

(*) Perceptron model

↓
first model with weights

↓
goal in ML in general: to find a function which satisfactorily maps input to output. These functions have parameters that need to learnt in order for the accurate mapping, for example weights and bias.

↓
MP neuron had simple function with only one parameter. (b)

↓
a function with a higher degree parameters have more freedom in the trajectory and are thus better able to generalize.

↓
MP neuron fails even on basic data when the the points are at $(0,0)$ and $(1,1)$ for same class and at $(0,1)$ and $(1,0)$ for the other class.

↓
Issues with MP neuron:

(i) linear model

(ii) fixed slope

(iii) few possible intercept values

data { (iv) boolean ips
(v) boolean ops

(vi) learning algorithm: brute force

(okay here, but is not scalable to real world problems)

non linear mappings are very hard to approximate

in perceptron \rightarrow real valued inputs

boolean outputs (binary classification)

weights for every i/p \rightarrow slopes can be varied.

linear model / function

$$\text{loss} = \sum_i \max(0, 1 - y_i * \hat{y}_i)$$

specific
to perceptron
model

sum over all
samples.

0 when correct prediction
1 when wrong prediction

(i) Data and Task



real valued inputs, $x \in \mathbb{R}^n$.



the feature's can take values from different ranges & scales, which might give some features (with bigger values), an inherent advantage of representation during aggregation, that can be avoided, by data preparation techniques like standardization, and normalization.

feature scaling

data scaling

$$z = \frac{x - \mu}{\sigma}$$

$\mu \rightarrow$ mean
 $\sigma \rightarrow$ standard deviation

Task: Binary classification task

(ii) Normalization

i) max (min-max scaling)
ii) min

$$x' = \frac{x - \min}{\max - \min}$$



all values b/w 0 & 1



max \rightarrow 1

min \rightarrow 0

difference b/w standardization & normalization.



Normalization \rightarrow min-max scaling

standardization \rightarrow z-score normalization



features rescaled to have mean = 0 and $\sigma = 1$.

(ii) Model \rightarrow weights associated to each i/p.

things to specify while describing a model.

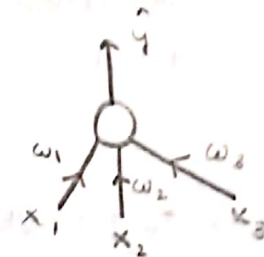
kind of i/p \rightarrow real values

kind of o/p \rightarrow boolean

parameters \rightarrow weights assigned to each feature

function:
$$\hat{y} = 1 \text{ if } \sum_{i=1}^n w_i x_i \geq b$$
$$\hat{y} = 0, \text{ otherwise}$$

this model is still linear.



weight is the importance or the weightage attached to different features to the aggregation which is used to achieve the desired o/p.

Desirable features can have +ve weights, varying in magnitude based on the extent of desirability, and similarly for undesirability.

the aggregation can be seen as the dot product of the feature vector & weight vector and then activation fn. can be applied.

$$\hat{y} = 1 \text{ (if } x \cdot w \geq \text{threshold})$$

$$\hat{y} = 0 \text{ (otherwise)}$$

since x 's can take real values

(-) geometric interpretation: since b can take

real values, \therefore the intercept is much freer to move, and similarly, as the slopes of the features are not fixed, the line can take more varied forms, allowing it to take more points

in its stride.



$$w_1 x_1 + w_2 x_2 - b = 0$$

$$x_2 = -\frac{w_1}{w_2} x_1 + \frac{b}{w_2}$$

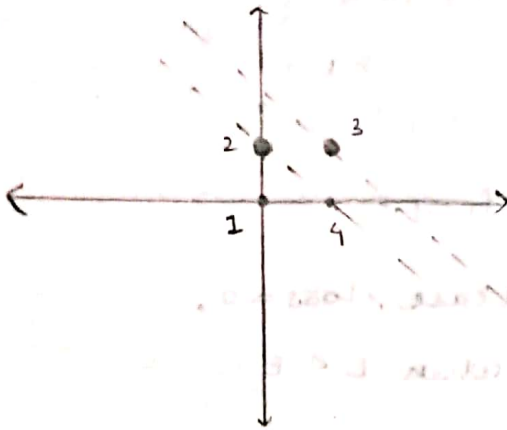
slope is
variable and
depends on w_1 & w_2

$$x_1 + x_2 - b = 0$$

$$x_2 = -x_1 + b$$

-1: slope

only
intercept
could
vary



If 2 and 3 are of same class and
1, and 4 from the other class.

