

## EE-559 – Deep learning

### 3.1. The perceptron

François Fleuret

<https://fleuret.org/ee559/>

Mar 4, 2020

The first mathematical model for a neuron was the Threshold Logic Unit, with Boolean inputs and outputs:

$$f(x) = \mathbf{1}_{\{w \sum_i x_i + b \geq 0\}}.$$

The first mathematical model for a neuron was the Threshold Logic Unit, with Boolean inputs and outputs:

$$f(x) = \mathbf{1}_{\{w \sum_i x_i + b \geq 0\}}.$$

It can in particular implement

$$\text{or}(u, v) = \mathbf{1}_{\{u+v-0.5 \geq 0\}} \quad (w = 1, b = -0.5)$$

$$\text{and}(u, v) = \mathbf{1}_{\{u+v-1.5 \geq 0\}} \quad (w = 1, b = -1.5)$$

$$\text{not}(u) = \mathbf{1}_{\{-u+0.5 \geq 0\}} \quad (w = -1, b = 0.5)$$

The first mathematical model for a neuron was the Threshold Logic Unit, with Boolean inputs and outputs:

$$f(x) = \mathbf{1}_{\{w \sum_i x_i + b \geq 0\}}.$$

It can in particular implement

$$\begin{aligned} \text{or}(u, v) &= \mathbf{1}_{\{u+v-0.5 \geq 0\}} & (w = 1, b = -0.5) \\ \text{and}(u, v) &= \mathbf{1}_{\{u+v-1.5 \geq 0\}} & (w = 1, b = -1.5) \\ \text{not}(u) &= \mathbf{1}_{\{-u+0.5 \geq 0\}} & (w = -1, b = 0.5) \end{aligned}$$

Hence, **any Boolean function can be build with such units.**

(McCulloch and Pitts, 1943)

The perceptron is very similar

$$f(x) = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

but the inputs are real valued and weights can be different (Rosenblatt, 1957).

The perceptron is very similar

$$f(x) = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

but the inputs are real valued and weights can be different (Rosenblatt, 1957).

It was originally motivated by biology, with  $w_i$  being the *synaptic weights*, and  $x_i$  and  $f$  firing rates. However, it is a (very) crude biological model.

To make things simpler we take responses  $\pm 1$ . Let

$$\sigma(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{otherwise.} \end{cases}$$



The perceptron classification rule boils down to

$$f(x) = \sigma(w \cdot x + b).$$

To make things simpler we take responses  $\pm 1$ . Let

$$\sigma(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{otherwise.} \end{cases}$$



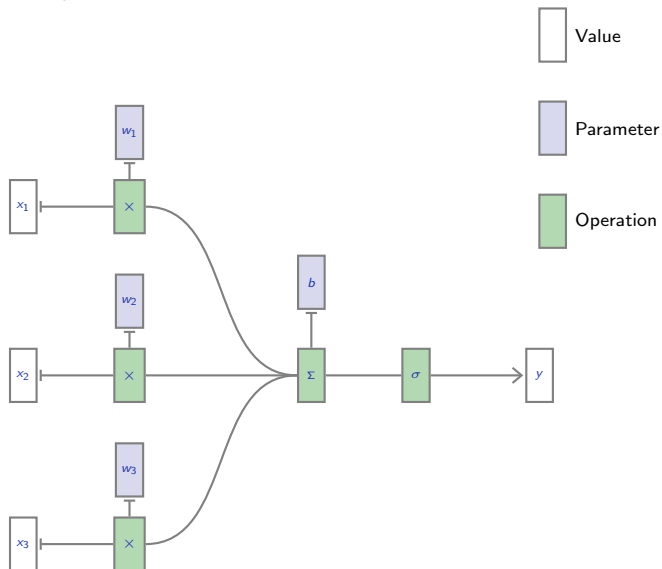
The perceptron classification rule boils down to

$$f(x) = \sigma(w \cdot x + b).$$

For neural networks, the function  $\sigma$  that follows a linear operator is called the **activation function**.

