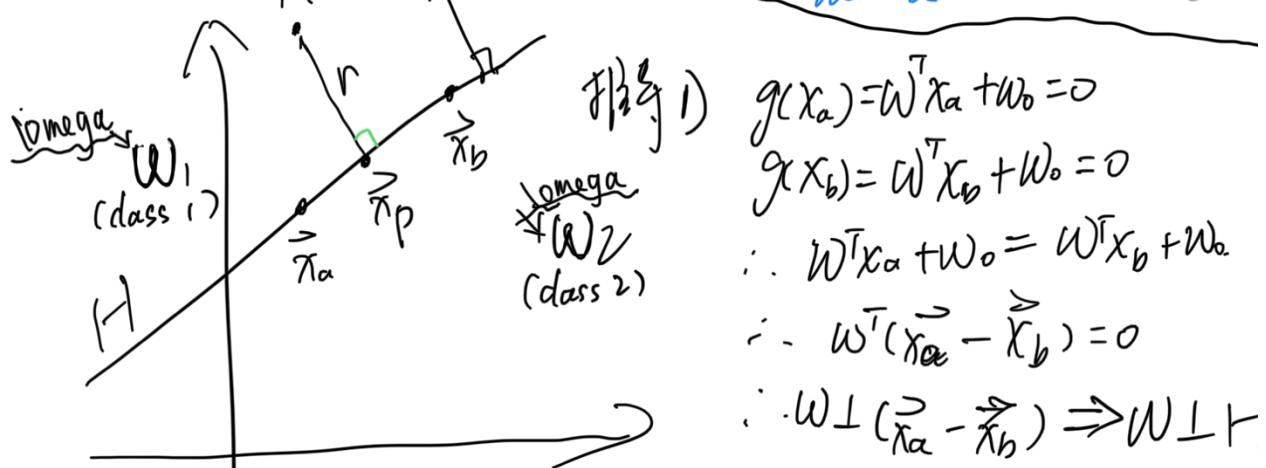
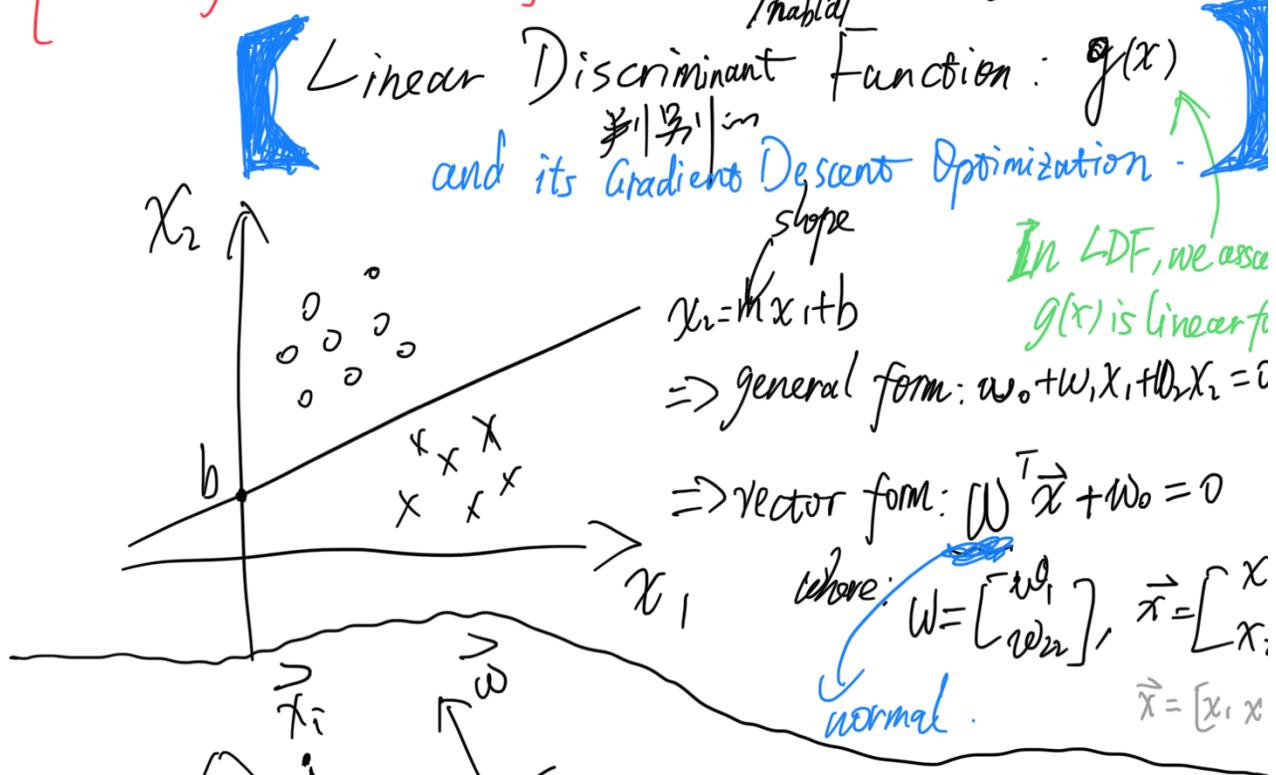


