

1. google 所提出的 BERT 就是一種自監督式 pre-train model，雖然可以有不同應用場景，但最常被應用在 NLP 模型上，特別是在訊號處理，語音辨識，或文字生成上，它所運用的技術主要是兩個層面上，一是做填空題，首先隨機的對一些 input 做 mask 處理，以文字為例便是隨機對個別某些字遮蔽(你也可以說是挖洞)，然後這些遮蔽字(token)透過 BERT 給定 embedding 與分類器做完線性變換後通過 softmax 輸出一個各個字分布機率的向量，再與 mask 前的 input 做 cross entropy 求最小值。另一則是預測一個句子是是另一個句子的下一句 next sentence prediction，將可能是兩個句子的文字組合，加入兩個 token SEP CLS 分別表示分隔與起點，最後透過 BERT 與分類器得出 YES 或 NO。
2. 部分同意，在走迷宮這類問題上，使用到相當多強化學習類型的模型，透過獎勵與懲罰機制來訓練 agent，舉例來說以 QLearning，模型內有一依據環境不斷更新的 Qtable 記錄了狀態 state 與動作 action 的對應表以及接下來各個決策預期所能得到的獎勵或懲罰，依照環境更新的刺激不斷地給予 agent 新的回饋，以此來訓練 agent 學習在當下的狀態做出何種動作的決策，這種類型的模型已經在走迷宮的研究行之有年，也多經常能解決設定好的迷宮，但我認為我們仍難以證明這種方式所找到的解必定會是最短路徑解，因此只能說部分同意。
3. 可能的，儘管 chatGPT 給了一個只能針對某些演算法或特定目標或特定資料最佳化，而且必須取決於問題的本質與演算法自身的回答。但我認為更精確地說只要是未被數學證明為最佳解的演算法，都仍有機會可以做進一步的最佳化，舉例來說圍棋作為人類最古老的遊戲，也能視為可以透過演算法處理的問題，惟其變化的可能性與計算量太高，人類有限的大腦必須透過固定的範式來對遊戲進行理解，這在圍棋術語中稱為定石，這些定石都是經過時間長河的棋士經驗與直覺累積形成的行之有年的思考模式代代傳授，我們數學知識體系與演算法也有異曲同工之妙，在那些相當依賴經驗與直覺的領域特別如此，充滿了人類思考模式的特質。然而 alphago 的誕生徹底的改寫了長久以來許多人類視之為圍棋常識的定石，正如同飛機不需要模仿鳥類振翅飛行，潛水艇不需要模仿魚類擺鰭游泳，人工智慧也不需要模仿人類思考，並且其上限也只受硬體與演算法限制，所以也許也很有機會能夠打破人類舊有的思考框架，拓展至更多人類思考的未竟之地。

4.  $w_0 = T_+ - T_-$

pt. Assume  $w_f^T$  is the solution, let  $w_0^T = 0$ , the zero-th component  $x_0 = 1$

$\forall x_i, i=1, 2, \dots, n, x_i \in D$ , known  $w_{t+1} \leftarrow w_t + \eta(x_t) x_{n(t)}$  where  $\forall t$

only update when  $\eta(x_t) w_t^T x_{n(t)} \leq 0$ , thus total number of "update" is  $T_+ + T_-$

$$\underbrace{w_f^T \leftarrow \eta_{f-1} x_{f-1} + w_{f-1}^T \quad \dots \quad w_i^T \leftarrow \eta_{i-1} x_{i-1} + w_{i-1}^T}_{(T_+ + T_-) \text{ times}}$$

where component  $w_0$  of  $w_f^T$ :  $w_0 = \underbrace{\eta_{f-1} x_{n(f-1)} + \dots + 0}_{\# \text{ of } \eta(x)=1: T_+, \# \text{ of } \eta(x)=-1: T_-}$ , where  $\eta(x) x_{n(t)} = \begin{cases} -x_{n(t)} & \text{if } \eta(x) < 0 \\ x_{n(t)} & \text{if } \eta(x) > 0 \end{cases}$

since  $x_0 = 1 \forall x_i \in D$ ,  $w_0 = \underbrace{1+1+1 \dots 1}_{\# \text{ of } 1: T_+} - \underbrace{(1+1+1 \dots 1)}_{\# \text{ of } -1: T_-} = T_+ - T_-$

5. pt.

