

---

# Comparison of performance of Standard Perceptron and Second order Perceptron

---

**Hardik Sahi**

David R. Cheriton School of Computer Science  
University of Waterloo

Waterloo, ON, N2L 3G1

h2sahi@uwaterloo.ca

CS 698: Intro to ML Fall 2017 Project Report

## Abstract

1 In this report, I perform empirical evaluation of performance of Second order  
2 perceptron algorithm and compare it with Standard perceptron algorithm. Perform-  
3 ance analysis is done based on mistake bound model for online learning. The  
4 shortcomings of standard perceptron and the way in which they are handled in  
5 second order perceptron are discussed. I try and replicate the results in paper [1]  
6 using synthetic data satisfying specific spectral or geometric properties.

## 7 1 Introduction

8 **Online methods** of machine learning are the ones in which the data is available to the learner  
9 sequentially and is used to update the best predictor for future data at each step. Such algorithms  
10 use a **mistake bound model** of learning which are discussed in [2]. This model is the one in which  
11 the learner is evaluated on the basis of number of mistakes it makes before converging to the correct  
12 hypothesis. For second order perceptron, the idea is that the algorithm under consideration must  
13 be able to sequentially classify the data points correctly ensuring that it does not make mistakes  
14 much greater than that made by the fixed/ referenced linear classifier on the same sequence  $S$   
15  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$  of input data. Advantage of the above approach of analysis is that the  
16 loss of learning algorithm, in terms of misclassified points with respect to reference classifier can be  
17 explained and nicely bound in terms of geometric properties of individual data sequences on which  
18 the algorithm is run.

19  
20 As will be shown in the next section, the Perceptron learning algorithm (**PLA**) is a gradient descent  
21 algorithm whose mistake bound is determined by the trace information (first order information). We  
22 get the second order perceptron algorithm by enhancing Standard Perceptron to become sensitive to  
23 second order information by incrementally computing sparse data correlation matrix. In this way,  
24 the enhanced algorithm is able to exploit certain geometric properties of data which are missed by  
25 standard Perceptron.

26  
27 A typical scenario under study in which the mistake bound of Second order perceptron is signif-  
28 icantly smaller than that of standard Perceptron is when the data lies on a flat ellipsoid so that the  
29 eigen values of correlation matrix is significantly different along different axis whereas the trace  
30 will be determined by the largest eigen value.

31  
32 In this report, I explain the properties of the enhanced perceptron and show how its behaviour  
33 depends on eigenvalues of data correlation matrix. Also, I will publish the results of application of

34 Second order perceptron on an artificial dataset and compare it with standard perceptron.  
 35

## 36 2 Standard Perceptron

37 A set of input data points  $S \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$  is **linearly separable** if there is a  
 38 hyperplane (which splits the input space into two half-spaces) such that all points of the first class  
 39 are in one half-space and those of the second class are in the other half-space. Hence the instance  
 40 vectors  $\mathbf{x}_i$  are consistently labeled according to whether they lie on positive ( $y_i = +1$ ) or the negative  
 41 side ( $y_i = -1$ ) of an unknown target hyperplane with normal vector  $\mathbf{u} \in \mathbb{R}^d$ .  
 42

43 The standard perceptron learning algorithm was discussed for the first time by **Rosenblatt** in [4]  
 44 Given the set  $S$  of data points and starting from an initial weight vector, the PLA is ultimately  
 45 able to yield a hyperplane (represented by weight vector which is normal to the hyperplane) that  
 46 classifies all the examples correctly after a finite number of iterations for a linearly separable dataset.  
 47

48 The prediction made by PLA at time  $t$  is given by  $\hat{y}_t = \text{SIGNUM}(\mathbf{w}^\top \mathbf{x}_t) \in \{-1, +1\}$ . We say that the  
 49 algorithm has made a mistake in making a prediction if  $\hat{y}_t \neq y_t$ . Hence, to correct its mistake, the  
 50 algorithm updates its internal state (maintained by weight vector  $\mathbf{w}$ ) following a simple additive  
 51 rule:  $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t$ . Geometrically speaking, the updated weight vector  $\mathbf{w}_t$  moves along the  
 52 direction of instance  $\mathbf{x}_t$  for which the prediction was wrong. Hence the number of updates made in  
 53 the weight vector is equal to the number of mistakes the algorithm makes in the process of learning.  
 54

55 As shown in [3] by **A. Novikoff**, the convergence of standard perceptron on linearly separable data  
 56 set is guaranteed after a finite number of iterations. That is, it makes at most  $(R/\gamma)^2$  mistakes on any  
 57 number  $t$  of examples where,

$$R(\text{Radius}) = \max_{1 \leq s \leq t} \|\mathbf{x}_s\|$$

$$\gamma(\text{Margin}) = \min_{1 \leq s \leq t} |\mathbf{u}^\top \mathbf{x}_t|$$

