

1. (20 points) Describe an application of self-supervised learning with 10-20 English or Chinese sentences.

The most famous application of self-supervised learning must be Bert, which can be used on masked token prediction and next sentence prediction.

To explain it briefly, we mask some words in sentences / change the sequence of sentences, Bert is able to fill the correct words into the blanks / whether the sentences are matched.

What can we do on these abilities? Bert can be applied to dozens of downstream tasks, even though they seem to have no relationship with those abilities: Bert can be used as pre-trained model.

With pre-trained model that has the ability of SSL, which learns from data without relying on explicitly annotated labels, many tasks can be done by AI with a greater accuracy, including Natural Language processing.

For example, Sentiment Analysis in text classification .We input sentences, while AI distinguishes the potentially feelings of those sentences, between positive or negative.

It can be used in companies, when customers ask for support about their products, AI need to decide the way it respond to the question or send warning to human making sure the customers problem would be solved in some better way.

Moreover, use those data in statistics and performance evaluation .

Natural Language Processing (NLP) applications can also benefit companies across various industries by automating tasks, improving customer interactions, gaining insights from text data, and enhancing decision-making processes.

To sum up, self-supervised learning is similar to stem cell in our body, it provides pre-trained model with cost-efficiency, better performance in realistic problems, and what's more important, transfer learning.

AI learns from unlabeled data , after we feed it with some labeled data, it becomes specific AI model which can be used in specific ways.

2. (20 points) The following answer was generated when asking chatGPT “*Can we use machine learning to search for the shortest path to exit a maze? Please make your answer brief.*” Do you agree or disagree with the answer? Provide your reasoning with 10-20 English or Chinese sentences.

Default (GPT-3.5)



Can we use machine learning to search for the shortest path to exit a maze? Please make your answer brief.



Yes, machine learning can be used to search for the shortest path to exit a maze. One common approach is to use reinforcement learning algorithms like Q-learning or Deep Q-Networks (DQN) to train an agent to navigate the maze and find the shortest path based on rewards and penalties associated with different actions and states in the maze.

Yes, I agree with the answer.

Reinforcement learning is used on partial/implicit information.

We can't tell AI about which path is the best ( $y_n=1$  or  $y_n=-1$ ), but we can give reward if it exit the maze by spending a shorter time or by less turns.

By reward, AI can learn from the process and keep trying to reach a better solution, depending on our reflection.

However, when to use Q-learning or DQN should depend on situation.

If the size of the maze is small, there won't be problem when using Q-learning/DQN.

If the size is extremely large, the use of DQN will be necessary.

Although both of the two methods learn from experiments that actions interact with the environment, update and keep finding better actions in states.

Q-learning stores the results of experiment by tables, recording state, action and rewards. By going through this table, it's able to decide the action it should take in specific state.

But this makes sense only when the combinations are finite.

When the maze getting larger, the combination of states and actions grow rapidly.

It won't be possible to save records by table.

While DQN, which uses a neural network to approximate the Q-values, is allowed to solve high-dimension problems.

Experience replay, target network, and sample efficiency of DQN are able to help us with more complex maze problems

3. (20 points) The following answer was generated when asking chatGPT “*Can Machine Learning be used to speed up any off-the-shelf algorithm? Please make your answer brief.*” Somehow Google DeepMind announced that they discovered faster sorting algorithms with Machine Learning in June 2023 (<https://www.nature.com/articles/s41586-023-06004-9>). The discovery seems contradictory to the answer from chatGPT. Why is that? Provide your reasoning with 10-20 English or Chinese sentences.

Default (GPT-3.5)



Can Machine Learning be used to speed up any off-the-shelf algorithm? Please make your answer brief.



No, not necessarily. Machine learning can improve certain algorithms by optimizing them for specific tasks or data, but it may not always speed up off-the-shelf algorithms, as it depends on the nature of the problem and the algorithm itself.

“Off-the-shelf” refers to pre-trained models or algorithms that are ready to use and don’t need additional training, which can save time and resources.

However, whether ml can improve off-the -shelf algorithm depends on problem complexity, algorithm adaptability...

For most of off-the-shelf algorithm, ml is unable to speed up.

In this case, the article focuses on two kind of small sort algorithm: fixed sort and variable sort.

Fixed sort is about how to sort in case of the same size of number, while variable sort is different size.