{ Bzrkz github ML notes }

$$\nabla = \left[ \frac{\partial}{\partial x} \frac{\partial}{\partial y} \frac{\partial}{\partial z} \right]^T$$



2)  $\vec{x}_i = \vec{x}_p + r \frac{\vec{w}}{\|\vec{w}\|_2}$   
     normalized, called unit vector

$\vec{x}_i - \vec{x}_p = r \frac{\vec{w}}{\|\vec{w}\|}$

3)  $\because g(x_p) = \vec{w}^T \vec{x}_p + w_0 = 0$ . ( $\because x_p$  在  $H$  上)

$$\begin{aligned} & \therefore \vec{w}^\top (\vec{x}_i - r \frac{\vec{w}}{\|\vec{w}\|}) + w_0 = 0 \Rightarrow \|\vec{w}\| \\ & \therefore \vec{w}^\top \vec{x}_i + w_0 = r \cdot \frac{\vec{w}^\top \vec{w}}{\|\vec{w}\|} \\ & \quad = r \|\vec{w}\| \\ & \therefore g(\vec{x}_i) = r \|\vec{w}\| \end{aligned}$$

$$\therefore \boxed{r = \frac{g(\vec{x}_i)}{\|\vec{w}\|}} \quad \begin{array}{l} \text{r: distance of } \vec{x}_i \text{ to the} \\ \text{decision boundary. } g(x) \text{ or it.} \end{array}$$

More confident if  $x$  is further away from decision boundary

Ex:  $\vec{w} = [1, 2]$ ,  $w_0 = 1$

① H (decision boundary)?

$$H: g(\vec{x}) = \vec{w}^\top \vec{x} + w_0 = [1 \ 2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + 1 = x_1 + 2x_2 + 1$$

$$\therefore \text{boundary } H: x_2 = -\frac{1}{2}x_1 - \frac{1}{2}$$

② Test two samples:  $(\omega: \text{omega is class.})$   
 $(w: \text{double-u is weight})$

