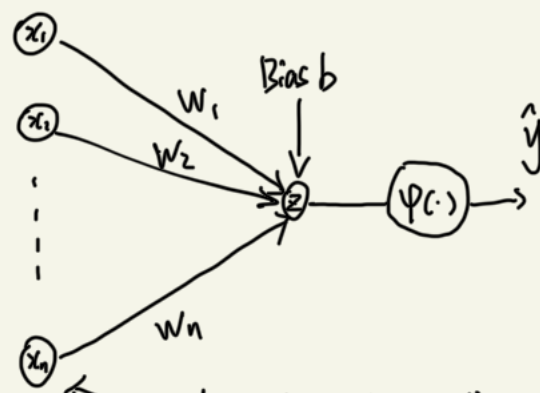


① Perceptron Model



Notation:

- I. w : weight, $w \in \mathbb{R}^m$
- II. b : bias, $b \in \mathbb{R}$
- III. x : input vector, $x \in \mathbb{R}^m$
- IV. y : label for each x , $y \in \mathbb{R}$
- V. \hat{y} : model predict label, $\hat{y} \in \mathbb{R}$
- VI. $\phi(\cdot)$: activation function, $\phi(\cdot): \mathbb{R}^m \rightarrow \mathbb{R}$

each input x has a "+1/-1" label named y .

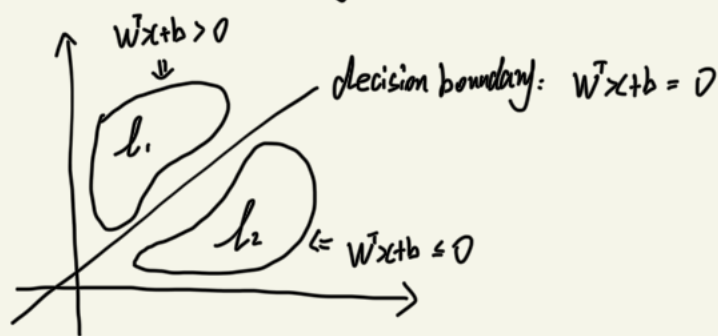
$$\begin{aligned} \text{i. } z &= w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b \\ &= \sum_{i=1}^n w_i x_i + b \end{aligned}$$

\Leftarrow Linear combination

$$= \langle w, x \rangle = w^T x$$

$$\text{ii. } \hat{y} = \phi(z), \quad \phi(\cdot): x \mapsto \text{sgn}(x) \quad \Leftarrow \text{linear to non-linear}$$

② Perceptron Training



i. weight update for perceptron

perceptron only update its weight when misclassify

- case 1: correctly classified, we have $w(n+1) = w(n)$
- case 2: incorrectly classified, i.e. $y(n+1)w(n)^T x(n+1) < 0$

we first initialize the model with $w(0)=0, b(0)=0$
and then we start input $x(1)$ and its label $y(1)$
we calculate and compared the result
if $y(1)w(0)^T x(1) < 0$
 $w(1) = w(0) + y(1)x(1)$
and then input the next, $x(2)$

e.g. consider input same data for $w(n+1)$

$$\begin{aligned} & y(n)w(n+1)^T x(n) \\ &= y(n)[w(n) + y(n)x(n)]^T x(n) \\ &= y(n)w(n)^T x(n) + y(n)^2 \|x(n)\|^2 \\ &= y(n)w(n)^T x(n) + \|x(n)\|^2 \end{aligned}$$

always positive,

which means the negative component $y(n)w(n)^T x(n)$ is turning to more positive.

P.S. where notation (n) means n^{th} input, w, b start from $w(0), b(0)$

P.S.2 sometimes we add learning rate η for update.
i.e. if $y(n)w(n)^T x(n) < 0$ {

$$\begin{aligned} w(n+1) &= w(n) + \eta y(n)x(n); \\ \} \end{aligned}$$

where $\eta \in (0, 1]$, which control the speed that model trains.

ii. perceptron convergence

In order to simplify the derivation, we add the bias term b and its coefficient "+1" to the w and x first place, respectively.

i.e. $y = \text{sgn}(w^T x)$
 $w = [b \ w_1 \ w_2 \ \dots \ w_n]^T$
 $x = [1 \ x_1 \ x_2 \ \dots \ x_n]^T$