In the experiment, they introduced AI called AlphaDev, it improved sort-3, sort-6, sort-7, except for sort-4, since it can't find a way to improve human code.

I think there are reasons:

First, machine learning comes up with the algorithm by learning from procedure. We designed a game for AlphaDev, in order to count the step it took. The research isn't completed by directly asking it the better solution.

Second, ChatGPT can't update the information after the year of 2021. Most of the news nowadays are out of the range. The research was done recently, which ChatGPT is unable to reach.

Third, the definition of off-the -shelf is about the algorithm that are used in many of human's research, commonly be seen as basic algorithm, and that's might be the reason of why AI thinks ML I can't improve it, some of which are too basic to be learn or too wonderful to be improve, for example, sort-4.

Last, just like the answer from ChatGPT, it depends on the nature of the problem, and the algorithm, some of off-the-shelf algorithms are unable to speed up, while some of which are improvable, just like fixed sort and variable sort.

4. (20 points) For the PLA algorithm introduced in class (page 8/22 of Lecture 2), assume that  $T_+$  mistakes happened during  $y_{n(t)} = 1$ , and  $T_-$  mistakes happened during  $y_{n(t)} = -1$ . Express  $w_0$ , the zero-th component of the PLA solution, in terms of  $T_+$  and  $T_-$ , and prove the result.

$$w_{t+1} = (w_t + y_{n(t)}x_{n(t)})$$

$$\begin{cases} w_0(t) + y_{n(t)}x_{n(t)} &= w_0(t+1) \\ w_0(t) + y_{n(t)}x_{n(t)} &= w_0(t-1) \end{cases}$$

since PLA updates only when mistakes occur,

$$\begin{cases} w_0(t) + 1 \rightarrow T^+ \text{ mistakes} \\ w_0(t) - 1 \rightarrow T^- \text{ mistakes.} \end{cases} \quad \begin{array}{c} \xrightarrow{\text{updates } T^+ \text{ times}} w_0(t) + T^+ \\ \xrightarrow{\text{updates } T^- \text{ times}} w_0(t) - T^- \end{array}$$

$$\text{By correcting all the mistakes} \rightarrow w_0(t) = T^+ - T^-$$

5. (20 points) Consider online spam detection with machine learning. We will represent each email  $\mathbf{x}$  by the distinct words that it contains. In particular, assume that there are at most  $m$  distinct words in each email, and each word belongs to a big dictionary of size  $d \geq m$ . The  $i$ -th component  $x_i$  is defined as  $\|\text{word } i \text{ is in email } \mathbf{x}\|$  for  $i = 1, 2, \dots, d$ , and  $x_0 = 1$  as always. We will assume that  $d_+$  of the words in the dictionary are more spam-like, and  $d_- = d - d_+$  of the words are less spam-like. A simple function that classifies whether an email is a spam is to count  $z_+(\mathbf{x})$ , the number of more spam-like words with the email (ignoring duplicates), and  $z_-(\mathbf{x})$ , the number of less spam-like words in the email, and classify by

$$f(\mathbf{x}) = \text{sign}(z_+(\mathbf{x}) - z_-(\mathbf{x}) - 0.5).$$

That is, an email  $\mathbf{x}$  is classified as a spam iff the integer  $z_+(\mathbf{x})$  is more than the integer  $z_-(\mathbf{x})$ .

Assume that  $f$  can perfectly classify any email into spam/non-spam, but is unknown to us. We now run an online version of Perceptron Learning Algorithm (PLA) to try to approximate  $f$ . That is, we maintain a weight vector  $\mathbf{w}_t$  in the online PLA, initialized with  $\mathbf{w}_0 = \mathbf{0}$ . Then for every email  $\mathbf{x}_t$  encountered at time  $t$ , the algorithm makes a prediction  $\text{sign}(\mathbf{w}_t^T \mathbf{x}_t)$ , and receives a true label  $y_t$ . If the prediction is not the same as the true label (i.e. a mistake), the algorithm updates  $\mathbf{w}_t$  by

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t \mathbf{x}_t.$$

Otherwise the algorithm keeps  $\mathbf{w}_t$  without updating

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t.$$

Prove or disprove that  $(4d+1)(m+1)$  upper bounds the number of mistakes that the online PLA can make for this spam classification problem.

