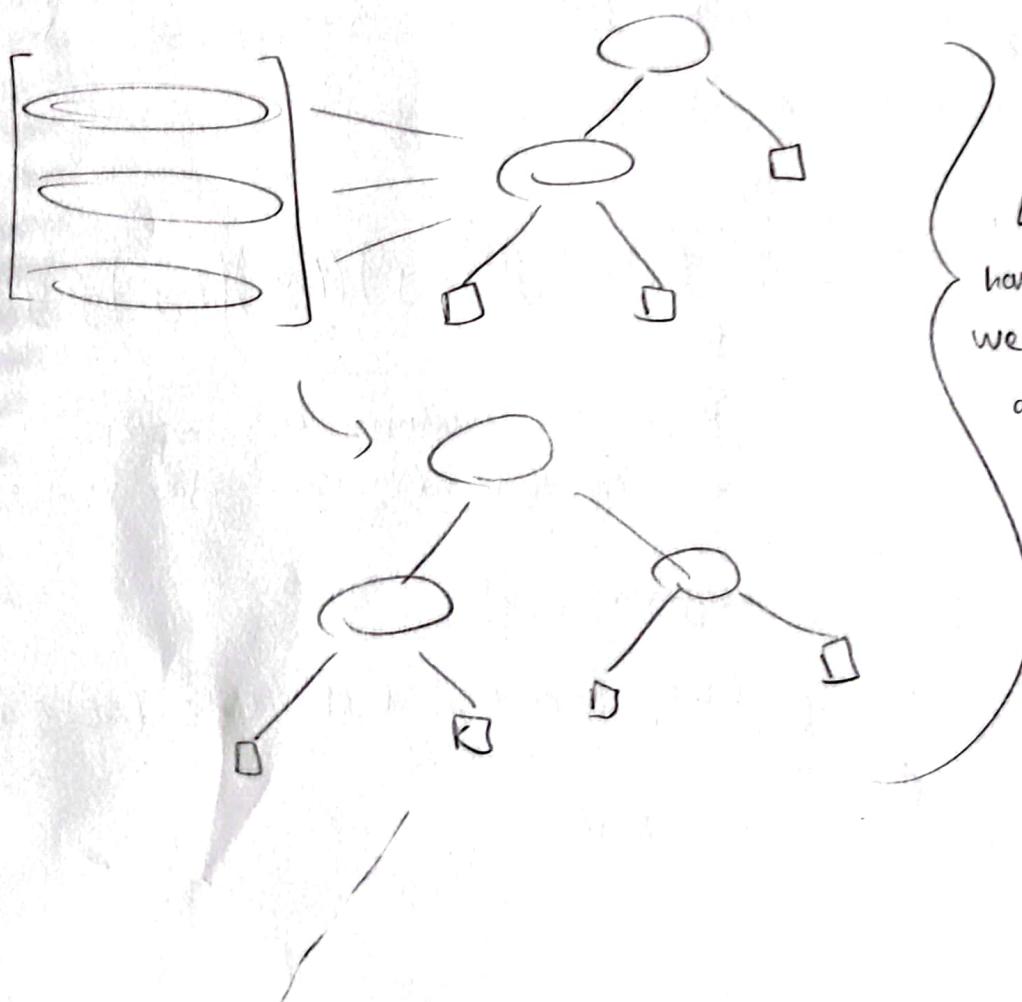
## Pruning trees

after constructing tree, remove subtrees for better generalization. (avoid overfitting)

decrease variance  $\rightarrow$  increase bias

## Random Forest

Randomizing features and creating a model with them.  
then combining all of them to get a good estimation.

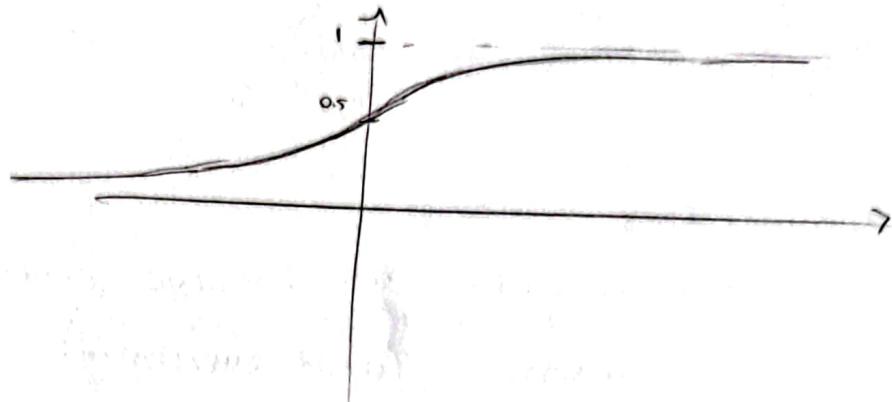


Let's say we have a new data we put our new data to all the decision trees and get the weighted result.

# NEW WEEK!

(Logistic Sigmoid function)

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



$$g'(s) = \frac{e^{-s}}{(1+e^{-s})^2} = \frac{1}{1+e^{-s}} \cdot \frac{e^{-s}}{1+e^{-s}} = g(s) \cdot (1-g(s))$$

$$\prod (y^t)^{r_t} (1-y^t)^{1-r_t}$$

Bernoulli  
remander

$$P(X=x) = \begin{cases} p & \text{for } x=1 \\ 1-p & \text{for } x=0 \end{cases}$$

$x_1$	$x_2$	$r$	$y$
0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0
1	1	1	1
1	1	1	1

$$l(\theta) = (1-y_1)(1-y_2)(1-y_3)(1-y_4)y_5^5y_6^6y_7^7y_8^8$$

If we can maximize this expression we can get a close classification.

$$\operatorname{argmax} l(\theta)$$

To do this we use gradient descent. But this time for maximize.

