

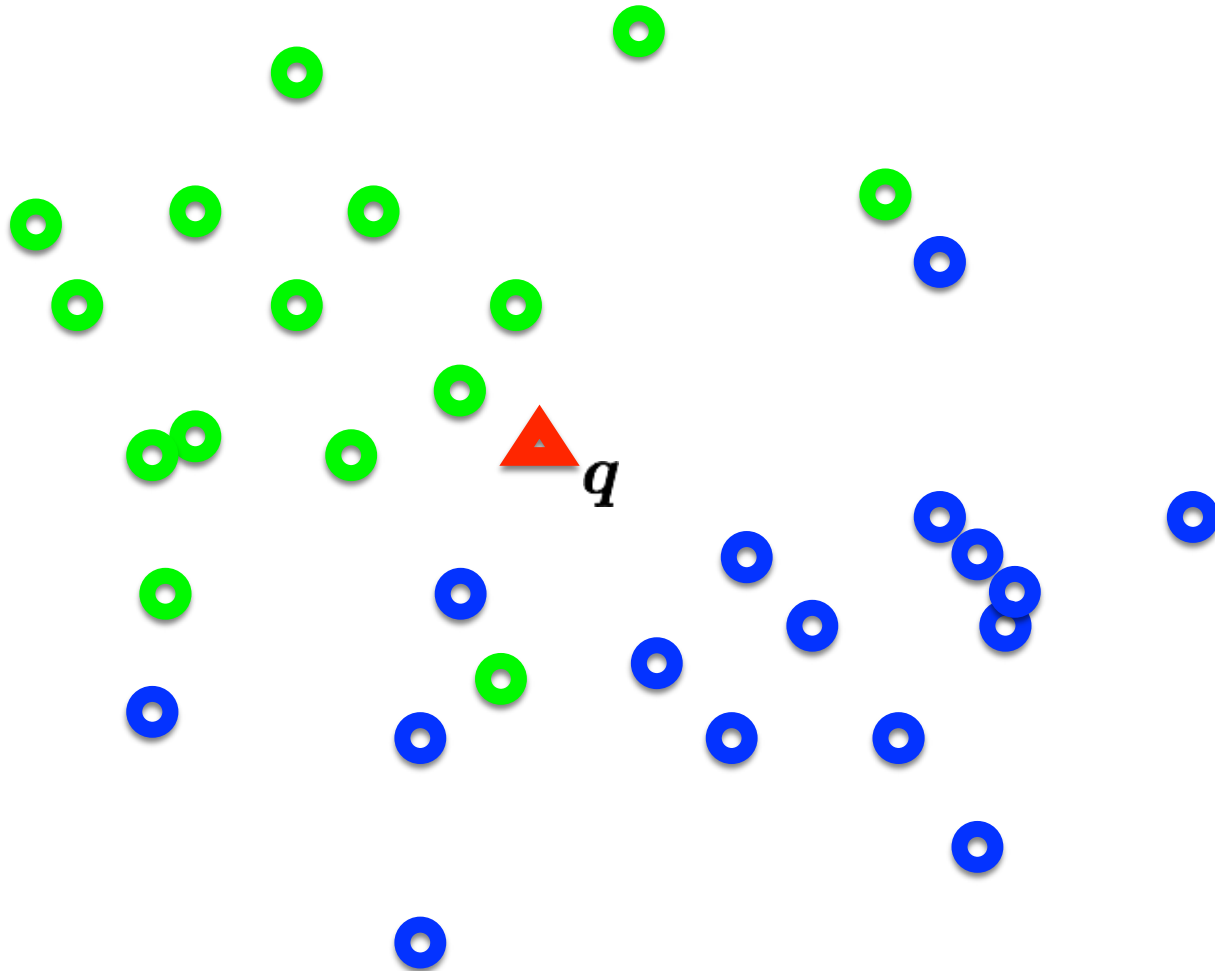
Lecture 6: Machine Learning 2

Sung In Cho

*Div. AISW
Dongguk Univ.*

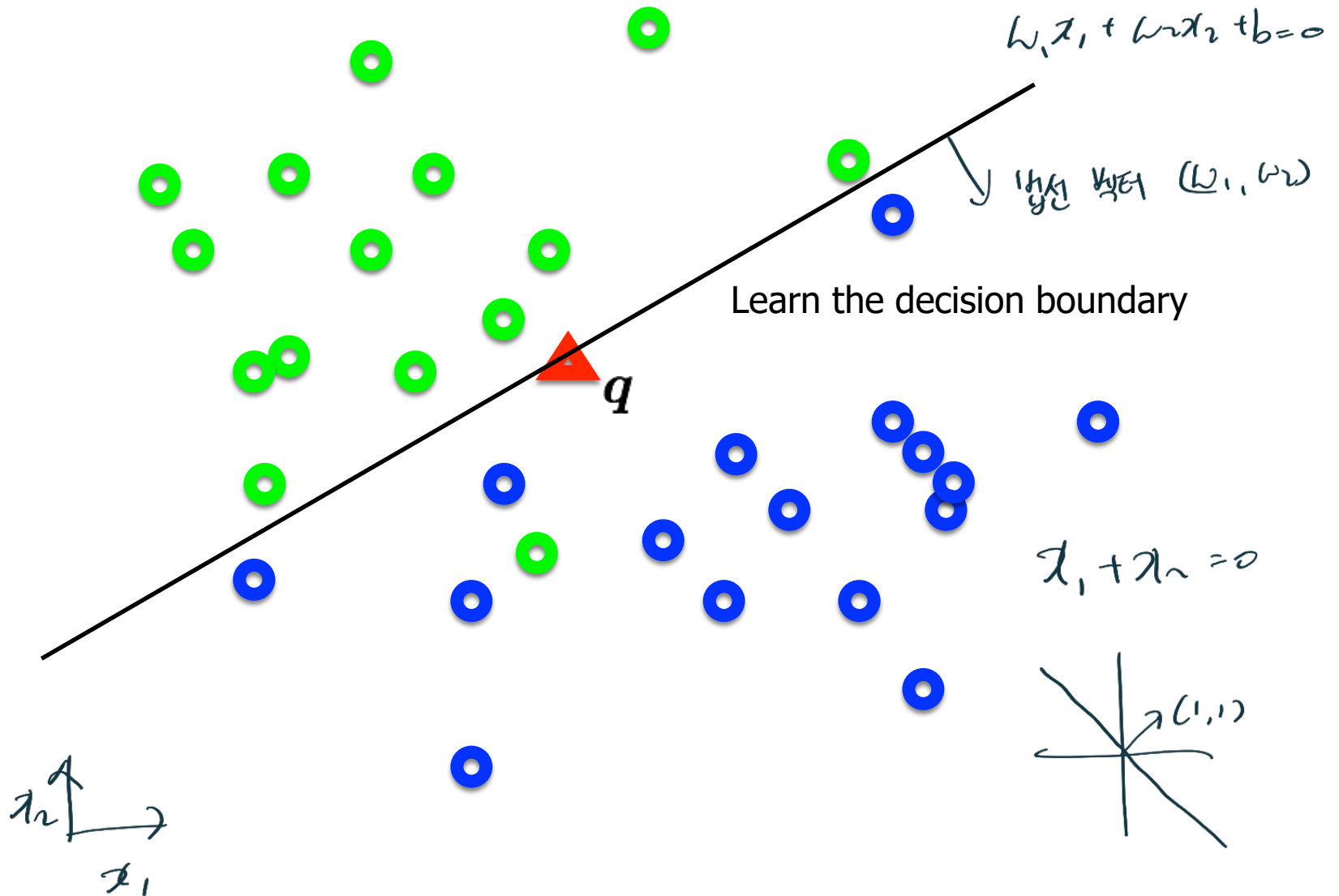
Concept of Support Vector Machine (SVM)

Distribution of data from two classes



Which class does q belong too?

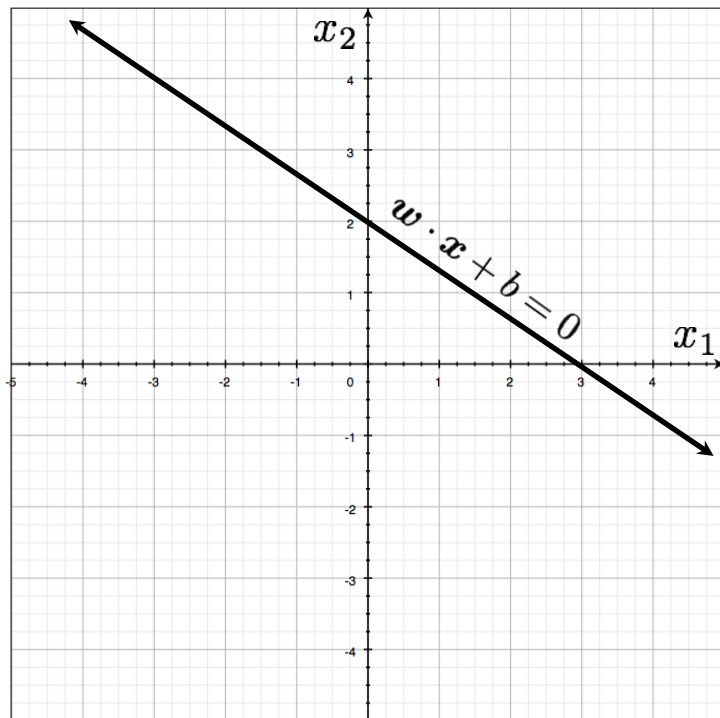
Distribution of data from two classes



First we need to understand hyperplanes...

Hyperplanes (lines) in 2D

$$w_1x_1 + w_2x_2 + b = 0$$



A line can be written as
dot product plus a bias

$$\begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b = 0$$

$$w \in \mathcal{R}^2$$

Another version, add a weight 1
and push the bias inside

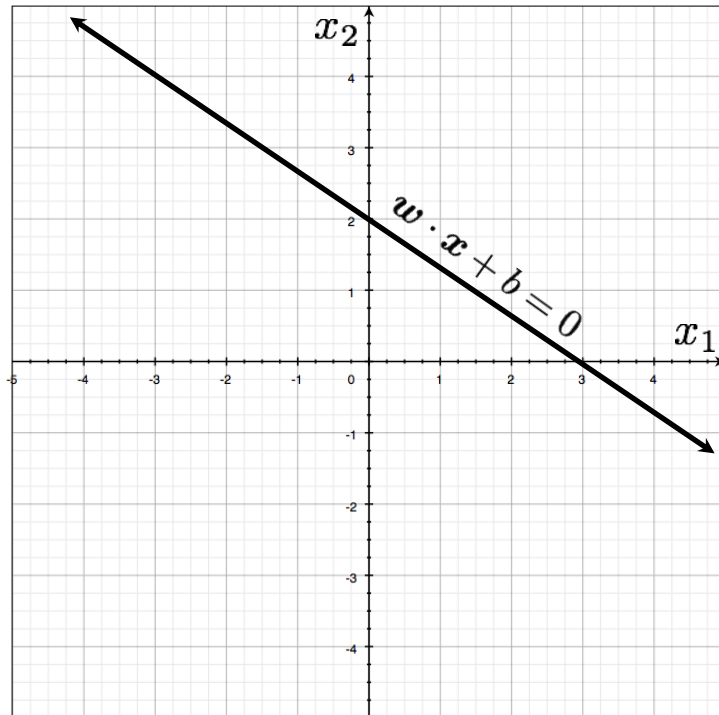
$$\begin{bmatrix} w_1 & w_2 & b \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix} = 0$$

$$w \in \mathcal{R}^3$$

Hyperplanes (lines) in 2D

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad (\text{offset/bias outside}) \quad \mathbf{w} \cdot \mathbf{x} = 0 \quad (\text{offset/bias inside})$$

