# Basic Perceptron

Sunday, March 6, 2022     8:58 AM

Rosenblatt's  Perceptron → 1958



Model: $v = \sum_{i=1}^{m} (w_i x_i) + b$

It is basically a pondered sum of the inputs.
The parameters to modify are only the weights (w)
and the bias (b).

→ The perceptron is a simple neural network algorithm used for classification
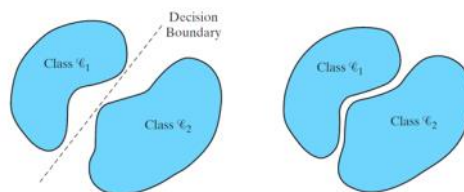problems. What it does is to generate a hyperplane:



Hyperplane: it is a plane
generalized to 1D, 2D, 3D and
4D...

the function
adjusts to a
hyperplane that
separates $C_1$ and
$C_2$ by updating
w (b is constant).

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

→ Training the perceptron
↳ there are 3 convergence criteria:
→ Convergence conditions
• The overall set of training vectors (input) is
well defined (not NaN data, etc).
• The sets are separable: we can decide on this
simply by training and if the output is constantly
mistaken or accumulated a big error, it is probably
because of non-separable data.

if it does
not con-
verge
=
NOT SE-
PARABLE

Data must be
separable
by a hyper
plane.

○ Training algorithm:
→ $\bar{w}_{(n+1)} = \bar{w}_{(n)} + \eta_{(n)} \bar{x}_{(n)}$ if $\bar{w}_{(n)} \bar{x}_{(n)} \leq 0$ and $\bar{x}_{(n)}$ is $C_1$
→ $\bar{w}_{(n+1)} = \bar{w}_{(n)} - \eta_{(n)} \bar{x}_{(n)}$ if $\bar{w}_{(n)} \bar{x}_{(n)} > 0$ and $\bar{x}_{(n)}$ is $C_2$



There are two different approaches for NN:

Algorithmic
approach

Bayesian
Approach

Brute Force
Gradient Descent

Mathematic
Input statistics

→ The Algorithmic Approach (Haykin algorithm)

*Variables and Parameters:*

$\mathbf{x}(n) = (m+1)$-by-1 input vector
$= [+1, x_1(n), x_2(n), ..., x_m(n)]^T$
$\mathbf{w}(n) = (m+1)$-by-1 weight vector
$= [b, w_1(n), w_2(n), ..., w_m(n)]^T$
$b$ = bias
$y(n)$ = actual response (quantized) → output
$d(n)$ = desired response → Supervised Learning
$\eta$ = learning-rate parameter, a positive constant less than unity

1. *Initialization.* Set $\mathbf{w}(0) = \mathbf{0}$. Then perform the following computations for time-step $n = 1, 2, ....$
2. *Activation.* At time-step $n$, activate the perceptron by applying continuous-valued input vector $\mathbf{x}(n)$ and desired response $d(n)$.
3. *Computation of Actual Response.* Compute the actual response of the perceptron as
$$y(n) = \text{sgn}[\mathbf{w}^T(n)\mathbf{x}(n)]$$
where sgn($\cdot$) is the signum function.
4. *Adaptation of Weight Vector.* Update the weight vector of the perceptron to obtain
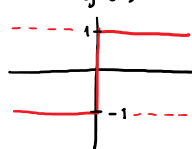$$\mathbf{w}(n+1) = \mathbf{w}(n) + \eta[\underbrace{d(n) - y(n)}_{error}]\mathbf{x}(n) \quad \rbrace \; n \text{ is the iteration}$$

where

$$d(n) = \begin{cases} +1 & \text{if } \mathbf{x}(n) \text{ belongs to class } \mathscr{C}_1 \\ -1 & \text{if } \mathbf{x}(n) \text{ belongs to class } \mathscr{C}_2 \end{cases}$$

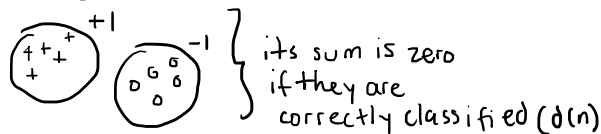5. *Continuation.* Increment time step $n$ by one and go back to step 2.

Signum Function sgn($\cdot$):



For any input, the signum function limits it to be either 1 (positive input) or $-1$ (negative input).

↓

$-1$ and $1$ end up being the classification of the data.