$$w = w + \nabla w$$

$$E = \log \prod_{t=1}^N (y^t)^{r^t} \cdot (1-y^t)^{1-r^t}$$

$$y^t = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2)}}$$

$$E = \sum_{t=1}^N r^t \log(y^t) + (1-r^t) \log(1-y^t) \rightarrow \text{maximize this} \\ (\log likelihood)$$

gradient descent

\* or minimize  $-\log \text{likelihood}$

$w_0, w_1, w_2 \leftarrow$  initialize small numbers around 0

$$w_0 = w_0 - \eta \frac{\partial E}{\partial w_0}$$

$$w_1 = w_1 - \eta \frac{\partial E}{\partial w_1}$$

$$\frac{\partial E}{\partial w_1} =$$

$$\frac{d(\log_a g(x))}{dx} = \frac{g'(x)}{g(x) \cdot \ln a}$$

$$\frac{d(a^{g(x)})}{dx} = a^{g(x)} \cdot g'(x) \cdot \ln a$$

HATIRLATMA

$$-\sum \frac{r^t \cdot (y^t)^{r^t}}{y^t} + (1-r^t) \cdot \frac{(1-y^t)^{1-r^t}}{1-y^t} = \frac{r^t y^t (1-y^t)}{y^t} + \frac{(1-r^t) \cdot (1-y^t) (y^t)}{1-y^t}$$

$$\sum (r^t - y^t) x_1^t$$

example

$x_1$	$x_2$	$r$
2	3	0
1	0	0
5	4	1
2	4	1

$$w_0 = 0.1 \quad w_1 = 0.1 \quad w_2 = 0.05$$

random initialize around 0

$$\eta = 0.1$$

$$y^1 = \frac{1}{1 + e^{-(0.1 + 0.1 \cdot x_1 + 0.05 \cdot x_2)}} = \frac{1}{1 + e^{-0.45}} = 0.61$$

$$y^2 = \frac{1}{1 + e^{-0.2}} = 0.54$$

$$y^3 = \frac{1}{1 + e^{-0.8}} = 0.68$$

$$y^4 = \frac{1}{1 + e^{-0.5}} = 0.62$$

$$w = w - \eta \nabla w$$

$$w = w - \eta \sum (r^t - y^t) \cdot x_1^t$$

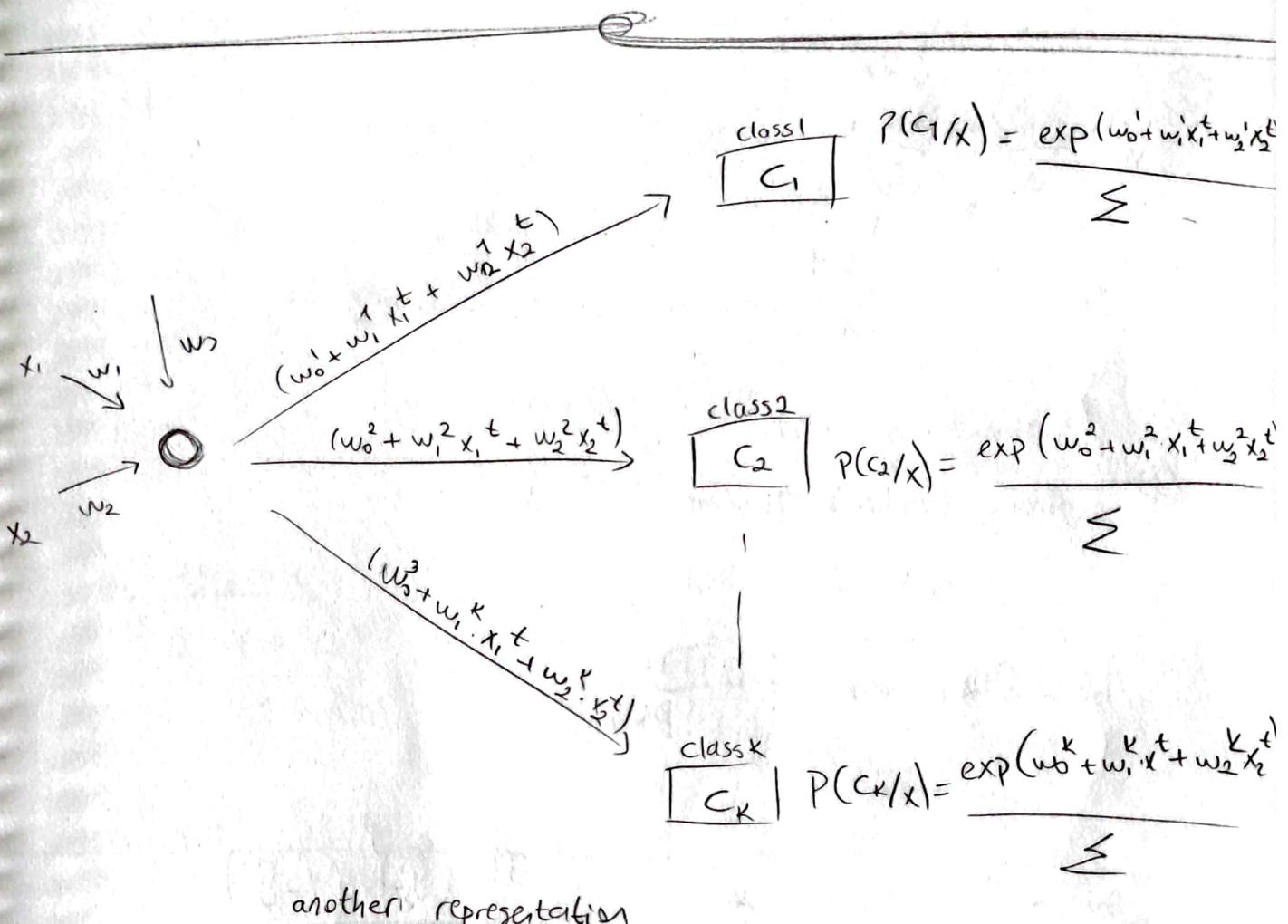
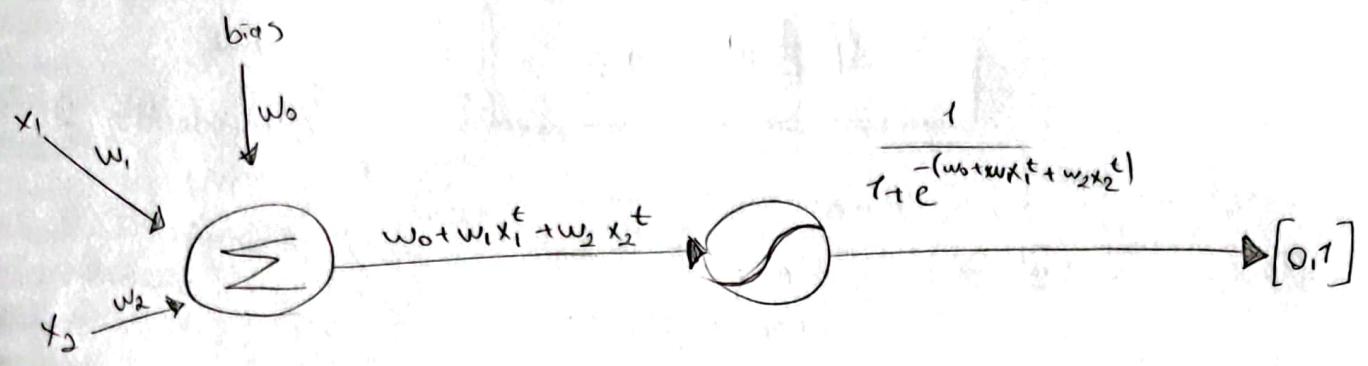
$$w_1 = w_1 - \eta \left[ (r^1 - y^1) x_1^1 + (r^2 - y^2) x_1^2 + (r^3 - y^3) x_1^3 + (r^4 - y^4) x_1^4 \right]$$

$$w_1 = w_1 - 0.1 \cdot [(0 - 0.61) \cdot 2 + (0 - 0.54) \cdot 1 + (1 - 0.68) \cdot 5 + (1 - 0.62) \cdot 2]$$