Assume perfect solution is  $w_f^T$ , known  $w_f^T w_{t+1} \geq w_f^T w_t + \min_n \eta_n w_f^T x_n$   
 $\geq w_f^T w_{t-1} + 2(\min_n \eta_n w_f^T x_n) \Rightarrow w_f^T w_T \geq w_f^T w_0 + (T) (\min_n \eta_n w_f^T x_n)$   
 Let  $w_0 = 0$ ,  $w_f^T w_T \geq T \cdot (\min_n \eta_n w_f^T x_n)$  where  $T$  is total # of "mistake"

And we know  $\|w_{t+1}\|^2 \leq \|w_t\|^2 + \max_n \|x_n\|^2 \Rightarrow \|w_f\| \leq \|w_0\| + \sqrt{T \max_n \|x_n\|^2}$

Thus  $\|w_t\| \leq \sqrt{T} (\max_n \|x_n\|)$

$$\therefore 1 \geq \frac{w_f^T w_T}{\|w_f\| \|w_T\|} \geq \frac{\min_n \eta_n w_f^T x_n \cdot T}{\|w_f\| \cdot \max_n \|x_n\| \cdot \sqrt{T}}$$

$\therefore \frac{w_f^T}{\|w_f\|}, \frac{w_T}{\|w_T\|}$  are unit vector,  $\frac{w_f^T w_T}{\|w_f\| \|w_T\|} \leq 1$

$$\therefore 1 \geq \frac{T \cdot \min_n y_n W_f^T X_n}{\sqrt{T} \cdot \max_n \|X_n\| \cdot \|W_f\|}, \quad T \leq \frac{(\max_n \|X_n\| \cdot \|W_f\|)^2}{(\min_n y_n W_f^T X_n)^2}$$

Since  $\max_n \|X_n\|$  is maximum Euclidean distance of all data from  $f$  there are  $m$  of  $d$  distinct words from  $X_n$  at most.

$$X_n = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ \vdots \end{bmatrix} \begin{matrix} \text{d-dimension} \\ \text{number of 1} = m \end{matrix} \quad \max_n \|X_n\| = \sqrt{\underbrace{1^2 + 1^2 + \dots + 1^2}_m} = \sqrt{1 + m}$$

and  $(\min_n y_n W_f^T X_n)^2$  is "closest" data to  $f$  which mean  $W_f$  and

$X_n$  have minimum inner product, since that  $Z_+(X) = Z_-(X)$ , where  $-0.5$

be  $W_0$  (threshold) of  $W_f$   $\therefore W_f^T X_n = -0.5$  where  $Z_+(X_n) = Z_-(X_n)$

As we know  $W_{t+1} \leftarrow W_t + y_t X_t$ , where  $y_t$  is sign and  $X_t$  is binary numbers 0 or 1, for every component  $W_i$  of  $W_f$ :

$$\cancel{W_i} \quad W_i \in \mathbb{Z} \quad \forall i=1, 2, \dots, n \quad \text{and} \quad W_0 = -0.5$$

$$\therefore \|W_f\| = \sqrt{\underbrace{(-0.5)^2 + (1)^2 + \dots + (-1)^2}_d} = \sqrt{\frac{1}{4} + d}$$

$d: \text{var-1 for minimize Euclidean distance}$



$$\therefore (\min_n y_n W_f^T X_n)^2 = \frac{1}{4} \quad \therefore T \leq \frac{(\max_n \|X_n\| \cdot \|W_f\|)^2}{(\min_n y_n W_f^T X_n)^2} = \frac{(m+1) \cdot (\frac{1}{4} + d)}{\frac{1}{4}}$$

$$= (4d+1)(m+1)$$

6. pt.