Note: For those who know the bag-of-words representation for documents, the representation we use is a simplification that ignores duplicates of the same word.

$$\max \|x_n\|^2 = R^2 = m+1 \quad (+1 \text{ 的项数} + (-1) \text{ 的项数} + 0 \text{ 的项数} + m+1)$$

$$x_1 \dots x_n \in \{-1, 0, 1\}$$



$$\|\mathbf{w}_t\| = \sqrt{(0.5)^2 + \sum (0^2 + \sum (-1)^2)} = \sqrt{\frac{1}{4} + d}$$

$$\min y_t x_n \mathbf{w}_t^T = \frac{1}{2}$$

since  $\mathbf{w}_t^T$  is the correct matrix  $\Rightarrow y_t \cdot \mathbf{w}_t^T x_n \geq 0 \quad \text{for } n \neq 0$

$$\min y_t \mathbf{w}_t^T x_n = 0 \quad \text{for } n \neq 0$$

$$\mathbf{w}_0 \cdot x_0 = +0.5$$

$$P = \frac{\min y_n x_n \mathbf{w}_t^T}{\|\mathbf{w}_t\|} = \frac{\frac{1}{2}}{\sqrt{\frac{1}{4} + d}} \quad . \quad R^2 = m+1$$

$$T \in \frac{R^2}{P^2} = \frac{m+1}{\frac{1}{4(d+1)}} = (4d+1)(m+1)$$

6. (20 points) Before running PLA, our class convention adds  $x_0 = 1$  to every  $\mathbf{x}_n$  vector, forming  $\mathbf{x}_n = (1, \mathbf{x}_n^{\text{orig}})$ . Suppose that  $x'_0 = -1$  is added instead to form  $\mathbf{x}'_n = (-1, \mathbf{x}_n^{\text{orig}})$ . Assume that running PLA on  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$  with a particular sequence  $n(t), t = 1, 2, \dots$  with  $\mathbf{w}_0 = \mathbf{0}$  returns  $\mathbf{w}_{\text{PLA}}$ , and running PLA on  $\{(\mathbf{x}'_n, y_n)\}_{n=1}^N$  with the same  $n(t)$  with  $\mathbf{w}_0 = \mathbf{0}$  returns  $\mathbf{w}'_{\text{PLA}}$ . Prove or disprove that  $\mathbf{w}_{\text{PLA}}$  and  $\mathbf{w}'_{\text{PLA}}$  are equivalent.

Prove.

$$W_{n+1} = W_n + \gamma_n y_n \mathbf{x}_n$$

$$x_0 = 1, x'_0 = -1$$

Since following the same sequence.  $n(t) . t=1,2,3 \dots$

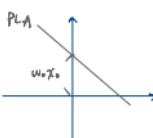
If catch a first data  $\mathbf{x}_1$  with  $y_1 = +1$

$$\rightarrow w_{0(1)} = w_0 + 1 = 1, w'_{0(1)} = w_0 - 1 = -1$$

$$\text{PLA} = \mathbf{w} \cdot \mathbf{x}_n = \mathbf{w}' \cdot \mathbf{x}'_n, (x_0 = 1, x'_0 = -1)$$

$$\begin{bmatrix} w_{0(1),0} \\ w_{0(1),1} \\ \vdots \\ w_{0(1),d} \end{bmatrix} = \begin{bmatrix} w_{0,0} \\ \vdots \\ w_{0,d} \end{bmatrix} + y_{1(1)} \begin{bmatrix} x_0 \\ x_{m(1),1} \\ \vdots \\ x_{m(1),d} \end{bmatrix}, x_0 = 1$$

$$\begin{bmatrix} w_{0(1),0} \\ \vdots \\ w_{0(1),d} \end{bmatrix} = \begin{bmatrix} w'_{0,0} \\ \vdots \\ w'_{0,d} \end{bmatrix} + y'_{1(1)} \begin{bmatrix} x'_0 \\ x'_{m(1),1} \\ \vdots \\ x'_{m(1),d} \end{bmatrix}, x'_0 = -1$$



$w_0 x_0 = w'_0 x'_0$ ,  $w$ 's elements are all the same  
for  $n \neq 0$ ,  $w_n = w'_n$  except for  $w_0$  and  $w'_0$

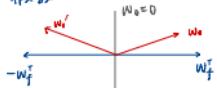
$\begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$  indicated the normal vector of the lines / hyperplanes.