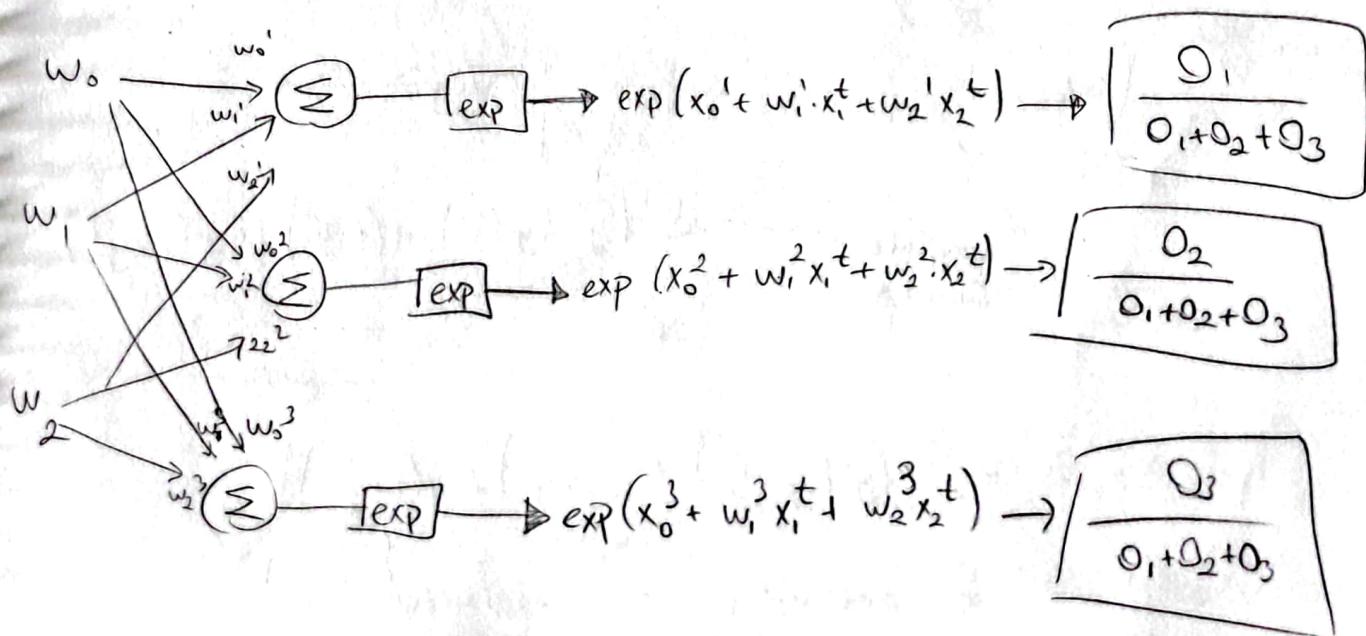
$$\frac{\text{new}-\text{old}}{0.04} = 0.1 - \frac{0.1 \cdot 0.6}{0.06}$$

and same goes for  $w_2$   
 $w_0$  goes same as well but  
 there is no ' $x$ ' term.

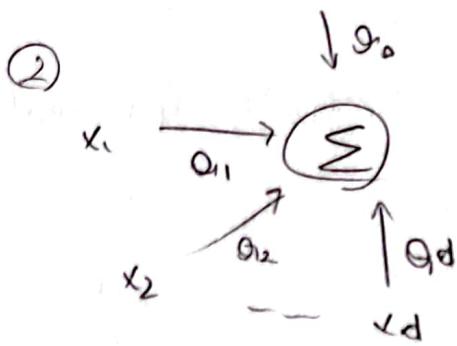
goes with converge



another representation



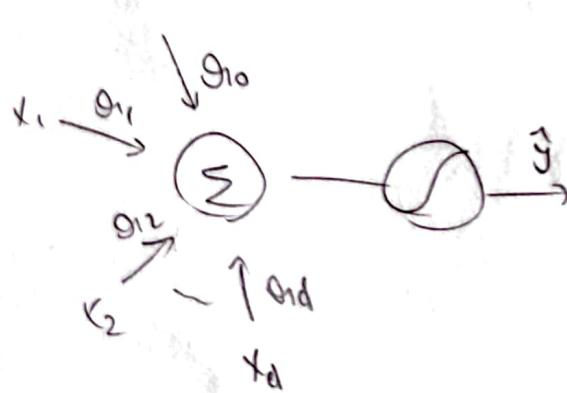
$$\textcircled{1} \quad \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \text{MSE} \quad (\text{error function})$$



\textcircled{3} Apply Gradient Descent

$$\theta_0, \theta_1, \theta_2, \dots, \theta_d$$

$$\theta_{k+1} = \theta_k - h \cdot \frac{\partial E}{\partial \theta_k}$$



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

$$\log$$

maximize likelihood  $\leftarrow \sum (y_i \cdot \log(\hat{y}_i) + (1-y_i) \cdot \log(1-\hat{y}_i))$

OR

$$-\sum (y_i \cdot \log(\hat{y}_i) + (1-y_i) \cdot \log(1-\hat{y}_i))$$

cross entropy

(minimize this)

$$O_1 = O_{101} + O_{111} \cdot x_1 + O_{121} \cdot x_2 + \dots + O_{1d1} \cdot x_d$$
  

$$O_2 = O_{202} + O_{212} \cdot x_1 + O_{222} \cdot x_2 + \dots + O_{2d2} \cdot x_d$$
  

$$O_3 = O_{303} + O_{313} \cdot x_1 + O_{323} \cdot x_2 + \dots + O_{3d3} \cdot x_d$$

to normalize these output values:

$$\hat{y}_k = \frac{\exp(O_k)}{\sum_{i=1}^N \exp(O_i)} \quad (\text{softmax function})$$