- \Rightarrow
- ① Hypothesis 1: All training datas have bounded Euclidean norms
i.e. $\|x(t)\|_{L^2} \leq \epsilon$ for all $t \in \mathcal{X}$
 - ② Hypothesis 2: Linear classifier with weight W_* exists
i.e. $y(t) W_*^T x(t) \geq \gamma$ for all $t \in \mathcal{X}$
and the classifier just right before W_* is $W(n)$
 - ③ Aim: using the constraint for $\cos(\cdot) \leq 1$,
set up $\cos(W_*, W(n))$ to find the boundary

a) Part 1. $\frac{W_*^T W(n)}{\|W_*\| \|W(n)\|} \leq 1$

we can first decomposition $W(n)$ using the "update rule"

$$\begin{aligned} W(n) &= W(n-1) + d(n) y(n) x(n) \\ &= W(n-2) + d(n-1) y(n-1) x(n-1) + d(n) y(n) x(n) \\ &\vdots \\ W(0) = 0 &= W(n-1) + d(1) y(1) x(1) + d(2) y(2) x(2) + \dots + d(n) y(n) x(n) \\ &= \sum_{i=1}^n d(i) y(i) x(i) \end{aligned}$$

$d(n) = \begin{cases} 1, & \text{if misclassify at } n^{\text{th}} \text{ input} \\ 0, & \text{otherwise} \end{cases}$

thus we can rewrite the inner product for $\langle W_*, W(n) \rangle$

$$W_*^T W(n) = W_*^T \left(\sum_{i=1}^n d(i) y(i) x(i) \right)$$

$$= \sum_{i=1}^n d(i) y(i) W_*^T x(i)$$

recall the hypothesis 2, we have
 $\min_i \{ y(i) W_*^T x(i) \} = \gamma$

not every input
cause update,
let's just assume
 k 's update exist

$$\begin{aligned} &\geq \sum_{i=1}^n d(i) \gamma \\ &\geq K \gamma \quad (\Leftrightarrow) \quad W_*^T W(n) \geq K \gamma \end{aligned}$$

At the same time we know
 $W_*^T W(n) \geq K \gamma$

thus we have

$$[W_*^T W(n)]^2 \geq K^2 \gamma^2$$

using the Cauchy-Schwarz Inequality.

$$\|W_*\|^2 \|W(n)\|^2 \geq [W_*^T W(n)]^2 \geq K^2 \gamma^2$$

\downarrow

$$\|W(n)\|^2 \geq \frac{K^2 \gamma^2}{\|W_*\|^2}, \text{ lower bound}$$

b) Part 2.

Calculate $\|W(n)\|^2$

$$\begin{aligned} &= \|W(n-1) + d(n) y(n) x(n)\|^2 \\ &= \|W(n-1)\|^2 + \|d(n) y(n) x(n)\|^2 + 2 d(n) y(n) W(n-1)^T x(n) \\ &\leq \|W(n-1)\|^2 + \|d(n) x(n)\|^2 \\ &\leq \|W(n-2)\|^2 + \|d(n-1) x(n-1)\|^2 + \|d(n) x(n)\|^2 \\ &\leq \sum_{i=1}^n \|d(i-1) x(i-1)\|^2 \end{aligned}$$

$\left\{ \begin{array}{l} \text{if no update at } n^{\text{th}} \text{ input:} \\ \quad d(n) = 0 \\ \text{if update at } n^{\text{th}} \text{ input:} \\ \quad y(n) W(n-1)^T x(n) < 0 \end{array} \right.$

