

Machine Learning HW1

B11901174 傅啟恩

1. One of the applications for self-supervised learning is “Speech-to-Speech Translation”, which use speech waveforms as input data (even without written form) and output as waveforms in another language. It utilizes SSL to train upstream model so that it can deal with the downstream tasks.

To be specific, the upstream model is trained through self-supervised learning, which means it leverage large amounts of unlabeled speech data to train models, and it can be fine-tuned on downstream tasks related to speech translation. For example, it can be further trained to perform tasks like speech recognition, machine translation, and speech synthesis.

2. Yes, I agree with the answer. Reinforcement learning can be used to deal with the shortest path in a maze problem by choosing a person/player as the agent, and the maze as the unchanged environment and properly setting rewards such as the fewer returning times, the higher the rewards, or the shortest the path, the higher the rewards.

For example, we can define a set of actions that the agent can take, such as moving up, down, left, or right. The agent starts at a specific position in the maze and explores the environment by taking different actions.

By using reinforcement learning algorithms, the agent can learn through trial and error to navigate the maze and optimize its actions to maximize the cumulative rewards. Over time, the agent can develop a policy that enables it to consistently find the shortest path to the exit.

3. It seems like ChatGPT does not disapprove any possibilities in this problem. I think the off-the-shelf algorithm, which is proved to be the fastest or most optimal, cannot be improved by machine learning. However, for problems that have not been proven yet, such as sorting problems, algorithms with fewer instructions than human benchmarks can be discovered using AlphaDev.

According to this article, AlphaDev has the potential to find algorithms that outperform human benchmarks and provide more efficient solutions. By utilizing AlphaDev and other similar machine learning techniques, we can push the boundaries of what is considered possible in algorithm design. With the

continuous advancements in machine learning, we can uncover innovative solutions that were previously unimaginable.

So, while the off-the-shelf algorithms may serve their purpose in many cases, I think AlphaDev and similar tools can have advancements in various domains, including sorting problems and beyond.

4.

4. Suppose $\hat{w}_{PNA} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$.

If "mistake" $\Leftrightarrow \text{sign}(w_t^T x_{n(t)}) \neq y_{n(t)}$,

$\hat{w}_{t+1} = \hat{w}_t + y_{n(t)} \begin{pmatrix} 1 \\ x_{n(t)} \end{pmatrix} \quad x_0 = 1$

and -threshold $= w_0 = 0$

\Rightarrow with T_+ mistakes as $y_{n(t)} = +1$,
 T_- " " " " $= -1$

$w_0 = 1 \times (T_+) \times (+1) + 1 \times (T_-) \times (-1) = (T_+) - (T_-)$

5.

5. As we proved in the course, the upperbound of the total steps

to $w_{PNA}(T)$ is bounded by $\frac{(\max_n \|x_n\|^2)}{(\min_n y_n \frac{w_f^T}{\|w_f\|} x_n)}^2$, and $X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}_{(d+1) \times 1}$

and composed of $\{0, 1\}$, which means $X = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 1 \end{bmatrix}$ at most has $m \times 1$

because in email x it contains at most m distinct words

$\Rightarrow \max_n \|x_n\|^2 \leq (m+1)^2 \quad x_0 = 1$
 $f(n) = \text{Sign}(\underbrace{z_+(x)}_{\text{spam-like}} - \underbrace{z_-(x)}_{\text{not spam-like}} - 0.5)$ ①

In each step we want to check $\text{Sign}(w_t^T x_t) \stackrel{?}{=} y_t$,

To approximate f , $w_f^T = [-0.5, \pm 1, \pm 1, \pm 1 \dots \pm 1]_{(d+1)}$
 ± 1 : spam-like, -1 : not spam-like

$\Rightarrow (\min_n y_n \frac{w_f^T}{\|w_f\|} x_n)$ happens if $(w_f$: final perfect w)
predict correctly

$$1. \quad z_+(x) = z_-(x) \overset{-0.5}{}, \text{sign}(w_f^T x_n) = y_n = -1$$

$$2. \quad z_+(x) = z_-(x) + 1 \overset{+0.5}{}, \text{sign}(w_f^T x_n) = y_n = +1$$

$$\text{And } \left(\min_n y_n w_f^T x_n \right)^2 = (0.5)^2 = 0.25$$

$$\|w_f\| = \sqrt{(0.5)^2 + d \times 1^2} = \sqrt{0.25 + d}$$

$$\Rightarrow \left(\min_n y_n \frac{w_f^T}{\|w_f\|} x_n \right)^2 = \frac{0.25}{0.25 + d} \quad \text{--- (2)}$$

$$\text{By (1), (2), (the upperbound of } T) = \frac{R^2}{\rho^2}$$

$$= \frac{\max_n \|x_n\|^2}{\left(\min_n y_n \frac{w_f^T}{\|w_f\|} x_n \right)^2} = \frac{(m+1)}{\frac{0.25}{0.25 + d}} = (4d+1)(m+1) \quad \#$$

$\Rightarrow (4d+1)(m+1)$ bounds the number of mistakes that PLA can make spam classification.

