

## Marginal Perceptron Convergence Proof

In this proof, the following assumptions are made

\* Fixed increment,  $\eta(i) = \eta = \text{constant} > 0$   
for ease we assume  $\eta = 1$  <sup>[without loss of generality]</sup> ①

\* All the data points are considered to be linearly separable

\* Sequential Gradient Descent technique is used

\*  $w(0)$  is arbitrary

$$w(i+1) = w(i) + \sum_i z_i x_i \eta; \text{ where } w(i)^T z_i x_i \leq B \quad \text{②}$$

where  $z_i x_i$  represents the reflected data points that are cyclically ordered.

and this eqn ② is iterated over many epochs.

from ① as  $\eta = 1$ , we can rewrite them as

$$w(i+1) = w(i) + \sum_i z_i x_i \quad \text{③}$$

Let us make another assumption,

$\sum_i z_i^2 x_i^2$  where  $i = 0, 1, 2, \dots, c$  be the subset of training data points that are misclassified

at each epoch's iteration.

∴ We can start the Marginal Perceptron algorithm as (i) p, iteration limit \*

$w(0) = \text{arbitrary value}$

$$w(i+1) = w(i) + \sum z^p x^p \quad \text{updates } i^{\text{th}} \text{ iteration}$$

$$\text{when } w(i)^T z^p x^p \leq B \quad \forall i$$

If  $\underline{w}$  is a (weight vector) solution,

$$\text{i.e. } \underline{w}^T z_n x_n > B \quad \forall n$$

then a  $\underline{w}$  is also a solution

$$\text{i.e. } a \underline{w}^T z_n x_n > B \quad \forall n \quad [if a > 0]$$

(4)

We know that the weight vector  $w$  has some error value that keeps decreasing in each iteration. let that error value be defined by

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

(5)

from equations (3) & (4), we can say that-

$$w(i+1) - a\underline{w} = w(i) + \sum z^p x^p - a\underline{w}$$

$\Rightarrow$

$$[w(i+1) - a \hat{w}] = [w(i) - a \hat{w}] + z^i x^i$$

for  $a > 0$ .

Squaring on both sides in term of norm

$$\begin{aligned} \|w(i+1) - a \hat{w}\|_2^2 &= \|w(i) - a \hat{w}\|_2^2 \\ &+ 2 (w(i) - a \hat{w})^T z^i x^i \\ &+ \|z^i x^i\|_2^2 \end{aligned}$$

$$\begin{aligned} \Rightarrow \|w(i+1) - a \hat{w}\|_2^2 &= \|w(i) - a \hat{w}\|_2^2 \\ &+ 2 w^T z^i x^i - 2 a \hat{w}^T z^i x^i \\ &+ \|z^i x^i\|_2^2 \end{aligned}$$

$\Rightarrow$

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

$$\|w(i+1) - a \hat{w}\|_2^2 \leq (\|w(i) - a \hat{w}\|_2^2 - 2 a \hat{w}^T z^i x^i + \|z^i x^i\|_2^2) + 0 \quad (6)$$

In order to maximize the equation at margin

we need to minimize  $\hat{w}^T z^i x^i$  and

maximize  $\|z^i x^i\|_2^2$  value such that

$$\hat{w}^T z^i x^i > B \quad \text{and} \quad \|z^i x^i\|_2^2 > 0$$



Let's assume variables to make eqn simpler

$$\text{let } b^2 = \max_j \|x_j\|_2^2$$

$$c = \min_j \{ \omega^T z_j x_j \} > B$$

$\therefore$  we can re-write eqn (6) as

$$\|w(i+1) - a\bar{w}\|_2^2 \leq \|w(i) - a\bar{w}\|_2^2 - 2ac + b^2 \quad (7)$$

$a > 0$

from (7) wkt  $\|w(i) - a\bar{w}\|_2^2$  is the error in weight vector.

$$\therefore (7) \Rightarrow E_w(i+1) \leq E_w(i) - 2ac + b^2$$

if we choose  $a$  such that  $a = \frac{b^2}{c}$

$$\Rightarrow E_w(i+1) \leq E_w(i) - \frac{2b^2}{c} c + b^2$$

$$\Rightarrow E_w(i+1) \leq E_w(i) - b^2$$

$\therefore$  we see that each iteration, error value is decreased by  $b^2$ .  $[b^2 > 0]$