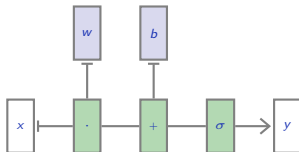


We can represent this “neuron” as follows:



We can also use tensor operations, as in

$$f(x) = \sigma(w \cdot x + b).$$



Given a training set

$$(x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}, \quad n = 1, \dots, N,$$

a very simple scheme to train such a linear operator for classification is the **perceptron algorithm**:

Given a training set

$$(x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}, \quad n = 1, \dots, N,$$

a very simple scheme to train such a linear operator for classification is the **perceptron algorithm**:

1. Start with  $w^0 = 0$ ,
2. while  $\exists n_k$  s.t.  $y_{n_k} (w^k \cdot x_{n_k}) \leq 0$ , update  $w^{k+1} = w^k + y_{n_k} x_{n_k}$ .

Given a training set

$$(x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}, \quad n = 1, \dots, N,$$

a very simple scheme to train such a linear operator for classification is the **perceptron algorithm**:

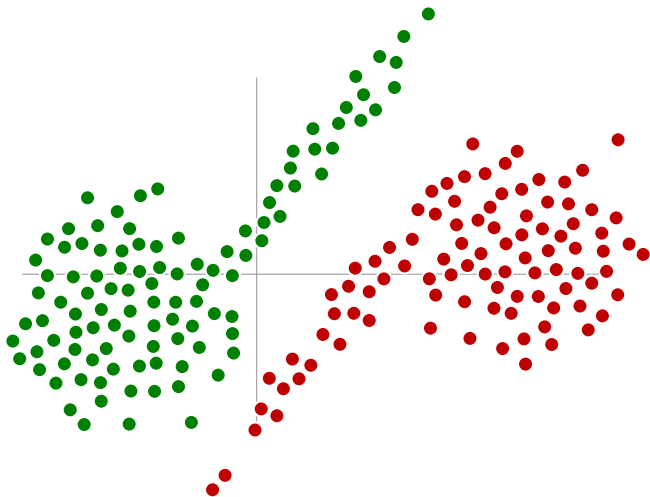
1. Start with  $w^0 = 0$ ,
2. while  $\exists n_k$  s.t.  $y_{n_k} (w^k \cdot x_{n_k}) \leq 0$ , update  $w^{k+1} = w^k + y_{n_k} x_{n_k}$ .

