Q.2  Perceptron with margin convergence proof



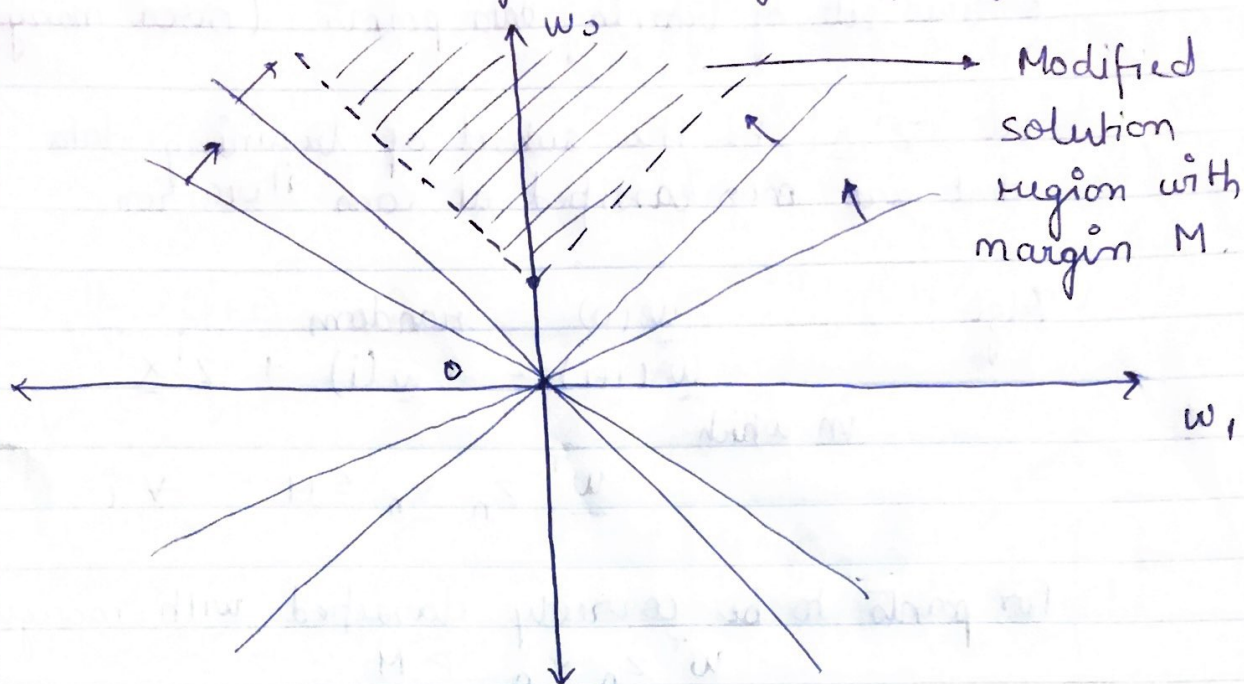Modified solution region with margin M.

ASSUMPTIONS:

- Fixed Increments $\eta(i) = \eta = $ constant $> 0$
- Sequential Gradient Descent
- Data points are linearly separable
- Use reflected data points $z_n \underline{x}_n$, $n = 1, 2, \cdots, N$

We can set $\eta = 1$, without loss of generality:
Let $z_n \underline{x}'_n = z_n \underline{x}_n$, $\eta > 0$
Then drop primes.

ALGORITHM:

- $\underline{w}(0) = $ arbitrary

- $\underline{w}(i+1) = \underline{w}(i) + z_i \underline{x}_i$ $\left[ \underline{w}^T z_i \underline{x}_i \leq M \right]$

in which $z_i \underline{x}_i$ , $i = 1, 2, \cdots$ are the cyclically ordered set of training data points (over many epochs).

Let $z^i \underline{x}^i$ be the subset of training data points that are misclassified at each iteration.

Algo :
$$\underline{w}(0) = \text{random}$$
$$\underline{w}(i+1) = \underline{w}(i) + z^i \underline{x}^i$$

in which
$$\underline{w}^T z_n \underline{x}_n \leq M \qquad \forall i$$

For points to be correctly classified with margin $M$,
$$\underline{w}^T z_n \underline{x}_n > M$$

If $\underline{\hat{w}}$ is a solution, then $a\underline{\hat{w}}$, $a > 1$ is also a solution

$$\therefore \underline{\hat{w}}^T z_n \underline{x}_n > M \qquad \forall n$$

$$\Rightarrow a \underline{\hat{w}}^T z_n \underline{x}_n > aM \qquad \forall n$$

- Need 'error measure' on $\underline{w}(i)$:

$$E_{\underline{w}}(i) = \| \underline{w}(i) - a\underline{\hat{w}} \|_2^2$$

- Show $E_{\underline{w}}(i)$ must decrease at each iteration $\Rightarrow$

$$\underline{w}(i+1) - a\underline{\hat{w}} = \underline{w}(i) - a\underline{\hat{w}} + z^i \underline{x}^i \qquad , a > 1$$

$$\Rightarrow \| \underline{w}(i+1) - a\underline{\hat{w}} \|_2^2 = \| \underline{w}(i) - a\underline{\hat{w}} \|_2^2 + 2 [\underline{w}(i) - a\underline{\hat{w}}]^T z^i \underline{x}^i + \| z^i \underline{x}^i \|_2^2$$