In this case

$$\hat{y}_1 = \frac{\exp(O_1)}{\sum_{i=1}^3 \exp(O_i)}$$

$$\prod_{t=1}^N \prod_{i=1}^3 (\hat{y}_i^{t_i}) = \text{likelihood}$$

$$mse = \sum (\hat{y}_i - y_i)^2$$

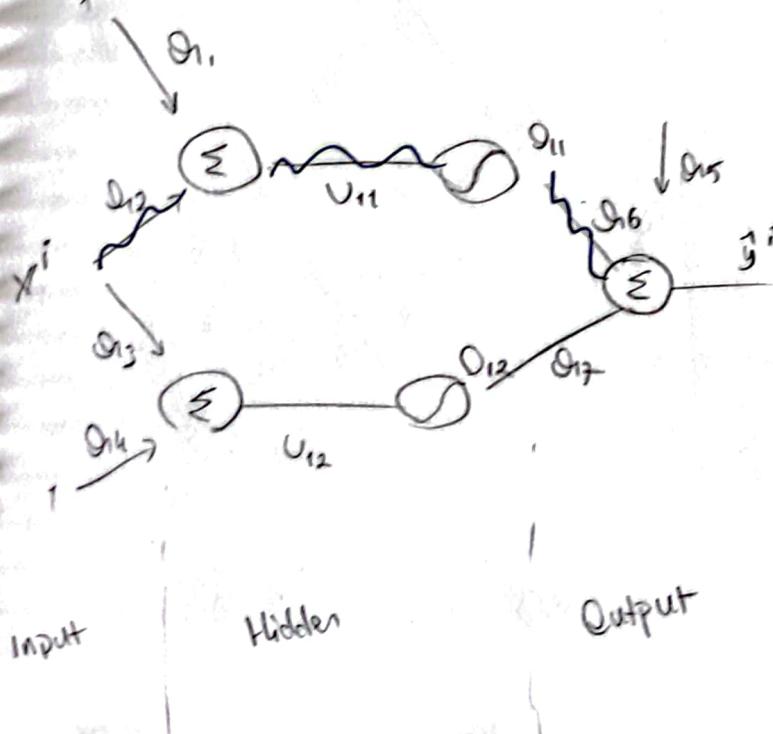
$$\hat{y} = O_{15} + O_{16} \cdot O_{11} + O_{17} \cdot O_{12}$$

$$O_{11} = \frac{1}{1 + e^{-u_{11}}}$$

$$O_{12} = \frac{1}{1 + e^{-u_{12}}}$$

$$u_{11} = O_{11} + O_{12} \cdot x^i$$

$$u_{12} = O_{14} + O_{13} \cdot x^i$$



$$\frac{\partial E}{\partial \theta_{12}} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial u_{11}} \cdot \frac{\partial u_{11}}{\partial \theta_{12}} \quad (\text{chain rule})$$

$$\sum (\hat{y}_i - y_i)$$

$$o_{16}$$

$$o_{11}, (1-o_{11})$$

$$x_i$$

REMINDER

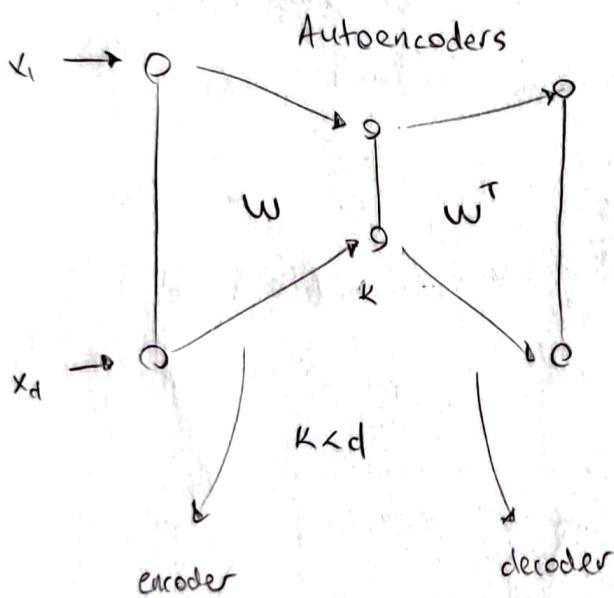
$$s \rightarrow g(s)$$

$$\frac{\partial g(s)}{\partial s} = g'(s)$$

$$g(s) = \frac{1}{1+e^{-s}} \quad g'(s) = \frac{e^{-s}}{(1+e^{-s})^2}$$

$$g(s) = g(s), (1-g(s))$$

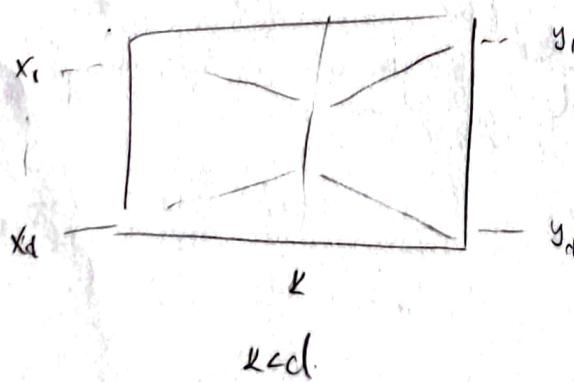
1 output  $\rightarrow$  linear regression  
( $\hat{y}$ )



$$[x]_{100,10}$$

$N \times d$

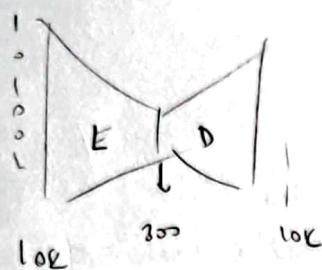
$$w^T \rightarrow k$$



$$\begin{array}{c|ccc} & y_1 & & \\ \hline x_1 & & y_2 & y_3 \\ x_2 & & x_3 & \end{array}$$