58 A typical scenario in which standard perceptron would be extremely slow even for linearly separable  
 59 data set would be when instance vector  $\mathbf{x}_t$  would have a very small projection over target hyperplane  
 60 (hence  $\mathbf{x}_t$  would be almost orthogonal to  $\mathbf{u}$ , making denominator extremely small ) and has a large  
 61 norm  $\|\mathbf{x}_t\|$  (numerator becomes huge) [Figure 2].  
 62

## 63 3 Whitened Perceptron: Adding sensitivity to second order information.

64 The second order perceptron algorithm is defined to exploit the aforementioned spatial orientation  
 65 of data points and hence uses second order information of data distribution as well. To be able to  
 66 use it, I first introduce an intermediate algorithm called Whitened perceptron algorithm.  
 67

68 Whitened perceptron is a non incremental version of standard perceptron which means that all the  
 69 instance vectors are known ahead of time, only the corresponding labels are unknown. Also, for  
 70 simplicity it is assumed that  $\text{span}\{\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_N\} = \mathbb{R}^d$ . This ensures that the data correlation matrix  
 71  $\mathbf{C} = \sum_{i=1}^T \mathbf{x}_i \mathbf{x}_i^\top$  is full rank and hence its inverse exists ( $\mathbf{C}^{-1}$  exists). Also, as  $\mathbf{C}$  is positive semi  
 72 definite matrix,  $\mathbf{C}^{-1/2}$  also exists. Given all the above, the whitened perceptron algorithm simply  
 73 involves using standard perceptron algorithm on the whitened sequence represented by  $\{(\mathbf{C}^{-1/2} \mathbf{x}_1,$   
 74  $y_1), (\mathbf{C}^{-1/2} \mathbf{x}_2, y_2) \dots (\mathbf{C}^{-1/2} \mathbf{x}_N, y_N)\}$ .  
 75

76 Now, in order to show the advantage of whitening the data sequence and applying standard  
 77 perceptron to it instead of simply applying the PLA directly to un-whitened data, I will show the

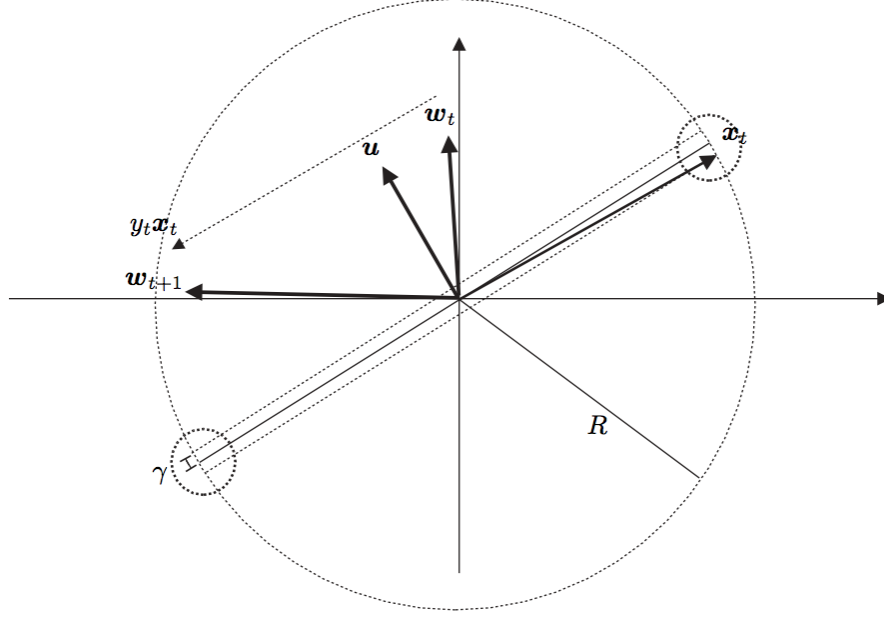


Figure 1: [1] A typical scenario in which perceptron convergence is very slow even for linearly separable dataset.  $\mathbf{u}$  is the target normal vector,  $\mathbf{w}_t$  is the weight vector that the algorithm maintains at the start of trial  $t$ . Because the instance vector  $\mathbf{x}_t$  has a very small projection on  $\mathbf{u}$  and length of  $\mathbf{x}_t$  is  $R$ , following the simple additive rule of perceptron:  $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t$  updates the weight vector to  $\mathbf{w}_{t+1}$  farther away from  $\mathbf{u}$  than it originally was. Hence, this slows down the convergence of perceptron algorithm

