# The Perceptron

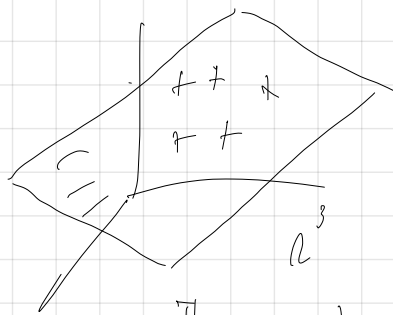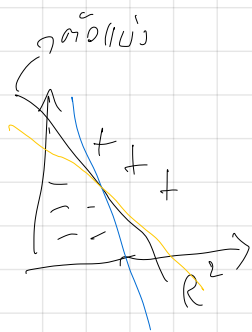- **Intro** : - The first consider learning algorithm

    - Assumption :

        - ทำงานกับ Binary Classification

        $y_i \in \{-1, +1\}$ #

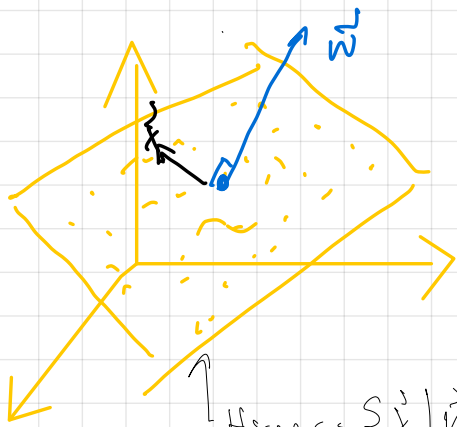        - Data is linearly separable.



- The perceptron will try to learn a hyperplane that separate between the data

# Hyper plane

   - a sub space whose dimension is one less than the dimesion of ambient space

     $\underline{\text{space ที่ใหญ่}}$

   $X = \mathbb{R}^d \qquad Hyper \subseteq \mathbb{R}^{d-1}$

   ### for high dimensional data จะทำงานได้ดี, เพราะ data point จะเป็นจำนวน สามารถหาเส้น hyper plane ได้
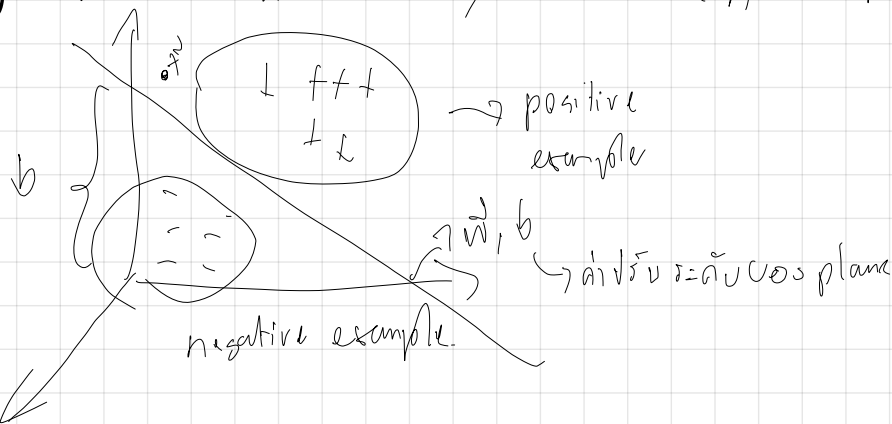


$\vec{w}^T \vec{x}_{+b} = 0$

$$\begin{bmatrix} w^1 \\ w^2 \\ w^3 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w^1x^1 + x^2x^2 + w^3x^3 \end{bmatrix} = 0$$

↑Hyper $= \{\vec{x} \mid \vec{w}^T \vec{x}_{+b} = 0\}$  bias    learn $\vec{w}$ (normal vector to hyperplane)

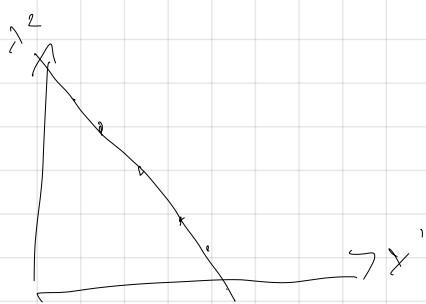**Traning** : learn normal vector $\vec{w}$, and bias term $b$ form date



