

NTUEE, WAN - CYUAN FAN 范萬泉

“no collaborators, with reference”

Homework 3

Problem 1 :

HW3 machine learning 范萬泉 BO4502/05

Decision tree.

1. Gini impurity : $1 - \sum_{k=1}^K M_k^2$, given $\sum_{k=1}^K M_k = 1$

$$= 1 - (M_1^2 + M_2^2 + M_3^2 + \dots + M_K^2)$$

求 Gini impurity $\geq \text{Max}$. 則利用 Cauchy-Schwarz Inequality

$$(M_1^2 + M_2^2 + \dots + M_K^2) \underbrace{\left(1^2 + 1^2 + \dots + 1^2 \right)}_{\text{K 1's}} \geq (M_1 + M_2 + \dots + M_K)^2$$

$$K(M_1^2 + M_2^2 + \dots + M_K^2) \geq (1)^2 = 1$$

則 $M_1^2 + M_2^2 + \dots + M_K^2$ 有 min = $\frac{1}{K}$

則 Gini impurity 有 Max = $1 - \frac{1}{K}$ *

Problem 2 :

2. Squared regression error

$$\begin{aligned}
 & \rightarrow M_+ (1 - (M_+ - M_-))^2 + M_- (1 - (M_+ - M_-))^2 \\
 & = M_+ (1 - M_+ + M_-)^2 + M_- (1 + M_+ - M_-)^2 \\
 & = M_+ (1 + M_+^2 + M_-^2 - 2M_+M_-) + M_- (1 + M_+^2 + M_-^2 - 2M_+M_-) \\
 & = \frac{1}{M_+ + M_-} + \frac{M_+^3 - M_+M_-^2}{M_+ + M_-} - \frac{2M_+^2 + 4M_+M_- - M_+^2M_-}{M_+ + M_-} \\
 & \quad + M_-^3 - 2M_-^2 \quad \text{given: } M_+ + M_- = 1 \\
 & = 1 + \frac{M_+ (M_+ - M_-) - M_- (M_+ - M_-) - 2(M_+ - M_-)^2}{M_+ + M_-} \\
 & = 1 + (M_+ - M_-)^2 - 2(M_+ - M_-)^2 \\
 & = 1 - (M_+ - M_-)^2 = 1 - (2M_+ - 1)^2 \\
 & = -4M_+^2 + 4M_+
 \end{aligned}$$

Gini impurity $\Rightarrow 1 - M_+^2 - M_-^2 \Rightarrow 1 - M_+^2 - (1 - M_+)^2$
 $= -2M_+^2 + 2M_+$

得證 $-4M_+^2 + 4M_+ = 2(-2M_+^2 + 2M_+)$

\Rightarrow Squared regression error = $2 \times (\text{Gini Impurity})$

得證 \Rightarrow 倍數關係

Problem 3:

Random Forest.

<ref: 講義 p.8/22 L10>

3. 挑出 $N - PN$ 個 samples

每個 sample 沒有被 sample 的機率為 $(1 - \frac{1}{N})^{PN}$ = 27%

$$\lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^{PN} = \lim_{N \rightarrow \infty} \left[\left(1 - \frac{1}{N}\right)^N\right]^P \text{ 根據 } e \text{ 的逼近}$$

$$\lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^N \approx e^{-1}$$

$$\text{因此 } \lim_{N \rightarrow \infty} \left[\left(1 - \frac{1}{N}\right)^N\right]^P \approx e^{-P}$$

沒 sample 的機率為 e^{-P}

$\therefore e^{-P}$ 將所有沒 samples 的點,

Problem 4:

4.

有 k 個 binary classification tree g_1, g_2, \dots, g_k .

若 G 輸出錯答案，其必至少有 $\frac{k+1}{2}$ 個 g_i 輸出錯答案。

$E_{\text{out}}(G)$ 的 worst case 為

(假設 sample 數量為 N 個)，即 $(\sum_{k=1}^k e_k)N$ 個 sample 全錯

$$G = \text{Uniform}(\{g_i\})$$

則 $E_{\text{out}}(G)$ 的 Upper bound 為

$$\frac{\left(\sum_{k=1}^k e_k\right)N}{\left(\frac{k+1}{2}\right)} \xrightarrow{\text{個點錯誤}} \text{至少 } \frac{k+1}{2} \text{ 個錯 } G \text{ 才全錯.}$$

$$\xrightarrow{\text{降 } N} E_{\text{out}}(G) = \frac{2}{k+1} \left(\sum_{k=1}^k e_k\right) \text{ 為 upper bound}$$

Problem 5:

Gradient Boosting.