$$g(\vec{x}) \begin{cases} \geq 0 \\ \leq 0 \end{cases}$$

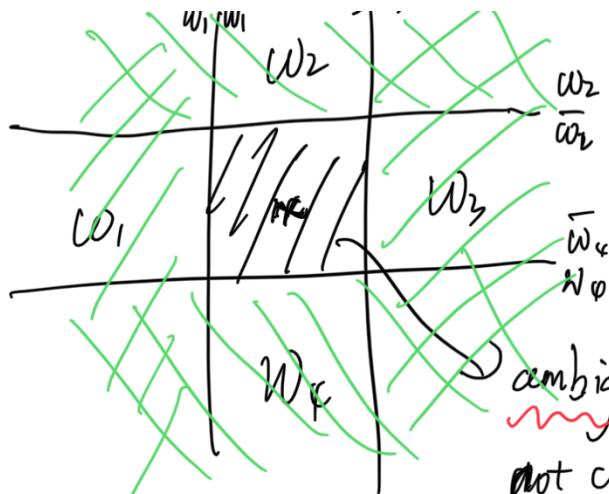
$$\vec{x}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \vec{x}_2 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$g(\vec{x}_1) = 0 + 2 \times 1 + 1 = 3 > 0 \Rightarrow \vec{x}_1 \in \omega_1$$

$$g(\vec{x}_2) = 0 + 2 \times (-1) + 1 = -1 < 0 \Rightarrow \vec{x}_2 \in \omega_2$$

The system is more confident of  $\vec{x}_1$  belonging to  $\omega_1$   
than  $\vec{x}_2$  belonging to  $\omega_2$  because  $|g(\vec{x}_1)| > |g(\vec{x}_2)|$

Multi-category LDF.  
 $\vec{w}_1, \vec{w}_2, \dots, \vec{w}_k$



ambiguity region (caused by  
not covering the whole space)

↓  
one solution of this issue is  
by Bayesian Decision Theory:

\* Assign  $x$  to  $w_i$  if  $g_i(x) \geq g_j(x), \forall j \neq i$

## Generalized LDF

### ① Linear DF:

$$g(x) = w_0 + \sum_{i=1}^d w_i x_i$$

### ② Higher order forms

- Quadratic DF.

$$g(x) = w_0 + \sum_{i=1}^d w_i x_i + \sum_{i=1}^d \sum_{j=1}^d w_{ij} x_i x_j$$

- Polynomial DF (cubic or higher)

$$g(x) = w_0 + \sum_i w_i x_i + \sum_{i,j} w_{ij} x_i x_j + \sum_{i,j,k} w_{ijk} x_i x_j x_k + \dots$$



So, how to find  $H$ , i.e. estimate decision boundary parameters

-  $w_0, w_1, \dots, w_d$  : free params of  $H$ .

e.g. 2 features, 2 classes.

$$g(\vec{x}) = w_1 x_1 + w_2 x_2 + w_0 = [w_0 \ w_1 \ w_2] \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

[augmenting]:  $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \xrightarrow{\vec{x}} \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \xrightarrow{\vec{y}}$  (开了一个维度)

means augmenting  
 $\vec{a} = [w_0 \ w_1 \ w_2], \vec{y} = [$

换变量:  $g(\vec{y}) = \vec{a}^\top \vec{y}$  (变量替换) PW: 应该是把Unknown写成

\* One of the good reasons to use augmenting is that:

\*  $\vec{a}$  is always locked to origin i.e.  $(0, 0, 0) \neq 0$   
 $\therefore (0 \ 0 \ 0) \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} = 0$

so the task now is to rotate  $\vec{a}$  via origin to find the best boundary curve.

⇒ Estimate  $\vec{a}$  or rotate  $\vec{a}$ .

+ We need good training sample set (Labeled).

If they are linearly separable, we can find the ideal  $\vec{a}$ , which is also called solution vector.

PW: now the goal is training not classification.

+ We need to integrate label information into an uniform function.

~~Normalization~~  $\hat{y} = \begin{cases} y_i & , y_i \in w_1 \\ -y_i & , y_i \in w_2 \end{cases}$

Suppose for  $w_1$  class:  $g(x) > 0$   
 for  $w_2$  class:  $g(x) < 0$

$$g(x) = \vec{a}^T \hat{y}_i < 0 \Rightarrow \vec{a}^T \hat{y}_i > 0$$

i.e. After normalization,

find the optimal  $\vec{a}$ , s.t.  $\vec{a}^T \hat{y}_i > 0, \forall \hat{y}_i (\forall x_i)$

? What if we get  $g(\hat{y}_i) < 0$ .  $\Rightarrow$  incorrect classification  
 for a training data 不对也不符合 label 不对.