Let  $\{(X_n, y_n)\}_{n=1}^N$  where  $\underline{x_0=1}$  forming  $X_n = (1, x_n^{\text{orig}})$  and

$\{(X'_n, y_n)\}_{n=1}^N$  where  $\underline{x'_0=-1}$  forming  $X'_n = (-1, x_n^{\text{orig}})$  have

same particular sequence  $n(t)$ ,  $t=1, 2, \dots$  and both with  $W_0=0$

Assume  $\cancel{W_{PLA}}$  and  $\cancel{W'_{PLA}}$  are solution of  $\cancel{DPLA} \{(X_n, y_n)\}_{n=1}^N$  and

$$W_1 = W_0 + y_n(0) \begin{bmatrix} 1 \\ x_{n(0)}^{\text{orig}} \end{bmatrix} \text{ base of } x_i, i=1,2,\dots, \quad W_2 = W_0 + \begin{bmatrix} x_0 + x_{n(1)} \\ x_0 x_{n(0)}^{\text{orig}} + y_{n(1)} x_{n(1)}^{\text{orig}} \end{bmatrix}, \dots$$

$$\therefore W_{PLA} = W_T + y_{n(T)} \begin{bmatrix} 1 \\ x_{n(T)}^{\text{orig}} \end{bmatrix} = \begin{bmatrix} y_{n(T)} + \dots + y_{n(0)} \\ y_{n(T)} x_{n(T)}^{\text{orig}} + \dots + y_{n(0)} x_{n(0)}^{\text{orig}} \end{bmatrix}$$

$$\therefore X'_n = \begin{bmatrix} -1 \\ x_{n(0)}^{\text{orig}} \end{bmatrix} \therefore W'_{PLA} = \begin{bmatrix} -(y_{n(T)} + \dots + y_{n(0)}) \\ x_{n(T)} x_{n(T)}^{\text{orig}} + \dots + y_{n(0)} x_{n(0)}^{\text{orig}} \end{bmatrix}$$

Since  $y_{n(t)} W_t^T X_{n(t)} = y_{n(t)} W_t^{T'} X'_{n(t)}$   $\forall t=1, 2, \dots$

$$\text{hence } \text{sign}(y_{n(t)} W_{PLA}^T X_{n(t)}) = \text{sign}(y_{n(t)} W_{PLA}^{T'} X'_{n(t)})$$

$\therefore W_{PLA}$  and  $W'_{PLA}$  are equivalent.



7.  
 pf. Known original upper bound is  $\frac{\max_n \|X_n\|^2}{\left(\min_n y_n \frac{W_f^T X_n}{\|W_f\|}\right)^2} \geq T$  (before normalization)

Let  $W_0 = 0$ ,  $T$  is total number of mistake corrections

$$\|W_T\|^2 \leq \max_n \left\| \frac{X_n}{\|X_n\|} \right\|^2 \cdot T = \max_n \|Z_n\|^2 \cdot T$$

we know  $1 \geq \frac{W_f^T W_T}{\|W_f\| \|W_T\|}$ , ~~since~~ and  $W_f^T W_{t+1} = W_f^T (W_t + y_{n(t)} \frac{X_{n(t)}}{\|X_{n(t)}\|})$

$$\therefore W_f^T W_{t+1} \geq W_f^T W_t + \min_n y_n W_f^T Z_n > W_f^T W_t$$

$$\text{since } W_f^T W_T \geq T \cdot \|W_f\| \rho_2, \quad 1 \geq \frac{W_f^T W_T}{\|W_f\| \|W_T\|} \geq T \cdot \frac{1}{\sqrt{T}} \rho_2 = \sqrt{T} \rho_2$$

$$\therefore 1 \geq \sqrt{T} \rho_2 \quad \therefore T \leq \frac{1}{\rho_2^2} \quad \text{new upper bound: } \frac{1}{\rho_2^2}$$

8.

pf. Assume  $D: \{(X_n, y_n)\}_{n=1}^N$  is linearly  $\gamma$ -separable:

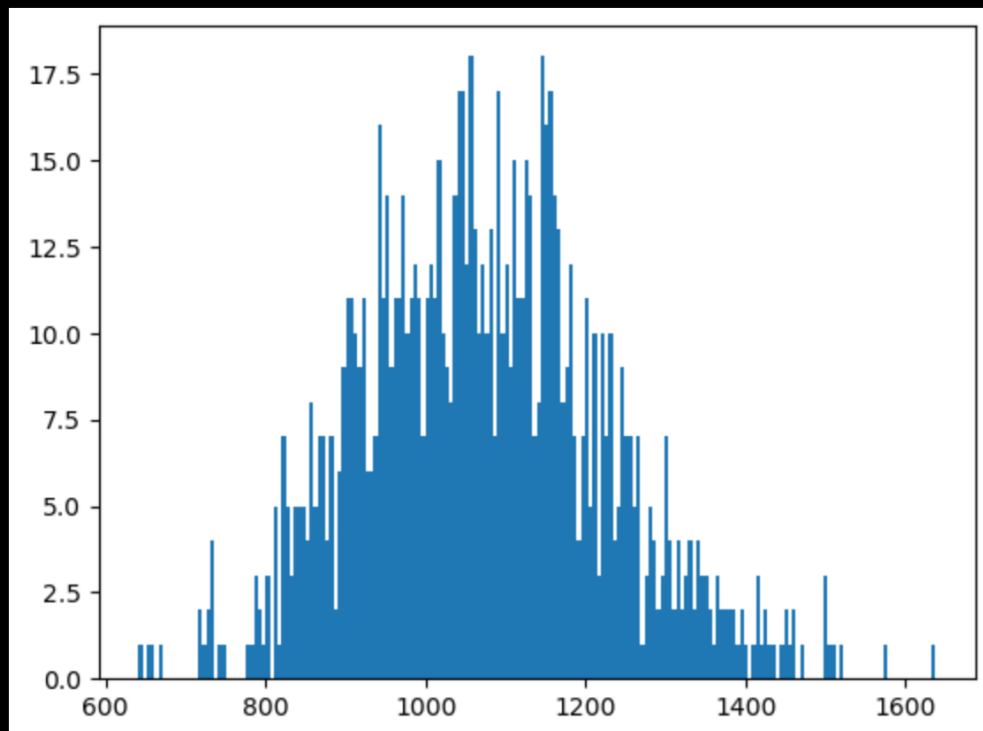
$\exists W_f$  s.t.  $\rho > \gamma$ , that is  $y_n \frac{W_f^T X_{n(t)}}{\|W_f\|} > \gamma$  (by considering PLA

is just a special case where  $\gamma=0$ ).  $\therefore W_f^T W_{t+1} > W_f^T W_t + \gamma$

$$\text{hence } \|W_{t+1}\|^2 \leq \|W_t\|^2 + 2\gamma + \|y_{n(t)} X_{n(t)}\|^2 \leq 2\gamma + \max_n \|X_n\|^2$$