78 mistake bound on such a whitened sequence under the assumption that the data sequence originally  
 79 is linearly separable, and hence will converge after a maximum of  $(R/\gamma)^2$  mistakes on any number  $t$   
 80 of examples.  
 81

82 From [1], it is realised that under the above conditions, whitened sequence is also linearly separable  
 83 with the hyperplane  $\mathbf{z} = C^{1/2} \mathbf{u}$ . Also,

$$\gamma_{\text{white}} = \frac{\gamma}{\|C^{1/2} \mathbf{u}\|}.$$

84 Hence, maximum number of mistakes that the perceptron makes on whitened data sequence before  
 85 convergence is given by,

$$\left(\frac{R_{\text{white}}}{\gamma_{\text{white}}}\right)^2 = \left(\frac{1}{\gamma^2}\right) \left(\max_t \|C^{-1/2} \mathbf{x}_t\|^2\right) \|C^{1/2} \mathbf{u}\|^2$$

86 Simplifying the above, we get,

$$\left(\frac{R_{\text{white}}}{\gamma_{\text{white}}}\right)^2 = \left(\frac{1}{\gamma^2}\right) \left(\max_t \mathbf{x}_t^T C^{-1} \mathbf{x}_t\right) (\mathbf{u}^T C \mathbf{u})$$

87 Clearly, the term  $\mathbf{x}_t^T C^{-1} \mathbf{x}_t$  can be associated with the extent of correlation in the input sequence  $S$ .  
 88 As can be seen from the equation above, the maximum number of mistakes made will be reduced  
 89 if the numerator is minimized. Consider the case in which the input sequence is highly correlated,  
 90 which implies that the the term  $\mathbf{x}_t^T C^{-1} \mathbf{x}_t$  is very small. Also, assume that the input instances have

91 very small projection over the target normal vector  $\mathbf{u}$ . In such a scenario the term  $\mathbf{u}^T \mathbf{C} \mathbf{u}$  is very  
 92 small. Hence, the numerator term of the above equation is minimized following the aforementioned  
 93 conditions. So, it can be concluded that the mistake bound for standard perceptron application  
 94 on whitened sequence is significantly smaller than that of standard perceptron application on  
 95 non-whitened sequence.  
 96

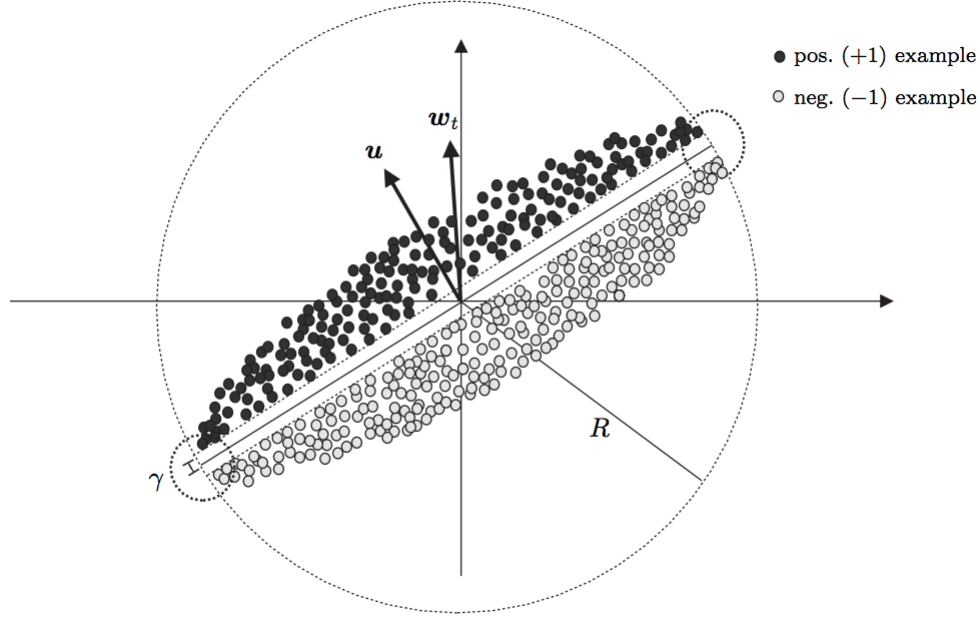


Figure 2: [1] The typical scenario which whitened perceptron can take advantage of but the standard perceptron algorithm struggles to converge.