$\Rightarrow$  After the previous two steps, i.e. augmenting and normalization  
 we can now find the optimal  $\vec{a}$ . Note: ~~After this step~~ ~~normalized~~

Before move on, we need to ask two questions:

① What are you optimizing?  $\rightarrow J$

② how do you prove it is optimal?  $\rightarrow$  global minimum

Find  $\vec{a}^*$   $\rightarrow$  Find criterion  $J(\vec{a})$ .  $\left[ \begin{array}{l} \text{目标函数} \\ \text{极值问题} \end{array} \right]$

$$\vec{a}^* = \underset{\vec{a}}{\operatorname{arg\,min}} J(\vec{a})$$

Solution vector: if differentiable  $\rightarrow \nabla J(\vec{a}) = 0$   
 continuous

Options of  $J(\vec{a})$

+ Count # of misclassifications(errors)

~~+ Error Correction methods~~ (perceptron) algorithm  
 + calculate measurable errors of  
 only from misclassified sample.

if linearly separable.

$$J_p(\vec{a}) = \sum_{y \in Y} (-\vec{a}^T y)$$

$$\frac{\partial J_p(\vec{a})}{\partial \vec{a}} : \nabla_{\vec{a}} J_p = \sum_{y \in Y} (-y)$$

misclassified samples

~~+ Squared Error methods~~

+ measure errors from all training samples, i.e.  
 misclassified samples and correctly classified samples.

Perceptron  $J_p(\vec{a})$

+ no samples misclassified:  $J_p(\vec{a}) < 0$

+ if  $y$  is misclassified,  $J_p(y) \leq 0$ , i.e.  $\vec{a}^T y \leq 0$ :  $J_p(\vec{a}) \geq 0$ .

$$\nabla_{\vec{a}} J_p = \sum_{y \in Y} (-y)$$

step ① Initialization: random  $\vec{a}(0)$   
 step ② Use gradient descent algorithm iteratively:

$$\vec{a}(k+1) = \vec{a}(k) - \eta(k) \nabla J(\vec{a}(k))$$

step ③ stop:  $|\nabla J(\vec{a}(k))| < \theta$

learning rate

Ex:	$x_1$	$x_2$	$w_1/w_2$ (class)	$\ w\  = 0.1$
	0.1	1.1	1	
	6.8	7.1	1	tolerance $\theta = 0$
	7.1	4.2	2	Goal: $a = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$
	-1.4	-4.3	2	s.t. $y_i \cdot a > 0, \forall i$ .

① Augmentation

$$y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0.1 & 6.8 & 7.1 & -1.4 \\ 1.1 & 7.1 & 4.2 & -4.3 \end{bmatrix}^T \quad (\text{to column vector})$$

$N \times (d+1)$  (把  $w_0$  从 samples 对齐过去)

② Normalization

$$\tilde{y} = \begin{bmatrix} \tilde{y}_1 & \tilde{y}_2 & \tilde{y}_3 & \tilde{y}_4 \end{bmatrix}^T$$

$$\tilde{y} = \begin{bmatrix} 0.1 & 6.8 & -7.1 & 1.4 \\ 1.1 & 7.1 & 4.2 & -4.3 \end{bmatrix}^T$$

now,  $y = \tilde{y}$ .