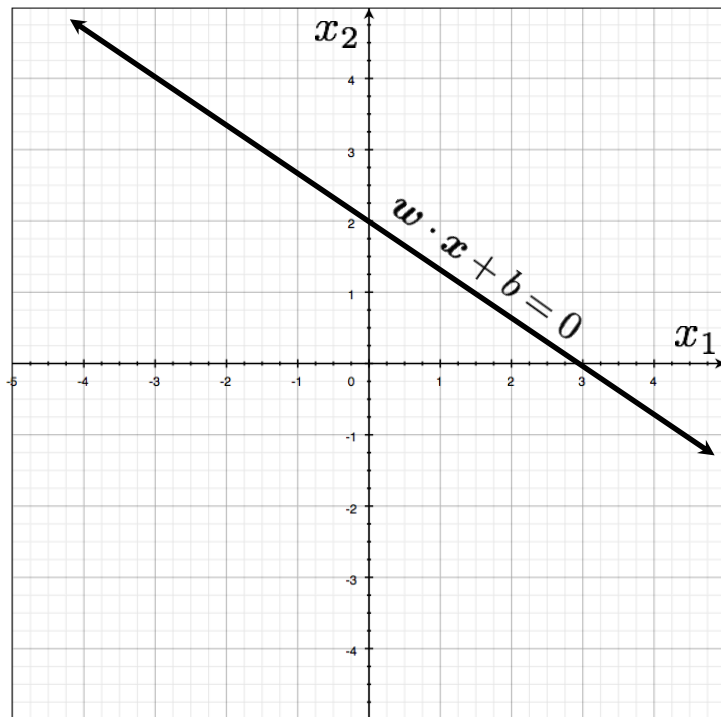
$$w_1x_1 + w_2x_2 + b = 0$$



Hyperplanes (lines) in 2D

$$w \cdot x + b = 0 \quad (\text{offset/bias outside}) \quad w \cdot x = 0 \quad (\text{offset/bias inside})$$

$$w_1 x_1 + w_2 x_2 + b = 0$$



Important property:
Free to choose any normalization of w

The line

$$w_1 x_1 + w_2 x_2 + b = 0$$

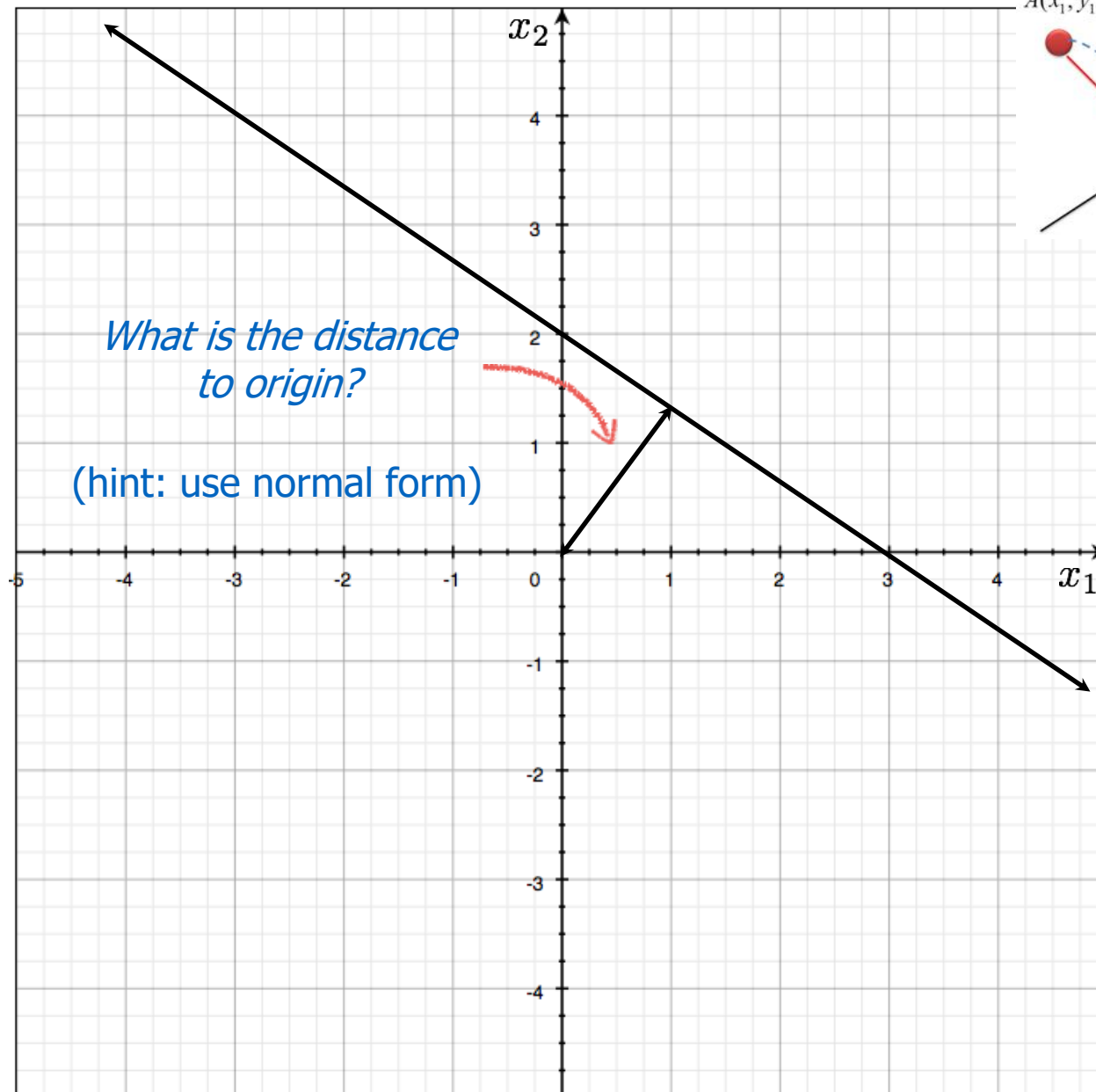
and the line

$$\lambda(w_1 x_1 + w_2 x_2 + b) = 0$$

Scale를 해도

define the same line

직접 바꾸지 않을 → 방향성도 중요!



*What is the distance
to origin?*

(hint: use normal form)

$A(x_1, y_1)$

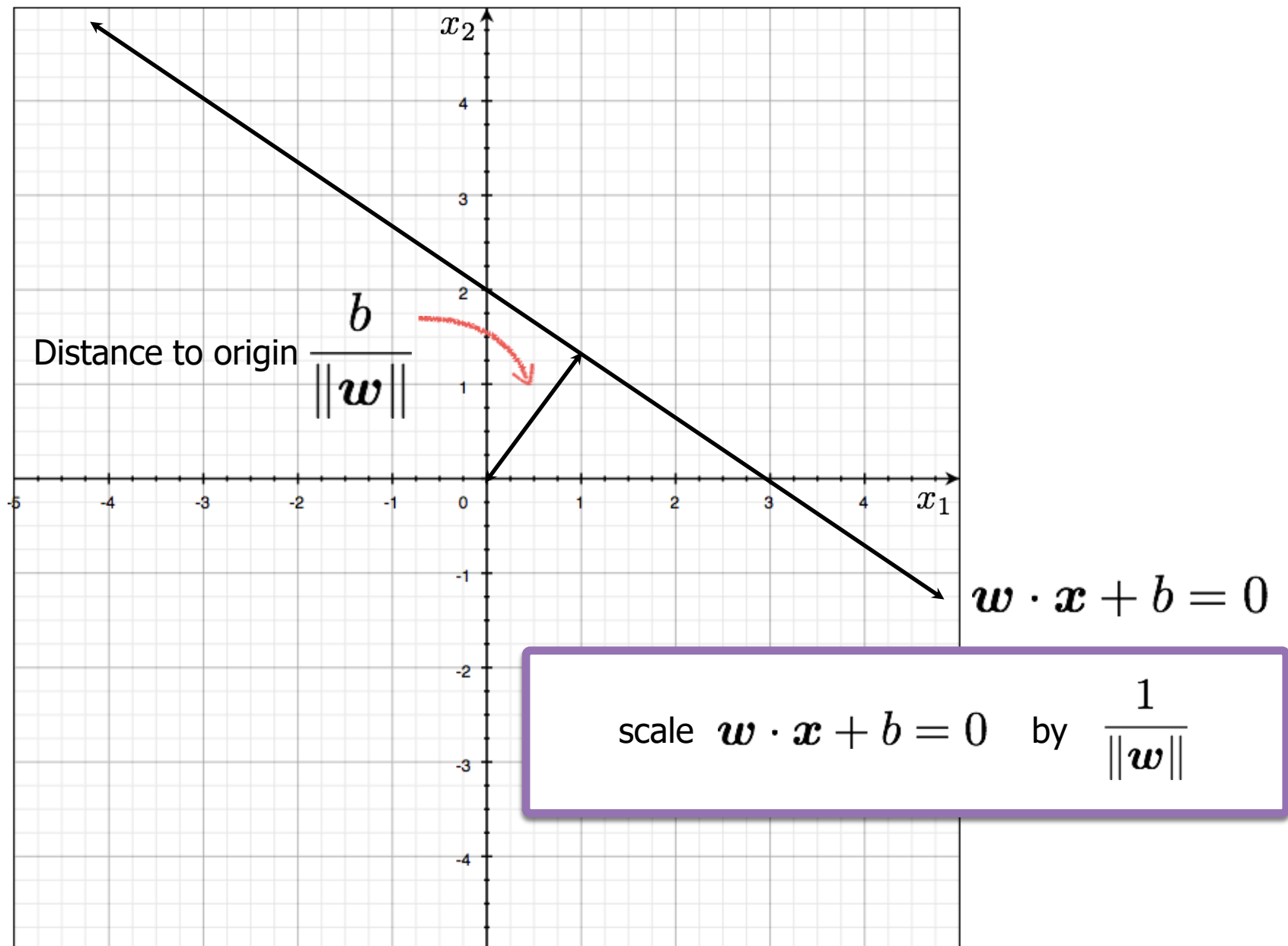
$$d = \frac{|ax_1 + by_1 + c|}{\sqrt{a^2 + b^2}}$$

