

Perceptron

By Frank Rosenblatt, Cornell, 1957.

"Linear classifier", one of the first ANN.

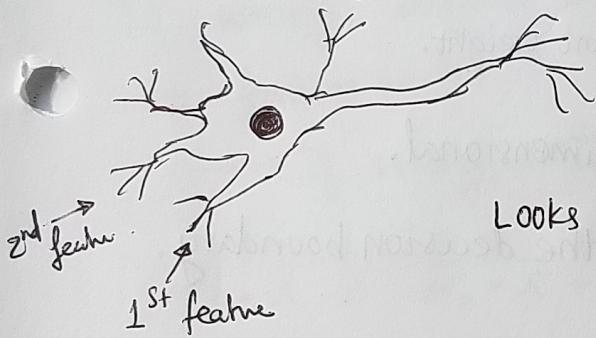
Linear classification :-

Feature vector $(\bar{x}_1, \dots, \bar{x}_k)$, for each \bar{x} (feature).

Ex/Trng :- $(\bar{x}(1), f(\bar{x}(1)))$, ...

↓
a k-dimensional vector.

Inspired by biology (very simplified).



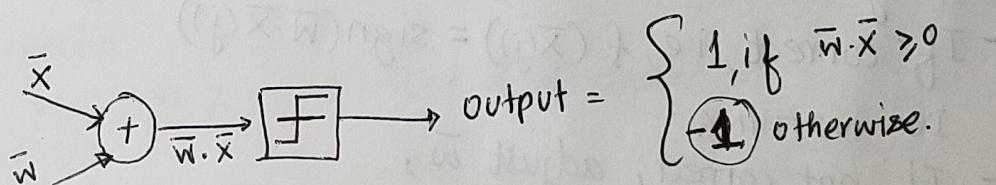
Get a bunch of inputs - if higher than a threshold, activate.

Looks like a thresholding circuit.

- ① Input:- feature values $(\bar{x}_1, \dots, \bar{x}_k)$.
- ② Each feature has a weight. $w_j \leftrightarrow$ feature j 's weight.

- ③ Sum The weighted summation is an activation function.

$$\text{activation}_w(\bar{x}) = \sum_{j=1}^k \bar{x}_j \cdot w_j = \bar{w} \cdot \bar{x}$$



Intercept/bias, add a dummy feature $\bar{x}_0 = 1$ for every example,
 w_0 is the weight of this feature

What does the perceptron classifier do?

1. Using the "training example" learn a

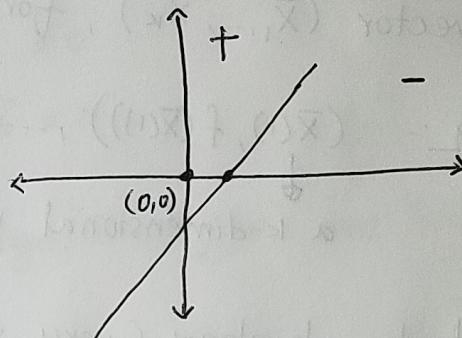
set of weight vectors (w_0, w_1, \dots, w_k)

2. For the new test \bar{x} , output $\text{sign}(\bar{w} \cdot \bar{x})$.

In 2D :- ($k=2$),

Is w_0 positive or negative?

Special feature, additional,
for the intercepts.



Large $w_j \rightarrow$ strong activation/important weight.

These are hyperplanes. very high dimensional.

⇒ Decision rule $\rightarrow \boxed{\bar{w} \cdot \bar{x} = 0}$ is the decision boundary.

How to learn the weights?

ALGORITHM :-

1. Start with weights = 0.

2. For iter = 1, ..., I

2(a) For each training example,

- classify with the current weight

- If correct, i.e., $f(\bar{x}(j)) = \text{sign}(\bar{w} \cdot \bar{x}(j))$

- DO NOTHING

- If not correct, adjust \bar{w} ,

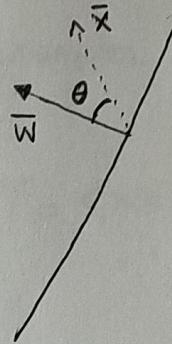
- change the \bar{w} toward the example.