Alg: ① Initialization:

$$a(0) = [0 \ 0 \ 0]^T$$

② 1st  $\xrightarrow{k \uparrow}$  iteration:  $a(0)^T y = [0 \ 0 \ 0]^T = [0 \ 0 \ 0]$  all classified

$$\nabla J_p(a) = \sum_i (-y_i) = \begin{bmatrix} 0 \\ -1.2 \\ -8.3 \end{bmatrix} \quad (y_i: \text{column vector})$$

$$a(1) = a(0) - \eta \nabla J_p(a(0))$$

$$= 0 - 0.01 \times \begin{bmatrix} 0 \\ -1.2 \\ -8.3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.012 \\ 0.083 \end{bmatrix}$$

2<sup>nd</sup> iteration:

$$\begin{array}{ccccccc} >0 & >0 & >0 \\ \nearrow & \searrow & \nearrow & \searrow & \nearrow & \searrow & \nearrow \\ -1 & 0 & -1 & 0 & -1 & 0 & 1 \end{array}$$

$\alpha(1) y = \dots = [0.01, -0.12, 0.15, 0.10]$

$\nearrow$   
 $\nwarrow$

$\text{<0: misclassified}$

$$\therefore \nabla_{\beta}(\alpha(k=1)) = -y_3 = \begin{bmatrix} 1 \\ 7.1 \\ 4.2 \end{bmatrix}$$

$$\therefore \alpha(2) = \alpha(k=1) - 0.01 \times \begin{bmatrix} 1 \\ 7.1 \\ 4.2 \end{bmatrix} = \begin{bmatrix} -0.0 \\ -0.06 \\ 0.04 \end{bmatrix}$$

3<sup>rd</sup> iteration ( $k=2$ ):

step 1: classification

$$\alpha^T(2) y = \dots = [0.028, -0.120, 0.15, 0.10]$$

$\nearrow$   
 $\nwarrow$

$\text{<0: misclassified}$

$$\text{Step 2: } \nabla_{\beta}(\alpha(2)) = -y_2 = \begin{bmatrix} 1 \\ -6.8 \\ -7.1 \end{bmatrix}$$

$$\text{Step 3: } \alpha(3) = \alpha(2) - 0.01 \times \begin{bmatrix} 1 \\ -6.8 \\ -7.1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0.008 \\ 0.11 \end{bmatrix}$$

4<sup>th</sup> iteration.

5<sup>th</sup> iteration.

$$H: g(x) = -0.01 - 0.063x_a + 0.069x_b = 0$$

(BW: but is this the best solution vector? e.g. equal margin to each class)



\* Problem formulated as that of minimizing  
or cost function.

↑ Gradient descend & perceptron criterion.

⇒ Batch and single sample perceptron criterion.

**Step 1** LDF: Augmentation

$$\square \text{LDF: } g(x) = w^T x + w_0$$

$$\square \text{Rewrite it: } g(x) = [w_0 \ a^T] \begin{bmatrix} 1 \\ x \end{bmatrix} = \underbrace{a^T}_{\substack{\text{new weight} \\ \text{Vector } a}} \underbrace{\begin{bmatrix} 1 \\ x \end{bmatrix}}_{\substack{\text{new feature} \\ \text{vector } y}} = a^T y = g[y]$$

□  $y$  is the Augmented Feature Vector  $\Rightarrow (d+1) \times 1$

**Step 2** LDF: Normalization (i.e. switch the  $y_i$  sign)

线性边  $\begin{cases} a^T y_i > 0 & \forall y_i \in C_1 \\ a^T y_i < 0 & \forall y_i \in C_2 \end{cases} \Rightarrow a^T (-y_i) > 0 \quad \forall y_i \in C_2$

$$\rightarrow n^T a \rightarrow n \cdot H u: \text{从上到下减 -it}$$

$\rightarrow$  If such  $a$  exists, it is called a solution vector.  
 $\underbrace{\text{If such } a \text{ exists, it is called a solution vector.}}$   
 $\quad (\rightarrow \text{linearly separable})$

**Step 3** LDF: criterion  $J(a) \rightarrow$  always be a scalar.

Particularly, we use perceptron criterion  $J_p(a)$   
 optimization method: gradient descent with  $J_p(a), \nabla J_p(a)$

$$J_p(a) = \sum_{y \in Y} (-a^T y), \quad \nabla J_p(a) = \sum_{y \in Y} (-y)$$

$\uparrow$  error samples : Initialization of Error Correct  
 methods.