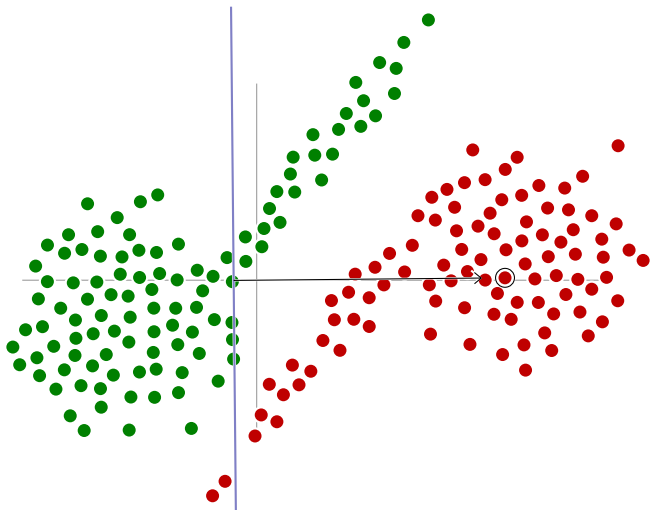
The bias  $b$  can be introduced as one of the  $w$ s by adding a constant component to  $x$  equal to  $1$ .

```
def train_perceptron(x, y, nb_epochs_max):
    w = torch.zeros(x.size(1))

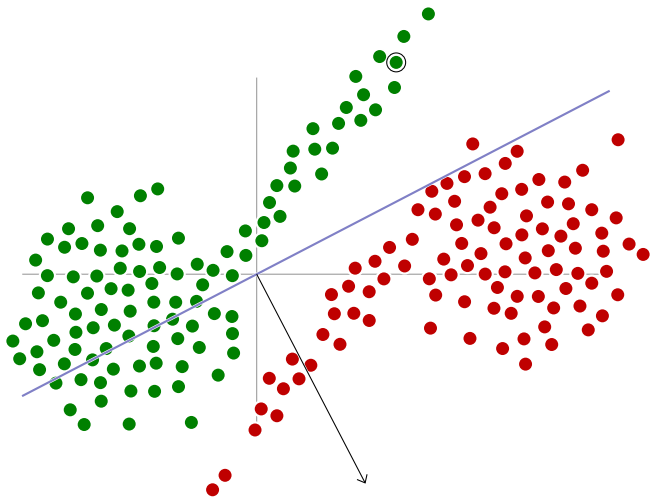
    for e in range(nb_epochs_max):
        nb_changes = 0
        for i in range(x.size(0)):
            if x[i].dot(w) * y[i] <= 0:
                w = w + y[i] * x[i]
                nb_changes = nb_changes + 1
        if nb_changes == 0: break;

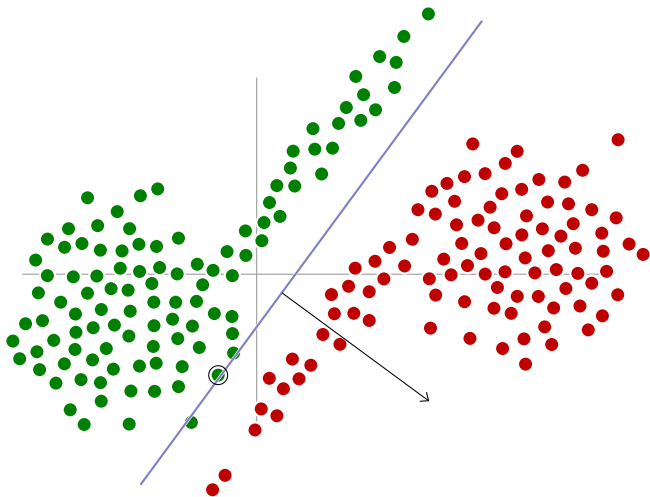
    return w
```