$\hat{\theta}$ = angle between \bar{w} & \bar{x}

determines the sign of $\bar{w} \cdot \bar{x}$.

$$\theta > 90^\circ \rightarrow -1$$

$$\theta \leq 90^\circ \rightarrow +1.$$



$$\text{output} = \begin{cases} +1 & \text{if } \bar{w} \cdot \bar{x} \geq 0 \\ -1 & \text{if } \bar{w} \cdot \bar{x} < 0. \end{cases}$$

If correct, i.e., $f(\bar{x}) = \text{sign}(\bar{w} \cdot \bar{x})$, i.e, +ve here, do nothing.

If wrong, adjust the weight vector by subtracting/adding the feature vector.

$$\rightarrow \text{If wrong, } \bar{w}_n = \bar{w}_o + f(\bar{x}) \cdot \bar{x}.$$

$$\bar{w}_n \cdot \bar{x} = \bar{w}_o \cdot \bar{x} + f(\bar{x}) \cdot \|\bar{x}\|_2^2, \text{ slowly make change towards } f(\bar{x}).$$

>Show the increase in the angle for the example above.

Multiclass Decision Rules.

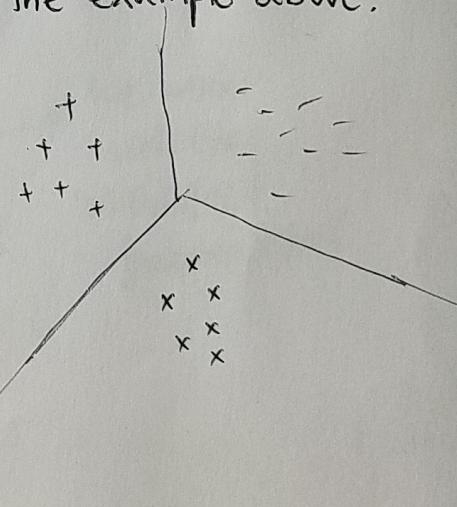
For each class we have a ~~selected~~ weight vector.

(w_+, w_-, w_x) , $C = \text{set of all classes}$.

Score of a class is $w \cdot \bar{x}$

for feature vector \bar{x} ,

$$c = \underset{c \in C}{\operatorname{argmax}} w_c \cdot \bar{x} \leftarrow \text{select class with highest score.}$$



Some notations, recap.

(45)

\bar{x} , $\bar{x}(1), \bar{x}(2), \dots, \bar{x}(n)$, are the features of examples.

$\bar{x}(i)$ is example 'i', therefore a $K-d$ vector.

$\bar{x}_j(i) \rightarrow$ jth feature of example 'i'; $T = \langle x(1), f(x(1)) \rangle$,

$\bar{x} \rightarrow$ new example/ one new feature vector.

$\bar{x}_j \rightarrow$ jth feature of \bar{x} .

$f(\bar{x}) \rightarrow$ label of \bar{x} , eg. $f(\bar{x}(i))$ - label of example 'i'.
 ↳ a scalar.

- Linear classifiers :- A number of classifiers with linear decision boundaries.

for example \bar{x}

\therefore decision boundary :- $w_1 \bar{x}_1 + \dots + w_K \bar{x}_K \geq t$

$\bar{x} = (\bar{x}_1, \dots, \bar{x}_K)$

scalar threshold.

$w \cdot \bar{x} \geq t$ ← decision rule.

can be written as $\bar{x}^* \cdot \bar{w}^* \geq 0$,

by making

$$\begin{cases} \bar{x}^* = (-1, \bar{x}_1, \dots, \bar{x}_K) \\ \bar{w}^* = (t, w_1, \dots, w_K) \end{cases}$$

We work with these modified features now!

'how to choose good linear boundaries'?