5. (from 韓義 18/35 L11)

$$g_t = g_1, g_1(x) = 11.26,$$

$$S_{n \text{ initial}} = 0.$$

→ compute $\alpha_t = \text{OneVarLinearRegression}(\{g_t(x_n), y_n - S_n\})$

$$\min_{\alpha} \frac{1}{N} \sum_{n=1}^N \text{err}(S_{n \text{ initial}} + \alpha g_1(x_n), y_n), \text{err}(s, y) = (s - y)^2$$

$$\frac{\partial}{\partial \alpha} \left(\frac{1}{N} \sum_{n=1}^N (11.26 \cdot \alpha - y_n)^2 \right) = 0$$

$$\rightarrow \frac{1}{N} \sum_{n=1}^N 2 \cdot 11.26 \cdot \alpha - 2 \cdot (11.26 \alpha - y_n) = 0$$

$$\rightarrow 11.26 \cdot N \cdot \alpha - \sum_{n=1}^N y_n = 0, \quad \alpha = \frac{\sum_{n=1}^N y_n}{11.26 N}$$

$$\alpha_1 = \frac{\sum_{n=1}^N y_n}{11.26 N}$$

Problem 6 :

6.

經過 $t=2$ iteration, steepest η as at 當 η_s .

則 from problem 5,

$$\frac{\partial}{\partial \eta} \left(\frac{1}{N} \sum_{n=1}^N (S_n + \eta_s g_t(x_n) - y_n)^2 \right) = 0$$

$$\rightarrow \frac{1}{N} \sum_{n=1}^N 2 \cdot g_t(x_n) \cdot (S_n + \eta_s g_t(x_n) - y_n) = 0, \quad b = 0$$

$$\rightarrow \sum_{n=1}^N g_t(x_n) (S_n + \eta_s g_t(x_n)) = 0$$

$$\rightarrow \sum_{n=1}^N S_n g_t(x_n) = \sum_{n=1}^N y_n g_t(x_n) - \eta_s g_t^2(x_n)$$

Problem 7:

7. gradient boosting , squared-error polynomial regression

$$\text{polynomial } g_1(x) = b_1 + w_1 x + w_2 x^2 + \dots + w_K x^K$$

則 $b_1, w_1, w_2, w_3 \dots w_K$ 為最佳解，則必滿足

$$\frac{\partial}{\partial b_1} \left(\sum_{n=1}^N ((y_n - S_n) - g_1(x_n))^2 \right) = 0$$

$$\Rightarrow \frac{\partial}{\partial b_1} \left(\sum_{n=1}^N (y_n - S_n - (b_1 + w_1 x_n + w_2 x_n^2 + \dots + w_K x_n^K))^2 \right) = 0$$

$$\Rightarrow \sum_{n=1}^N (y_n - S_n - (b_1 + w_1 x_n + \dots + w_K x_n^K)) = 0$$

$$\Rightarrow b_1 = \frac{1}{N} \left(\sum_{n=1}^N (y_n - S_n) - \sum_{n=1}^N (w_1 x_n + w_2 x_n^2 + \dots + w_K x_n^K) \right) \quad \text{--- ①}$$

$$\frac{\partial}{\partial w_1} \left(\sum_{n=1}^N ((y_n - S_n) - g_1(x_n))^2 \right) = 0$$

$$\Rightarrow \frac{\partial}{\partial w_1} \left(\sum_{n=1}^N (y_n - S_n - (b_1 + w_1 x_n + w_2 x_n^2 + \dots + w_K x_n^K))^2 \right) = 0$$

$$\Rightarrow \sum_{n=1}^N (y_n - S_n - (b_1 + w_1 x_n + \dots + w_K x_n^K)) \cdot x_n = 0$$

$$\Rightarrow w_1 = \frac{\sum_{n=1}^N (y_n - S_n - (b_1 + w_2 x_n^2 + \dots + w_K x_n^K)) x_n}{\sum_{n=1}^N (x_n)^2}$$

⋮

同理 $w_2, w_3 \dots w_K$.