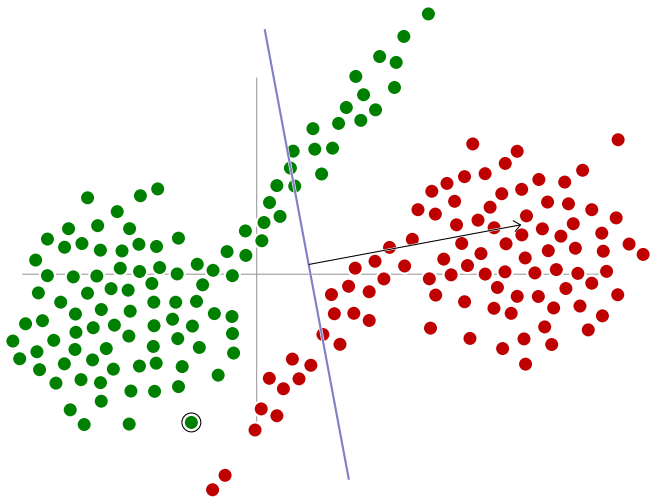


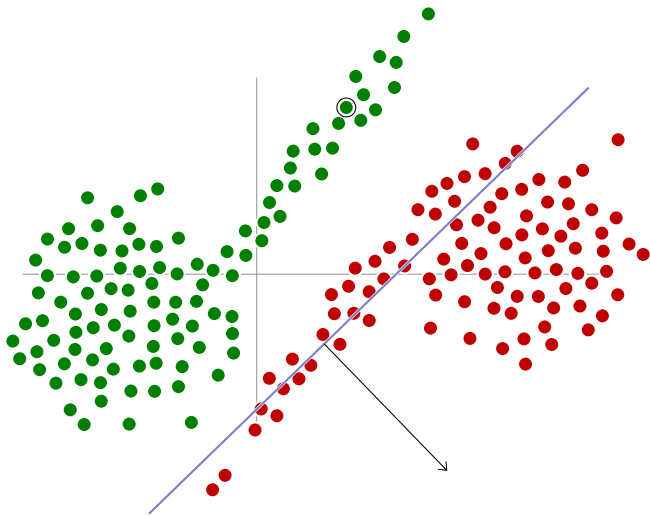


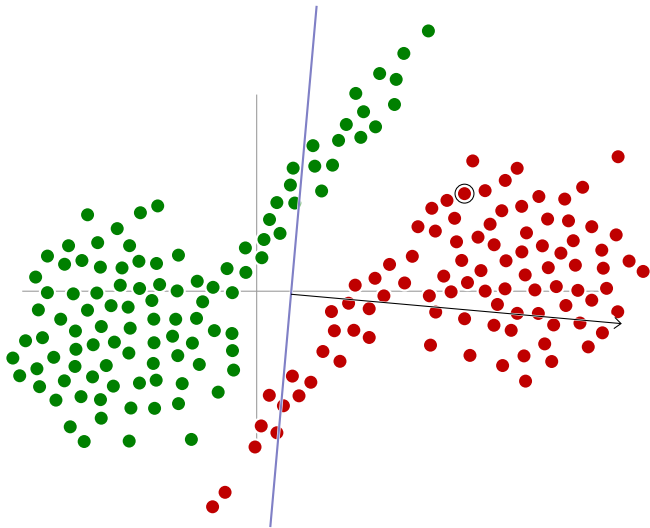


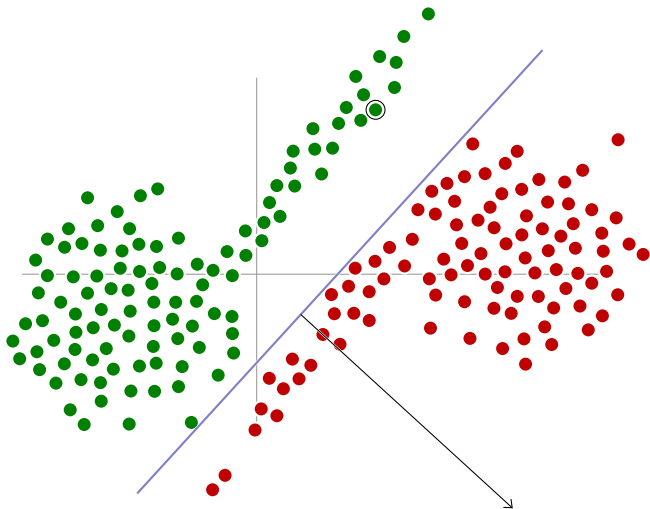




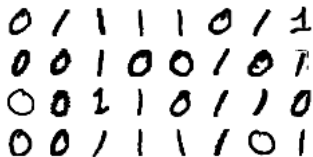




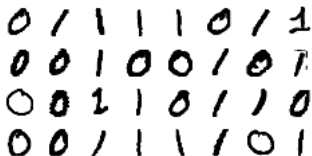




This crude algorithm works often surprisingly well. With MNIST's "0"s as negative class, and "1"s as positive one.



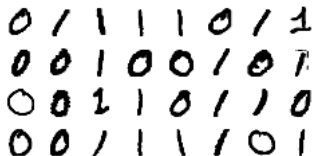
This crude algorithm works often surprisingly well. With MNIST's "0"s as negative class, and "1"s as positive one.



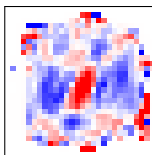
```
epoch 0 nb_changes 64 train_error 0.23% test_error 0.19%
epoch 1 nb_changes 24 train_error 0.07% test_error 0.00%
epoch 2 nb_changes 10 train_error 0.06% test_error 0.05%
epoch 3 nb_changes 6 train_error 0.03% test_error 0.14%
epoch 4 nb_changes 5 train_error 0.03% test_error 0.09%
epoch 5 nb_changes 4 train_error 0.02% test_error 0.14%
epoch 6 nb_changes 3 train_error 0.01% test_error 0.14%
epoch 7 nb_changes 2 train_error 0.00% test_error 0.14%
epoch 8 nb_changes 0 train_error 0.00% test_error 0.14%
```



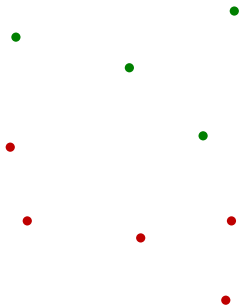
This crude algorithm works often surprisingly well. With MNIST's "0"s as negative class, and "1"s as positive one.