6.

$$b. \quad w_0 = 0 \rightarrow \boxed{X_n = \begin{pmatrix} 1 \\ \vdots \\ x_{n,orig} \end{pmatrix}, \quad X_n' = \begin{pmatrix} -1 \\ \vdots \\ x_{n,orig} \end{pmatrix}} \rightarrow w_{PLA} \rightarrow w_{PLA'}$$

with same $n(t)$, $t=0, 1, 2, \dots$,

Each updates,

$$w_1 = w_0 + y_{n(0)} \frac{y_{n(0)} w_0^T x_{n(0)}}{\|w_0\|^2} \quad , \quad w_1' = w_0 + y_{n(0)} \begin{pmatrix} -1 \\ \vdots \\ x_{n(0),orig} \end{pmatrix}$$

$$w_2 = w_1 + y_{n(1)} \begin{pmatrix} 1 \\ \vdots \\ x_{n(1),orig} \end{pmatrix} \quad w_2' = w_1' + y_{n(1)} \begin{pmatrix} -1 \\ \vdots \\ x_{n(1),orig} \end{pmatrix}$$

$$= \begin{pmatrix} y_{n(0)} + y_{n(1)} \\ y_{n(0)} x_{n(0),orig} + y_{n(1)} x_{n(1),orig} \\ \vdots \end{pmatrix} = \begin{pmatrix} -y_{n(0)} - y_{n(1)} \\ y_{n(0)} x_{n(0),orig} + y_{n(1)} x_{n(1),orig} \\ \vdots \end{pmatrix}$$

$$w_{PLA} = w_T + y_{n(T)} \begin{pmatrix} 1 \\ \vdots \\ x_{n(T),orig} \end{pmatrix} \quad w_{PLA'} = w_T' + y_{n(T)} \begin{pmatrix} -1 \\ \vdots \\ x_{n(T),orig} \end{pmatrix}$$

$$= \begin{pmatrix} \sum_i y_{n(i)} \\ \sum_i y_{n(i)} x_{n(i),orig} \\ \vdots \end{pmatrix} = \begin{pmatrix} -\sum_i y_{n(i)} \\ \sum_i y_{n(i)} x_{n(i),orig} \\ \vdots \end{pmatrix}$$

And for every t , $y_{n(t)} w_t^T x_{n(t)} = y_{n(t)} w_t'^T x_{n'(t)}$!

\Rightarrow w_{PLA} is not necessarily equal to $w_{PLA'}$ unless $\sum_i y_{n(i)} = 0$ ✗

But without considering $w_{PLA}[0]$, $w_{PLA'}[0]$, They will both

return the same $\text{sign}(w_{PLA} x_n)$ and $\text{sign}(w_{PLA'} x_n')$, which

$$\text{means } \text{sign} \left(\begin{pmatrix} \sum_i y_{n(i)} \\ \sum_i y_{n(i)} x_{n(i),orig} \end{pmatrix} \begin{pmatrix} 1 \\ x_{n(T),orig} \end{pmatrix} \right)$$

$$= \text{sign} \left(\begin{pmatrix} -\sum_i y_{n(i)} \\ \sum_i y_{n(i)} x_{n(i),orig} \end{pmatrix} \begin{pmatrix} -1 \\ x_{n'(T),orig} \end{pmatrix} \right) \quad \neq$$

7.