- dimensionality reduction

Word 2 Vec  $\rightarrow$  [ ]  
300 vector

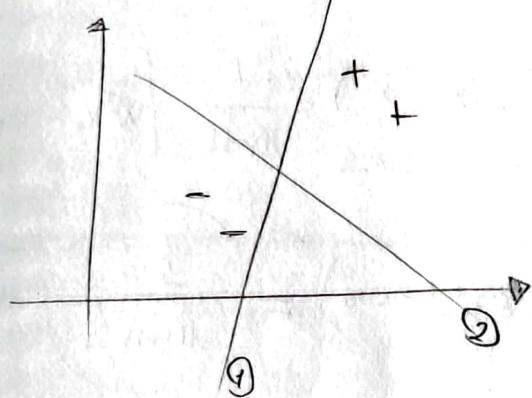


I have lost the will to live --

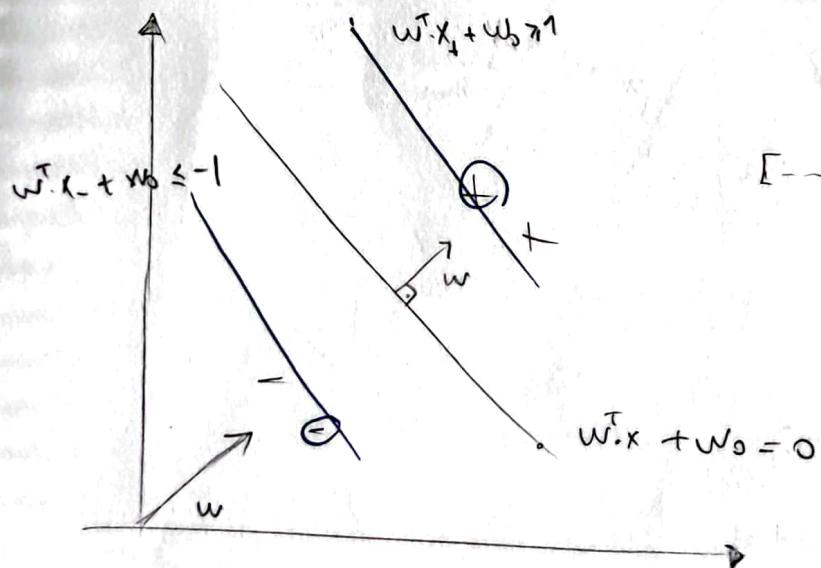
(one hot  
encoding attention)

translation  $\rightarrow$  RNN

### Support Vector Machines -SVM-



2 is better  
we try to maximize the  
distance between boundary  
and data.



$$w^T x + w_0 \geq 1$$

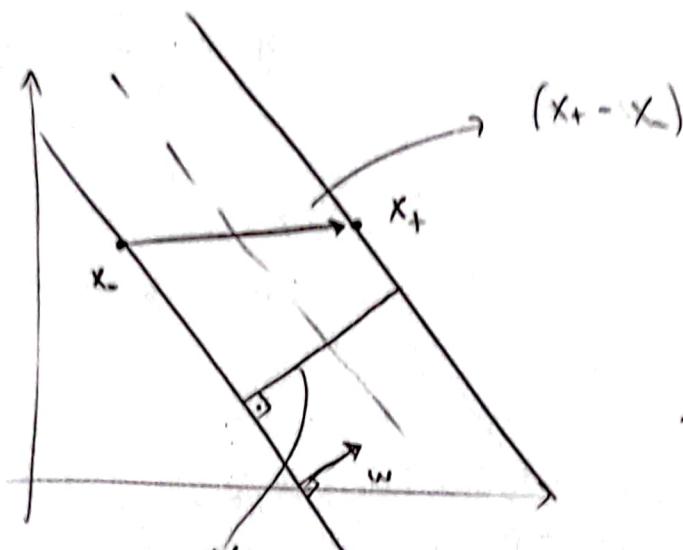
$$[-\dots]. \begin{bmatrix} 1 \\ 1 \end{bmatrix} + b$$

$w$  vector will be  
perpendicular to decision  
boundary.

$$\begin{aligned} w^T x_+ + w_0 &\geq 0 \\ w^T x_- + w_0 &\leq 0 \end{aligned}$$

by definition our dataset is separable.

$$y^i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$$



$$\mathbf{w}^T \mathbf{x}_+ + w_0 = 1$$

$$\mathbf{w}^T \mathbf{x}_- + w_0 = -1$$

project  $(\mathbf{x}_+ - \mathbf{x}_-)$  vector onto  $\mathbf{w}$  vector

maximize DISTANCE

$$\begin{array}{c} \text{if } \\ \mathbf{a} \\ \text{is} \\ \mathbf{w} \end{array} \Rightarrow \frac{\mathbf{a}^T \mathbf{w}}{\|\mathbf{w}\|}$$

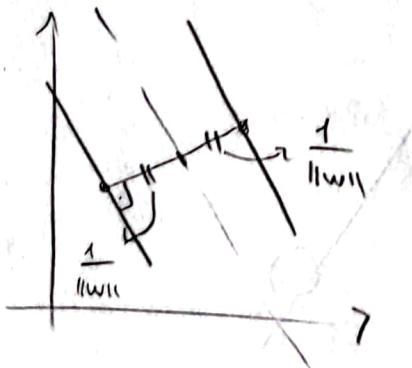
$$\text{Distance} = \frac{\mathbf{w}^T (\mathbf{x}_+ - \mathbf{x}_-)}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|} \cdot [\mathbf{w}^T \mathbf{x}_+ - \mathbf{w}^T \mathbf{x}_-]$$

$$= \frac{1}{\|\mathbf{w}\|} \cdot (1 - w_0 + 1 + w_0) = \left( \frac{2}{\|\mathbf{w}\|} \right)$$

$$\text{so maximize } \frac{1}{\|\mathbf{w}\|} = \frac{1}{\sqrt{\mathbf{w}^T \mathbf{w}}}$$

or

$$\text{minimize } \mathbf{w}^T \mathbf{w}$$



$$\frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$$

to get rid of  
2 when taken derivative

Lagrange

Logrange

$$L_p = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [y_i (w^T x_i + w_0) - 1]$$

$$y_i (w^T x_i + w_0) \geq 1$$

$$\alpha_i > 0$$

$$\text{minimize } \frac{1}{2} w^T w$$

w

$$\text{subject to } y_i (w^T x_i + w_0) \geq 1$$

$$-\sum_i (\alpha_i y_i w^T x_i + \alpha_i y_i w_0 - \alpha_i)$$

$$[w_0 - w_d] \begin{bmatrix} x_0 \\ x_d \end{bmatrix} \xrightarrow{\text{NO } w}$$

$$w_0 = w + w_1 x_1 + \dots + w_d x_d$$

derivative with respect to w  
 $\Rightarrow x^i$

$$\frac{\partial L}{\partial w} = 0 = w - \sum_i \alpha_i y_i x_i$$

$$i) w = \sum_i \alpha_i y_i x_i$$

$$ii) \frac{\partial L}{\partial w_0} = \sum_i \alpha_i y_i = 0$$

LAST WEEK :)