97 For verifying the point that the whitened perceptron algorithm has significantly smaller mistake  
 98 bound as compared to that of standard perceptron algorithm for the input sequence with starkly  
 99 different eigen values along different axis such that the data is maximally squashed along the  
 100 direction of  $\mathbf{u}$  (degenerate scenario), we perform the following analysis:  
 101

102 Because of symmetric distribution of data, we have the following eigenstructure of matrix  $\mathbf{C}$ :  $\lambda_{\min}$   
 103 and  $\lambda_{\max}$  are the minimal and maximal eigenvalues of  $\mathbf{C}$  respectively. As per our aforementioned  
 104 degenerate scenario,  $\lambda_{\min}$  is aligned along the target normal vector  $\mathbf{u}$  and  $\lambda_{\max}$  is aligned orthogonal  
 105 to  $\mathbf{u}$  (represented by  $\mathbf{u}^\perp$ ).

$$\begin{aligned}
 \lambda_{\min} &= \mathbf{u}^T \mathbf{C} \mathbf{u} \\
 &= \mathbf{u}^T \left( \sum_{i=1}^T \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{u} \\
 &= \sum_{i=1}^T (\mathbf{u}^T \mathbf{x}_i)^2 \\
 &= \gamma^2 T
 \end{aligned}$$

106 Since  $\lambda_{\min} + \lambda_{\max} = \sum_{i=1}^T \|\mathbf{x}_i\|^2 \implies (\mathbf{u}^T \mathbf{x})^2 + ((\mathbf{u}^\perp)^T \mathbf{x})^2 = \|\mathbf{x}\|^2$

$$\begin{aligned}\mathbf{x}_t^\top C^{-1} \mathbf{x}_t &= \frac{(\mathbf{u}^\top \mathbf{x}_t)^2}{\lambda_{\min}} + \frac{((\mathbf{u}^\perp)^\top \mathbf{x}_t)^2}{\lambda_{\max}} \\ &= \frac{\gamma^2}{\gamma^2 T} + \frac{\|\mathbf{x}_t\|^2 - \gamma^2}{\sum_{i=1}^T \|\mathbf{x}_i\|^2 - \gamma^2 T}\end{aligned}$$

Hence,

$$\begin{aligned}\left(\frac{R_{\text{white}}}{\gamma_{\text{white}}}\right)^2 &= \left(\frac{1}{\gamma^2}\right) \left(\max_t \mathbf{x}_t^\top C^{-1} \mathbf{x}_t\right) (\mathbf{u}^\top C \mathbf{u}) \\ &= 1 + \frac{R^2 T - \gamma^2 T}{\sum_{i=1}^T \|\mathbf{x}_i\|^2 - \gamma^2 T}\end{aligned}$$

It can be clearly seen that the above bound for maximum number of mistakes approaches 2 (which is a constant) as norm of instance vectors  $\mathbf{x}_t$  approaches  $R$ , which is independent of both margin  $\gamma$  and radius  $R$  of the ball containing data. Hence it proves that the mistake bound is extremely small in the degenerate scenario as shown above.

#### 4 Second order perceptron.

Second order perceptron algorithm is viewed as an incremental version of Whitened perceptron discussed above which involves iterative computation of sparse correlation matrix based only on a small subset of data points observed so far. The pseudocode for Second order perceptron is described below.

**Initialize:**  $X_0 = \emptyset$ ,  $\mathbf{v}_0 = \mathbf{0}$ ,  $k=1$ .

**Repeat**  $t = 1, 2, \dots$

1. Access instance vector  $\mathbf{x}_t$ ;
2. Set  $S_t = [X_{k-1} \ \mathbf{x}_t]$ ;
3. Predict  $\hat{y}_t = \text{SIGNUM}(\mathbf{w}^\top \mathbf{x}_t) \in \{-1, +1\}$ , where  $(S_t S_t^\top)^+ \mathbf{v}_{k-1}$ ;
4. Get actual label corresponding to input,  $y_t$ ;
5. Update the internal state of the algorithm if  $y_t \neq \hat{y}_t$  as follows:

$$\begin{aligned}\mathbf{v}_k &= \mathbf{v}_{k-1} + y_t \mathbf{x}_t \\ X_k &= S_t \\ k &= k+1\end{aligned}$$

Clearly we can see that the number of columns in  $X_k$  matrix is equal to the number of mistakes that the algorithm makes on the dataset. Hence the second order perceptron is also mistake bound like the standard perceptron algorithm.

Clearly the algorithm iterates over the data points over and over again. It converges in case the system of data points is linearly separable but stays in loop if the set is not linearly separable, just like the case with standard perceptron algorithm. But the number of mistakes it makes before convergence is smaller than that made by standard perceptron algorithm.