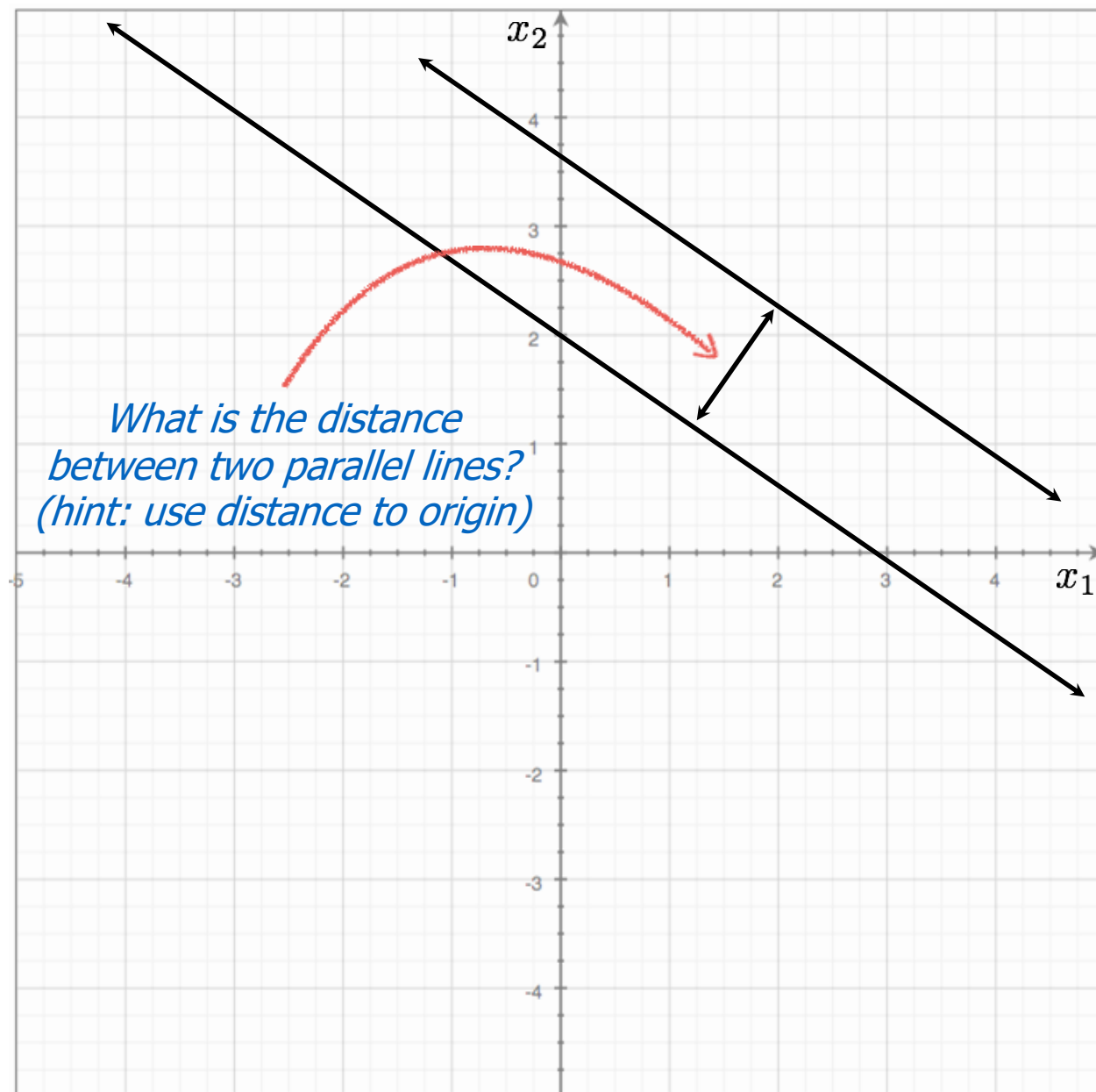
$$w_1 \quad w_2 \quad b$$

$$ax + by + c = 0$$

$$(0,0) \text{ or } \frac{|b|}{\|w\|}$$

$$w \cdot x + b = 0$$



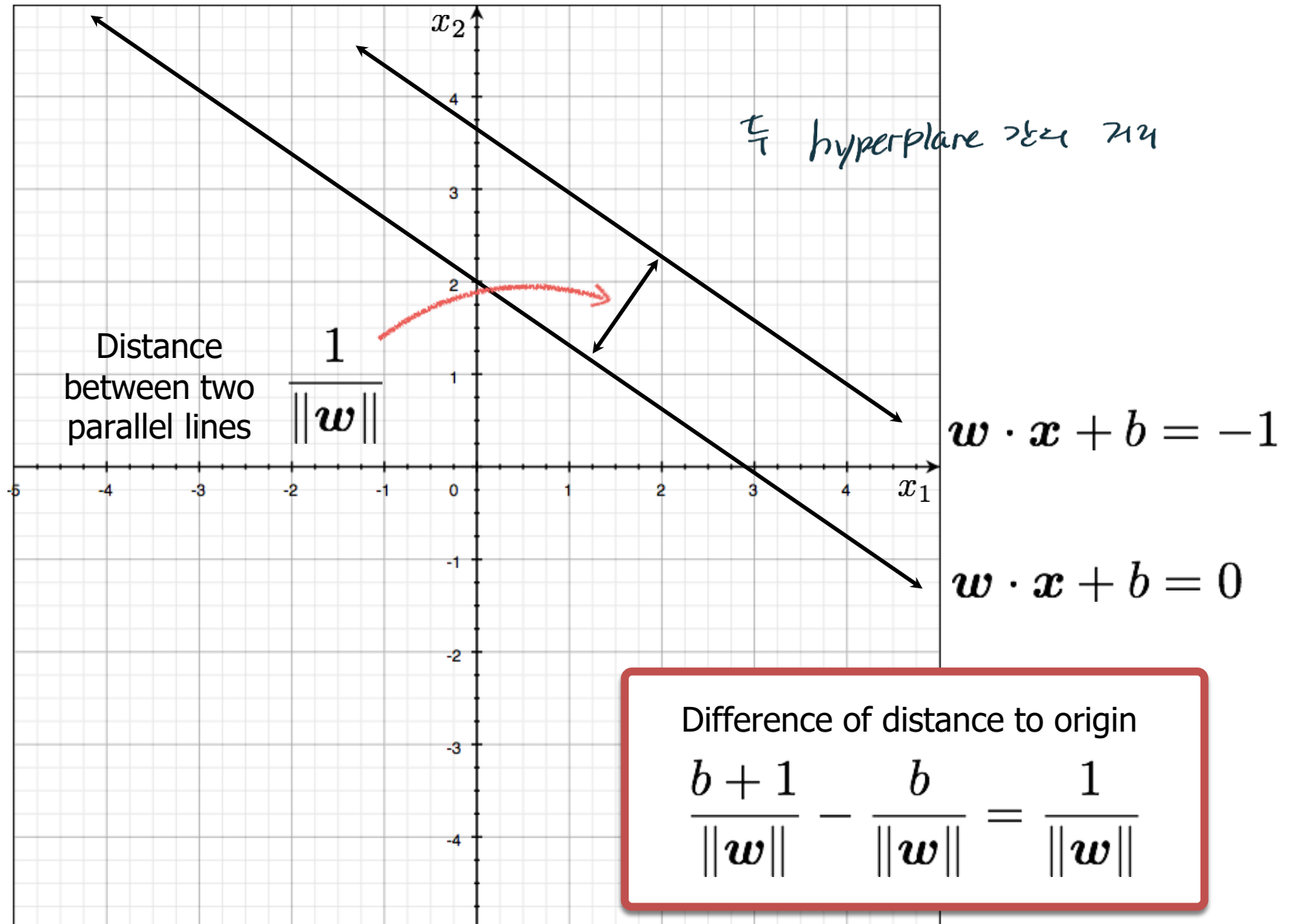


*What is the distance
between two parallel lines?
(hint: use distance to origin)*

$$w \cdot x + b = -1$$

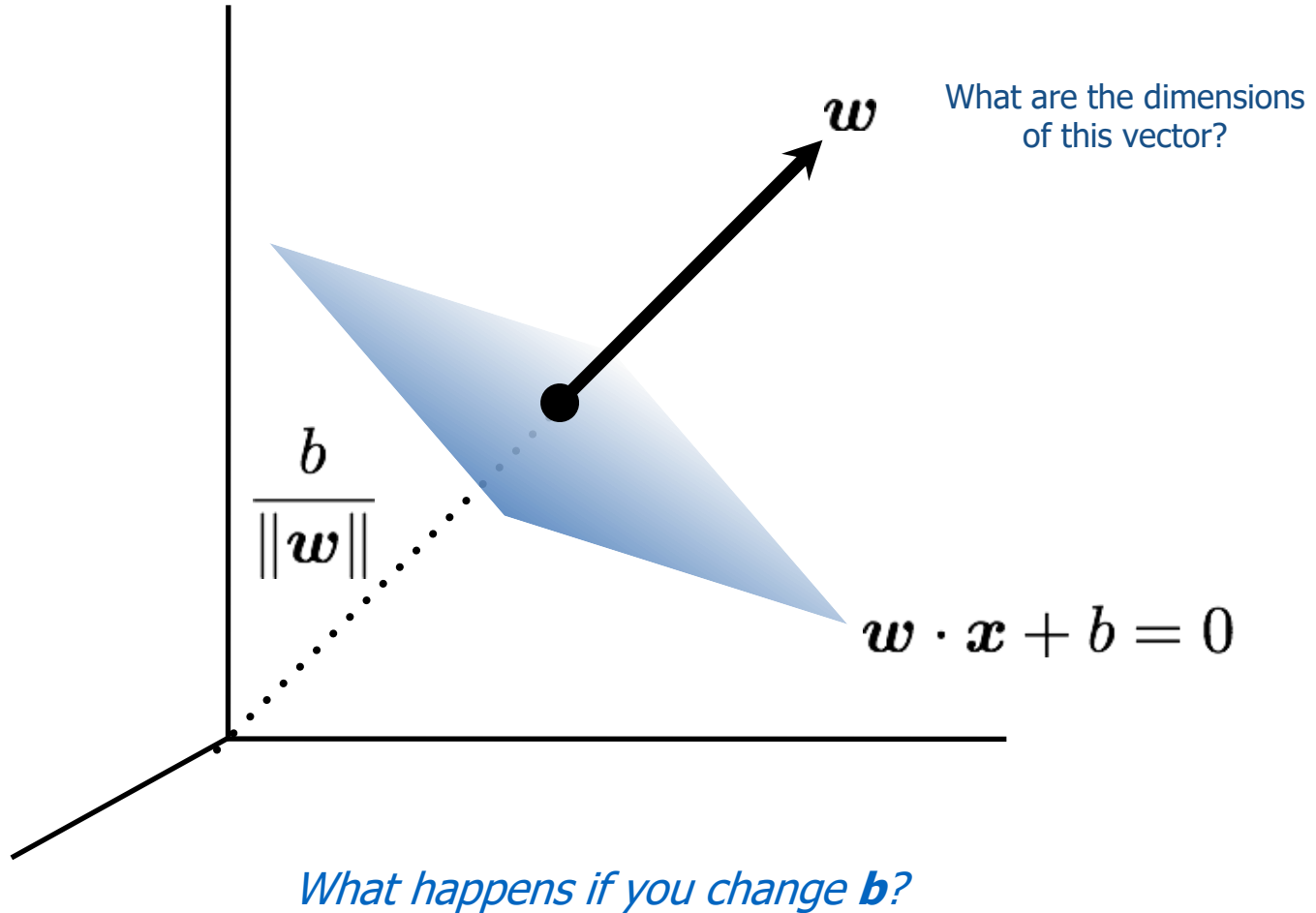
$$w \cdot x + b = 0$$

3D, 4 Dimension으로 바뀌어도 수식은 변하지 않는다.

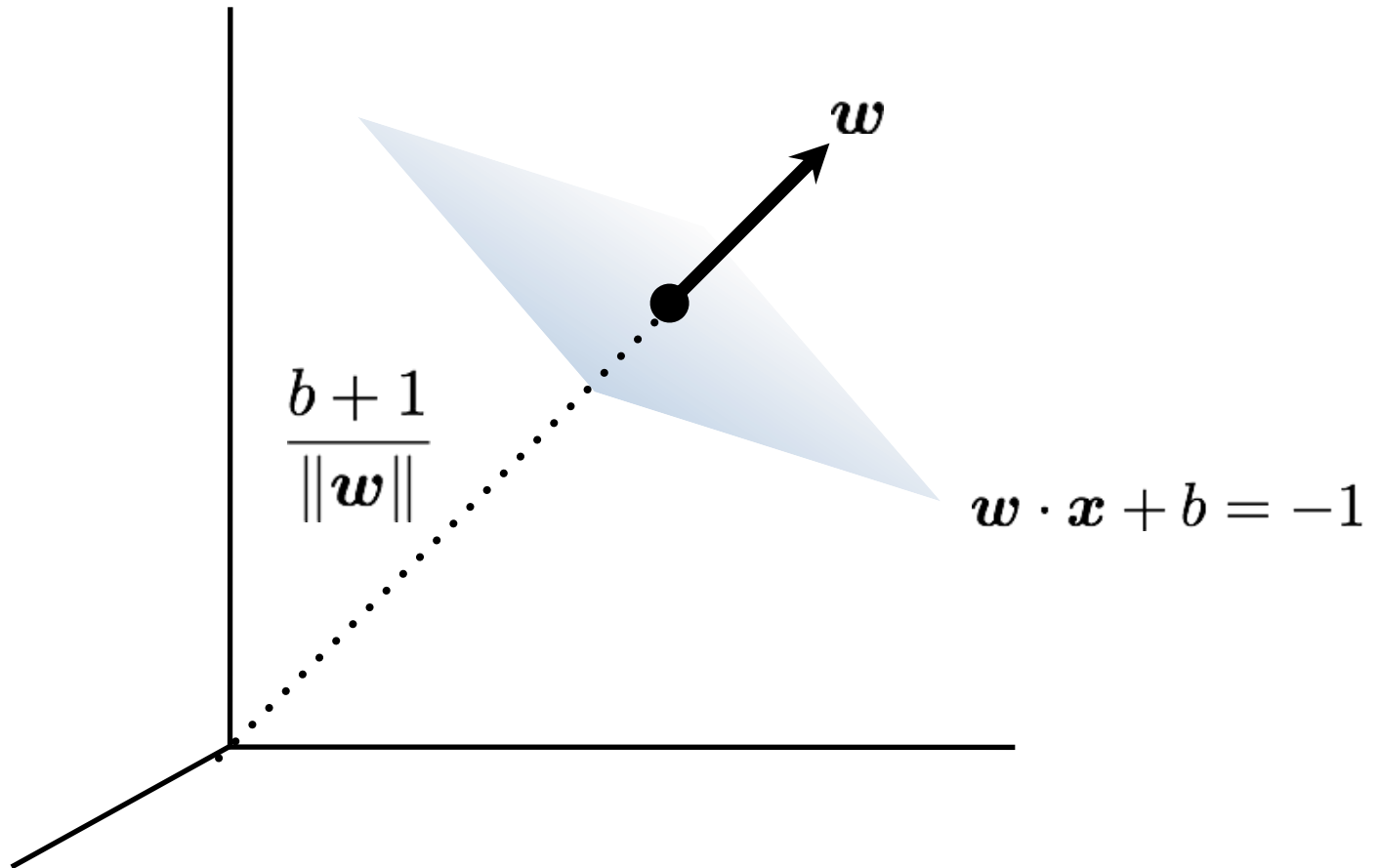


Now we can go to 3D ...

