

1. $(2, 3)$

x_1	x_2	class
1	1	+1
-1	-1	-1
0.5	0.5	-1
0.1	0.5	-1
0.2	0.2	+1
0.9	0.5	+1

$$W^T x = 0 \quad \text{and} \quad W = [1, 1]$$

Step 1 iter 1

(i) $y = 1 \times 1 + 1 \times 1 = 2$, class = +1 \rightarrow correctly classified.

(ii) $y = 1 \times (-1) + 1 \times (-1) = -2$, class = -1 \rightarrow correctly classified.

(iii) $y = 1 \times 0.5 + 1 \times 0.5 = 1$, class = -1 \rightarrow wrongly classified.

hence $W = W - x = [1 - 0.5, 1 - 0.5]$ (as wrongly classified class is -1)

$$W = [1, 0.5]$$

(iv) $y = 1 \times 0.1 + 0.5 \times 0.5 = 0.35$, class = -1 \rightarrow wrongly classified.

hence $W = W - x = [1 - 0.1, 0.5 - 0.5]$ (as wrongly classified class is -1)

$$W = [0.9, 0.0]$$

(v) $y = 0.9 \times 0.2 + 0 \times 0.2 = 0.18$, class = +1 \rightarrow correctly classified.

(vi) $y = 0.9 \times 0.9 + 0 \times 0.5 = 0.81$, class = +1 \rightarrow correctly classified.

iter - 2

(i) $y = 0.9 \times 1 + 0 \times 1 = 0.9$, class = +1 \rightarrow correctly classified.

(ii) $y = 0.9 \times (-1) + 0 \times (-1) = -0.9$, class = -1 \rightarrow correctly classified.

(iii) $y = 0.9 \times 0.5 + 0 \times 0.5 = 0.45$, class = -1 \rightarrow wrongly classified.

hence $W = W - x = [0.9 - 0.5, 0 - 0.5]$ (as wrongly classified class is -1)

$$W = [0.4, -0.5]$$

(iv) $y = 0.9 \times 0.1 + (-0.5) \times 0.5 = -0.16$, class = -1 \rightarrow correctly classified.

(v) $y = 0.9 \times 0.2 + (-0.5) \times 0.2 = 0.08$, class = +1 \rightarrow correctly classified.

(vi) $y = 0.9 \times 0.9 + (-0.5) \times 0.5 = 0.56$, class = +1 \rightarrow correctly classified.

iter - 3

(i) $y = 0.9 \times 1 + (-0.5) \times 1 = 0.4$, class = +1, correctly classified.

(ii) $y = 0.9 \times (-1) + (-0.5) \times (-1) = 0.4$, class = -1, correctly classified.

(iii) $y = 0.9 \times 0.5 + (-0.5) \times 0.5 = -0.25$, class = -1, correctly classified.

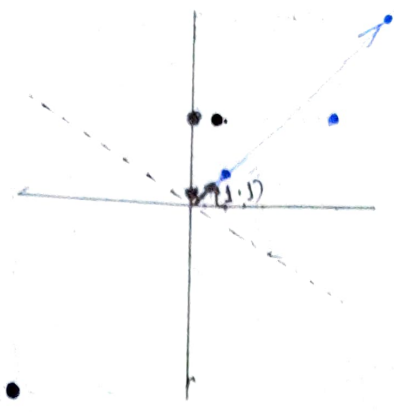
(iv) $y = 0.9 \times 0.1 + (-0.5) \times 0.5 = -0.16$, class = -1, correctly classified.

(v) $y = 0.9 \times 0.2 + (-0.5) \times 0.2 = 0.08$, class = +1, wrongly classified.

(vi) $y = 0.9 \times 0.9 + (-0.5) \times 0.5 = 0.56$, class = +1, correctly classified.