#### 5 Pros and Cons of Second order perceptron algorithm.

As will be clear from the simulations on artificial data set in the following section, the second order perceptron algorithm exploits the spectral properties of data distribution in the sense that it takes advantage of the maximally squashed data distribution, something that the standard perceptron fails to and hence converges with fewer mistakes as compared to standard perceptron.

On the downside, the implementation of second order perceptron involves taking pseudo-inverse of the data correlation matrix iteratively which makes the algorithm slower as compared to standard perceptron algorithm.

## 144 6 Running the algorithms on artificial datasets.

145 I have tried to replicate the results provide in paper [1]. In order to do so, I tried creating an artificial  
 146 dataset that is linearly separable and also satisfies the spectral properties that enhanced perceptron  
 147 can take advantage of.

148  
 149 In order to create an artificial dataset I sampled data from a multivariate Gaussian distribu-  
 150 tion making sure that the correlation matrix has single dominant eigen vector perpendicular to  
 151 the direction of the maximally squashed data distribution. As will be evident by the result of the  
 152 simulation, the number of mistakes made by second order perceptron would be smaller than those  
 153 made by standard perceptron algorithm.

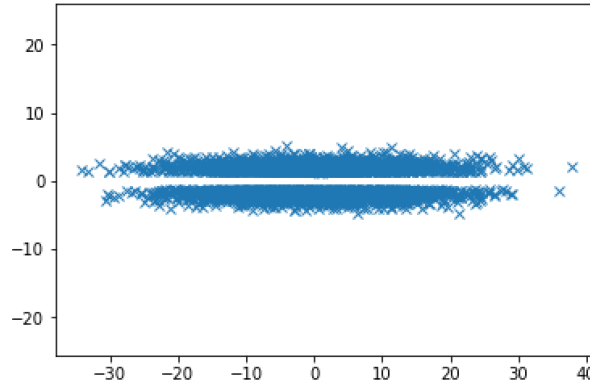
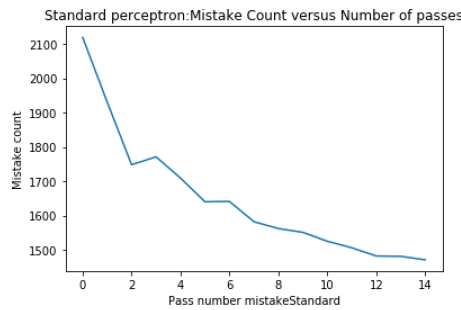
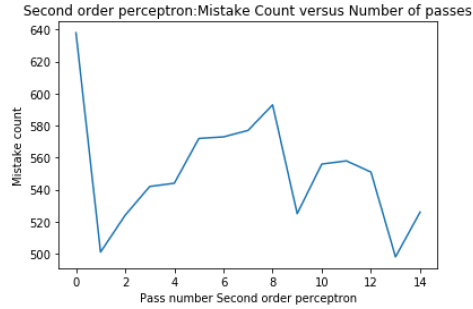


Figure 3: Distribution of the artificial dataset corresponding to the 2 major attributes represented along X and Y axis. The data is sampled from multivariate Gaussian distribution with significantly different eigen values along the 2 axis.



(a) Running standard Perceptron algorithm



(b) Running second order perceptron algorithm on data set

Figure 4: Comparison of mistake count of the standard and second order perceptron. Note the significantly lower mistake count made by second order perceptron.

Algorithm	Mistake percent
Standard Perceptron	35.81%
Second order perceptron	11.99%

Table 1: Mistake count comparison for the two algorithms.

## 154 7 Conclusion.

155 As evident by the results of the simulation above, it is clear that second order perceptron algorithm  
156 has significantly lower bound in case of specific geometric distribution of data. Having said that,  
157 the second order perceptron is slower as compared to standard one as it involves taking the pseudo-  
158 inverse of the incrementally computed data correlation matrix.

## 159 References

- 160 [1] Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. A second-order perceptron algorithm.  
161 In *Proceedings of the 15th Annual Conference on Computational Learning Theory*, COLT '02,  
162 pages 121–137, London, UK, UK, 2002. Springer-Verlag.
- 163 [2] Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold  
164 algorithm. *Mach. Learn.*, 2(4):285–318, April 1988.
- 165 [3] A. B. Novikoff. On convergence proofs on perceptrons. In *Proceedings of the Symposium on*  
166 *the Mathematical Theory of Automata*, volume 12, pages 615–622, New York, NY, USA, 1962.  
167 Polytechnic Institute of Brooklyn.
- 168 [4] F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization  
169 in the brain. *Psychological Review*, 65(6):386–408, 1958.