→ positive example

→ $\vec{w}, b$ → ค่าปรับระดับของ plane

negative example

**testing** : $h(x) = sign(\vec{w}^T \vec{x} + b)$

(ดูสัญลักษณ์= ถ้า = 0 อยู่ใน hyperplane (-1,+1)
  >0  positive
  <0  Negative

# Point of hyperplane

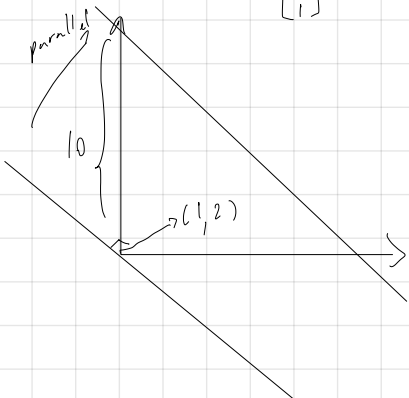$$x^2 = -2x^1 + 10$$

$$y = -2x + 10$$



$$\vec{w}^T \vec{x} + b = 0$$

inner product

$$[w_1 \; w_2] \begin{bmatrix} x^1 \\ x^2 \end{bmatrix} + b = 0$$

$$w_1 x^1 + w_2 x^2 + b = 0 \qquad = x^2 = -2x^1 + 10$$

$$= x^2 + 2x^1 - 10$$

$$w_1 = 2 \quad w_2 = 1 \quad b = -10 \qquad \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$



2) $\vec{x}_{test} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$ ; $h(\vec{x}_{test}) = ?$

$$h(\vec{x}_{test}) = sign(\vec{w}^T \vec{x}_{test} + b)$$



$$= sign\left([2 \; 1] \begin{bmatrix} 2 \\ 3 \end{bmatrix} + (-10)\right)$$
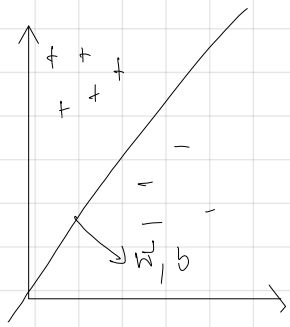
$$= sign(10 + (-10)) = 0 = Negative$$

$$sign(4 + 3 + (-10)) = neg.$$

# Perception

- $y_i \in \{-1, +1\}$
- D must be linearly separable
- Testing : $sign(\vec{w}^T \vec{x} + b)$
- Traing: learn to find $\vec{w}, b$ from D

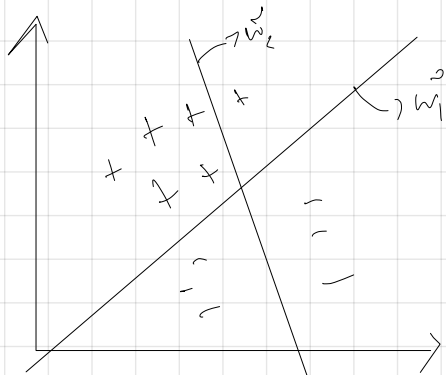## How to learn $\vec{w}, b$? ★ $\vec{w}'\vec{x}'$

$\vec{w} = \begin{bmatrix} w^1 \\ w^2 \\ \vdots \\ w^d \end{bmatrix}$ ; $\vec{x} = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^n \end{bmatrix}$ $= \vec{w}^T \vec{x} + b$ , $w^1 x^1 + w^2 x^2 \ldots w^d x^d + b(1)$

$\vec{w}' = \begin{bmatrix} w^1 \\ w^2 \\ w^d \\ b \end{bmatrix}$ $\vec{x}' = \begin{bmatrix} x^1 \\ \vdots \\ x^d \\ 1 \end{bmatrix}$

— $\vec{w}_1$ is the right normal vector but $\vec{w}_2$ isn't

$\vec{w}^T \vec{x} \longrightarrow > 0 \oplus$ predict
$\longrightarrow < 0 \ominus$ predict

$y(\vec{w}^T \vec{x})$ ถ้าถูก prediction $\oplus$, $y \oplus > 0$
↑ predict
↑ true label $\longmapsto 1 \oplus, y \ominus \leq 0$ ✗ เดา $\ominus \neq \oplus$ : misclassification
$\longmapsto -1 \ominus, y\ominus > 0$