7. Start from $w_0 = 0$

$$\textcircled{1} \|w_T\|^2 \leq \max_n \left\| \frac{x_n}{\|x_n\|} \right\|^2 \times T$$

$$\begin{aligned} \textcircled{2} w_f^T w_{t+1} &= w_f^T (w_t + y_n(t) \frac{x_n(t)}{\|x_n(t)\|}) \\ &\geq w_f^T w_t + \min_n y_n w_f^T x_n \end{aligned}$$

$$w_f^T w_T \geq T \cdot \|w_f\| \rho_z$$

$\textcircled{1} + \textcircled{2}$

$$1 \geq \frac{w_f^T}{\|w_f\|} \frac{w_T}{\|w_T\|} \geq T \rho_z \frac{1}{\sqrt{T}} = \sqrt{T} \rho_z$$

$$\sqrt{T} \rho_z \leq 1, \quad T \leq \frac{1}{\rho_z^2} \quad \#$$

8. τ -separable

\Leftrightarrow exist perfect w_f such that $y_n w_f^T x_n > \tau$

$$y_n(t) w_f^T x_n(t) > \min_n y_n w_f^T x_n(t) > \tau$$

$$\textcircled{1} w_f^T w_{t+1} = w_f^T (w_t + y_n(t) x_n(t)) > w_f^T w_t + \tau$$

$$\textcircled{2} \|w_{t+1}\|^2 \leq \|w_t\|^2 + 2\tau + \|y_n(t) x_n(t)\|^2$$

$$\leq \|w_t\|^2 + \underbrace{2\tau + \max_n \|y_n x_n\|^2}_{= R^2} = R^2$$

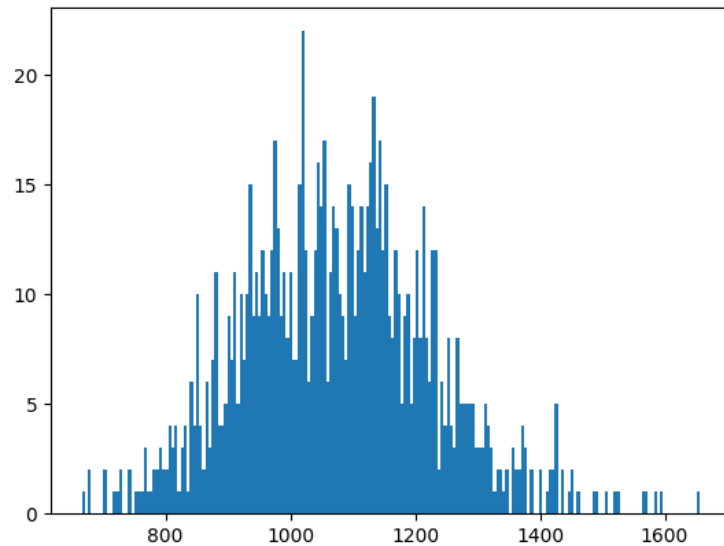
\Rightarrow like PLA, PAM has upper bound for updates time and "lines" are more and more aligned with w_f

\Rightarrow PAM will halt in finite number $T \leq \frac{R^2}{\tau^2}$ of steps $\#$

$$1 \geq \frac{w_f^T w_T}{\|w_f\| \|w_T\|} \geq \frac{T \cdot \tau}{\sqrt{T} \cdot R}$$

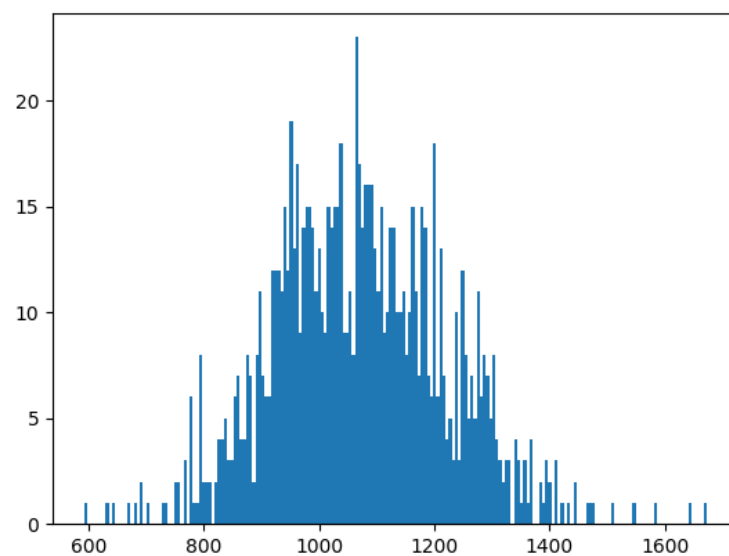
9. (Original)

Median number of updates: 1073.5 (updates)