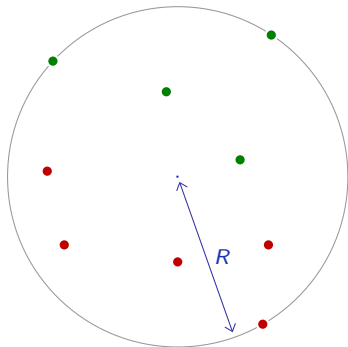
```
epoch 0 nb_changes 64 train_error 0.23% test_error 0.19%
epoch 1 nb_changes 24 train_error 0.07% test_error 0.00%
epoch 2 nb_changes 10 train_error 0.06% test_error 0.05%
epoch 3 nb_changes 6 train_error 0.03% test_error 0.14%
epoch 4 nb_changes 5 train_error 0.03% test_error 0.09%
epoch 5 nb_changes 4 train_error 0.02% test_error 0.14%
epoch 6 nb_changes 3 train_error 0.01% test_error 0.14%
epoch 7 nb_changes 2 train_error 0.00% test_error 0.14%
epoch 8 nb_changes 0 train_error 0.00% test_error 0.14%
```



We can get a convergence result under two assumptions:



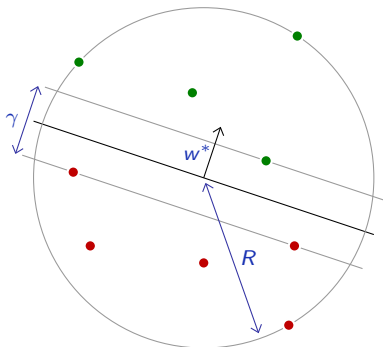
We can get a convergence result under two assumptions:



1. The  $x_n$  are in a sphere of radius  $R$ :

$$\exists R > 0, \forall n, \|x_n\| \leq R.$$

We can get a convergence result under two assumptions:



1. The  $x_n$  are in a sphere of radius  $R$ :

$$\exists R > 0, \forall n, \|x_n\| \leq R.$$

2. The two populations can be separated with a margin  $\gamma$ :

$$\exists w^*, \|w^*\| = 1, \exists \gamma > 0, \forall n, y_n (x_n \cdot w^*) \geq \gamma/2.$$

To prove the convergence, let us make the assumption that there still is a misclassified sample at iteration  $k$ . We have

$$\begin{aligned}w^{k+1} \cdot w^* &= (w^k + y_{n_k} x_{n_k}) \cdot w^* \\&= w^k \cdot w^* + y_{n_k} (x_{n_k} \cdot w^*) \\&\geq w^k \cdot w^* + \gamma/2 \\&\geq (k+1)\gamma/2.\end{aligned}$$

To prove the convergence, let us make the assumption that there still is a misclassified sample at iteration  $k$ . We have

$$\begin{aligned}w^{k+1} \cdot w^* &= (w^k + y_{n_k} x_{n_k}) \cdot w^* \\&= w^k \cdot w^* + y_{n_k} (x_{n_k} \cdot w^*) \\&\geq w^k \cdot w^* + \gamma/2 \\&\geq (k+1)\gamma/2.\end{aligned}$$

Since

$$\|w^k\| \|w^*\| \geq w^k \cdot w^*,$$

we get

$$\begin{aligned}\|w^k\|^2 &\geq (w^k \cdot w^*)^2 / \|w^*\|^2 \\&\geq k^2 \gamma^2 / 4.\end{aligned}$$

And

$$\begin{aligned}\|w^{k+1}\|^2 &= w^{k+1} \cdot w^{k+1} \\ &= (w^k + y_{n_k} x_{n_k}) \cdot (w^k + y_{n_k} x_{n_k}) \\ &= w^k \cdot w^k + 2 \underbrace{y_{n_k} w^k \cdot x_{n_k}}_{\leq 0} + \underbrace{\|x_{n_k}\|^2}_{\leq R^2} \\ &\leq \|w^k\|^2 + R^2 \\ &\leq (k+1) R^2.\end{aligned}$$

Putting these two results together, we get

$$k^2 \gamma^2 / 4 \leq \|w^k\|^2 \leq k R^2$$

hence

$$k \leq 4R^2 / \gamma^2,$$

hence no misclassified sample can remain after  $\lfloor 4R^2 / \gamma^2 \rfloor$  iterations.