$w'_0 x'_0 / w_0 x_0$  stands for the intercept of PLA,  $w_0 x_0 = w'_0 x'_0$

Thus,  $W_{\text{PLA}}$  and  $W'_{\text{PLA}}$  refer to the same lines / hyperplanes.

Although  $W_{\text{PLA}(0)} \neq W'_{\text{PLA}(0)}$ .

解答



$w_0 = -w'_0$   
 $w_f^T x_n = w_f^T x'_n \rightarrow$  Both refers to the same ideal PLA.

7. (20 points) Dr. Norman thinks PLA will be highly influenced by very long examples, as  $w_t$  changes drastically if  $\|x_n\|$  is large. Hence, ze decides to preprocess the training data by normalizing each input vector i.e.,  $\tilde{x}_n = \frac{x_n}{\|x_n\|}$ . What is PLA's upper bound on page 16/22 of Lecture 1 change with this preprocessing procedure in terms of  $\rho_x = \min_n \frac{y_n w_f^T x_n}{\|w_f\|}$ . Prove your answer.

$$\begin{aligned} \|W_{t+1}\|^2 &= \|W_t + y_{n(t)} x_{n(t)}\|^2 \leq \|W_t\|^2 + \max_n \|x_n\|^2 \rightarrow \|W_t\|^2 \leq \|W_{t+1}\|^2 + 1 = T \\ W_t^T W_t &\geq W_t^T W_t + \min_n y_n w_f^T x_n \rightarrow \frac{W_t^T W_t}{\|W_t\|} \geq \min_n \frac{y_n w_f^T x_n}{\|w_f\|} \cdot T \\ \sqrt{T} \cdot C &\leq \frac{W_t^T \cdot W_t}{\|W_t\|^2 + \|W_t\|} = p \cdot T \cdot \frac{1}{\sqrt{T}} = p \cdot \sqrt{T} \cdot C = p \\ \therefore \sqrt{T} \cdot C &\leq 1, \quad T \leq \frac{1}{C^2} = \frac{1}{p^2} \end{aligned}$$

8. (20 points) In PLA, we defined  $\rho = \min_n y_n w_f^T x_n$ , and  $\rho$  is often called the unnormalized **margin** of  $w_f$ . Margin is related to the minimal distance between  $x_n$  and hyperplane  $w_f$ , and will be a main concept when we introduce the Support Vector Machine (SVM) later in this class. Before that, let us play with a variant of PLA, the Perceptron Algorithm using Margins (PAM). The difference between PAM and the original PLA is that PAM updates on  $x_n$  if

$$y_n w_f^T x_n \leq \tau, < \rho$$

where  $\tau \geq 0$  is the margin we would like to achieve.

For any prove  $\tau$ , assume that the data set  $\{(x_n, y_n)\}_{n=1}^N$  is linearly  $\tau$ -separable. That is,  $w_f$  satisfies  $\rho > \tau$ . Prove that PAM halts in a finite number of steps (Hint: PLA is just a special case with  $\tau = 0$ ).

$$\begin{aligned} \|W_{t+1}\|^2 &\leq \|W_t + y_{n(t)} X_{n(t)}\|^2 \\ &= \|W_t\|^2 + 2y_{n(t)} W_t^T X_{n(t)} + \|y_{n(t)} X_{n(t)}\|^2 \\ &\leq \|W_t\|^2 + 2\tau + \|y_{n(t)} X_{n(t)}\|^2 \\ W_f^T W_{t+1} &\geq W_f^T (W_t + y_{n(t)} X_{n(t)}) \\ &= W_f^T W_t + y_{n(t)} X_{n(t)} W_f^T \geq W_f^T W_t + \min_n y_n w_f^T x_n W_f^T > W_f^T W_t + \tau \\ \|W_t\|^2 &\leq \|W_{t+1}\|^2 + 2\tau + R^2 \\ \|W_t\|^2 &\leq \|W_{t+1}\|^2 + 2\tau + R^2 \\ \Rightarrow \|W_t\|^2 &\leq T(\tau + R^2) \quad W_f^T W_t > W_f^T W_{t-1} + \tau \\ \Rightarrow \|W_t\| &\leq \sqrt{T(\tau + R^2)} \quad \Rightarrow W_f^T W_t > T\tau \\ \therefore \|W_t\| &\leq \sqrt{T} \cdot \sqrt{\tau + R^2} \end{aligned}$$