若 $\alpha \neq 1$ 時，則

$$\frac{\partial}{\partial b} \left(\sum_{n=1}^N (y_n - S_n - \alpha g_1(x_n))^2 \right) = 0 \quad \text{亦為 optimal solution.}$$

$$\frac{\partial}{\partial w_1} \left(\sum_{n=1}^N (y_n - S_n - \alpha g_1(x_n))^2 \right) = 0$$

則可求出。

$$\frac{\partial}{\partial b_1} \left(\sum_{n=1}^N (y_n - S_n - \alpha (b_1 + w_1 x_n^1 + \dots + w_k x_n^k))^2 \right) = 0$$

$$\Rightarrow \sum_{n=1}^N (y_n - S_n - \alpha (b_1 + w_1 x_n + \dots + w_k x_n^k)) \cdot \alpha = 0$$

$$\Rightarrow b_1 = \frac{1}{\alpha N} \sum_{n=1}^N (y_n - S_n - \alpha (w_1 x_n + \dots + w_k x_n^k))$$

亦為 optimal solution.

此時 $b_1 \neq$ ①式中的 b_1 ，但皆為 optimal solution.

i. 矛盾

因此，在 polynomial regression 替代 decision trees 時，

$\alpha < 1$ 成立。

Problem 8:

Neural Network.

8. OR($x_1, x_2 \dots x_d$)

代表其中只要有-個 true (+1) 則輸出(+1)

$$w_i = \begin{cases} +1 & , x_i = +1 \\ 0 & , x_i = -1 \end{cases}$$

→ 若用以上 w_i ，則所有 x_i 為 -1 時均為 0。

而只要有一個 x_i 為 +1， $\sum_{i=0}^d w_i x_i > 0$

$g_A(x)$ 則有 OR 的 function.

Problem 9:

9. generally ($1 \leq l < L$) (from L12 14/23)

$$\epsilon_n = (y_n - S_j^{(l)})^2 = (y_n - \sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)})^2$$

$$\frac{\partial \epsilon_n}{\partial w_{ij}^{(l)}} = \frac{\partial \epsilon_n}{\partial S_j^{(l)}} \cdot \frac{\partial S_j^{(l)}}{\partial w_{ij}^{(l)}} = \delta_j^{(l)} \cdot (x_i^{(l-1)})$$

$$\begin{aligned} * \delta_j^{(l)} &= \frac{\partial \epsilon_n}{\partial S_j^{(l)}} = \sum_{k=1}^{d^{(l+1)}} \frac{\partial \epsilon_n}{\partial \delta_k^{(l+1)}} \cdot \frac{\partial \delta_k^{(l+1)}}{\partial S_j^{(l)}} \cdot \boxed{\frac{\partial x_j^{(l)}}{\partial S_j^{(l)}}} \\ &= \sum_{k=1}^{d^{(l+1)}} (\delta_k^{(l+1)}) (w_{jk}^{(l)}) \cdot (\tanh'(\delta_j^{(l)})) \end{aligned}$$

<Note>

$$\tanh'(x) = \operatorname{sech}^2(x)$$

$$\operatorname{sech}^2(x) = 1 - \tanh^2(x)$$

$$\frac{\partial \epsilon_n}{\partial w_{ij}^{(l)}} = \sum_{k=1}^{d^{(l+1)}} [(\delta_k^{(l+1)}) (w_{jk}^{(l)}) (1 - \tanh^2(\delta_j^{(l)}))] (x_i^{(l-1)})$$

$\therefore w_{ij}^{(l)}$: all the initial weight are set to 0

$$\therefore \frac{\partial \epsilon_n}{\partial w_{ij}^{(l)}} = 0$$

for i, j, l , $\frac{\partial \epsilon_n}{\partial w_{ij}^{(l)}} = 0$
each

Problem 10 :

10.