∴ we can write as

$$0 \leq \cancel{E_w(i+1)} E_w(i+1) \leq E_w(i) - b^2 \quad \text{--- (8)}$$

at some iteration  $(i_0)$ ,  $E_w(i_0) < b^2$

∴ eqn (8) becomes implausible

therefor the algorithm's iteration stops at  $(i_0 - 1)^{\text{th}}$  iteration, and the corresponding weight vector is the converged weight vector.

$$3) \quad \Delta w(i) = w(i+1) - w(i) \quad \text{--- (1)}$$

$$J(w) = \sum_{n=1}^N J_n(w) \quad \text{--- (2)}$$

$$\eta(i) = \eta = \text{const} > 0$$

- a) Using Stochastic Gradient Descent, variant 2, a training data point is randomly picked with replacement and single sample update is done.

In Gradient Descent,

$$w(i+1) = w(i) - \eta \sum_{n=1}^N \nabla J_n(w) \quad \text{--- (3)}$$

for single sample update;

$$(3) \Rightarrow w(i+1) = w(i) - \eta \nabla J_n(w) \quad \text{--- (4)}$$

And the expected value of  $\Delta w(i)$

$$E\{\Delta w(i)\} = \sum_{n=1}^N \Delta w(i) p(n)$$

where  $p(n)$  = probability with which points are chosen and  $n = 1, 2, \dots, N$  with uniform probability.

$$P(n) = \frac{1}{(N-1)+1} = \frac{1}{N}$$

$$\therefore E\{\Delta w(i)\} = \sum_{n=1}^N \Delta w(i) \times \frac{1}{N}$$

①  $\Rightarrow$

$$E\{\Delta w(i)\} = \sum_{n=1}^N (w(i+1) - w(i)) \times \frac{1}{N}$$

④  $\Rightarrow$

$$= \sum_{n=1}^N -\eta \nabla J_n(w) \times \frac{1}{N}$$

$$= -\frac{\eta}{N} \sum_{n=1}^N \nabla J_n(w)$$

$$E\{\Delta w(i)\} = -\frac{\eta}{N} \nabla \sum_{n=1}^N J_n(w)$$

(10)

②  $\Rightarrow$

$$E\{\Delta w(i)\} = -\frac{\eta}{N} \nabla J(w)$$



b) Find  $E \left\{ \sum_{n=0}^{N-1} \Delta \omega(t) \right\}$

where each  $\Delta \omega(t)$  is IID.

from (a)  $E \{ \Delta \omega(t) \} = -\frac{\eta}{N} \nabla \sum_{n=1}^N J_n(\omega)$   
(a)

$$= -\frac{\eta}{N} \nabla J(\omega)$$

$$E \left\{ \sum_{n=0}^{N-1} \Delta \omega(t) \right\} = \sum_{n=0}^{N-1} E \{ \Delta \omega(t) \}$$

$$[ E[ax+b] = aE(x) + b ]$$

$$= (N-1) + 1 E \{ \Delta \omega(t) \}$$

$$= N \cdot -\frac{\eta}{N} \nabla J(\omega)$$

$$E \left\{ \sum_{n=0}^{N-1} \Delta \omega(t) \right\} = -\eta \nabla J(\omega)$$

(b)

$$-\eta \nabla \sum_{n=1}^N J_n(\omega)$$



c) Batch Gradient Descent :

$$w(i+1) = w(i) + \eta(i) \sum_n z_n x_n - \eta_n$$

$$= w(i) - \eta(i) \nabla_w J(w)$$

$$= w(i) - \eta(i) \nabla_w \sum_{n=1}^N J_n(w)$$

—————> ①

$$w(i+1) - w(i) = -\eta(i) \nabla_w \sum_{n=1}^N J_n(w)$$

$$= E \left\{ \sum_{n=0}^{N-1} \Delta w(i) \right\}$$

[∴ ⑥]

part ③ result is similar to part ⑥

i.e. in batch gradient descent, we consider all the misclassified data points

but in stochastic gradient, we consider single sample and repeat for 'N' samples.  
It is computationally stable because it can be used for a set with more minima/maxima.