$$\begin{aligned} 1 &\geq \frac{W_f^T \cdot W_t}{\|W_f\| \cdot \|W_t\|} \geq \sqrt{T} \cdot C \\ \sqrt{T} \cdot \sqrt{\tau + R^2} &\leq \\ 1 &\geq \frac{\tau}{\sqrt{\tau + R^2}} \cdot \sqrt{T} \\ \sqrt{T} &\leq \frac{\sqrt{\tau + R^2}}{\tau} \\ T &\leq \frac{\tau + R^2}{\tau^2} \end{aligned}$$

## Experiments with Perceptron Learning Algorithm

Next, we use an artificial data set to study PLA. The data set with  $N = 256$  examples is in

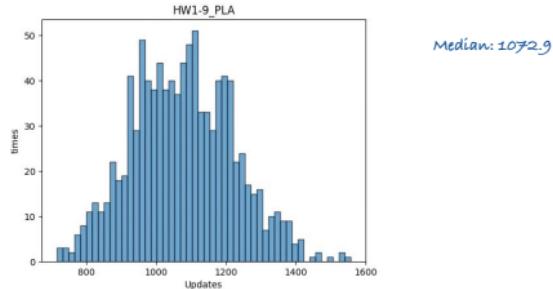
[http://www.csie.ntu.edu.tw/~htlin/course/ml23sfall/hw1/hw1\\_train.dat](http://www.csie.ntu.edu.tw/~htlin/course/ml23sfall/hw1/hw1_train.dat)

Each line of the data set contains one  $(\mathbf{x}_n, y_n)$  with  $\mathbf{x}_n \in \mathbb{R}^{12}$ . The first 10 numbers of the line contains the components of  $\mathbf{x}_n$  orderly, the last number is  $y_n$ . Please initialize your algorithm with  $\mathbf{w} = \mathbf{0}$  and take  $\text{sign}(0)$  as  $-1$ .

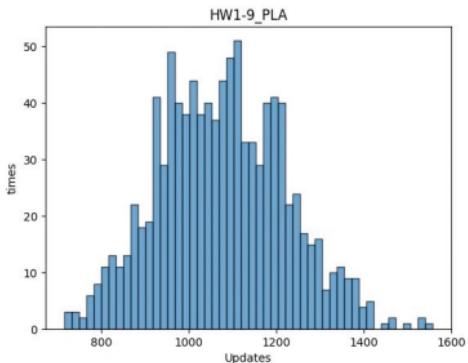
9. (20 points, \*) Please first follow page 4/22 of Lecture 2, and add  $x_0 = 1$  to every  $\mathbf{x}_n$ . Implement a version of PLA that randomly picks an example  $(\mathbf{x}_n, y_n)$  in every iteration, and updates  $\mathbf{w}_t$  if and only if  $\mathbf{w}_t$  is incorrect on the example. Note that the random picking can be simply implemented with *replacement*—that is, the same example can be picked multiple times, even consecutively. Stop updating and return  $\mathbf{w}_t$  as  $\mathbf{w}_{\text{PLA}}$  if  $\mathbf{w}_t$  is correct consecutively after checking  $5N$  randomly-picked examples.

*Hint: (1) The update procedure described above is equivalent to the procedure of gathering all the incorrect examples first and then randomly picking an example among the incorrect ones. But the description above is usually much easier to implement. (2) The stopping criterion above is a randomized, more efficient implementation of checking whether  $\mathbf{w}_t$  makes no mistakes on the data set.*

Repeat your experiment for 1000 times, each with a different random seed. Plot a histogram to visualize the distribution of the number of updates needed before returning  $\mathbf{w}_{\text{PLA}}$ . What is the median number of updates?