$\forall (x_i, y_i) \in D$ , $y_i(\vec{w}_1^T \vec{x}_i) \geqslant 0$ $\longrightarrow \vec{w}_1^T \vec{x}_i$ มากกว่า $\vec{w}_1^T \vec{x}_i \oplus$ เป็อลากตำแหน่ง
 เอาทุก $x_i y_i$ จาก D

classify correctly for        สมการ train
any data set in D

## Misclassification : $\exists (x_i, y_i) \in D, y_i(\vec{w}^T \vec{x}_i) \leq 0$

$\Rightarrow \vec{w}$ is not the right normal vector

## Perceptron algorithm:

```
0: w' ← 0        ; |w'| = d+1
1: while true:
2:     m ← 0    // counter to count # of misclassif..
3:     for (xi yi) ∈ D:
4:         if yi(w'ᵀxi) ≤ 0: // if miss
5:             m ← m+1    // update counter m
6:             w' ← w' + yi xi  // update w'
7:     if m=0:
8:         break
```

$y_i(w_1^T x_i) \leq 0$ ; miss
$< 0 \quad > 0$

$\vec{w}_{t+1} \leftarrow \vec{w}_t + (-1)(\vec{x}_1)$

$\underbrace{\vec{w}_t - \vec{x}_1}$

Update (rotate)

เกิดที่ไปโดนสายไปเหรอ?
OK
ชูสองนิ้ว = ไม่โกรธ

สรุป:

- Assume D is linearly separable
- The perceptron algoritm will find
  a hyperplane that separates the data
- $y \Rightarrow$ binary classification

$(x_i, y_i) \in D \rightarrow \boxed{\text{learn}} \rightarrow \left(\vec{w}, b\right)$   $\vec{w} \in \mathbb{R}^{d+1}$

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

( bias ใช้เป็นได้

$\rightarrow \vec{w}$

$\vec{x}_i$ ใหม่ $= \begin{bmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^d \\ 1 \end{bmatrix}$

Algorithm : try to adjust the vector $\vec{w}$ to satisfy

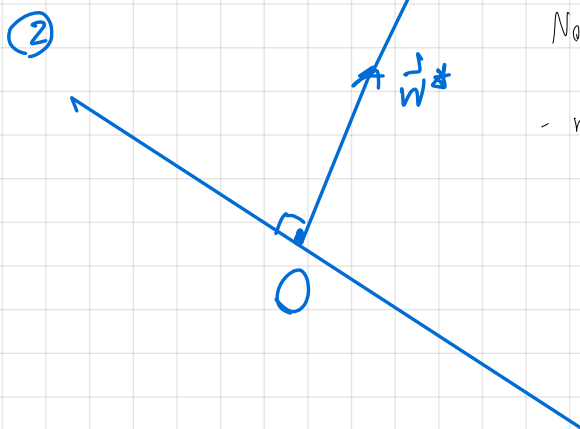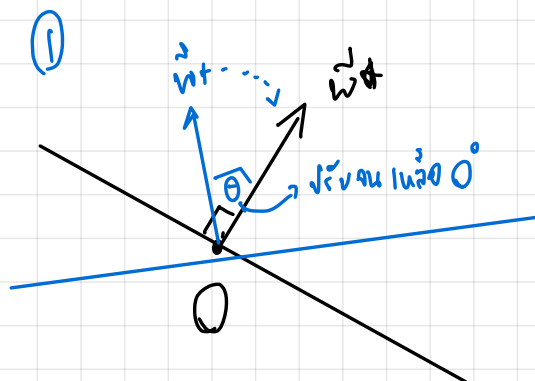$$\boxed{\forall (x_i, y_i) \in D, \; y_i(\vec{w}^T \vec{x}_i) > 0}$$

- update (when mis classifying): $y_i(\vec{w}^T \vec{x}_i) \leq 0$
  - $\vec{w} \leftarrow \vec{w} + y_i \vec{x}_i$

Q: When the algorithm terminates
   how can we be sure that $\vec{w}$ defines
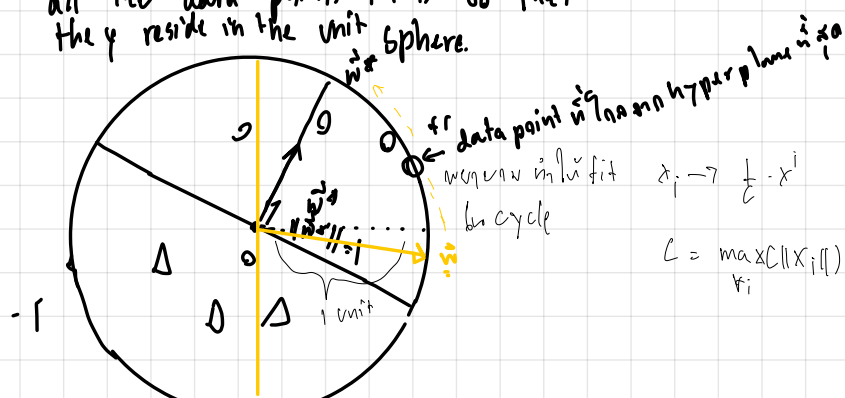   the right hyperplane?