i. Relaxation procedures (smoothen the error change)

-  $J_q(a)$ : squared error

// -  $J_r(a)$ : Normalized squared error

$J_p(a)$  abrupt change (cons)

### Basic Gradient Descent Algorithm

begin initialize  $a$ , threshold  $\theta$ ,  $\eta(\cdot)$ ,  $k \leftarrow 0$

do  $k \leftarrow k+1$

$a \leftarrow a - \eta(k) \nabla J(a)$

until  $|\eta(k) \nabla J(a)| < \theta$

return  $a$ .

$$J(a) = \begin{cases} J_p(a) \\ J_q(a) \\ J_r(a) \end{cases}$$

$$\begin{cases} J_p(a) \\ J_q(a) \\ J_r(a) \end{cases}$$

choose one of the above  
 when designing an algorithm

end

## Relaxation Procedures

\* Squared error criterion:

$$J_q(\alpha) = \sum_{y \in Y} (\alpha^T y)^2$$

- It is the generalization of perceptron criterion.

- Error surface is continuous and smooth.

$$\nabla J_q = 2 \cdot \sum_{y \in Y} (\alpha^T y) \cdot \vec{y} \quad ; \quad \nabla J_p = \sum (-\vec{y})$$

↳ dominated by the longest sample vectors. Can use GDI  
 (largest  $\|y\|$ )

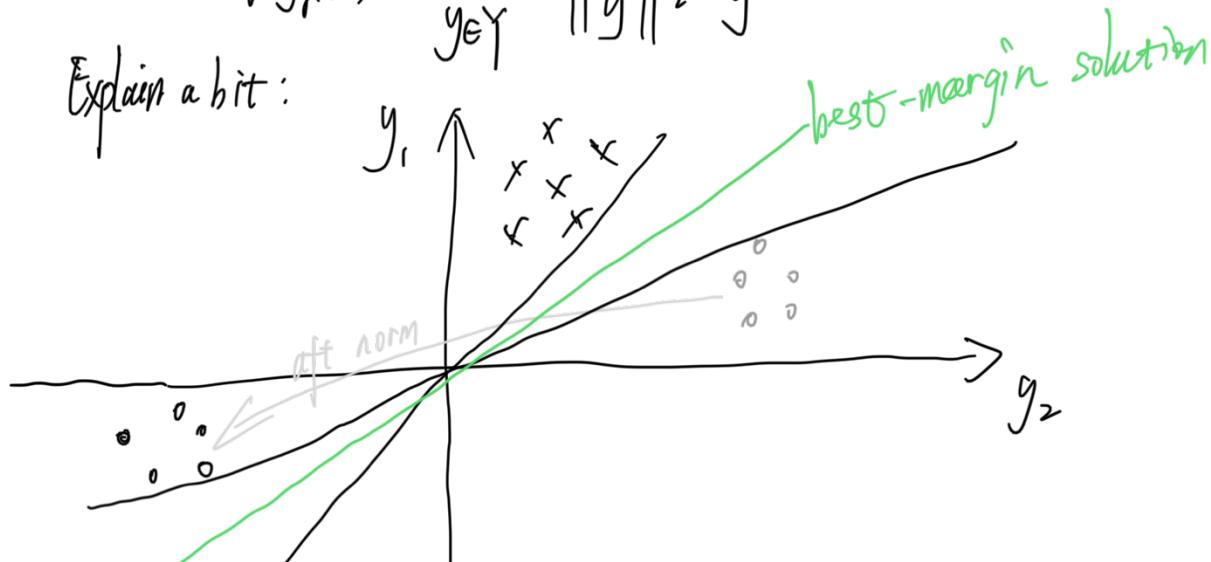
\* Normalized squared-error criterion

$$J_r(\alpha) = \frac{1}{2} \sum_{y \in Y} \frac{(\alpha^T y - b)^2}{\|y\|^2} \quad b: \text{margin } \leftarrow \begin{array}{l} \text{maximize} \\ \text{as well} \end{array}$$