Instead of minimizing  $L_p$ , people maximize  $L_D$

$$L_D = \frac{1}{2} \left( \underbrace{\sum_i \alpha_i y^i x^{(i)}}_{a^T} \right)^T \left( \underbrace{\sum_j \alpha_j y^j x^{(j)}}_{a^T} \right) - \underbrace{\sum_i \left[ \alpha_i y^i x^{(i)} \right]^T}_{a^T} \underbrace{\sum_j \left[ \alpha_j y^j x^{(j)} \right]}_{a^T}$$

$$+ \sum_i \alpha_i y^i w_0 - \sum_i \alpha_i$$

$$= \boxed{\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^i y^{(i)} (x^{(i)})^T x^{(j)} + \sum_i \alpha_i}$$

exp/

$$y=1, \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\alpha_1 \cdot 1 \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \alpha_2 \cdot 1 \cdot \begin{bmatrix} -1 \\ 1 \end{bmatrix} + \alpha_3 \cdot (-1) \cdot \begin{bmatrix} -1 \\ -1 \end{bmatrix} = a$$

$$\sum_i \alpha_i y^i x^{(i)} = a$$

$$y=-1, \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$(\alpha_1 [1]^\top + \alpha_2 [-1]^\top - \alpha_3 [1]^\top) \cdot (\alpha_1 [1] + \alpha_2 [-1] - \alpha_3 [1])$$

$$= \alpha_1^2 [1 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \alpha_1 \alpha_2 [1 1] \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \alpha_1 \alpha_3 [1 1] \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ + \alpha_2 \alpha_1 [-1 1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \alpha_2^2 [-1 1] \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \alpha_2 \alpha_3 [1 1] \begin{bmatrix} -1 \\ -1 \end{bmatrix} \\ - \alpha_3 \alpha_1 [-1 -1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \alpha_2 \alpha_3 [-1 -1] \begin{bmatrix} -1 \\ 1 \end{bmatrix} + \alpha_3^2 [-1 -1] \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$= (2\alpha_1^2) + 0 \cancel{\alpha_1 \alpha_2} - 2\alpha_1 \alpha_3 + 0 \cancel{\alpha_1 \alpha_2} + (2\alpha_2^2) + 0 \cancel{\alpha_2 \alpha_3} + 2\alpha_2 \alpha_3 + 0 \cancel{\alpha_2 \alpha_3} \\ + \sum_i \alpha_i$$

$$L_D = 2\alpha_1^2 + 2\alpha_2^2 + 2\alpha_3^2 + 4\alpha_1 \alpha_3 + \alpha_1 \alpha_2 \alpha_3 \rightarrow \alpha_1 + \alpha_2 - \alpha_3 = 0$$

to find maximum

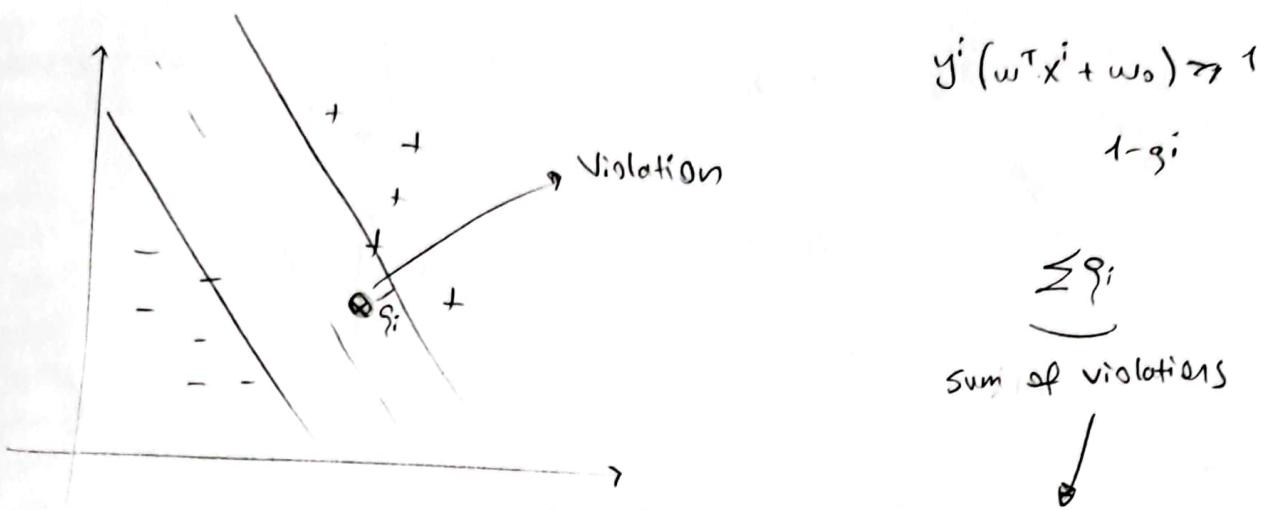
$$\alpha_3 = \alpha_1 + \alpha_2$$

coordinate ascent

\* For most of the points  $\alpha$ 's will be zero,  
non-zero  $\alpha$  values  $\rightarrow$  support vectors

another example on the book ML Making sense of Data  
"Peter Flach"

example 7.5, p. 215



$$\text{minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w} + c \cdot \sum_i g_i$$

a constant to make  
fine tuning on violations

Now we have this:

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$\sum_{i=1}^n \alpha_i y^{(i)} = 0$$

$0 \leq \alpha_i \leq C$

new constraint

$$\frac{\lambda^2}{2} \left[ \begin{matrix} -1 & 2 \\ 2 & 2 \end{matrix} \right] \left[ \begin{matrix} -1 \\ 2 \end{matrix} \right] + \dots$$

2-dim

$$x = \left[ \begin{matrix} -1 \\ 2 \end{matrix} \right] \begin{matrix} x_1 \\ x_2 \end{matrix}$$

$$-\frac{1}{2} \leq \sum_j \alpha_j y^{(j)} y^{(j)} (x^{(j)})^T x^{(j)} + \sum_i \alpha_i$$

$$(z^{(i)})^T (z^{(i)})$$

Let's rename  
 $(\phi(x))^T \cdot (\phi(x)) \rightarrow \phi^T(x) \cdot \phi(y)$

$$\phi(x) = \begin{bmatrix} 1 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \\ \sqrt{2}x_1 x_2 \\ x_1^2 \\ x_2^2 \end{bmatrix} = z$$

$$\phi^T(x) \cdot \phi(y) = 1 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2 + x_1^2 y_1^2 + x_2^2 y_2^2 \rightarrow \text{polynomial}$$

$$K(x, y) = (x^T y + 1)^2 = \phi^T(x) \cdot \phi(y)$$

Aim: Data may not be separable in lower dim-space  
 Therefore, we try higher dim-spaces but efficient

Here  $\phi(x)$  is a function that takes a 2-dim matrix and maps it into a 6-dim one. Dim++

$$-\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(x^i, x^j) + \sum_i \alpha_i$$

↓

Kernel      mercer conditions

## Online Session

Artificial neural networks take their inspiration from brain

1960 - Multilayer Perceptron



1986 - Back Propagation

1989 - Convolutional Neural Network

1990-1994 - Neural Nets in the wild (Support Vector Machines ?)