Perceptron's covergence Proof:

- provide strong formal guarantee of
  the separating hyperplane by the algorithm

- setup : - Assumption: D is linearly separable

  - Equivalently, there exists $\vec{w}^*$ such that
    $\hookrightarrow$ define right hyperplane (ขอบเขตที่ใช่จริง)

    $\vec{w}$ vector ในalgorithms ที่ converge
    เข้าใกล้ $\vec{w}^*$

    $\forall (x_i, y_i) \in D, \; y_i(\vec{w}^{*T} \vec{x}_i) > 0$

①



②



Note: there are infinitely many such $\vec{w}^*$
- we will focus on $\vec{w}^*$ with $\|\vec{w}^*\| = 1$
  (unique)

- Futher more, we consider the rescaling of
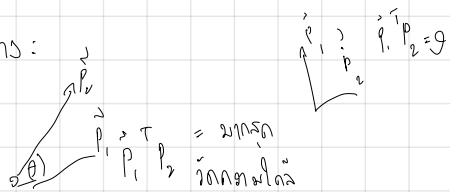  all the data points in D so that
  they reside in the unit sphere.



+r data point ที่ไกลจาก hyperplane ที่สุด
  เพราะงั้น เกิดเป็น fit $x_i \rightarrow \frac{1}{c} \cdot x^i$
  bi cycle

  $c = \max_{\forall i} (\|X_i\|)$

Effect caused by an update:

- Each update is made in hope of tuning $\vec{w}$ toward $\vec{w}^*$

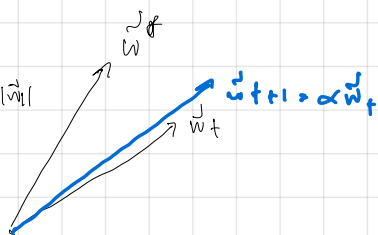- to measure how much the vector $\vec{w}$ chaning towards $\vec{w}^*$,

we consider two terms:

① $\vec{w}^T \vec{w}^*$



$\vec{p}_1, \vec{p}_1^T \vec{p}_2 = 9$

$\vec{p}_1 , \vec{p}_1^T \vec{p}_2 = มากสุด$
เวกเตอร์ใกล้

(closeness between $\vec{w}$ and $\vec{w}^*$, when $\vec{w} = \vec{w}^*$
$\Rightarrow \vec{w}^T \vec{w}^*$ is at maximum)

② $\vec{w}^T \vec{w}$
$\propto \sqrt{\vec{w}^T \vec{w}} = \|\vec{w}\|$



$\vec{w}_{t+1} \propto \vec{w}_t$

① เริ่ม ไฟθ= ② เริ่มน้อย ⇒ เปลี่ยน Direction ใกล้ $\vec{w}^*$

# analyze ① $\vec{w}^T \vec{w}^*$     $\vec{w} \leftarrow \vec{w} + y_i x_i$

$\vec{w}^T \vec{w}^* \rightarrow (\vec{w} + y\vec{x})^T \vec{w}^*$

$\vec{w}^T \vec{w}^* \Rightarrow \vec{w}^T \vec{w}^* + y\vec{x}^T \vec{w}^*$   หน่วยตรงโครงเข้าด้ม
ตก่ม่อ ขึ้นไม่เปลี่ยน
$\underbrace{}_{\geq 0}$ เท่า

$\vec{w}^T \vec{w}^* = \vec{w}^T \vec{w}^* + y\vec{x}^T \vec{w}^* \Rightarrow \vec{w}^T \vec{w}^* \geq \vec{w}^T \vec{w}^* + \gamma$
$\underbrace{}_{\geq \gamma}$

② $\vec{w}^T \vec{w} \rightarrow (\vec{w} + y\vec{x})^T (\vec{w} + y\vec{x})$

$\vec{w}^T \vec{w} \rightarrow \vec{w}^T \vec{w} + \underbrace{2y\vec{x}^T \vec{w}}_{\leq 0} + \underbrace{y^2 \vec{x}^T \vec{x}}_{1 \cdot \square \atop 1 \cdot \leq 1}$