using the same
assume as part 1:
 k^{th} update

$$\leq K \epsilon^2$$

recall the hypothesis 1:
 $\min_i \|x(i)\| = \epsilon$

c) Conclusion

a. conclusion 1: upper / lower bound for $\|W_{cn}\|$

$$\frac{k^2 \gamma^2}{\|W_{*}\|^2} \leq \|W_{cn}\|^2 \leq k \epsilon^2$$

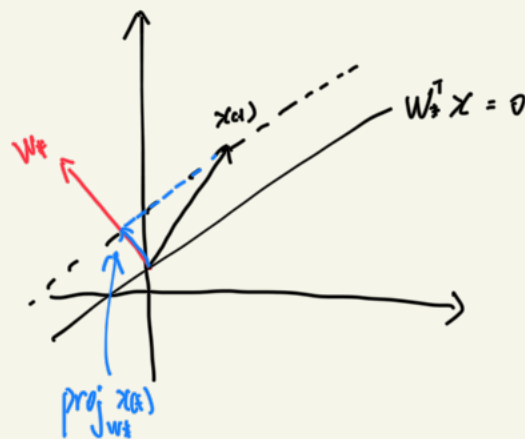
b) conclusion 2: update time k

$$\cos(W_{*}, W_{cn}) = \frac{W_{*}^T W_{cn}}{\|W_{*}\| \|W_{cn}\|} \stackrel{\text{part 1}}{\geq} \frac{k \gamma}{\|W_{*}\| \|W_{cn}\|} \geq \frac{k \gamma}{\|W_{*}\| \sqrt{k \epsilon^2}} \stackrel{\text{part 2}}{\geq}$$

$$\text{that is } 1 \geq \frac{k \gamma}{\|W_{*}\| \sqrt{k \epsilon^2}} \Leftrightarrow k \leq \frac{\epsilon^2 \|W_{*}\|^2}{\gamma^2}$$

d) Advanced derivation

Set x_{ct} be the data sample most close to decision boundary (i.e. $y_{ct} W_{*}^T x_{ct} = \gamma$)



$$\left\| \text{proj}_{W_*} x_{ct} \right\| = \left\| \frac{x_{ct}^T W_*}{W_*^T W_*} \cdot W_* \right\|$$

$$\begin{aligned} y_{ct} W_*^T x_{ct} &= \gamma \\ |W_*^T x_{ct}| &= \gamma = \frac{|W_*^T x_{ct}|}{\|W_*\|^2} \cdot \|W_*\| \end{aligned}$$

$$d_{\min} = \frac{\gamma}{\|W_*\|} \quad \Leftarrow \text{the min distance for the data set to the decision boundary}$$

rewrite the $k \leq \frac{\epsilon^2 \|W_{*}\|^2}{\gamma^2}$ to

$$k \leq \frac{\epsilon^2}{d_{\min}^2} \quad \leftarrow \text{the more small the geometric margin that separates the training set, the more difficult for the training.}$$

③ Batch Perceptron : changing the size we input

i. Premise : We input a Set of sample x named $\mathcal{H} = \{x_{c1}, x_{c2}, \dots, x_{cn}\}$

Instead of input one by one like we used to do.

And for the dataset \mathcal{H} , a batch of misclassified samples \mathcal{Q}

is used to compute the adjustment. and $\mathcal{Q} \in \mathcal{H}$.

ii. Error : For those misclassified samples, we compute the total error

$$J(w) = \sum_{x_{ct} \in \mathcal{Q}} (1 - y_{ct} w^T x_{ct})$$

and how w changes effect to the total error

$$\nabla J(w) = \frac{\partial J(w)}{\partial w} = \sum_{x_{ct} \in \mathcal{Q}} (-y_{ct} x_{ct})$$

Using projection and norm to calculate distance is invented by my own, inspired by the lecture notes from MIT.

iii. Update:

(Method of Gradient Descent)

$$W' := W - \eta \nabla J(W)$$

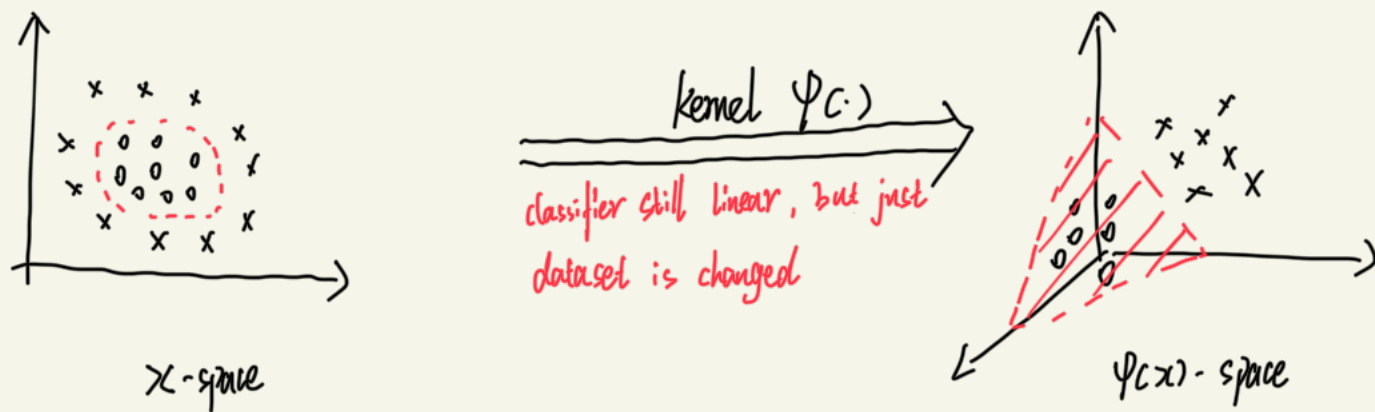
Recall the one by one input case.