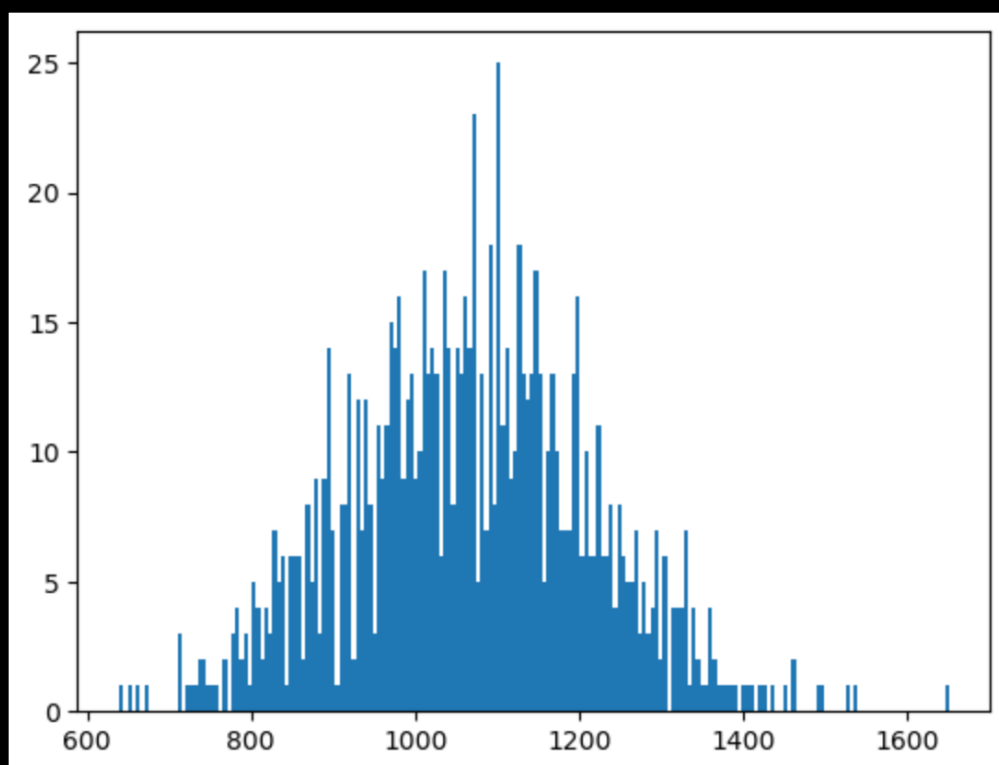
Define  $R^2 = 2\gamma + \max_n \|y_n X_n\|^2$   $\therefore$  PAM halts in  $T$  steps,  
 (considering  $\gamma > 0$ ) where  $T \leq \frac{R^2}{\gamma^2}$  #