There is no way for MCP to classify
this accurately, whereas perceptron
owing to its freedom on the slope
side can.



But perceptron still can't completely operate on the
XOR example, due to its linear nature.



perceptron model can only deal with linearly
separable data.

(iii) Loss function

$$L = 1(y \neq \hat{y}) \rightarrow \text{Indicator function}$$



whenever the condition is
satisfied, the fn. takes value = 1.
else takes value = 0.



squared error loss is same for MCP
(binary 0/ps)

(iv) Learning algorithm

typical recipe: a) initialize w_1, w_2, b

b) iterate over data, get predictions associated with current parameters of the model, and accordingly calculate loss associated with the predictions.

c) update (w, b) based on the loss

(update depends on the learning algorithm)

d) repeat b, c, till satisfied

↓
in ideal case, loss = 0,

but practically, when $L < \epsilon$

↓
some small quantity.

↓
or the loss seems to have hit a minima (cannot be decreased further)

•) The equations of the model can be generalized over the threshold, such that

$$w_n x_n + w_{n-1} x_{n-1} + \dots + w_1 x_1 + \underbrace{w_0 x_0}_{-b} \geq 0$$

$$\text{i.e. } w \cdot x \geq 0$$

$$\text{or } \sum_{i=0}^n w_i x_i \geq 0$$

where x_0 is fixed and equal to 1.

perceptron learning algorithm :

$P \leftarrow$ inputs with label 1 ;

$N \leftarrow$ inputs with label 0 ;

initialize ~~w~~ randomly ;

while ! convergence do :

when all P are classified correctly .

Pick random $x \in P \cup N$;

if $x \in P$ and $\sum w_i x_i < 0$, then

$$w = w + x$$

end

if $x \in N$ and $\sum w_i x_i > 0$, then

$$w = w - x$$

end

end

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} - \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

•) why it works?

$$w = [w_1, w_2, \dots, w_n]$$

$$x = [x_1, x_2, \dots, x_n]$$

$$\cos \theta = \frac{w \cdot x}{\|w\| \cdot \|x\|} \rightarrow \text{sign of } \cos \theta \text{ depends on the dot product}$$

↓

for $x \in P$ if $w \cdot x < 0$ then it means that the angle (θ) between this x and the current w is greater than 90° , and we need it to be less than 90°

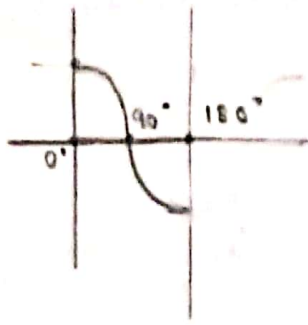
↓

$$w_{\text{new}} = w + x$$

↓

$$\text{if } \text{angle}(w, x) > \text{angle}(w_{\text{new}}, x)$$

the angle b/w the weights & x should dec. if it can't be made acute straightaway.



$$\cos(\alpha_{\text{new}}) \propto w_{\text{new}}^T \cdot x$$

$$\propto (w + \Delta w)^T \cdot x$$

$$\propto (w^T \cdot x) + \Delta w^T \cdot x$$

the always

$$\propto \cos \alpha + x^T \cdot x$$

$$\text{if } \cos(a) < \cos(b)$$

$$a > b$$

(monotonically decreasing)

$$\therefore \alpha_{\text{new}} < \alpha$$

$$\text{since } \cos(\alpha_{\text{new}}) > \cos(\alpha)$$

$$\cos \alpha + x^T \cdot x > \cos(\alpha) \quad \checkmark$$

↓

movement to the required condition.

↓

similar argument can be made for the negative or the 0 class.

* The perceptron algorithm will only converge if the data is linearly separable

* Defn. of perceptron for linearly separable: Two sets P & N of points in an N dimensional space are called absolutely linearly separable, if $n+1$ real numbers w_0, w_1, \dots, w_n exist such that every point $(x_1, x_2, \dots, x_n) \in P$ satisfies $\sum_{i=1}^n w_i \cdot x_i \geq w_0$ and every point $(x_1, x_2, \dots, x_n) \in N$ satisfies $\sum_{i=1}^n w_i \cdot x_i < w_0$.

(v) evaluation : checking the performance of the model .

↓

accuracy of the model on the test data .

↓

$$\text{accuracy} = \frac{\text{no. of correct predictions}}{\text{total no. of predictions}}$$