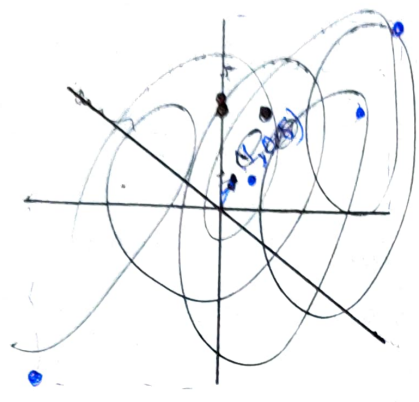
Therefore we can conclude in 3 steps the perceptron algorithm will converge.

1.2

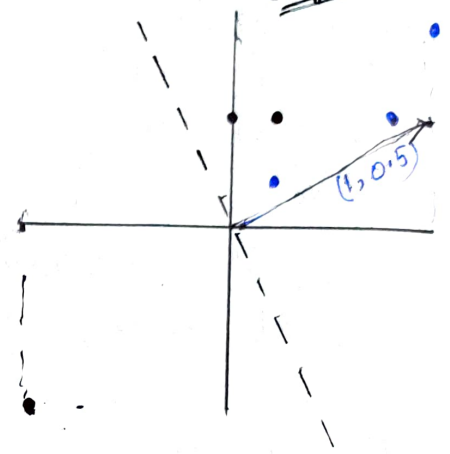


initial.

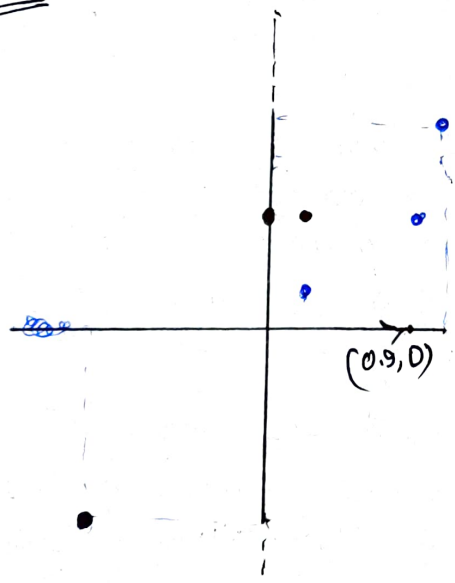
Step 1



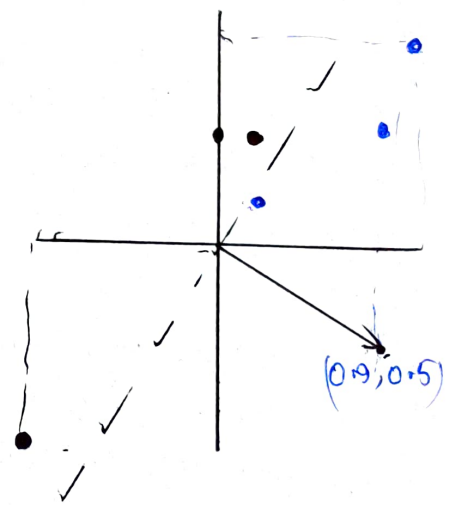
Step 1



Step-2



Step-3



The final decision boundary is $w^T x = 0$

$$\therefore [x_1 \ x_2] \begin{bmatrix} 0.9 \\ 0.5 \end{bmatrix} = 0 \quad \therefore 0.9x_1 + 0.5x_2 = 0$$

$$\therefore x_1 = 0.5, x_2 = -0.9$$

$$\text{OR } x_1 = -0.5, x_2 = 0.9$$

1.3

Let's consider the sets P and N are finite and linearly separable.
 Let's assume ~~the~~ we get the updated weighted vector w_t from
 perceptron learning algorithm in t no. of steps which is finite.
 So if we prove t is finite then ~~the~~ we can conclude
 the perceptron learning algorithm will converge in finite
 no. of steps.

we can make three simplification without losing any generality.

- (i) the sets P and N can be join in a single set $P' = P \cup N^-$
 where N^- set contains negated element of N .
- (ii) the vectors in P' can be normalized, because if a weight
 vector w is found so that $w \cdot x > 0$ is also valid for any
 vector γx where γ is a constant
- (iii) now we can say w^* is normalized vector of w .