9



Median number of updates is 1066.0

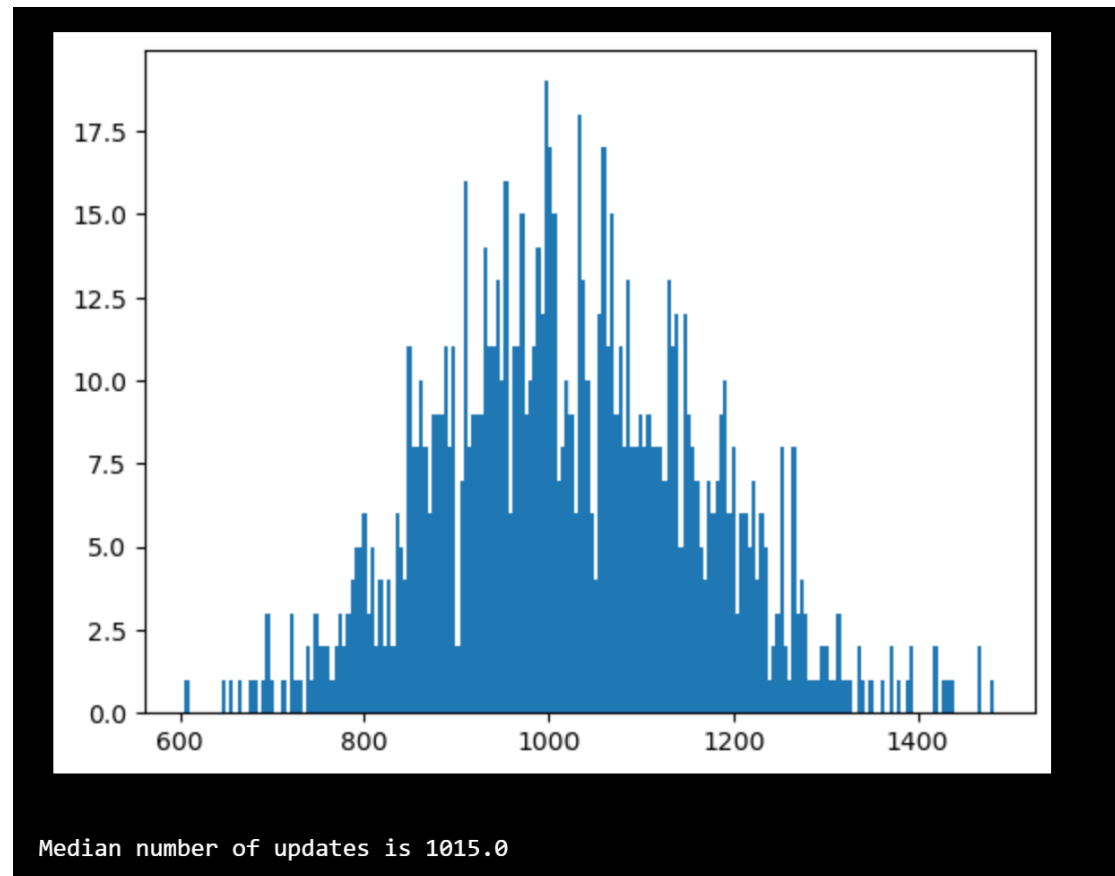
10



Median number of updates is 1068.0

中位數幾乎與 9. 沒有太多差異，可能是因為即使都乘 11.26，但相對尺寸仍差不多所經過的  $t$  差不多的關係

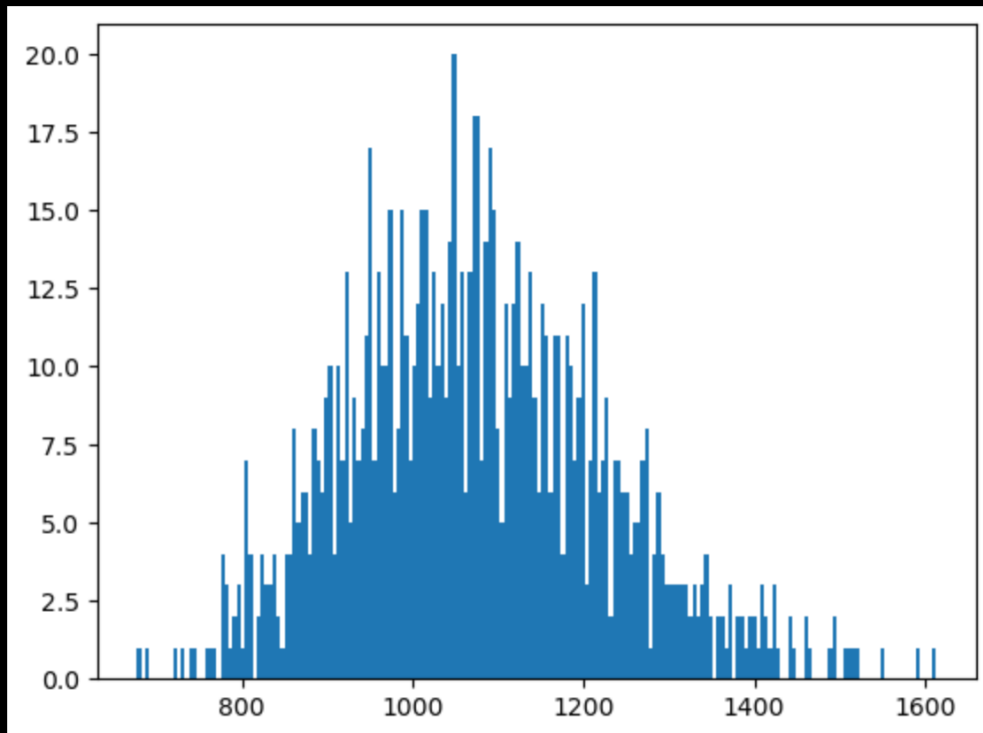
11



到這裡中位數就有很顯著的下降，這有可能是因為我們提高了  $w_0$ ，無意識中影響到它逼近的步驟了。

12

與 9. 10. 幾乎一樣，表示即使重複同樣的步驟，改變仍幾乎沒有。



Median number of updates is 1069.0

13

不一定 視情況而定