Hyperplanes (planes) in 3D

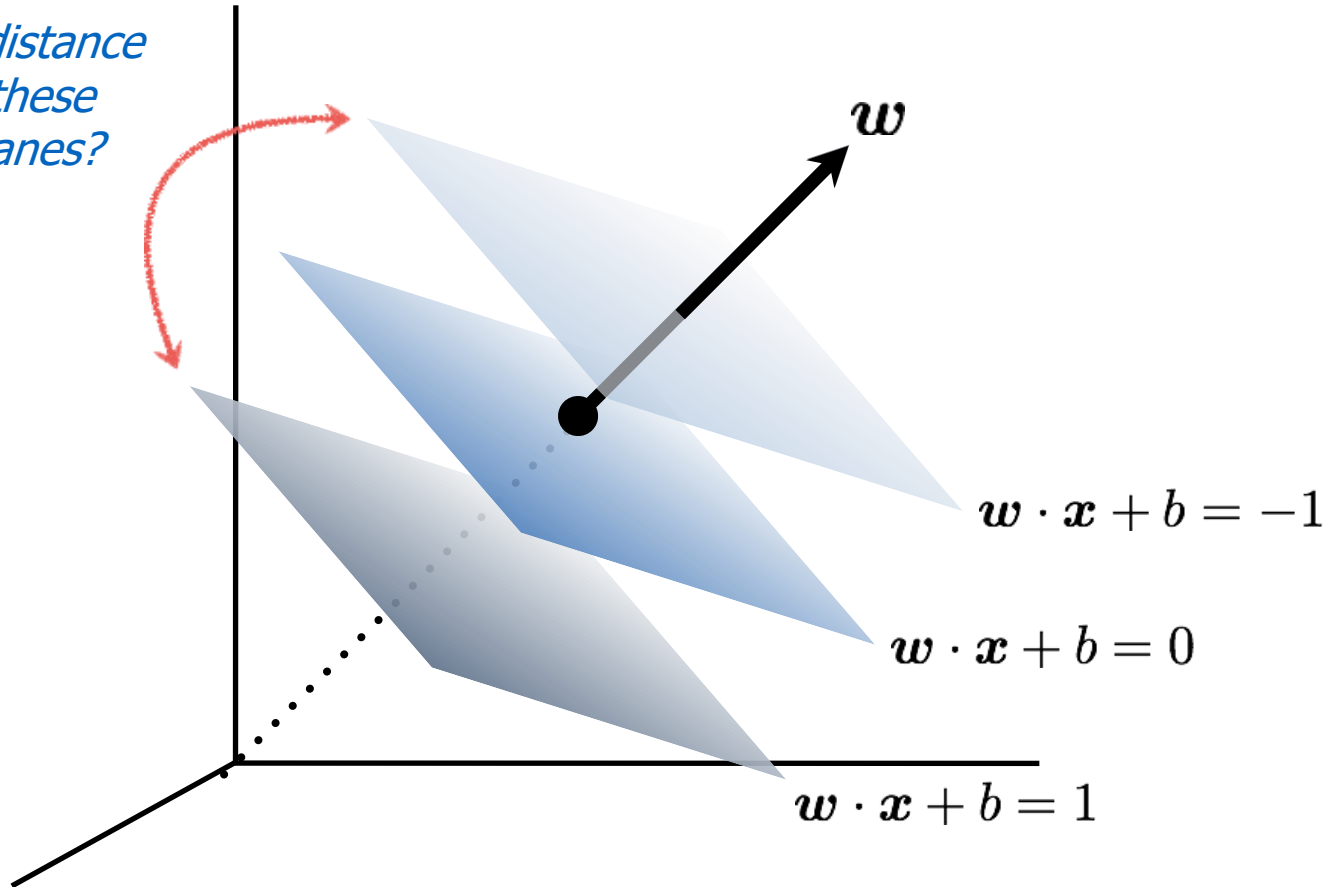


Hyperplanes (planes) in 3D

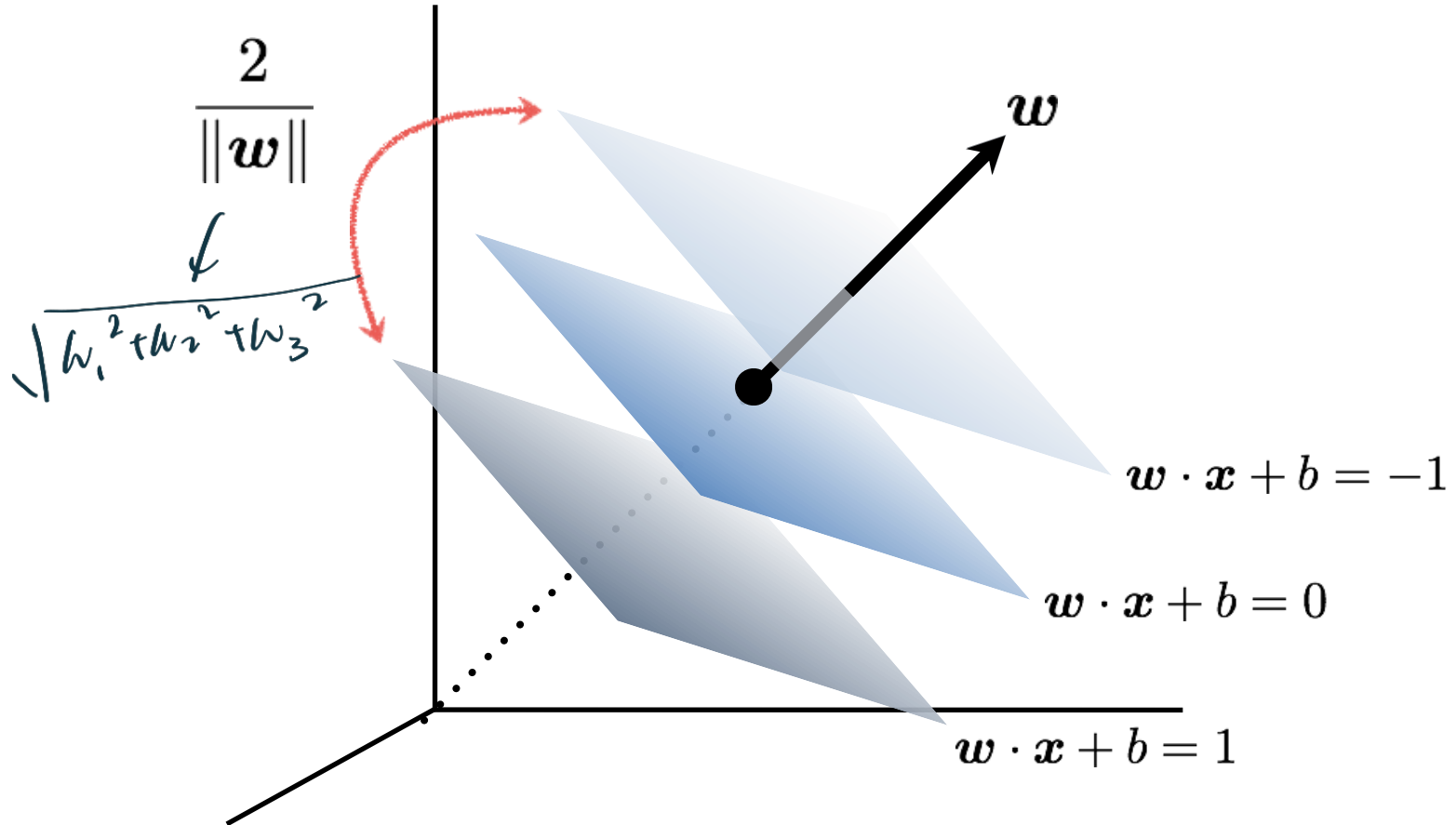


Hyperplanes (planes) in 3D

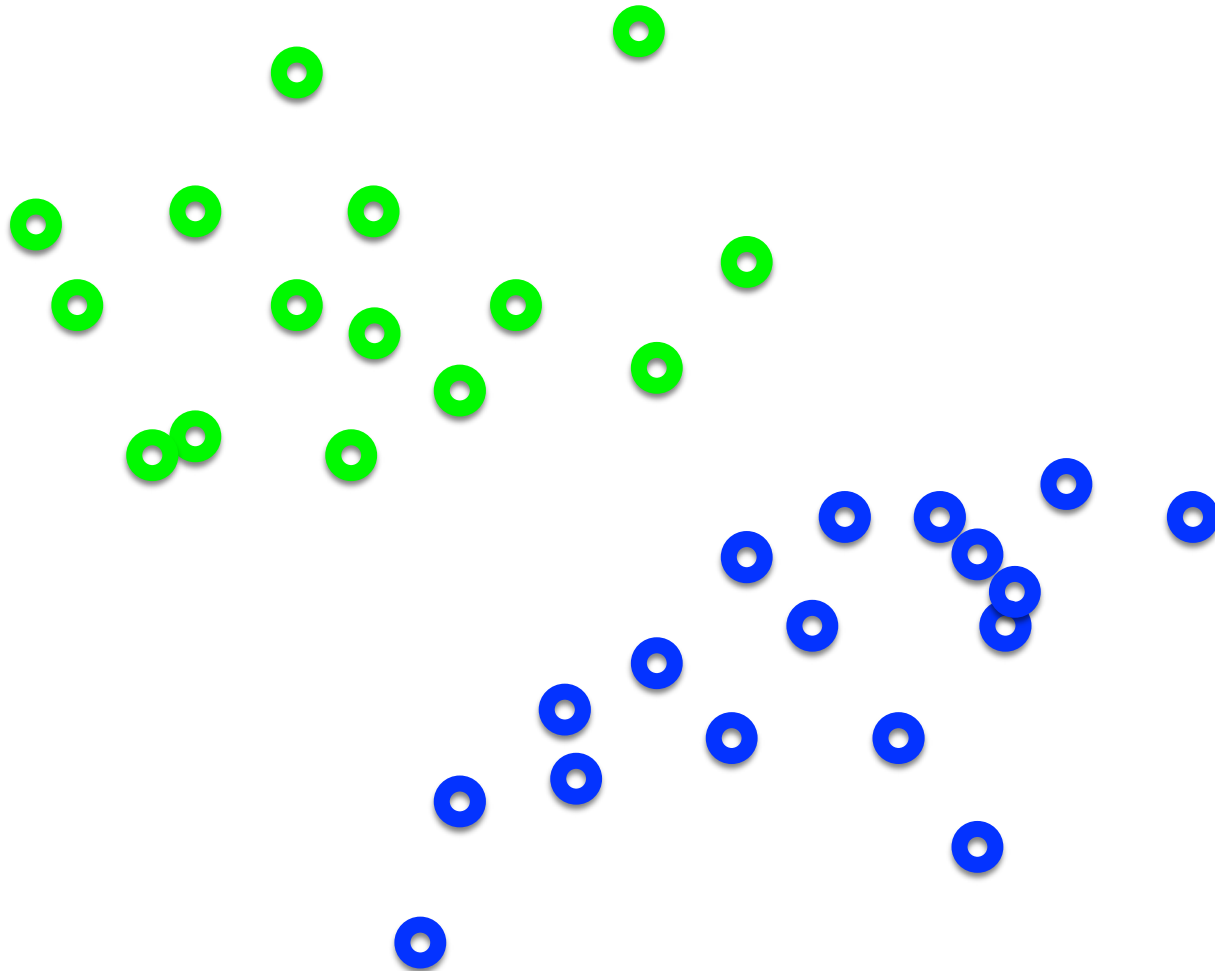
*What's the distance
between these
parallel planes?*