Loss Function: to minimize the loss since the ideal case is



its sum is zero if they are correctly classified ($d(n)$)

$$w(n+1) = w(n) - \eta(n) \nabla J(w) \quad \text{where } n \text{ is the iteration}$$
↳ so that the hyperplane decreases along the axis
$$= w(n) + \eta(n) \sum x(n) d(n)$$
↳ makes it oscilate to approach the convergence point

→ Bayes classifier → no a priori information (Algorithmic approach)

○ We have a cost function (Loss function)
$$J(\bar{w}) = \sum_{\bar{x}(n) \in \mathcal{H}} (-\bar{w}^T \bar{x}(n) d(n))$$

○ We aim for minimizing the cost function, by looking for the opposite direction of the gradient:

$J$ depends on all $\bar{w}$ values ← $\nabla J(\bar{w}) = \sum_{\bar{x}(n) \in \mathcal{H}} (-\bar{x}(n) d(n))$   $\nabla = \left[\dfrac{\partial}{\partial w_1}, \dfrac{\partial}{\partial w_2}, ..., \dfrac{\partial}{\partial w_m}\right]^T$

where $n$ is the iteration.    $\underbrace{\qquad}_{\text{since we got } m \text{ inputs}}$

this jump depends on $\eta$ (learning constant)

○ Thus, the algorithm for training the network as:

$$\bar{w}_{(n+1)} = \overbrace{\bar{w}_{(n)}}^{vector} \ominus \overbrace{\eta_{(n)}}^{constant} \overbrace{\nabla J(\bar{w})}^{vector}$$

The hyperplane oscilates towards the final solution ⎤ Until convergence near $J=0$

$$= \bar{w}_{(n)} \oplus \eta_{(n)} \sum_{\bar{x}(n) \in \mathcal{H}} (\bar{x}_{(n)} d_{(n)})$$
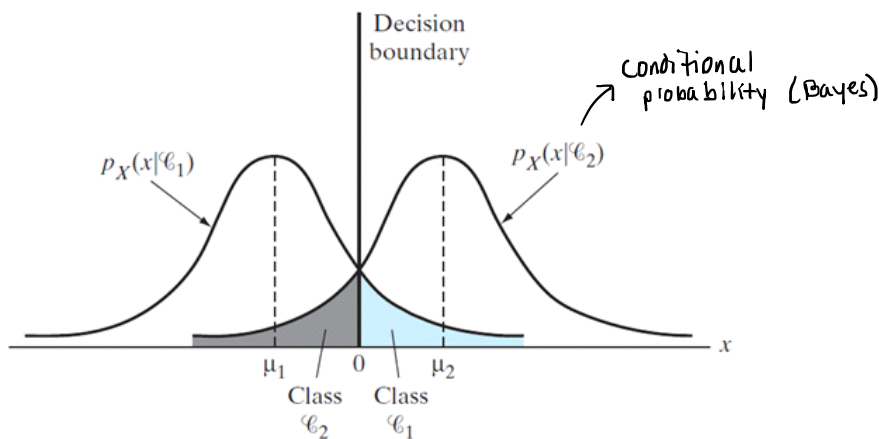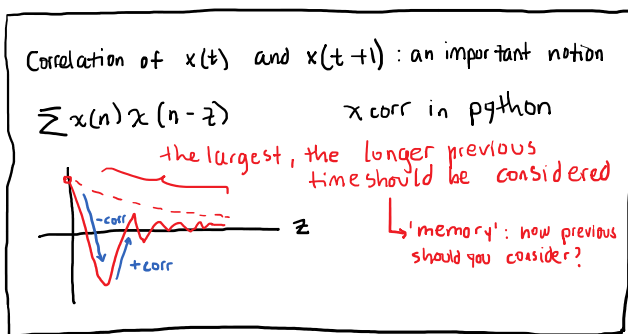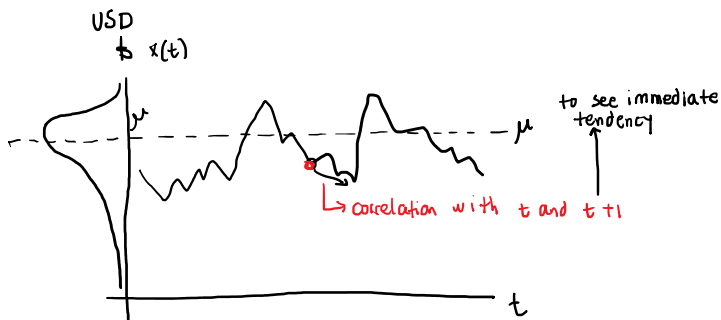
 cost function

$$= \bar{w}_{(n)} \oplus \eta_{(n)} \sum_{\bar{x}(n) \in \mathcal{H}} (\bar{x}_{(n)} d_{(n)})$$

towards the final solution

near $J = 0$



cost function

→ **Gaussian Classifier (Bayes approach)**

What happens since we cannot know the input? Statistics help us to describe the input to the algorithm.