Let's consider. after $t+1$ steps. weight vector w_{t+1} has been
 computed hence we can say. at time t , a vector p_i was
 incorrectly classified by weight vector w_t , hence $w_{t+1} = w_t + p_i$

$$\therefore \text{angle between } w_{t+1} \text{ and } w^* \Rightarrow \cos \theta = \frac{w^* \cdot w_{t+1}}{\|w_{t+1}\|}$$

$$\begin{aligned} \text{we know } w^* \cdot w_{t+1} &= w^* \cdot (w_t + p_i) \\ &= w^* \cdot w_t + w^* \cdot p_i \\ &\text{which is } > w^* \cdot w_t + \Delta \end{aligned}$$

with $\Delta = \min \{ w^* \cdot p_i \mid p_i \in P' \}$
 as w^* creates an absolute linear separation between P
 and N , $\Delta > 0$

$$\therefore w^* \cdot w_{t+1} > w^* \cdot w_0 + (t+1)\Delta$$

$$\begin{aligned} \text{again } \|w_{t+1}\|^2 &= (w_t + p_i) \cdot (w_t + p_i) \\ &= w_t^2 + 2p_i \cdot w_t + p_i^2 \end{aligned}$$

as $w_t \cdot p_i \geq 0$ or negative.

$$\begin{aligned} \|w_{t+1}\|^2 &\leq \|w_t\|^2 + \|p_i\|^2 \\ &\leq \|w_t\|^2 + 1 \end{aligned}$$

$$\text{again } \|w_{t+1}\|^2 = \|w_0\|^2 + (t+1) \text{ as } p_i \text{ is normalized}$$

$$\therefore \text{wsp} = \frac{W^* W_0 + (t+1)A}{\sqrt{(W_0)^2 + (t+1)}}$$

$$\frac{W^* W_0 + (t+1)A}{\sqrt{(W_0)^2 + (t+1)}} \text{ grows proportionally to } \sqrt{t}$$

As A is positive, it can be arbitrarily large, but wsp is bound $\& \text{wsp} \leq 1$, so t must be bounded by a maximum value which is finite.

3.
(i)

Sigmoid:

$$\textcircled{a} \sigma(x) = \frac{1}{1+e^{-x}}$$

$$\begin{aligned} \frac{d}{dx}(\sigma(x)) &= \frac{d}{dx}(1+e^{-x})^{-1} \\ &= (1+e^{-x})^{-2} \left[\frac{d}{dx}(-e^{-x}) + \frac{d}{dx}(1) \right] \\ &= (1+e^{-x})^{-2} \cdot (-e^{-x} \cdot -1) = \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \frac{1}{1+e^{-x}} \cdot \left(1 - \frac{1}{1+e^{-x}}\right) \\ &= \textcircled{a} \sigma(x) (1 - \sigma(x)). \end{aligned}$$

(ii) tanh:

$$\sigma(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$\begin{aligned} \frac{d}{dx} \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right) &= \frac{e^x + e^{-x}}{(e^x + e^{-x})^2} \frac{d}{dx}(e^x - e^{-x}) - \frac{e^x - e^{-x}}{(e^x + e^{-x})^2} \frac{d}{dx}(e^x + e^{-x}) \\ &= \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2} \\ &= 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2 = 1 - (\tanh(x))^2 \end{aligned}$$

ReLu:

$$f(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ x & \text{for } x \geq 0 \end{cases}$$

$$\therefore f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

3.2 Strategies to avoid over-fitting in neural network:

Strategies are

(i) Simplifying the model or the decrease the complexity by reducing hidden layers and by creating smaller network

(ii) Early stopping

(iii) By using data augmentation

(iv) ~~By~~ by using regularization

(v) by using dropouts.

3.3.

$$f(z) = z_1 x + z_2 y, \quad x = [1, 1]^T, \quad y = [1, 1]^T$$

$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^1 \quad f(z) = z_1 x + z_2 y$$

$$z = g(x) = [x^2, x^3]$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial z} \times \frac{\partial z}{\partial x} = \begin{bmatrix} x \\ y \end{bmatrix} [2x^2, 3x^2]$$

$$= \begin{bmatrix} 2x^2 x_1 & 3x^2 y_1 \\ 2x^2 x_2 & 3x^2 y_2 \end{bmatrix} \bigg|_{\substack{x=[1,1] \\ y=[1,1] \\ x=2}} = \begin{bmatrix} 8 & 12 \\ 8 & 12 \end{bmatrix}$$

3.4.