Distance is normalized with margin.

$$\nabla J_r(\alpha) = \sum_{y \in Y} \frac{\alpha^T y - b}{\|y\|^2} \vec{y}$$

Explain a bit:



$\checkmark / \checkmark$

---

$J_q, J_r \rightarrow$  batch/single sample mode

① Test misclassification

② Update of  $a$  (if any misclassification)

batch  $\begin{cases} J_q: a \leftarrow a - \eta(\cdot) \sum_{y \in Y} (a^T y) \vec{y} \\ J_r: a \leftarrow a - \eta(\cdot) \sum_{y \in Y} \frac{b - a^T y}{\|y\|^2} \vec{y} \end{cases}$

$Y: \text{misclassified samples}$

Simple sample  $\begin{cases} J_q: a \leftarrow a - \eta(\cdot) (a^T y_k) \vec{y}_k \\ J_r: a \leftarrow a - \eta(\cdot) \frac{b - a^T y_k}{\|y_k\|^2} \vec{y}_k \end{cases}$

↑ ↑ ↑

All is error-correction.

---

If Not linearly separable,

\* Error correcting procedures not so effective.

\* Why?

- modify weight only when error is encountered because there is always error.

---

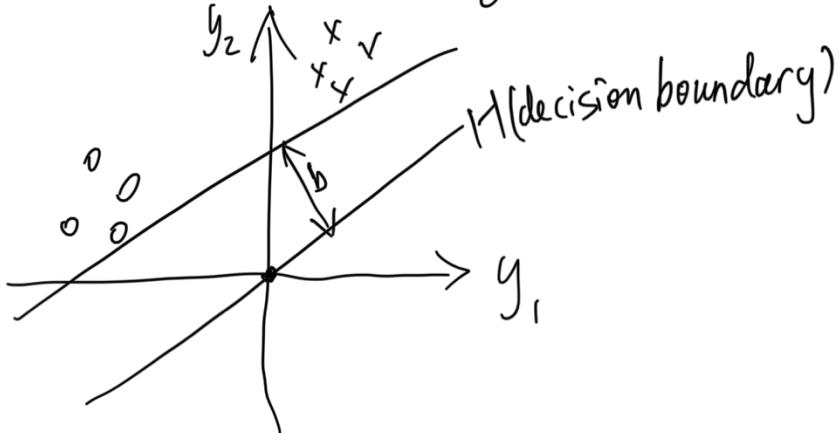
learning rate  $\eta(k)$  is a variable as well.

---

Now it's not error correction anymore since we're considering misclassified and correctly classified

samples, i.e. all sample are considered.

# Minimum Squared-Error (MSE) Procedures



$$a_i^T y_i = b_i > 0$$

formulate the problem

$$N \text{ Samples} \begin{bmatrix} y_{10} & y_{11} & \cdots & y_{1d} \\ y_{20} & y_{21} & \cdots & y_{2d} \\ \vdots & \vdots & & \vdots \\ y_{N0} & y_{N1} & \cdots & y_{Nd} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_d \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix} \quad \text{A.K.A. } \boxed{\vec{Y} \vec{a} = \vec{b}}$$

Can we find such a  $\mathbf{S}$

If  $Y$  is (square matrix and) nonsingular,  $\vec{a} = Y^{-1}\vec{b}$ .

But in practice,  $N \gg d+1$ ,  $Y$  is a tall matrix,  $\vec{a} = ?$

$\Rightarrow$  Answer: no solution of  $\vec{a}$ .

$\Rightarrow$  Answer: no solution.

matrix inverse.)