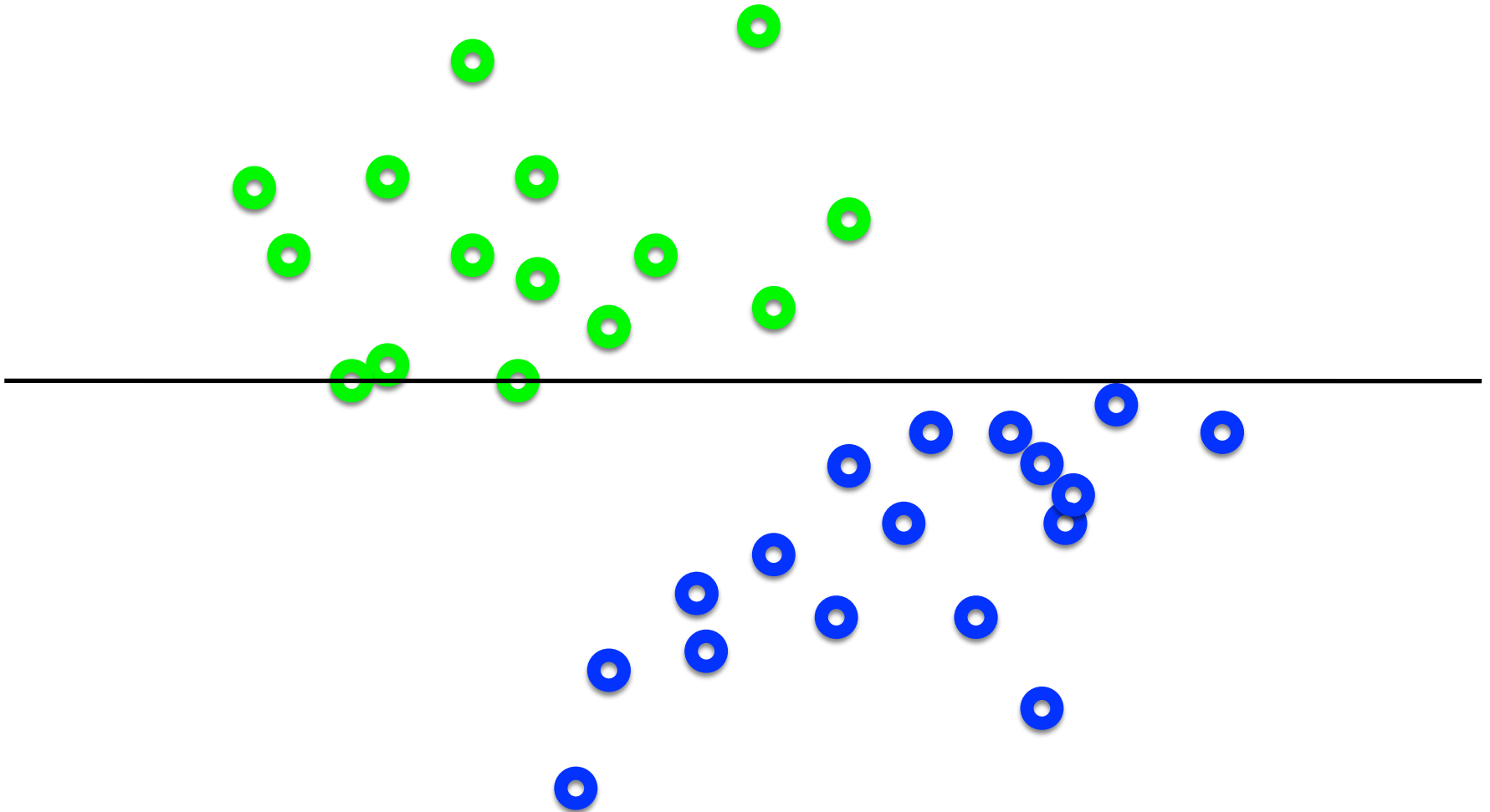
Hyperplanes (planes) in 3D



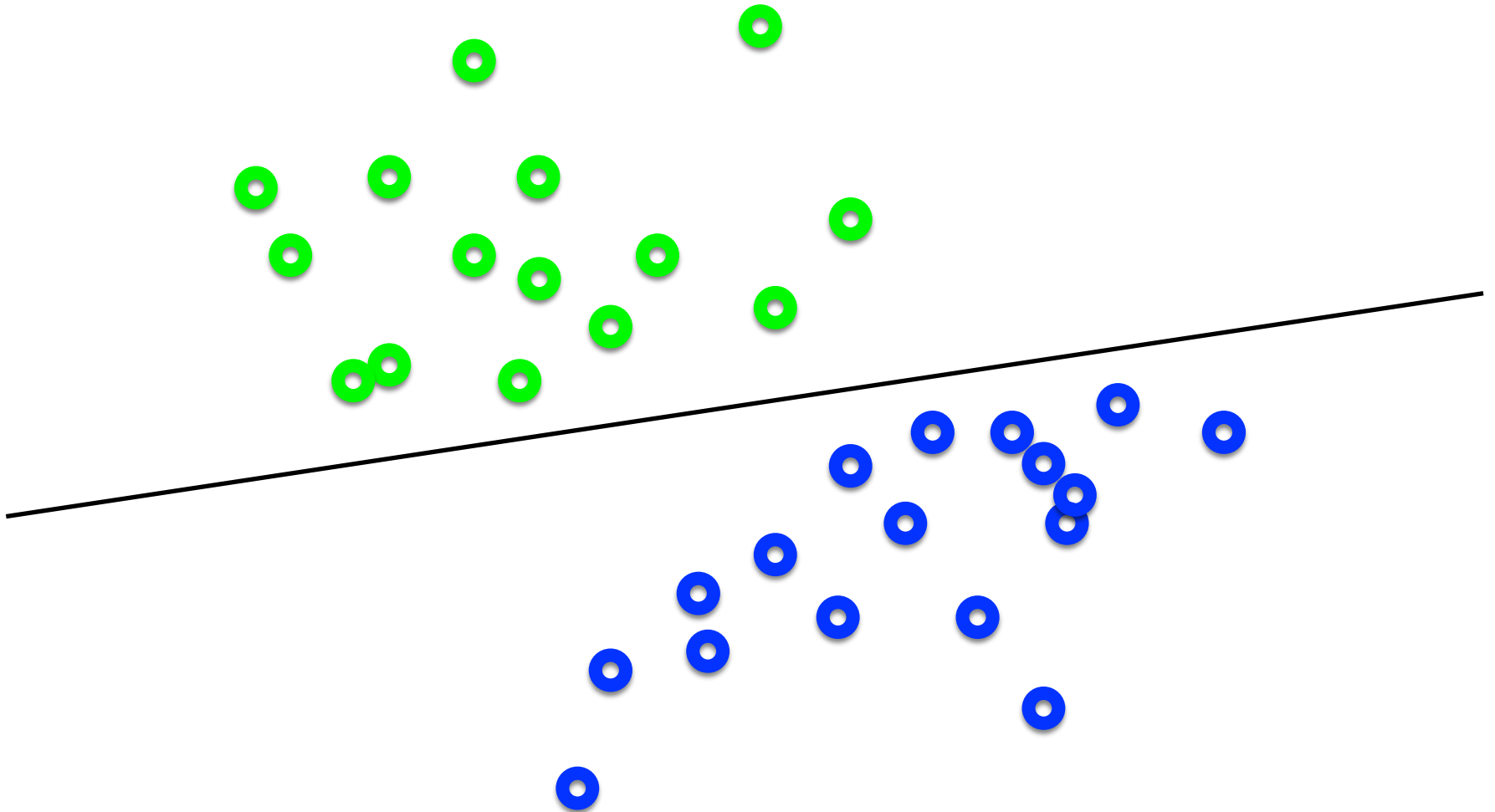
What's the best **w**?



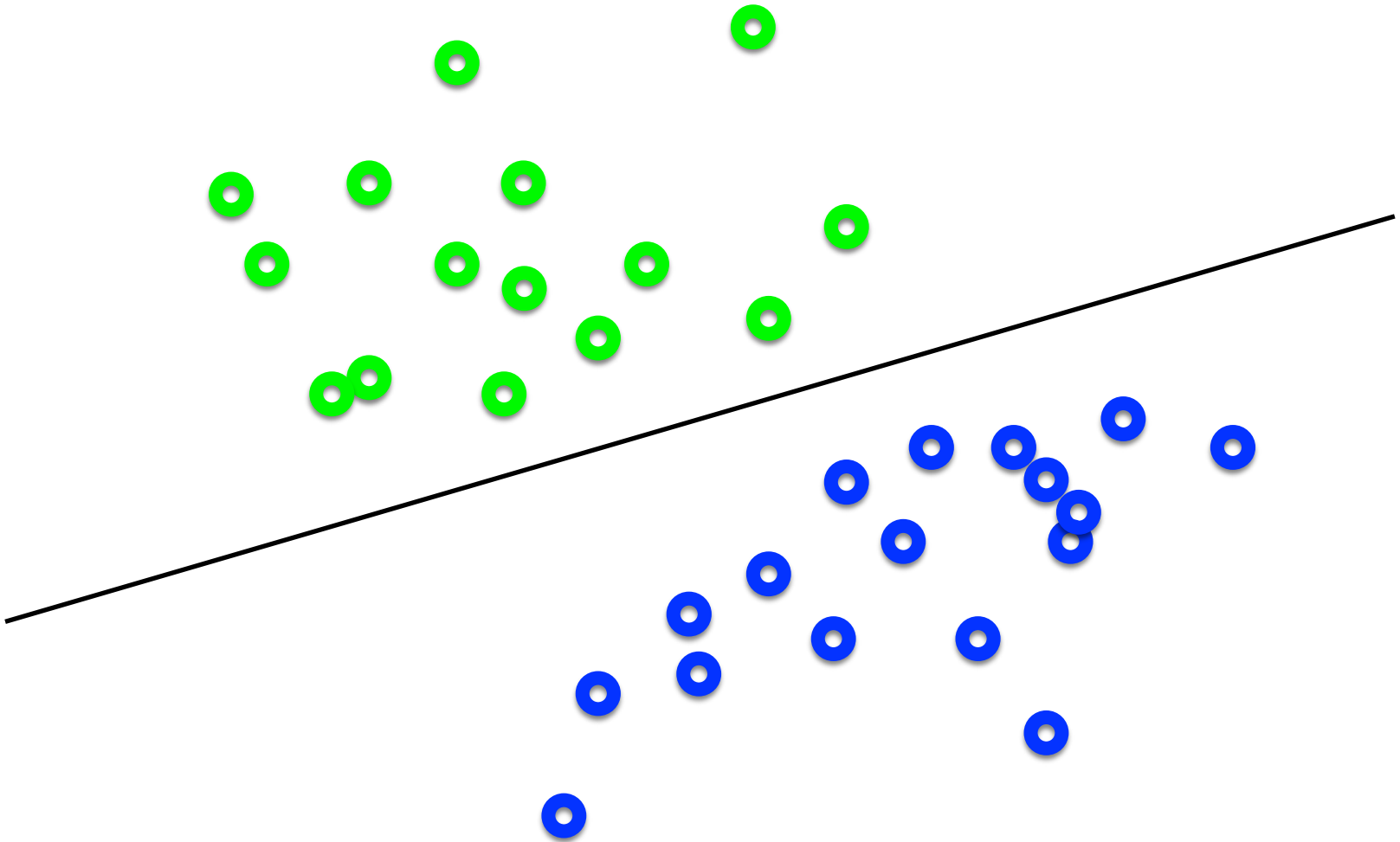
What's the best \mathbf{w} ?