$y \in \{-1, 1\}$

$\|x\| \leq 1$

$\|x\| = \sqrt{x^T x} \leq 1$
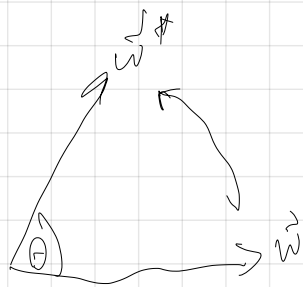
(w must misclassify on x)

$\Rightarrow \vec{w}^T \vec{w} \leq \vec{w}^T \vec{w} + 1$

for each update:

① $\vec{w}^T \vec{w}^* \geq \vec{w}^T \vec{w}^* + \gamma$
  $\vec{w}^T \vec{w} \leq \vec{w}^T \vec{w} + 1$



margin:

$\gamma = \min_{\forall i} |x_i^T \vec{w}^*|$

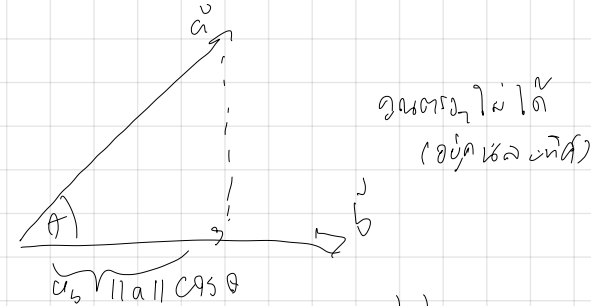$\gamma = ความ น้อยสุด vector
$\vec{w}^*$ กับ $x_i$ ตัว น้อยสุด

# Inner product / dot product

Inner: $\vec{a}, \vec{b} \in \mathbb{R}^d$

$$\vec{a} \circ \vec{b} \;/\; \vec{a}^T \vec{b} = \sum_{i=1}^{d} a_i b_i$$

$\underset{\text{scalar}}{}$

Geometrically the inner product $\vec{a}^T \vec{b}$ can be used to measure how close the vector $\vec{a}$ to vector $\vec{b}$, and vice versa



ฉายเงาไปไว้ (ฉากต้องเหมือนกัน)

$\|a\| \cos\theta$

$\theta$ is the angle between $\vec{a}, \vec{b}$

$\vec{a}^T \vec{b} = (\|\vec{a}\| \cos\theta) \vec{b}$

$\underset{\text{product}}{\rightsquigarrow} \;= (\|\vec{a}\|)(\|\vec{b}\|) \cos\theta$

$$\boxed{\vec{a}^T \vec{b} = \|\vec{a}\| \|\vec{b}\| \cos\theta}$$
Const cosnt

$\theta = 0$
$\theta = 180° (\pi)$
$-1 \le \cos\theta \le 1$
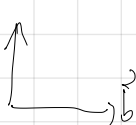
when $\vec{a}^T \vec{b}$ is at max
$\Rightarrow \cos\theta = 1 \Rightarrow \theta = 0$
$\vec{a}, \vec{b}$ are in the same direction.



- when $\vec{a}^T \vec{b}$ is at min
$\Rightarrow \cos\theta = 0 \Rightarrow \theta = \pi$
$\vec{a}, \vec{b}$ are in opposite direction



- when $\vec{a}^T \vec{b}$ is zero
$\Rightarrow \cos\theta = 0 \Rightarrow \theta = \frac{\pi}{2} \Rightarrow \vec{a}, \vec{b}$ are perpendicular / orthogonal
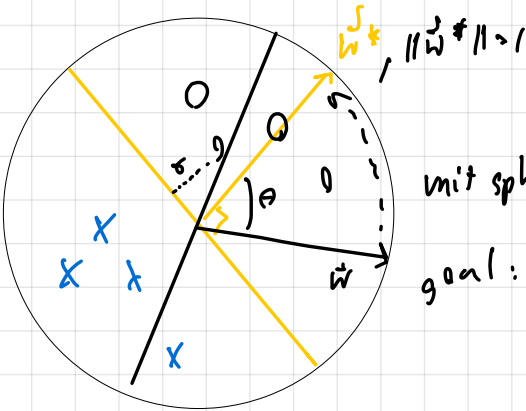


$*\ \theta = 0 \Rightarrow \cos\theta = 1 \Rightarrow \vec{a}^T \vec{b}$ is maximum A

- Assumption: $\exists \vec{w}^*$ by $\forall (x_i, y_i) \in D, \; y_i(\vec{w}^{*T} \vec{x}_i) > 0$
and $\|\vec{w}^*\| = 1$

- Transformation to Unit sphere:
$\forall (x_i, y_i) \in D, \; \|x_i\| \le 1$

- Define margin $\gamma = \min_{\forall i} |\vec{x}_i^T \vec{w}^*|$ ($\gamma$ is the distance from the close point to the hyperplane.)