\*<sup>pseudoinverse</sup>  $\Rightarrow J_s(a) = \|\vec{y}\vec{a} - \vec{b}\|^2 = (\vec{y}\vec{a} - \vec{b})^T(\vec{y}\vec{a} - \vec{b})$

$$\therefore a = \arg \min_a J_s(a)$$

$$\nabla J_s(a) = \sum_{i=1}^N 2(a^T y_i - b_i) g_i \\ = 2 \vec{Y}^T (\vec{y}\vec{a} - \vec{b}) = 2(\vec{Y}^T \vec{y}\vec{a} - \vec{Y}^T \vec{b})$$

Let  $\nabla J_s(a) = 0$ , then

$$\vec{Y}^T \vec{y}\vec{a} = \vec{Y}^T \vec{b}$$

$\uparrow \quad \uparrow$   
 $(d+1) \times N \quad N \times (d+1)$   
 $\underbrace{\quad \quad}_{(d+1) \times (d+1)}$

If nonsingular which most likely is true because all samples are independent to each other (covariance matrix), then

$$\vec{a} = \frac{(\vec{Y}^T \vec{Y})^{-1} \vec{Y}^T}{(d+1) \times (d+1)} \vec{b} = \frac{\vec{Y}^+ \vec{b}}{(d+1) \times N \quad N \times 1}$$

called pseudoinverse:  $\vec{Y}^+$

$(d+1) \times N$

Example:

$$X = \begin{bmatrix} 1 & 2 & 2 & 3 \\ 2 & 0 & 3 & 1 \end{bmatrix}^T \xrightarrow{\text{Aug.}} Y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 3 \\ 2 & 0 & 3 & 1 \end{bmatrix}^T \xrightarrow{\text{Norm.}} \begin{bmatrix} 1 & 1 & -1 & - \\ 1 & 2 & -2 & - \\ 2 & 0 & -3 & - \end{bmatrix}$$

$N \times d$  columns  $b = [1, 1, 1, 1]^T$   $\leftarrow$  ~~for linear equation have linear~~

assume  $U = L^{-1}$

= the inverse matrix you assume.

$$\therefore Y^T Y = \begin{bmatrix} 1 & 1 & -1 & -1 \\ 1 & 2 & -2 & -3 \\ 2 & 0 & -3 & -1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 0 \\ -1 & -2 & -3 \\ -1 & -3 & -1 \end{bmatrix} = \begin{bmatrix} 4 & 8 & 6 \\ 8 & 18 & 11 \\ 6 & 11 & 14 \end{bmatrix} \rightarrow |Y^T Y| =$$

$$\therefore (Y^T Y)^{-1} = \frac{1}{36} \begin{bmatrix} 131 & -46 & -22 \\ -46 & 20 & 4 \\ -20 & 4 & 8 \end{bmatrix}$$

$$\therefore Y^+ = (Y^T Y)^{-1} \cdot Y^T = \dots$$

$$\therefore a = Y^+ b = \begin{pmatrix} 11/3 \\ -4/3 \\ 2/3 \end{pmatrix}$$

$$\therefore \text{Decision boundary } H: \frac{11}{3} - \frac{4}{3}x_1 - \frac{2}{3}x_2 = 0$$

\* This pseudoinverse problem is a convex problem, which means we do have a global optimal solution.

□ If it's not batch procedure, that is, for single sample, this pseudoinverse method is call Least Mean Squared-error (LMS) algorithm. (又称单样本 MSE procedure using gradient descent method).

~~SVM~~ : kernels

SVM is fit for small amount of data, not good for deep learning.

## Summary

$\bar{J}_p(\vec{a}) \rightarrow$  two important steps. Why?

$$\bar{J}_q(\vec{a})$$

$$\bar{J}_r(\vec{a})$$

$\bar{J}_s(\vec{a}) \rightarrow$  Pseudoinverse, LMS(single sample)

Margin  $r, b$