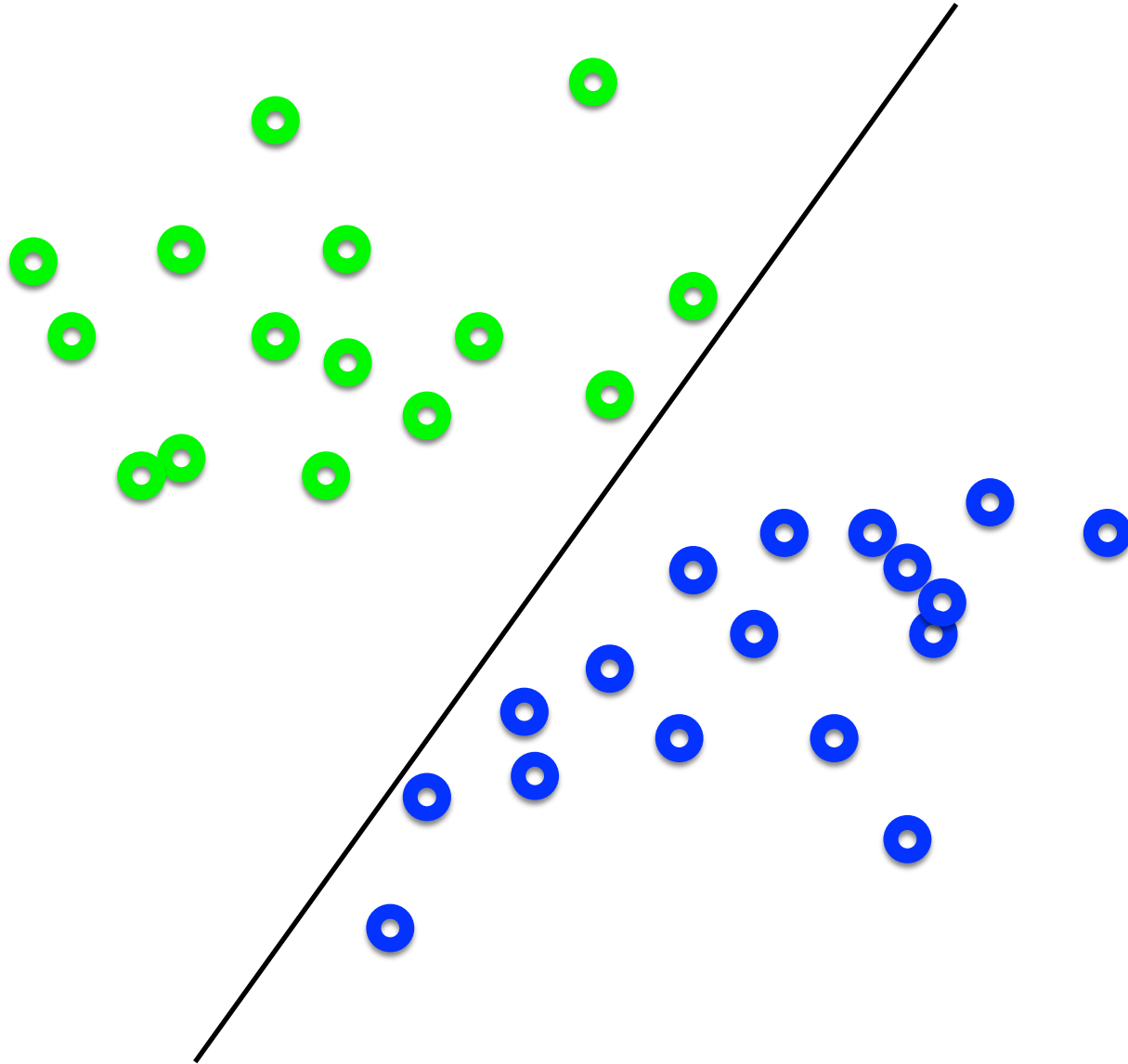
What's the best \mathbf{w} ?



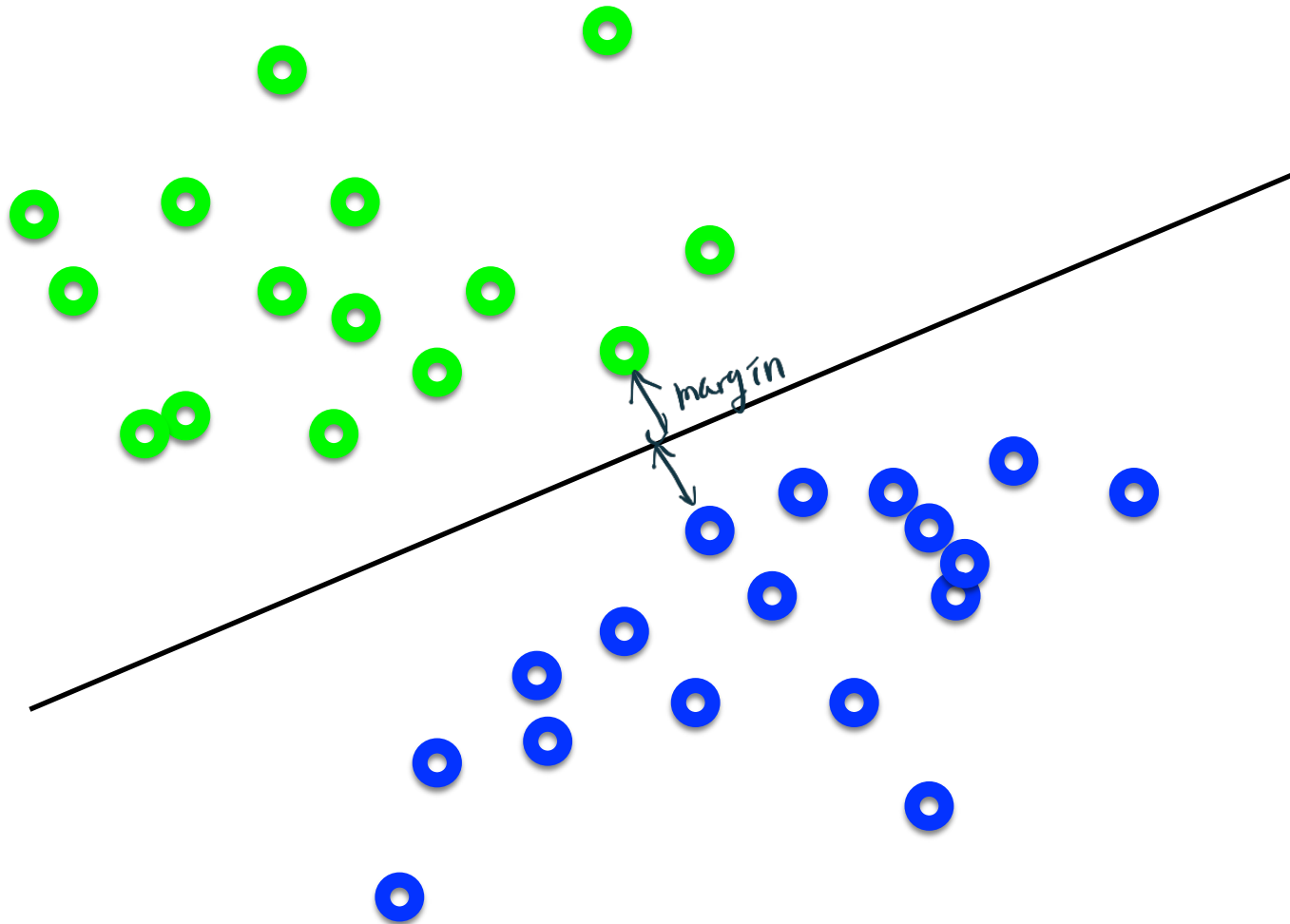
What's the best \mathbf{w} ?



What's the best \mathbf{w} ?

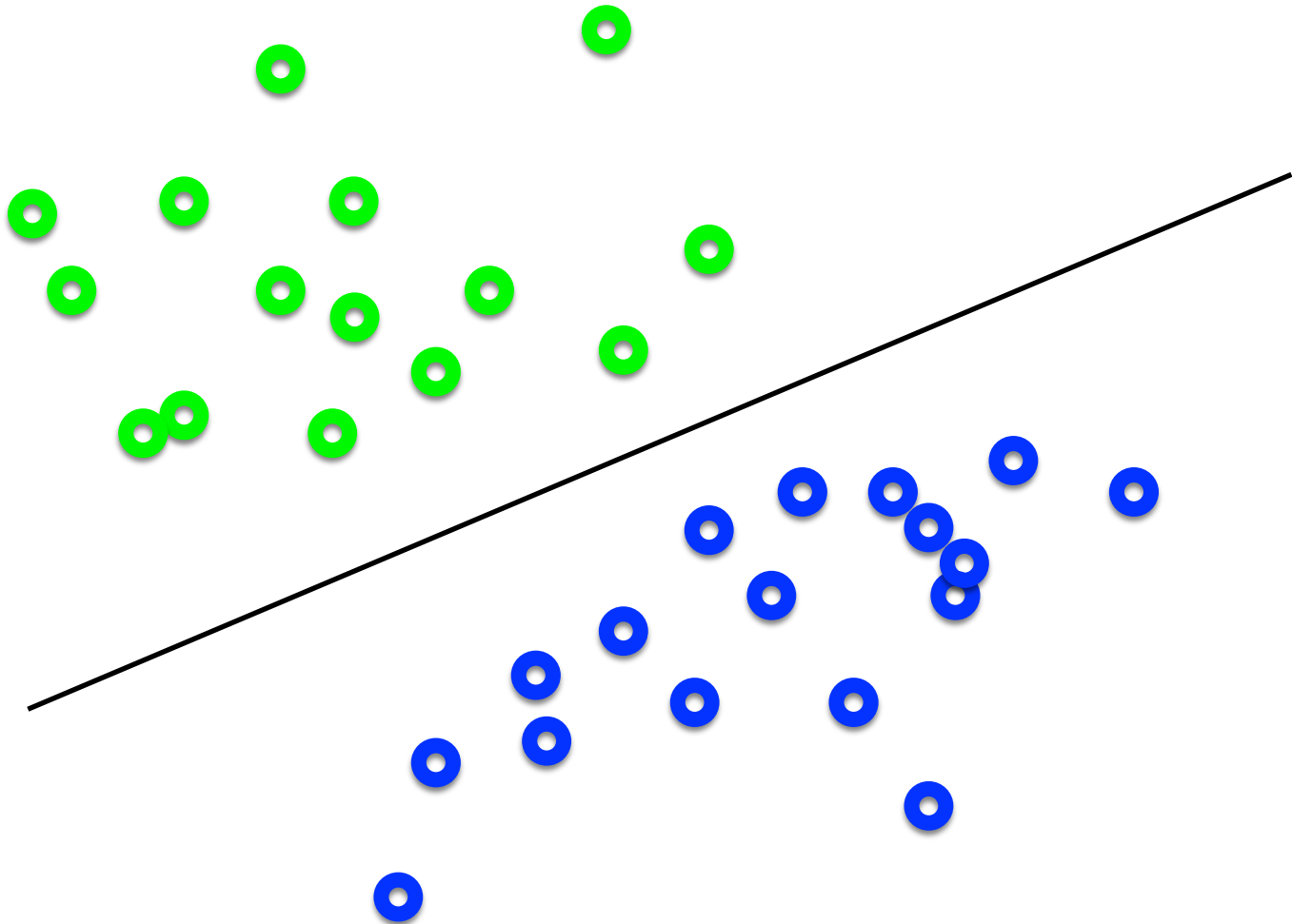


What's the best \mathbf{w} ?