USD $\$ \, x(t)$

$\mu$

to see immediate tendency

→ correlation with $t$ and $t+1$

$t$

---

Correlation of $x(t)$ and $x(t+1)$: an important notion

$$\sum x(n) x(n-z) \qquad x \text{ corr in python}$$

the largest, the longer previous time should be considered

-corr

+corr

$z$

'memory': how previous should you consider?

---

Decision boundary

conditional probability (Bayes)

$p_X(x|\mathscr{C}_1)$

$p_X(x|\mathscr{C}_2)$

$\mu_1$   $0$   $\mu_2$

Class $\mathscr{C}_2$   Class $\mathscr{C}_1$

$x$

The Bayesian approach thus informs the input data statistics to the algorithm also, so that the classifier improves.

→ **Bayes Classifier — minimum average risk**

prob of $\bar{x}$ given $C_1$          prob of $\bar{x}$ given $C_2$

$$\mathbb{R} = C_{11} P_1 \int_{\mathcal{H}_1} P_{\bar{x}}(\bar{x}|C_1) d\bar{x} + C_{22} P_2 \int_{\mathcal{H}_2} P_{\bar{x}}(\bar{x}|C_2) d\bar{x}$$

correct ones

→ Regions (4)

$$+ C_{21} P_1 \int_{\mathcal{H}_2} P_{\bar{x}}(\bar{x}|C_1) d\bar{x} + C_{12} P_2 \int_{\mathcal{H}_1} P_{\bar{x}}(\bar{x}|C_2) dx$$

false negative/positives

$\bar{x}$ being $C_1$ but in region 2

$$\Lambda = \frac{P_{\bar{x}}(\bar{x}|C_1)}{P_{\bar{x}}(\bar{x}|C_2)}$$

If this is bigger than this → $\Lambda \geq 1$ → Classified as $C_1$, etc.

→ "The probability that an $\bar{x}$ value belongs to Class $C_2$"

• Probability Density Function for the Gaussian Case

↪ "The probability that an X value belongs to Class $C_2$"

- Probability Density Function for the Gaussian Case

$$P_{\bar{x}}(\bar{x}|C_i) = \frac{1}{(2\pi)^{m/2}(\det\mathbb{C})^{1/2}} e^{\wedge\left(-\frac{1}{2}(\bar{x}-\mu_i)^T \mathbb{C}^{-1}(\bar{x}-\mu_i)\right)} , i=1,2$$

  * Note $\mathbb{C}$ is the correlation matrix

- Applying the logarithm of the function of verisimilitude

$$\log \Lambda(\bar{x}) = -\frac{1}{2}(\bar{x}-\mu_1)^T \mathbb{C}^{-1}(\bar{x}-\mu_1) + \frac{1}{2}(\bar{x}-\mu_2)\mathbb{C}^{-1}(\bar{x}-\mu_2)$$

$$= (\mu_1-\mu_2)^T \mathbb{C}^{-1}\bar{x} + \frac{1}{2}(\mu_2^T\mathbb{C}^{-1}\mu_2 - \mu_1^T\mathbb{C}^{-1}\mu_1)$$

- Leaving the neuron with the following weights (linear problem)

$$y = \bar{w}^T\bar{x} + b$$
$$\bar{y} = \log \Lambda(\bar{x})$$

$$\boxed{\begin{array}{l} \bar{w} = \mathbb{C}^{-1}(\mu_1-\mu_2) \\ b = \frac{1}{2}(\mu_2^T\mathbb{C}^{-1}\mu_2 - \mu_1^T\mathbb{C}^{-1}\mu_1) \end{array}}$$
No iterations
$\mathbb{C}^{-1}$: correlation matrix of all input variables

  - Do you have the relationship of all input variables?
    ↪ YES: Bayesian approach (compute the corr matrix)
    ↪ NO: Algorithmic approach