$\vec{w}^*$ , $\|\vec{w}^*\| > 1$

unit sphere

goal: $\vec{w} \rightarrow \vec{w}^*$

Note:
$\vec{w}^T \vec{w}^* = \|\vec{w}\| \|\vec{w}^*\| \cos\theta$

$\cos\theta = \dfrac{\vec{w}^T \vec{w}^*}{\|\vec{w}\| \|\vec{w}^*\|}$

We consider 2 terms

① $\vec{w}^T \vec{w}^*$ (to be large)

② $\vec{w}^T \vec{w}$ (to be small)

$\cos\theta = \dfrac{\vec{w}^T \vec{w}^*}{\sqrt{\vec{w}^T \vec{w}}} \approx 1 \quad -1 \le \cos\theta \le 1$

for each update

① $\vec{w}^T \vec{w}* \rightarrow \geq \vec{w}^T \vec{w}* + \gamma$

(each update increase $\vec{w}^T \vec{w}*$ by at least $\gamma$)

② $\vec{w}^T \vec{w} \rightarrow \leq \vec{w}^T \vec{w} + 1$

(each update increase $\vec{w}^T \vec{w}$ by at most 1)

Suppose after M updates, we have $\cos\theta < 1$

Then, $- \vec{w}^T \vec{w}* \geq \gamma M$
$\quad\quad - \vec{w}^T \vec{w} \leq M$

$\theta = 0 \Rightarrow \cos\theta = 1 = \dfrac{\vec{w}^T \vec{w}*}{\sqrt{\vec{w}^T \vec{w}}} \geq \dfrac{\gamma M}{\sqrt{\vec{w}^T \vec{w}}} \geq \dfrac{\gamma M}{\sqrt{M}}$
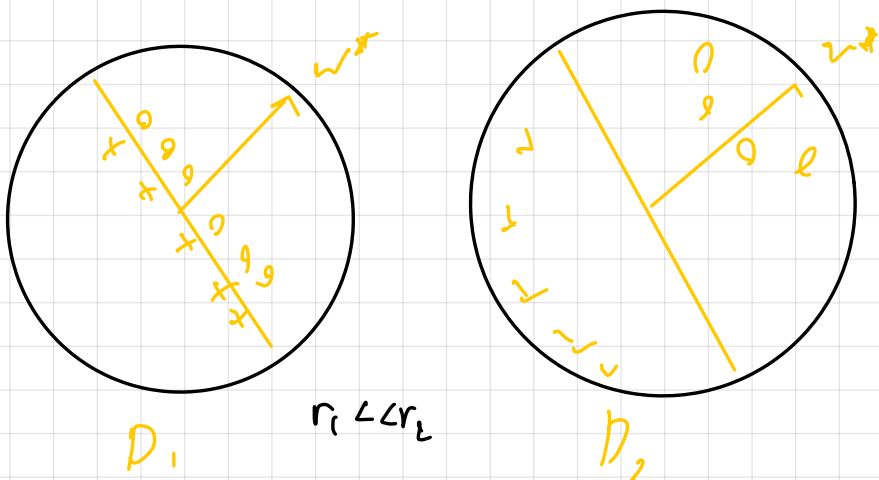
$\Rightarrow 1 \geq \dfrac{\gamma M}{\sqrt{M}}$

$\sqrt{M} \geq \gamma M$

$M \geq \gamma^2 M^2$

$1 \geq \gamma^2 M \Rightarrow M \leq \dfrac{1}{\gamma^2}$

Theorem : The perceptron will find $\vec{w}*$ within at most $\dfrac{1}{\gamma^2}$ updates.



$r_1 < < r_2$

$D_1$ $\quad\quad$ $D_2$

Q : ถ้ารัน perceptron บน $D_1, D_2$

กรณีไหนจะหา $\vec{w}*$ เจอใน $D_1 / D_2$ ก่อน

: ได้ออกมา $\gamma$ น้อย กว่า hyper plane
จะเข้าใกล้กันกว่า กาก $\frac{1}{\gamma^2}$

k-NN vs Perceptron

- for low-dimensional data, k-NN is very efficient.

for high-dimensional data, Perceptron is more suitable

- Limitations