1995 - Long Short term memory

1998-2005 - More N.N in the wild

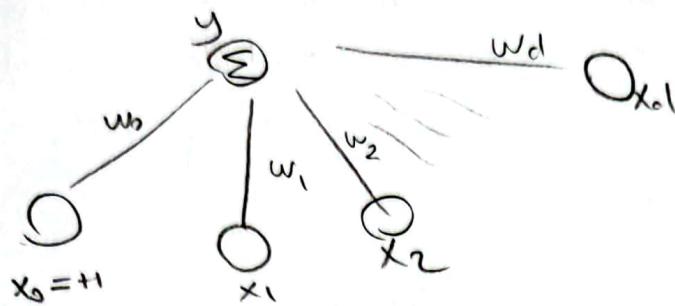
2005-2010 - Bottleneck (not enough data and CPU power)

2010 - Renarr : Deep Learning

2013 - CNN + GPU + ImageNet

2015 - AlphaGo

## Perception



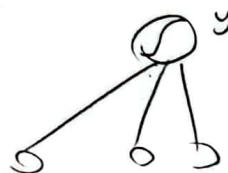
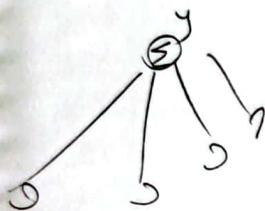
$$y = \sum_{j=1}^d w_j x_j + w_0 = w^T x$$

(1962)

what perceptron does

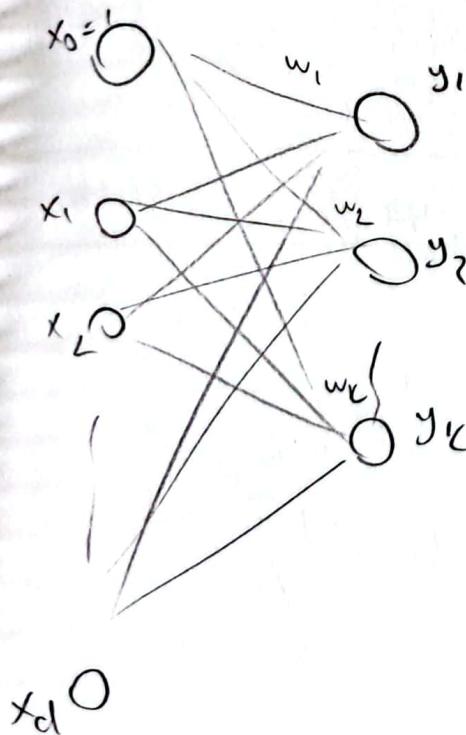
regression:  $y = w x + w_0$

classification  $y = 1(w x + w_0 > 0)$



$$y = \text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

K outputs



regression:  $y_i = \sum_{j=1}^d w_{ij} \cdot x_j + w_{i0} = w_i^T x$

classification

$$o_i = w_i^T x$$

$$y_i = \frac{\exp o_i}{\sum_k \exp o_k}$$

## Training

online

to train our weights ( $w$ )

batch

$\rightarrow$  instances seen one by one  
 $\rightarrow$  instances seen as whole

gradient descent

$$w = w - \eta \frac{\partial E}{\partial w}$$

stochastic gradient descent  $\rightarrow$  update after  
a single pattern

## Regression

$$E = \sum_{t=1}^T (r^t - y^t)^2 \quad \rightarrow \Delta w_j^t = \eta (r^t - y^t) x_j^t$$

↓  
actual      ↓ predicted

for every

$$E^t(w | x^t, r^t) = \frac{1}{2} \cdot (r^t - y^t)^2 = \frac{1}{2} [r^t - (w^t x^t)]^2$$

## Classification

$E \rightarrow$  cross entropy

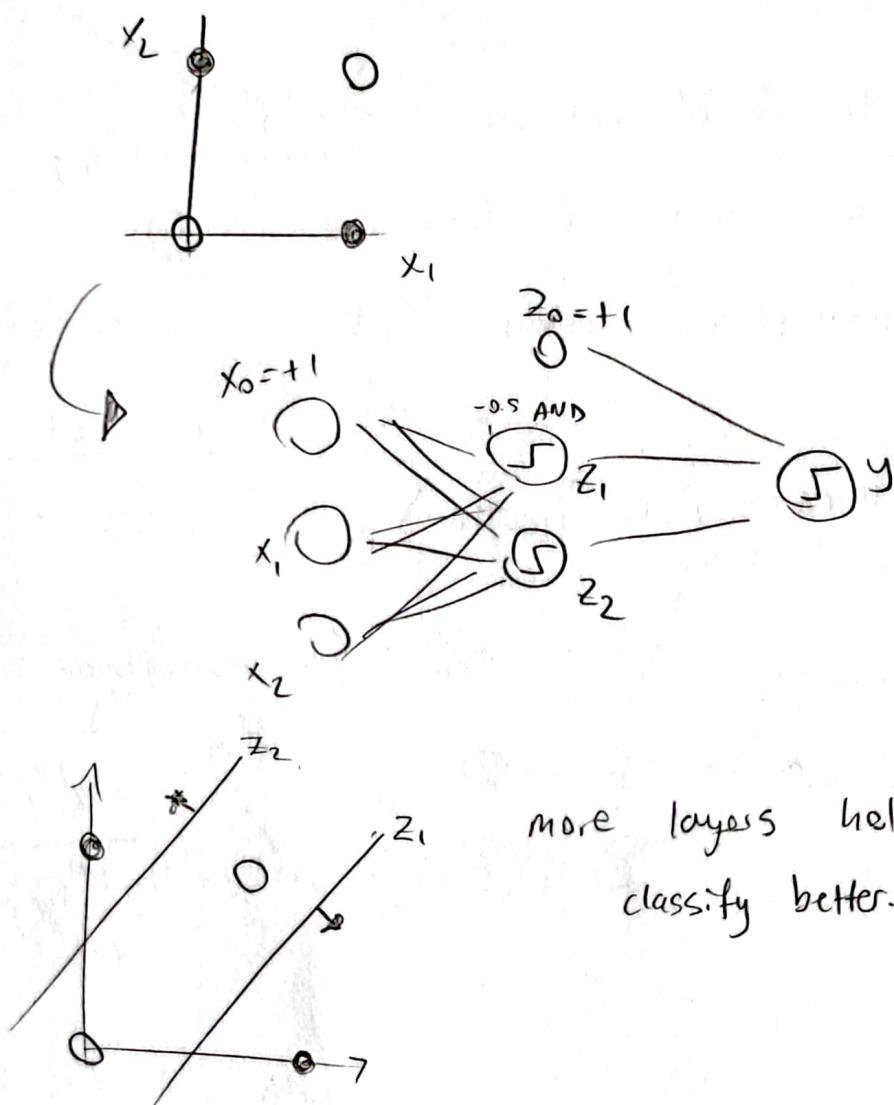
$$\log \pi(y^t)^{r^t} \cdot (1-y^t)^{1-r^t}$$