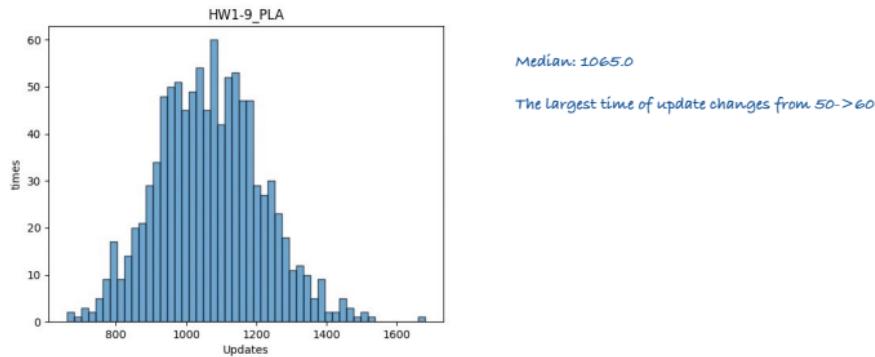


10. (20 points, \*) Scale up each  $\mathbf{x}_n$  by 11.26, including scaling each  $x_0$  from 1 to 11.26. Then, run PLA on the scaled examples for 1000 experiments, each with a different random seed. Plot a histogram to visualize the distribution of the number of updates needed before returning  $\mathbf{w}_{\text{PLA}}$ . What is the median number of updates? Compare your result to that of Problem 9. Describe your findings.

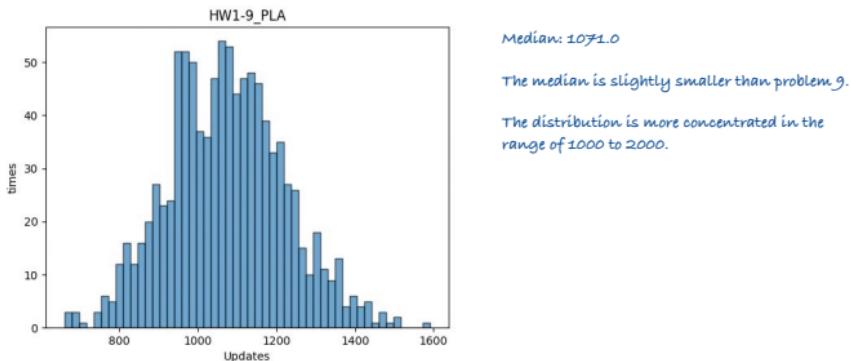


Median: 1072.5  
The whole distribution is the same as problem 9.

11. (20 points, \*) Set  $x_0 = 11.26$  to every  $\mathbf{x}_n$  instead of  $x_0 = 1$ , and do not do any scaling. Repeat the 1000 experiments above. Plot a histogram to visualize the distribution of the number of updates needed before returning  $\mathbf{w}_{\text{PLA}}$ . What is the median number of updates? Compare your result to that of Problem 9. Describe your findings.



12. (20 points, \*) Set  $x_0 = 1$  to every  $\mathbf{x}_n$ , and do not do any scaling. Modify your PLA above to a variant that keeps correcting the same example until it is perfectly classified. That is, when selecting an incorrect example  $(\mathbf{x}_{n(t)}, y_{n(t)})$  for updating, the algorithm keeps using that example (that is,  $n(t+1) = n(t)$ ) to update until the weight vector perfectly classifies the example (and each update counts!). Repeat the 1000 experiments above. Plot a histogram to visualize the distribution of the number of updates needed before returning  $\mathbf{w}_{\text{PLA}}$ . What is the median number of updates? Compare your result to that of Problem 9. Describe your findings.



### Bonus: Does Normalization Help PLA?

13. (Bonus 20 points) In Problem 7, you showed that the upper bound of PLA changes when running on the normalized data instead of the original data. Does it mean that normalization can speed up PLA? Why or why not? Note that this is a bonus problem and TAs can set a high standard on your logical arguments when deciding whether to give you the points.

Yes  
In problem 7, the result shows that the maximum update time  $T$  change from  $(R/p)^2$  to  $1/(p^2)$   
 $R$  refers to  $\max \|x\|$ , while  $p$  indicates the minimum distance from data to PLA

After we normalized the data,  $R$  will reduce to 1, however,  $p$  would also reduce from  $p$  to  $p'$ .

The maximum update time depends on the ratio of  $R$  and  $p$ .

If the data itself were scattered,  $p > p'$ ,  $T' < T$ , normalization did speed up.  
If were concentrated, normalization may not improve PLA. (but won't be worse than origin)

We take four situation for example.

In conclusion, though the maximum update time reduce, it doesn't mean it will definitely improve.  
After normalization, update time won't be larger, but whether it will be smaller depends on data.