Perceptron :-

(46)

Start with $\bar{w}^* = \mathbf{0}$ For iter = 1, ..., \blacksquare ITER,For $i = 1, \dots, n$,if $g_0(\bar{x}^*(i) \cdot \bar{w}^*) \neq f(\bar{x}(i))$ (namely, $f(\bar{x}(i)) \times (\bar{w}^* \cdot \bar{x}(i)) < 0$)

$$\bar{w}^* \leftarrow \bar{w}^* + f(\bar{x}(i)) \cdot \bar{x}^*(i)$$

Output $\bar{w}^* = (w_0, w_1, \dots, w_k)$ Decision rule = $w_0 + w_1 \bar{x}_1 + \dots + w_k \bar{x}_k \geq 0$ Example 'i', $\underset{\text{+ve}}{\text{sgn}}(\bar{w}^* \cdot \bar{x}(i)) \neq \underset{\text{-ve}}{\text{f}}(\bar{x}(i))$,

$$(\bar{w}^* + f(\bar{x}(i)) \cdot \bar{x}^*(i)) \cdot \bar{x}^*(i) = \underset{\text{+ve}}{\bar{w}^* \cdot \bar{x}^*(i)} + \underset{\text{-ve}}{f(\bar{x}(i)) \cdot |\bar{x}^*(i)|^2}$$

'makes it small.'

Recall :-

For a plane $\bar{w}^* \cdot \mathbf{0} = 0$, and any point \bar{x}^* , $d(x^*, \bar{w}^*) \left| \frac{\bar{w}^* \cdot \bar{x}^*}{|\bar{w}^*|} \right|$ is the distance of \bar{x}^* from the plane \perp to \bar{w}^* .

Margin :-

Separability :- \bar{w}^* separates $\bar{x}(1), \dots, \bar{x}(n)$ if it gets all the training examples correct.

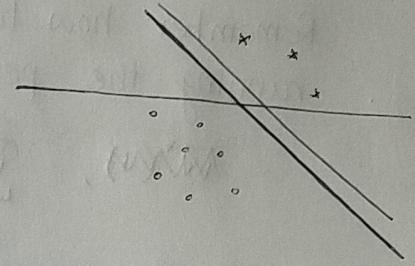
$$\text{Margin}(\bar{w}^*, T) = \begin{cases} \min_{i=1, \dots, n} d(\bar{x}^*(i), \bar{w}^*) & \text{if } \bar{w}^* \text{ separates } T \\ -\infty & \text{o/w} \end{cases}$$

Margin of the data, = maximum margin of all separating \bar{w}^* . (47)

' γ ' → denotes the margins.

$$\gamma(T) = \max_{\bar{w}^*} \text{margin}(T, \bar{w}^*)$$

* Larger margin, 'higher wiggle room'



Theorem ⇒ Suppose we run the perceptron algorithm, on a dataset 'T' with margin γ , and $|\bar{x}(i)|_1 \leq 1$. Then after at most $\frac{4}{\gamma^2}$ updates to \bar{w}^* perceptron converges

\bar{w}^* = all zeros initially,

$$\bar{w}_{\text{opt}}^* \rightarrow |\bar{w}_{\text{opt}}^* \cdot \bar{x}(i)| \geq \gamma$$

$$\bar{w}^*(j+1) \cdot \bar{w}_{\text{opt}}^* = \bar{w}^*(j) \cdot \bar{w}_{\text{opt}} + \underbrace{f(\bar{x}(i)) \times (\bar{w}_{\text{opt}}^* \cdot \bar{x}(i))}_{> \gamma}$$

$$\geq (j+1) \cdot \gamma$$

$$|\bar{w}^*(j+1)|_2^2 \leq 2(j+1)$$

$$\text{By induction } |\bar{w}^*(j+1)|_2^2 = |\bar{w}^*(j)|_2^2 + |\bar{x}(i)|_2^2 + 2 \underbrace{f(\bar{x}(i)) \cdot (\dots)}_{< 0}$$