$$X^{(L)} = q \quad , \quad q_k = X_k^{(L)} = \frac{\exp(S_k^{(L)})}{\sum_{i=1}^K \exp(S_i^{(L)})}$$

$$\frac{\partial e}{\partial S_k^{(L)}} = \frac{\partial e}{\partial q_k} \frac{\partial q_k}{\partial S_k^{(L)}} = -\frac{v_k}{q_k} \left(\frac{\partial q_k}{\partial S_k^{(L)}} \right)$$

$$\frac{\partial q_k}{\partial S_k^{(L)}} = \frac{\exp(S_k^{(L)}) \cdot \sum_{i=1}^K \exp(S_i^{(L)}) - \exp(S_k^{(L)}) \exp(S_k^{(L)})}{\left[\sum_{i=1}^K \exp(S_i^{(L)}) \right]^2}$$

< 說明 >

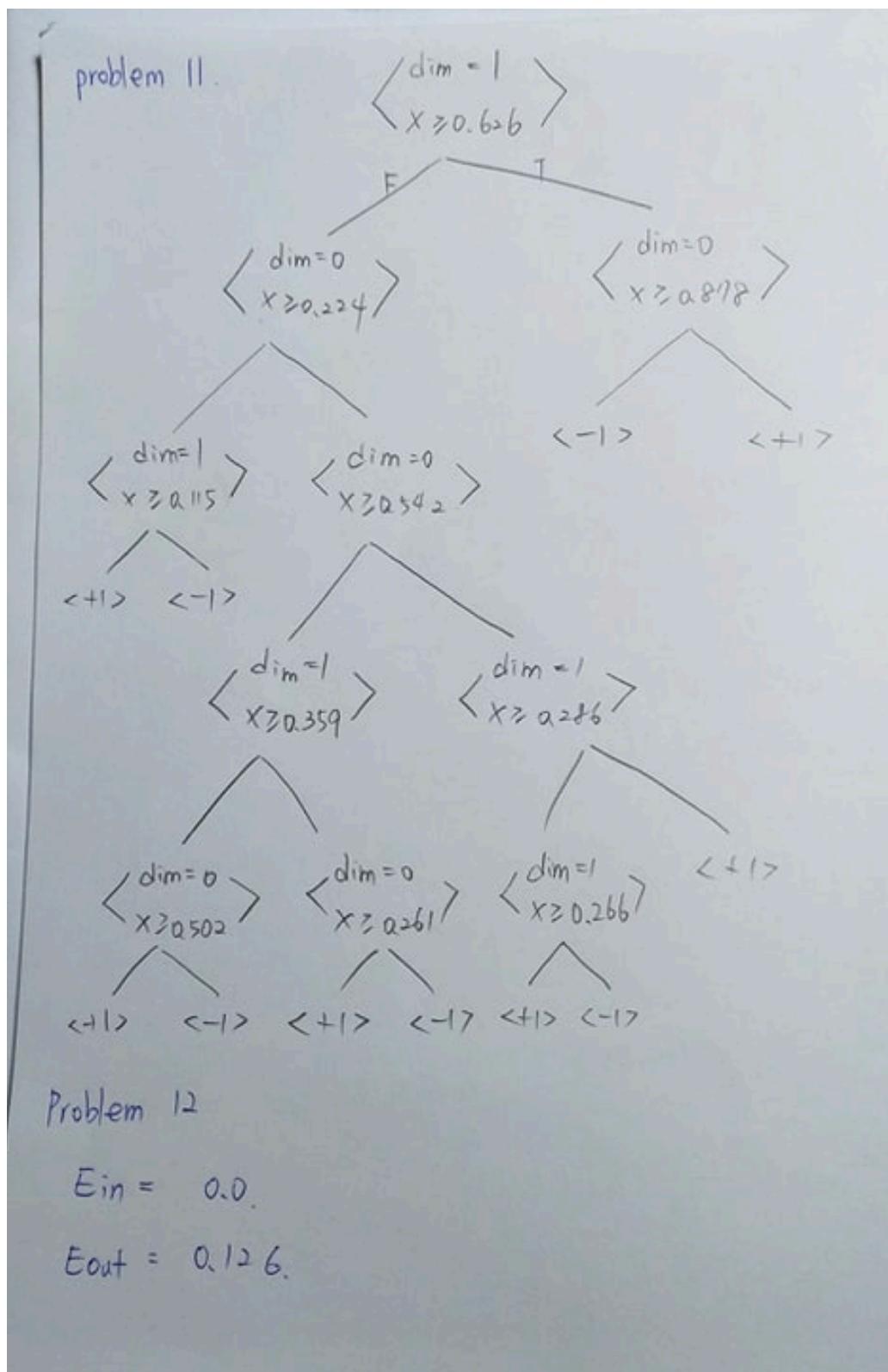
此處對 $S_k^{(L)}$ 微分，注意 $\sum_{i=1}^K \exp(S_i^{(L)})$ 其中有一項有 $S_k^{(L)}$
因此代入微分除法公式。