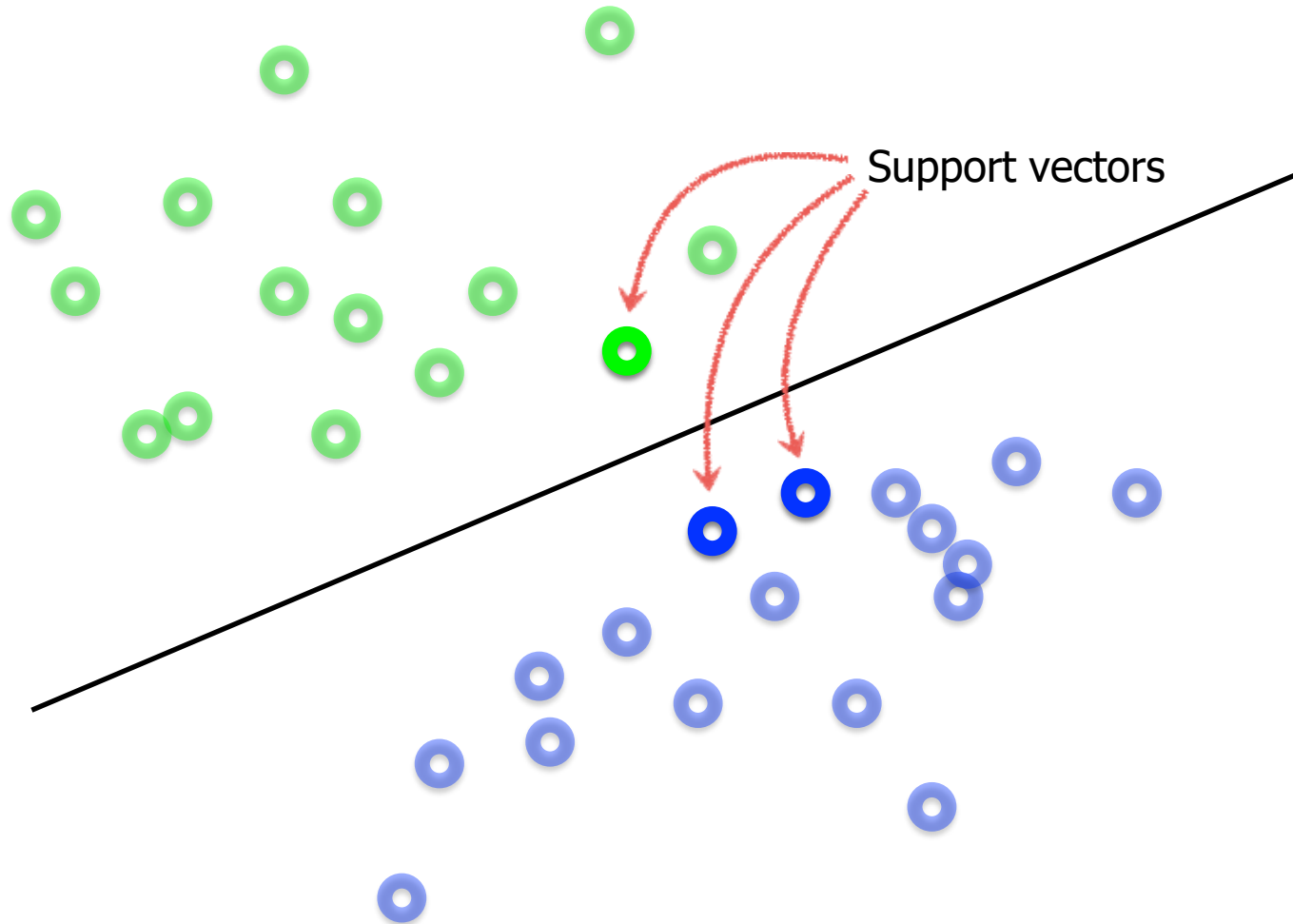
Intuitively, the line that is the farthest from all interior points

What's the best \mathbf{w} ?



Maximum Margin solution:
most stable to perturbations of data

What's the best \mathbf{w} ?



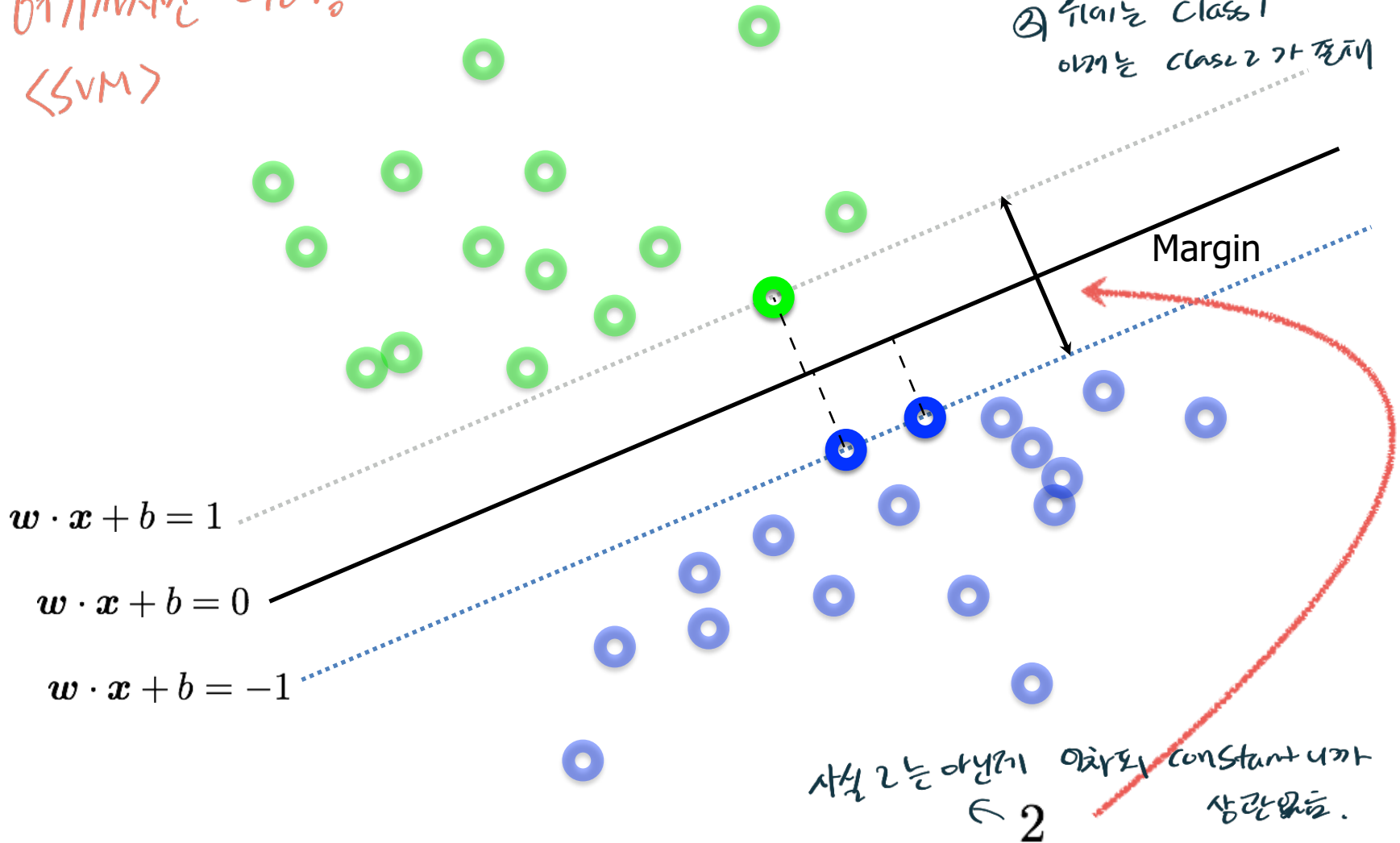
Want a hyperplane that is far away from 'inner points'

Find hyperplane w such that ...

017기까지만 아예!
<SVM>

(constraints : ① Support vector를 지나야 함.)

② Support는 class 1
이제는 class 2가 존재



The gap between parallel hyperplanes $\frac{2}{\|w\|}$ is maximized
 $\|w\|$ is minimized

Another explanation

선 위에 x_p 라는 점이 있다면 가정

Let's say

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$

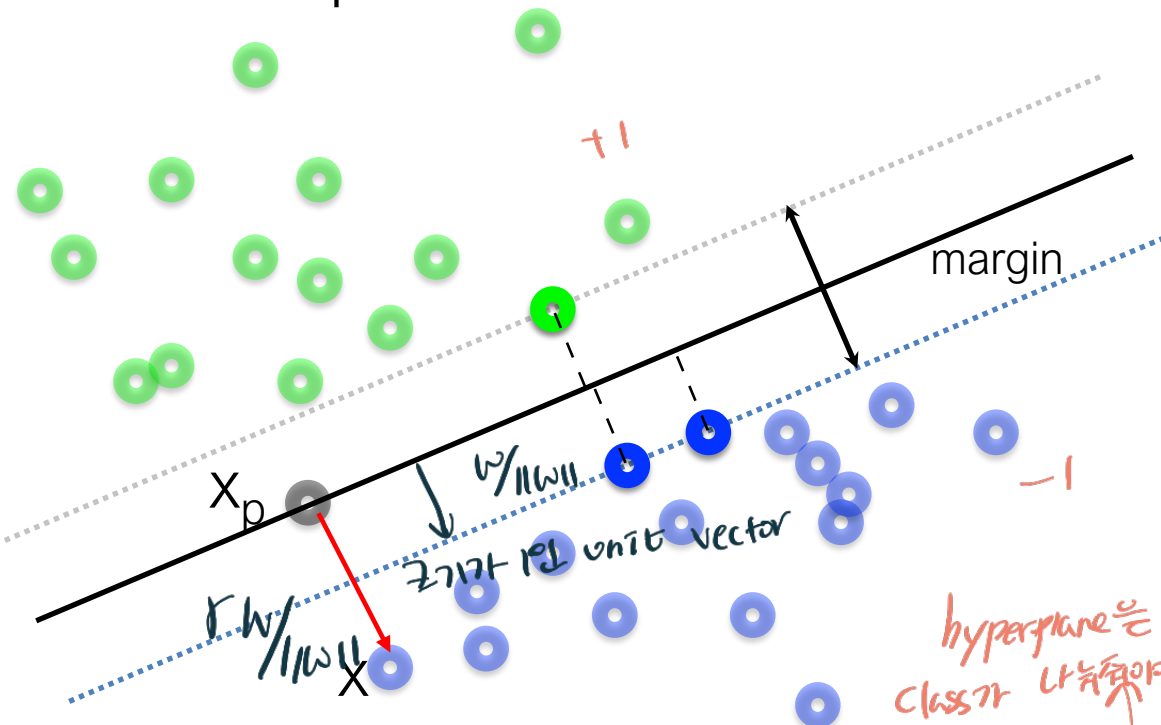
우리가 찾는 hyper plane

A point x on the boundary has

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b = 0$$

A positive point x has

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b = a, a > 0$$



hyperplane은 기준으로
class가 나뉘어야 함.

a 는 constant라 1, 2, or 2.2
나눠줘도 됨.

Maximize margin

$$r = a / \|\mathbf{w}\|$$

We are going to measure the distance
between an arbitrary point \mathbf{x} and a point \mathbf{x}_p on the boundary
and on the perpendicular line from \mathbf{x} to the boundary

unit vector

$$\mathbf{x} = \mathbf{x}_p + r \cdot \mathbf{w} / \|\mathbf{w}\|, f(\mathbf{x}_p) = 0$$

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b = \mathbf{w} \cdot (\mathbf{x}_p + r \cdot \mathbf{w} / \|\mathbf{w}\|) + b = \mathbf{w} \cdot \mathbf{x}_p + b + r \cdot \mathbf{w} \cdot \mathbf{w} / \|\mathbf{w}\| = r \|\mathbf{w}\|$$

The distance is $r = f(\mathbf{x}) / \|\mathbf{w}\|$

$$\max_{\mathbf{w}, b} 2r = 2a / \|\mathbf{w}\|$$

$$\text{s.t. } (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq a, \forall j$$

label ± 1

a is an arbitrary number
and can be normalized

$$\min_{\mathbf{w}, b} \|\mathbf{w}\|$$

$$\text{s.t. } (\mathbf{w} \cdot \mathbf{x}_j + b) y_j \geq 1, \forall j$$

$$\frac{w_1^2 + w_2^2}{\sqrt{w_1^2 + w_2^2}}$$

Can be formulated as a maximization problem

$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

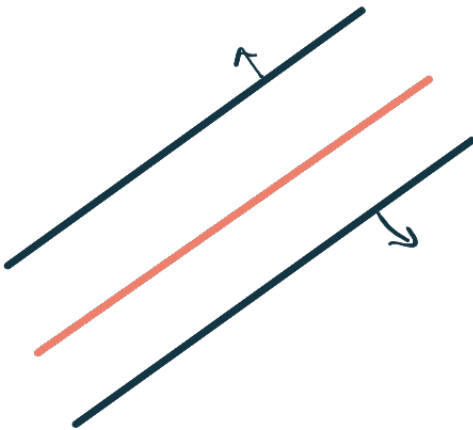
$$\text{subject to } \mathbf{w} \cdot \mathbf{x}_i + b \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \text{ for } i = 1, \dots, N$$

What does this constraint mean?



Label of the data point

Why is it +1 and -1?



Can be formulated as a maximization problem

$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

$$\text{subject to } \mathbf{w} \cdot \mathbf{x}_i + b \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \text{ for } i = 1, \dots, N$$

Equivalently,

Where did the 2 go?