In this, $\qquad 2\underline{w}^{(i)T} z^i \underline{x}^i \leq 2M.$ $\qquad\qquad >aM$

$$\Rightarrow \| w(i+1) - a\hat{\underline{w}} \|_2^2 \leq \| w(i) - a\hat{\underline{w}} \|_2^2 - 2a\overbrace{\hat{w}^T z^i \underline{x}^i}$$
$$-\cdot 2\underline{w}^{(i)T} z^i \underline{x}^i \qquad\qquad + \| z^i \underline{x}^i \|_2^2$$

$$\Rightarrow \| w(i+1) - a\hat{\underline{w}} \|_2^2 \leq \| w(i) - a\hat{\underline{w}} \|_2^2 + 2M - 2aM$$
$$+ \| z^i \underline{x}^i \|_2^2$$
$$\searrow >0$$

Let $\qquad b^2 \overset{\Delta}{=} \max\limits_{j} \| x_j \|_2^2 = \left[ \begin{array}{c} \text{length of largest} \\ \text{data point (vector)} \end{array} \right]^2$

$$\therefore \quad \| w(i+1) - a\hat{\underline{w}} \|_2^2 \leq \| w(i) - a\hat{\underline{w}} \|_2^2 + 2M - 2aM + b^2$$

Now choose $\qquad a = \dfrac{b^2 + M}{M} > 1$, we get,

$$\| w(i+1) - a\hat{\underline{w}} \|_2^2 \leq \| w(i) - a\hat{\underline{w}} \|_2^2 - b^2$$
$$\Rightarrow \quad E\,\underline{w}(i+1) \leq E\underline{w}(i) - b^2$$

$\Rightarrow$ so each iteration reduces $E\underline{w}$ by at least $b^2$.

· Applying forcing arguement

$$0 \leq E_{\underline{w}}(i+1) \leq E_{\underline{w}}(i) - b^2 \qquad\qquad \forall i$$

For some $i_0$, we would have $E\underline{w}(i_0) < b^2$
so that,
$$0 \leq E\,\underline{w}(i+1) \leq E\underline{w}(i) - b^2 < 0$$
which is impossible

$\Rightarrow$ Iteration must cease at $i = i_0$ (or sooner)

$\therefore$ Algorithm converges at a solution weight vector at $(i_0 - 1)^{th}$ iteration or sooner.

**Q3.** a) $\quad \Delta \underline{w}(i) = \dfrac{\underline{w}(i+1) - \underline{w}(i)}{N}$

$\quad\quad\quad J(\underline{w}) = \sum\limits_{n=1}^{N} J_n(\underline{w})$

$\quad\quad\quad \eta(i) = \eta$

$\therefore \quad E\{\Delta \underline{w}(i)\} = E\{\underline{w}(i+1) - \underline{w}(i)\}$

$\quad\quad\quad = E\{\underline{w}(i+1)\} - E\{\underline{w}(i)\}$

$\quad\quad\quad = E\{\underline{w}(i) + \eta \, z_n \underline{x}_n^{(i)}\} - E\{\underline{w}(i)\}$

$\quad\quad\quad = E\{\underline{w}(i)\} + E\{\eta \, z_n \underline{x}_n^{(i)}\} - E\{\underline{w}(i)\}$

$\quad\quad\quad = \eta \, E\{z_n \underline{x}_n^{(i)}\} = \eta \{z_n \underline{x}_n^{(i)} \cdot p(z_n \underline{x}_n^{(i)})\}$

$\quad\quad\quad = -\dfrac{\eta}{N} \, \nabla_{\underline{w}} J_{\bullet}(\underline{w}) = -\dfrac{\eta}{N} \sum\limits_{n=1}^{N} \nabla_{\underline{w}} J_n(\underline{w})$

b) $\quad E\left\{\sum\limits_{i=0}^{N-1} \Delta \underline{w}(i)\right\} \quad\quad\quad \Delta \underline{w}(i)$ are iid

$\quad\quad = E\{N(\Delta \underline{w}(i))\}$

$\quad\quad = N \, E\{\Delta \underline{w}(i)\}$

$\quad\quad = -\eta \cdot \dfrac{N}{N} \cdot \nabla_{\underline{w}} J_{\bullet}(\underline{w}) \quad\quad\quad$ from (a)

$\quad\quad = -\eta \, \dfrac{\nabla_{\underline{w}} J_{\bullet}(\underline{w})}{N}$

$\quad\quad = -\eta \sum\limits_{n=1}^{N} \nabla_{\underline{w}} J_n(\underline{w})$

c) $\Delta \underline{w}(i) = \eta \sum_{\substack{\text{all }\text{ } \text{s.t.} \\ X_n \in X}} z_n X_n$

$= -\eta \nabla_w J(\underline{w})$

$= -\eta \sum_{n=1} \nabla_w J_n(\underline{w})$ .

Thus comparing (b) and (c), we can see that the expected value of the sum of weight difference for stochastic gradient descent, variant 2 is equal to the difference in weights for batch gradient descent.

In Stochastic Gradient Descent - Variant 2, we randomly pick a training data point (with replacement) and perform single sample update on the weight, whereas in block gradient descent, we update weight for all misclassified datapoints in one iteration.

Batch gradient descent is great for convex, or relatively smooth error manifolds. Whereas, stochastic gradient descent is better for error manifolds that have lots of local maxima / minima. Computationally, stochastic gradient descent is much faster. ~~Adding~~ This computational ~~advantage~~ advantage is leveraged by performing many more steps than conventional batch gradient descent, resulting in a close model to that found by the latter.