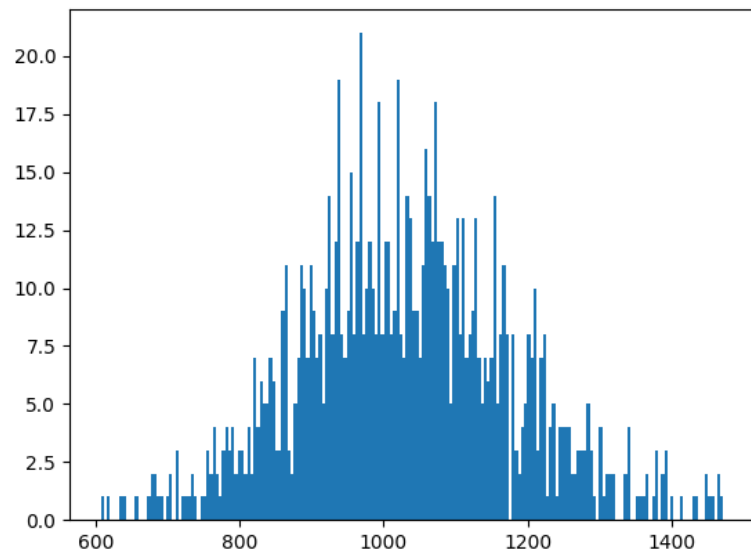
10. (Scaling by 11.26)

Median number of updates: 1067.0 (updates), compared to problem 9, the number of updates decreases a little bit or even the same, which means we need similar steps to get w_{PLA} .



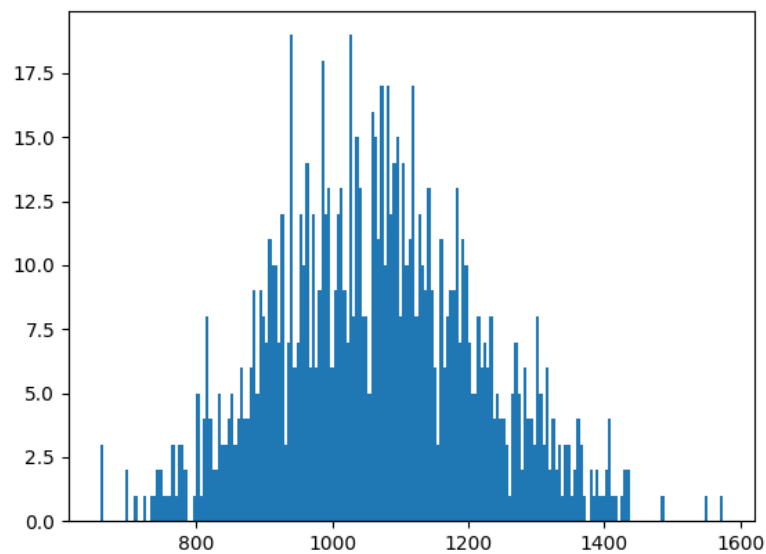
11. (Set $x_0 = 11.26$)

Median number of updates: 1024 (updates), compared to problem 9, the number of updates decreases about 5%, which means we need less steps to get w_{PLA} .



12. (Keep correcting the same example until it is perfectly classified)

Median number of updates: 1064.5 (updates), compared to problem 9, the number of updates decreases a little bit or even the same, which means we need similar steps to get w_{PLA} .



13.

13. Give one example that normalization actually speed up PLA:

$$\begin{aligned} x_1 &= (1, 50) & y_1 &= -1 & \text{and } w_0 &= (0, 0) \\ x_2 &= (1, 100) & y_2 &= +1 \end{aligned}$$

Go through PLA: (By code)

and choosing randomly

It needs ^{on} average 13825 updates to get w_{PLA}

$$w_{PLA} \approx (-302, 100)$$

While after normalization, and given $x_0 = 1$

$$\hat{x}_1 = \begin{matrix} \hat{x}_0 & x_1 & x_2 \end{matrix} (0.0199, 0.0199, 0.9996)$$

$$\hat{x}_2 = (0.0099, 0.0099, 0.9999)$$

It takes on average 2 updates to get

$$w_{PLA} \approx (-0.0099, 0.0003)$$

\Rightarrow In this case normalization seems speed up PLA

While for another case

$$\begin{aligned} x_1 &= (1, 1) & y_1 &= -1 & w_0 &= 0 \\ x_2 &= (1, -1) & y_2 &= +1 \end{aligned}$$

By running code (randomly choose x)

It takes 2 updates to get $W_{pin} = [0, -2]$

after normalization

$$\hat{x}_1 = [\overset{\hat{x}_0}{0.577}, 0.577, 0.577]$$

$$\hat{x}_2 = [0.577, 0.577, -1.577]$$

It also on average takes 2 steps to get W_{pin}

$$W_{pin} = [0, -1.577]$$

\Rightarrow Normalization cannot speed up more on simple data

In conclusion, normalization can speed up in some case, especially for complicated data \neq