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

$$\text{subject to } y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, N$$

What happened to the labels?

'Primal formulation' of a linear SVM

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

Objective Function

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$ for $i = 1, \dots, N$

Constraints

This is a convex quadratic programming (QP) problem
(a unique solution exists)

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

Objective Function

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$ for $i = 1, \dots, N$

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 = 0 \text{ for } \mathbf{x}_i \in \text{support vectors}$$

Constraints

Lagrange Multiplier Method

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1]$$

minimize w.r.t. \mathbf{w} and b

maximize w.r.t. $\alpha_i \geq 0 \forall i$

α : Lagrange multiplier

- Finding Solution

$$\begin{aligned}\nabla_{\mathbf{w}} L &= \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 & \mathbf{w} &= \sum_i \alpha_i y_i \mathbf{x}_i & \textcircled{1} \\ \nabla_b L &= -\sum_i \alpha_i y_i = 0 & -\sum_i \alpha_i y_i &= 0 & \textcircled{2}\end{aligned}$$

If we know only the value of α , we can find \mathbf{w}

Organize formulas into functions for α

$$\begin{aligned}& \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i + b) - 1] \\&= \sum_{i=1}^N \alpha_i + \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w} \cdot \mathbf{x}_i)] = \sum_{i=1}^N \alpha_i + \frac{1}{2} \mathbf{w}^T \mathbf{w} - \mathbf{w}^T \mathbf{w} \\&= \sum_{i=1}^N \alpha_i - \frac{1}{2} \mathbf{w}^T \mathbf{w} = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j = L(\alpha)\end{aligned}$$

Since the coefficient of the α is negative, the problem of finding the minimum value is replaced with the problem of finding the maximum value.

Find α that maximizes L
 \rightarrow Dual Form

- Finding Solution

$$\max L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j =$$

From ②

$$\sum_i \alpha_i y_i = 0, \alpha_i \geq 0, i = 1, \dots, n$$

From KKT condition (KKT complimentary slackness)

$$\alpha_i = 0 \quad or \quad y_i (w^T x_i + b) - 1 = 0$$

$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 = 0 \rightarrow \alpha_i \neq 0, \mathbf{x}_i$ is support vector

otherwise, \mathbf{x}_i is not support vector \rightarrow it cannot influence to find the \mathbf{w}

Concept of Adaboost

Classifier ensemble

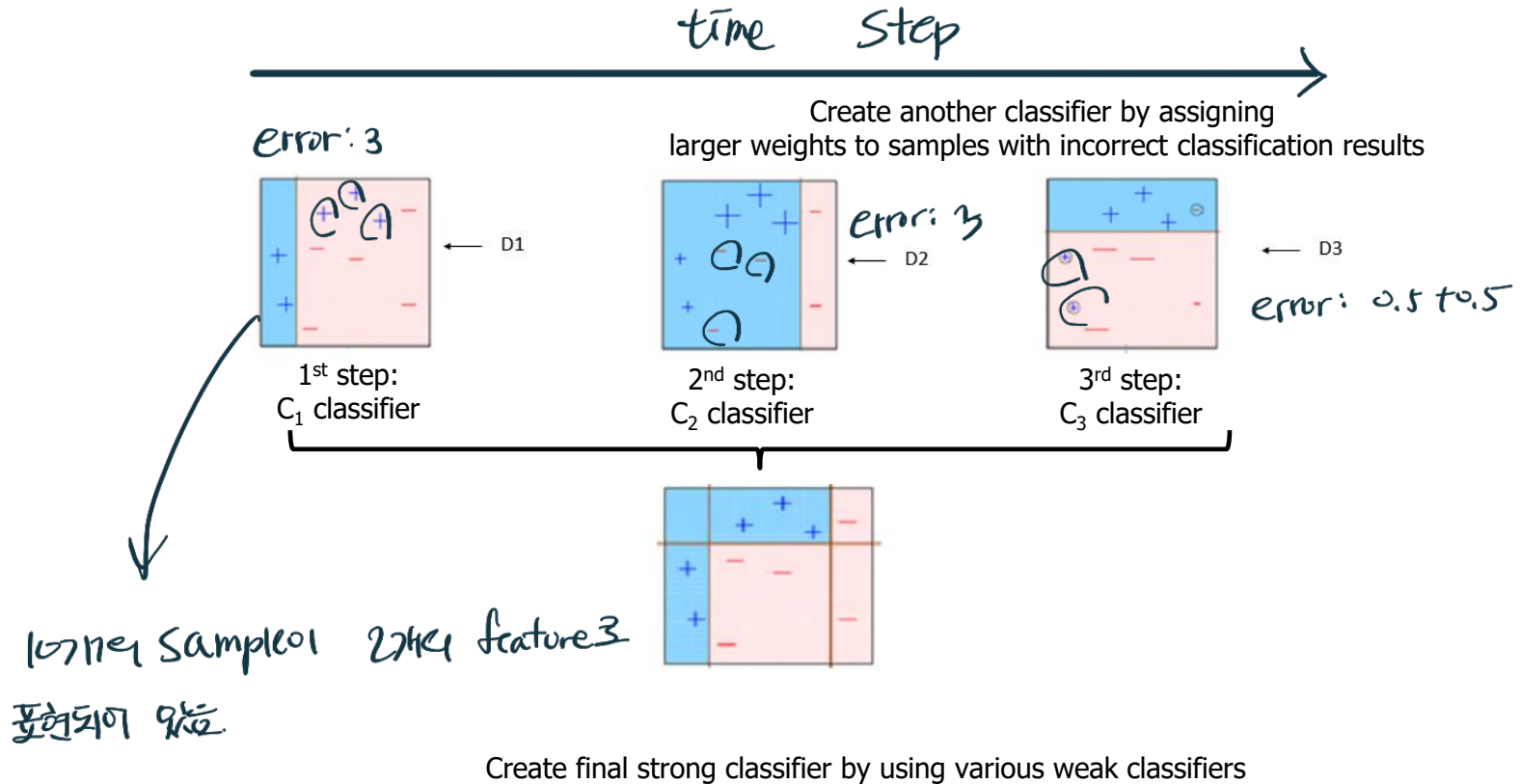
- ◆ The importance of collaboration (ensemble)
 - Collection of individual abilities
 - Better results if you cooperate → Motivation of classifier ensemble

다양한 관점 (집단지성)

Adaboost

◆ Boosting

- Most popular boosting: Adaboost



Adaboost

◆ Boosting

- Most popular boosting: Adaboost
- Using effective resampling
 - The boosting algorithm is created so that the k th classifier c_k and the next c_{k+1} are related
- Samples in training set X are separated by correct or incorrect classification results by c_k
 - The correct sample lowers the weight, and the wrong sample is still a 'tricky' opponent, increasing the weight
 - c_{k+1} is learned considering the above weight

Adaboost

앙상블 기법인 $f(x_i)$ 를 하는 과정을 고려하지 않음.

vector: 1x feature 수

label: scalar (0 or 1)

Input: Training sample, $X = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$, K is the number of classifier

Output: Classifier ensemble, $C = \{(c_k, \alpha_k), 1 \leq k \leq K\}$ // α_k is the Confidence of classifier, c_k

```

1  C = ∅;
2  for(j=1 to N)  $w_j = 1/N$ ; // Use the same weights initially
3  for(k=1 to K) {
4      Learning classifier  $c_k$  considering  $w_1, w_2, \dots, w_N$  // More weighted samples are more important
5       $\varepsilon = 0$ ;
6      for(j=1 to N) if( $c_k(x_j) \neq t_j$ )  $\varepsilon = \varepsilon + w_j$ ; // Calculate error (sum of weights of wrong samples)
7      if( $\varepsilon < 0.5$ ) { // Classifier  $c_k$  is taken only when the error is less than 0.5
8           $\alpha_k = \frac{1}{2} \log\left(\frac{1-\varepsilon}{\varepsilon}\right)$ ; // Confidence of classifier  $c_k$ 
9          for(j=1 to N)
10             if( $c_k(x_j) \neq t_j$ )  $w_j = w_j \times e^{\alpha}$ ; // Increase the weight of wrong samples
11             else  $w_j = w_j \times e^{-\alpha}$ ; // decrease the weight of correct samples
12             Normalize the sum of the  $w_1, w_2, \dots, w_N$  to be 1.
13              $C = C \cup (c_k, \alpha_k)$ ;
14         }
15     else {  $c_k = \text{Nil}$ ;  $C = C \cup (c_k, 0)$ ; }
16 }
```

normalize
한 값의 scale이
다양함

k 가 일어서 많은만 한지

weight가 크면, 틀렸을 때 error가 커짐.
더 정확하게 보겠다.

prediction과 실제 label이랑 다르면,

Classifier 1의 경우

$\varepsilon = 0.3$

0.5보다 크면, 버린다.

$N=3$
 $(0,0) \rightarrow 0 \quad w_1 = \frac{1}{3}$
 $(1,0) \rightarrow 0 \quad w_2 = \frac{1}{3}$
 $(1,1) \rightarrow 1 \quad w_3 = \frac{1}{3}$
 $K=2$
 $(1,1)$ 은 1과 0으로만 가지면
 0.5로 구분 $w_1 = \frac{1}{4}$
 $w_2 = \frac{1}{4}$
 $w_3 = \frac{1}{4}$

C 은 2개
 $0.5 \neq$

Adaboost

Input: Training sample, $X = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}$, K is the number of classifier, t : ground truth

Output: Classifier ensemble, $C = \{(c_k, \alpha_k), 1 \leq k \leq K\}$ // α_k is the Confidence of classifier, c_k

```
1  C = ∅;
2  for(j=1 to N)  $w_j = 1/N$ ;
3  for(k=1 to K) {
4      Learning classifier  $c_k$  considering  $w_1, w_2, \dots, w_N$ 
5       $\varepsilon = 0$ ;
6      for(j=1 to N) if( $c_k(\mathbf{x}_j) \neq t_j$ )  $\varepsilon = \varepsilon + w_j$ ;
7      if( $\varepsilon < 0.5$ ) {
8           $\alpha_k = \frac{1}{2} \log\left(\frac{1 - \varepsilon}{\varepsilon}\right)$ ;
9          for(j=1 to N)
10             if( $c_k(\mathbf{x}_j) \neq t_j$ )  $w_j = w_j \times e^{\alpha}$ ;
11             else  $w_j = w_j \times e^{-\alpha}$ ;
12             Normalize the sum of the  $w_1, w_2, \dots, w_N$  to be 1.
13              $C = C \cup (c_k, \alpha_k)$ ;
14         }
15     else { $c_k = \text{Nil}$ ;  $C = C \cup (c_k, 0)$ ; }
16 }
```

Adaboost

Input: Training sample, $X = \{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_N, t_N)\}$, K is the number of classifier, t : ground truth

Output: Classifier ensemble, $C = \{(c_k, \alpha_k), 1 \leq k \leq K\}$ // α_k is the Confidence of classifier, c_k

```
1  C = ∅;
2  for(j=1 to N)  $w_j = 1/N$ ; // Use the same weights initially
3  for(k=1 to K) {
4      Learning classifier  $c_k$  considering  $w_1, w_2, \dots, w_N$  // More weighted samples are more important
5       $\varepsilon = 0$ ;
6      for(j=1 to N) if( $c_k(\mathbf{x}_j) \neq t_j$ )  $\varepsilon = \varepsilon + w_j$ ; // Calculate error (sum of weights of wrong samples)
7      if( $\varepsilon < 0.5$ ) { // Classifier  $c_k$  is taken only when the error is less than 0.5
8           $\alpha_k = \frac{1}{2} \log\left(\frac{1-\varepsilon}{\varepsilon}\right)$ ; // Confidence of classifier  $c_k$ 
9          for(j=1 to N)
10             if( $c_k(\mathbf{x}_j) \neq t_j$ )  $w_j = w_j \times e^{\alpha}$ ; // Increase the weight of wrong samples
11             else  $w_j = w_j \times e^{-\alpha}$ ; // decrease the weight of correct samples
12             Normalize the sum of the  $w_1, w_2, \dots, w_N$  to be 1.
13              $C = C \cup (c_k, \alpha_k)$ ;
14         }
15     else { $c_k = \text{Nil}$ ;  $C = C \cup (c_k, 0)$ ; }
16 }
```