stochastic      gradient      descent

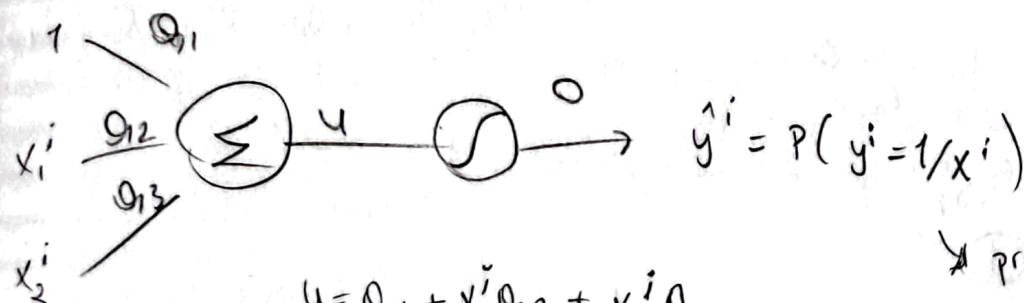
$$E^t(w | x^t, r^t) = r^t \log y^t + (1-r^t) \log (1-y^t)$$

$$\Delta w_{ij}^t = \eta (r^t - y_i^t) x_j^t$$

a linear model, may not correctly classify 4 data points  $\rightarrow$  XOR problem



### MLP - 1 neuron



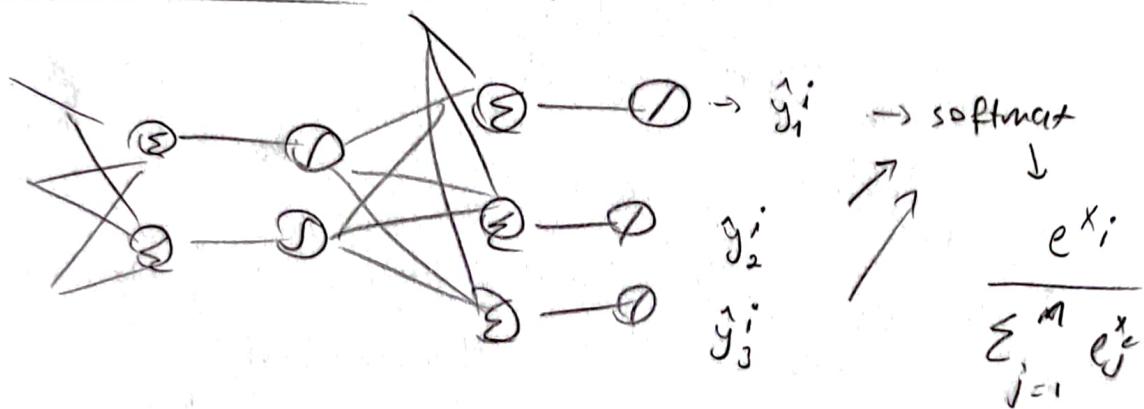
probability of  
 $y^i = 1$  given that

$$o = \frac{1}{1+e^{-u}} = \hat{y}^i$$

if MLP ends with  $\rightarrow$  a sigmoid function  $\rightarrow$  classification  
 $\downarrow$   
 $a \leq \dots \rightarrow$  regression

- Organize the neural network  $\rightarrow$  MSE
- Define an error function  $\rightarrow$  cross entropy
- Update parameters using G.D.

### MLP - Multi class (classification)

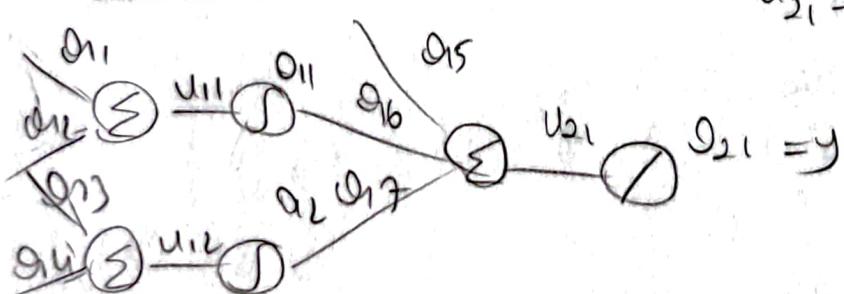


### Back propagation

regression

$$E = \frac{1}{2}(y - r)^2 \quad y = \theta_{21} = u_{21}$$

$$u_{21} = \theta_{15} + \theta_{16} \cdot \theta_{11} + \theta_7 \theta_{12}$$



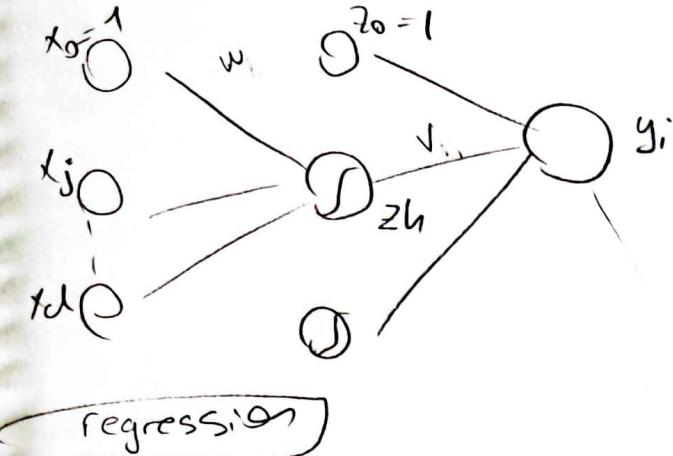
USE  
CHAIN RULE

$$\Delta \theta_{16} = \frac{\partial E}{\partial \theta_{16}} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial \theta_{21}} \cdot \frac{\partial \theta_{21}}{\partial u_{21}} \cdot \frac{\partial u_{21}}{\partial \theta_{16}}$$

$$(y - r) \theta_{11} \dots$$

$$y - r \quad \bar{1} \quad \bar{1} \quad \bar{1} \quad \bar{\theta_{11}}$$

## Multilayer Perceptrons



## General Formula

$$\Delta V_h = \sum (r^t - y^t) z_h^t$$

$$\Delta w_{hj} = -\eta \frac{\partial E}{\partial w_{hj}}$$

$$\Delta w_{hj} = \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t$$

## Multi output MLP

MLP with 2 hidden layers

## Summary

① define problem

## regression

msc

## Classification

A diagram consisting of two circles, one on the left and one on the right, connected by a single horizontal line segment that passes through the center of each circle.

$$\textcircled{2} \quad \oint \frac{1}{2}(r-y)^2$$

$$\Pi y'(1-y)^{1-v}$$

$$E(w, v) = \frac{1}{2} (r - \bar{y})^2$$

unknowns

$$\hat{y} = v \cdot z \quad \frac{\partial E}{\partial v} = (r - \hat{y}) \cdot z$$

$$\frac{\partial E}{\partial v} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial v} \quad \rightarrow$$

$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$z \rightarrow \text{sigmoid}$$

$$z' = z(1-z)$$