$$= |\bar{w}^*(j)|_2^2 + \sqrt{2}$$

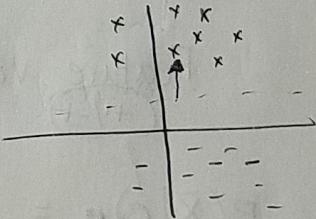
$$(j+1) \cdot \gamma < |\bar{w}^*(j+1)|_2 < \sqrt{2} \cdot \sqrt{j+1} \Rightarrow \boxed{j+1 \leq \frac{1}{\gamma^2}}$$



Problems

(1) One nasty example.

Flip everything around



(2) Weak generalization.

→ SVM later.

(3) Does not "learn from correct samples".

4. Multiclass even separable cannot be learnt.

• Voted Perceptron.

(48)

Remember how long each hyperplane survives, as we keep running the perceptron algorithm.

$$\text{VW}^*(\bar{w}), \quad \begin{bmatrix} \bar{w}^{*(1)} \\ \vdots \\ \bar{w}^{*(C)} \end{bmatrix}, \quad t_1, \quad t_C$$

$$\text{Sign} \left(\sum_{j=1}^C t_j \cdot \text{Sign}(\bar{w}^{*(j)} \cdot \bar{x}) \right).$$

Works well, huge storage!, slow prediction.

• Averaged perceptron: $\text{sign} \left(\sum_{j=1}^C t_j \cdot (\bar{w}^{(j)} \cdot \bar{x}) \right)$

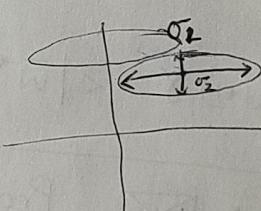
<Can maintain running averages> \rightarrow Faster.

~~Gaussian NB~~ Gaussian NB as a linear model.

Two classes, $\{+1, -1\}$, / $\{c_1, \dots, c_L\}$.

$\{\hat{\mu}_{j|+1}, \hat{\mu}_{j|-1}\}$ Each feature conditioned on class is Gaussian.

$\{\hat{\sigma}_{j|+1}^2, \hat{\sigma}_{j|-1}^2\}$ Suppose, $\underbrace{\hat{\sigma}_{j|+1}^2 = \hat{\sigma}_{j|-1}^2 = \hat{\sigma}_j^2}$



Feature variances are same given class.

$P(\bar{x}|c) = \prod_{j=1}^K P(\bar{x}_j|c) = \prod_{j=1}^K \frac{1}{\sqrt{2\pi\hat{\sigma}_j^2}} \exp \left(-\frac{(\bar{x}_j - \hat{\mu}_{j|c})^2}{2\hat{\sigma}_j^2} \right).$

Decision Rule :-

(49)

$$P(C=+1) \cdot \prod_{j=1}^K \frac{1}{\sqrt{2\pi\sigma_j^2}} \cdot \exp\left(-\frac{(\bar{x}_j - \hat{\mu}_{j|+1})^2}{2\sigma_j^2}\right) \geq P(C=-1) \dots$$

$$\Leftrightarrow \log(P(+1)) - \sum_{j=1}^K \frac{(\cancel{x_j^2} - 2\bar{x}_j \cdot \hat{\mu}_{j|+1} + \hat{\mu}_{j|+1}^2)}{2\sigma_j^2} \geq \log P(-1) - \sum_{j=1}^K \left(\cancel{x_j^2} \right)$$

$$\Rightarrow \alpha \times \sum_{j=1}^K \bar{x}_j \cdot \frac{(\hat{\mu}_{j|+1} - \hat{\mu}_{j|-1})}{\sigma_j^2} + \boxed{\log \frac{P(+1)}{P(-1)} - \sum_{j=1}^K \frac{(\hat{\mu}_{j|+1}^2 - \hat{\mu}_{j|-1}^2)}{\sigma_j^2}} \geq 0.$$

A linear classifier.