$$= W + \eta \left(\sum_{x(n) \in \mathcal{X}} y(n) x(n) \right)$$

$$W' = W + \eta y(n) x(n)$$

Accumulate all updates one by one into a single update

④ Kernel Method



kernel: Mapping an original feature vector x to an expanded version $\psi(x)$
 \uparrow
 information in x is expanded

i. Feel hard to explicitly compute $\psi(\cdot)$?

Solution: Noticed that the main calculate in learning model is "inner product".

Use $k(x, z)$ to replace $\psi(x)^T \psi(z)$
 define $k(\dots)$ to skip define $\psi(\cdot)$

example: define $\psi(x) = [x_1^2, \sqrt{2}x_1x_2, x_2^2]$

calculate $\psi(x)^T \psi(z)$, $x = [x_1, x_2]^T$, $z = [z_1, z_2]^T$

$$\psi(x)^T \psi(z) = [x_1^2, \sqrt{2}x_1x_2, x_2^2] \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix}$$

$$= x_1^2 z_1^2 + 2x_1x_2 z_1z_2 + x_2^2 z_2^2$$

instead we can define a function have same result

$$\begin{aligned} \underline{k(x, z)} &= (x^T z)^2 = x_1^2 z_1^2 + 2x_1x_2 z_1z_2 + x_2^2 z_2^2 \\ &= \psi(x)^T \psi(z) \end{aligned}$$

ii. Use kernel method in perceptron

Rewrite the weight in "dec-type".

step1: Initialize weight $w = 0$, $b = 0$

$$w(n-1) = w(n-2) + \lambda(n-1) y(n-1) \psi(x(n-1))$$

Step 2: At the n^{th} input, we have weights represented by

$$W_{(n-1)} = \sum_{i=1}^{n-1} d(i) y(i) \phi(x(i))$$

$$= \sum_{i=1}^{n-1} d(i) y(i) \phi(x(i))$$

Step 3: Calculate output before activated, using kernel

$$\begin{aligned} Z(n) &= W_{(n-1)}^T \phi(x(n)) + b \\ &= \left[\sum_{i=1}^{n-1} d(i) y(i) \phi(x(i)) \right]^T \phi(x(n)) + b \\ &= \sum_{i=1}^{n-1} d(i) y(i) \underbrace{[\phi(x(i))^T \phi(x(n))]}_{\text{inner product}} + b \\ &= \sum_{i=1}^{n-1} d(i) y(i) k(x(i), x(n)) + b \end{aligned}$$

Step 4: Mark the update. Execute update in next iteration instead of changing the value of $W(n)$

We simply mark $d(n) = 1$, if the n^{th} input is misclassified

Then the update will take effect during the next calculation output i.e. in $n+1^{th}$ input

$$Z = \left[\sum_{i=1}^n d(i) y(i) \phi(x(i)) \right]^T x(n+1)$$

$$= \left[\sum_{i=1}^{n-1} d(i) y(i) \phi(x(i)) + \underbrace{d(n) y(n) \phi(x(n))}_{\text{update}} \right]^T x(n+1)$$

update is not execute \Rightarrow
when identifying a misclassification (i.e. if \dots)
but execute when calculate the output
before the activating function.

$$W(n) = W(n-1) + d(n) y(n) \phi(x(n))$$

$$= \sum_{i=1}^{n-1} d(i) y(i) k(x(i), x(n+1)) + d(n) y(n) k(x(n), x(n+1))$$

$$= \sum_{i=1}^n d(i) y(i) k(x(i), x(n+1))$$