(i) Mean Squared Error:

Mean squared error is defined by the following equation.

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2.$$

(ii) Binary Cross Entropy:

BSE

Binary cross entropy is defined by the below equation

$$BSE = -\frac{1}{n} \sum_{i=1}^n y_i \cdot \log(\hat{p}(y_i)) + (1 - y_i) \cdot \log(1 - \hat{p}(y_i)).$$

where $\hat{p}(y)$ is the probability of expected result

(iii) Categorical Cross Entropy:

Categorical cross entropy is defined by the following equation.

$$CC = -\sum_{x \in \mathcal{X}} p(x) \log(\hat{p}(x))$$

where $p(x)$ is the probability distribution of x .

3.5.

(i) Batch Gradient Descent:

Batch gradient descent takes all the training data into consideration to take a single step. We take the average of the gradients of all training examples and then use that mean gradient to update our parameters. So that's one step of gradient descent is one epoch.

8.5.

(ii)

Stochastic Gradient Descent:

In Stochastic Gradient Descent we consider just one example at a time to take a single step. We do the following steps in one epoch for SGD.

- (i) Take an example
- (ii) Feed it to neural network.
- (iii) Calculate its gradient
- (iv) Use the calculated gradient above to update the weights.
- (v) Repeat all steps above for all examples in the training set to training dataset.

(iii) Mini-Batch Gradient Descent:

~~For mini batch gradient descent~~

For mini Batch Gradient descent we take a batch of a fixed no. of training examples which is less than the actual dataset called a mini-batch. and do the following steps.

- (i) Pick a mini batch
- (ii) Feed it to the neural network.
- (iii) Calculate the mean gradient of the mini batch.
- (iv) Use the ~~mean~~ calculated mean gradient to update the weights.
- (v) Repeat all above steps for all minibatches.

3.6.

(iv) Momentum based Gradient Descent:

Momentum based gradient descent is the calculate the exponentially weighted average of the gradients and use them instead of updating the weights. which works faster than regular algorithm for gradient descent.

3.5

(v) Adam Solver:

Adam solver can be defined as the combination of RMSprop and ~~Stochastic~~ Stochastic Gradient descent with momentum. It uses squared gradients to scale the learning rate like RMSprop also takes advantage of momentum by moving average of the gradient descent instead of gradient.

2.2.

(i) Change Loss function:

As we change the loss function from simplest ~~to the sum~~ of Mean squared error to the complex of cross entropy loss we can observe that we get a better accuracy for less no. of epochs as it provide a better back propagation and ~~converge~~ the network converges at a less no. of steps. where ~~as the~~ as simpler loss function model underfits and we get less accuracy. but with increase in epochs we can observe that the simpler loss function work better than the complex loss function because that time it overfits the model.

2.2: (ii) Change in learning rate:

Keeping other parameters fixed if we ~~work~~ change the learning rate we can observe ~~the~~ that for the less no. of epochs if learning rate is high it provide a decent accuracy because with higher learning rate the ~~the~~ network converges at a higher speed and fits the model but with higher no. of ~~accura~~ epochs if we increase the learning rate initially ~~it~~ the accuracy increases and after a ~~cer~~ value with increase in the learning rate accuracy decreases as the model misses its local minima ~~and~~ overfits.

2.2: (iii) Change in Number of hidden layers:

~~Q~~ Keeping other parameters fixed if we ~~increase~~ with less no. of hidden layers or minimal no. of hidden layers we can observe less accuracy ~~at~~ its a very simple neural network but as we ~~increase~~ increase the hidden layers we can observe ~~an~~ improvement in accuracy as the complexity of the network increases. But again after a certain no. of layers if we increase the no. of layers accuracy decreases as its too complex and information flow doesn't reach till the end layers and we didn't get optimum result.