Putting these two results together, we get

$$k^2 \gamma^2 / 4 \leq \|w^k\|^2 \leq k R^2$$

hence

$$k \leq 4R^2 / \gamma^2,$$

hence no misclassified sample can remain after  $\lfloor 4R^2 / \gamma^2 \rfloor$  iterations.

This result makes sense:

- The bound does not change if the population is scaled, and
- the larger the margin, the more quickly the algorithm classifies all the samples correctly.

The perceptron stops as soon as it finds a separating boundary.

Other algorithms maximize the distance of samples to the decision boundary, which improves robustness to noise.

The perceptron stops as soon as it finds a separating boundary.

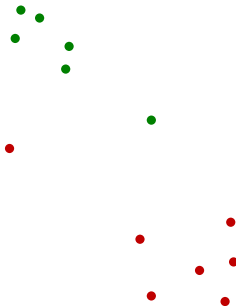
Other algorithms maximize the distance of samples to the decision boundary, which improves robustness to noise.

Support Vector Machines (SVM) achieve this by minimizing

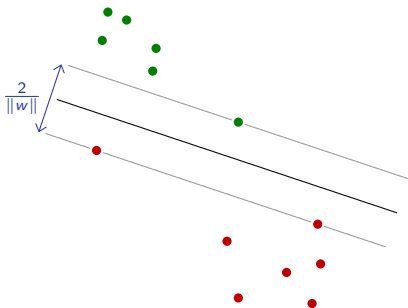
$$\mathcal{L}(w, b) = \lambda \|w\|^2 + \frac{1}{N} \sum_n \max(0, 1 - y_n(w \cdot x_n + b)),$$

which is convex and has a global optimum.

$$\mathcal{L}(w, b) = \lambda \|w\|^2 + \frac{1}{N} \sum_n \max(0, 1 - y_n(w \cdot x_n + b))$$

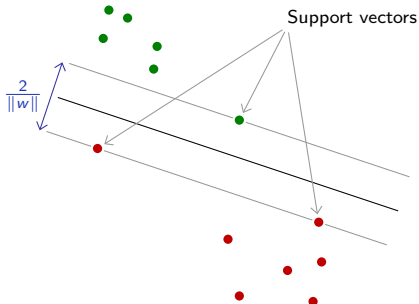


$$\mathcal{L}(w, b) = \lambda \|w\|^2 + \frac{1}{N} \sum_n \max(0, 1 - y_n(w \cdot x_n + b))$$



Minimizing  $\max(0, 1 - y_n(w \cdot x_n + b))$  pushes the  $n$ th sample beyond the plane  $w \cdot x + b = y_n$ , and minimizing  $\|w\|^2$  increases the distance between the  $w \cdot x + b = \pm 1$ .

$$\mathcal{L}(w, b) = \lambda \|w\|^2 + \frac{1}{N} \sum_n \max(0, 1 - y_n(w \cdot x_n + b))$$



Minimizing  $\max(0, 1 - y_n(w \cdot x_n + b))$  pushes the  $n$ th sample beyond the plane  $w \cdot x + b = y_n$ , and minimizing  $\|w\|^2$  increases the distance between the  $w \cdot x + b = \pm 1$ .

At convergence, only a small number of samples matter, the “support vectors”.

The term

$$\max(0, 1 - \alpha)$$

is the so called “hinge loss”



The end



## References

- W. S. McCulloch and W. Pitts. **A logical calculus of the ideas immanent in nervous activity**. The bulletin of mathematical biophysics, 5(4):115–133, 1943.
- F. Rosenblatt. **The perceptron—A perceiving and recognizing automaton**. Technical Report 85-460-1, Cornell Aeronautical Laboratory, 1957.