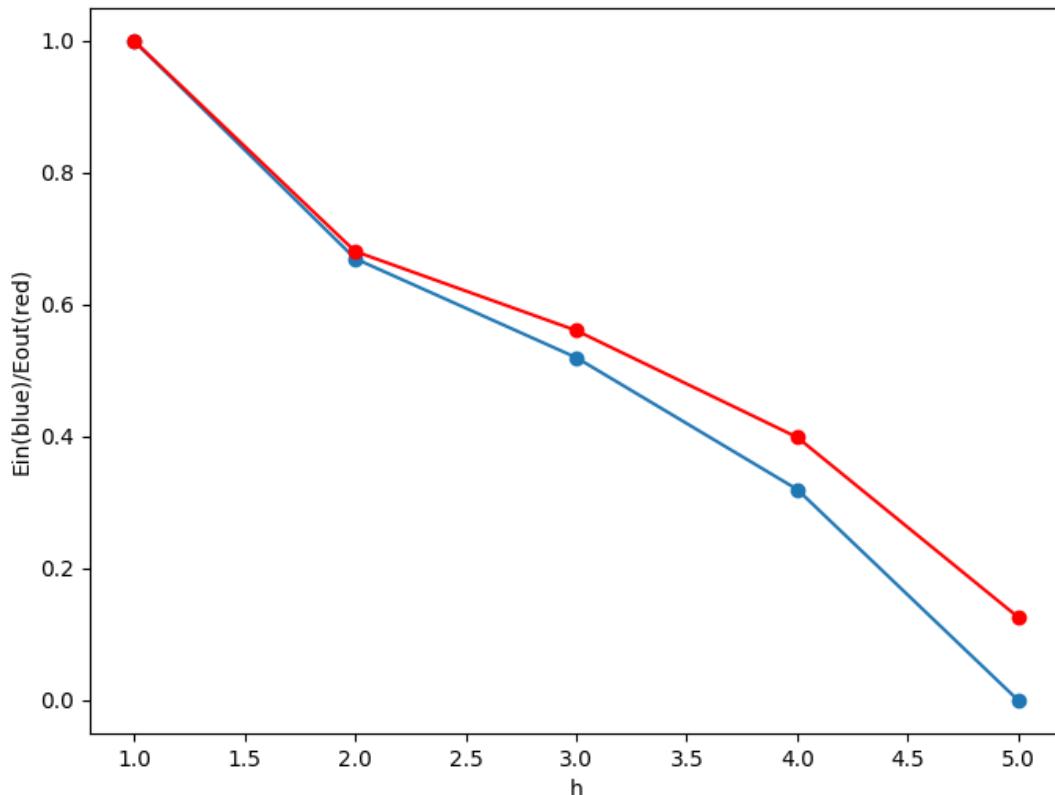
$$= \frac{\exp(S_k^{(L)}) \left[\sum_{i=1}^K \exp(S_i^{(L)}) - \exp(S_k^{(L)}) \right]}{\left[\sum_{i=1}^K \exp(S_i^{(L)}) \right]^2}$$

$$= q_k (1 - q_k) \quad v_k \text{ 值為 1.}$$

$$\frac{\partial e}{\partial S_k^{(L)}} = -\frac{v_k}{q_k} (q_k (1 - q_k))$$

$$= q_k - v_k \quad \text{得證}$$

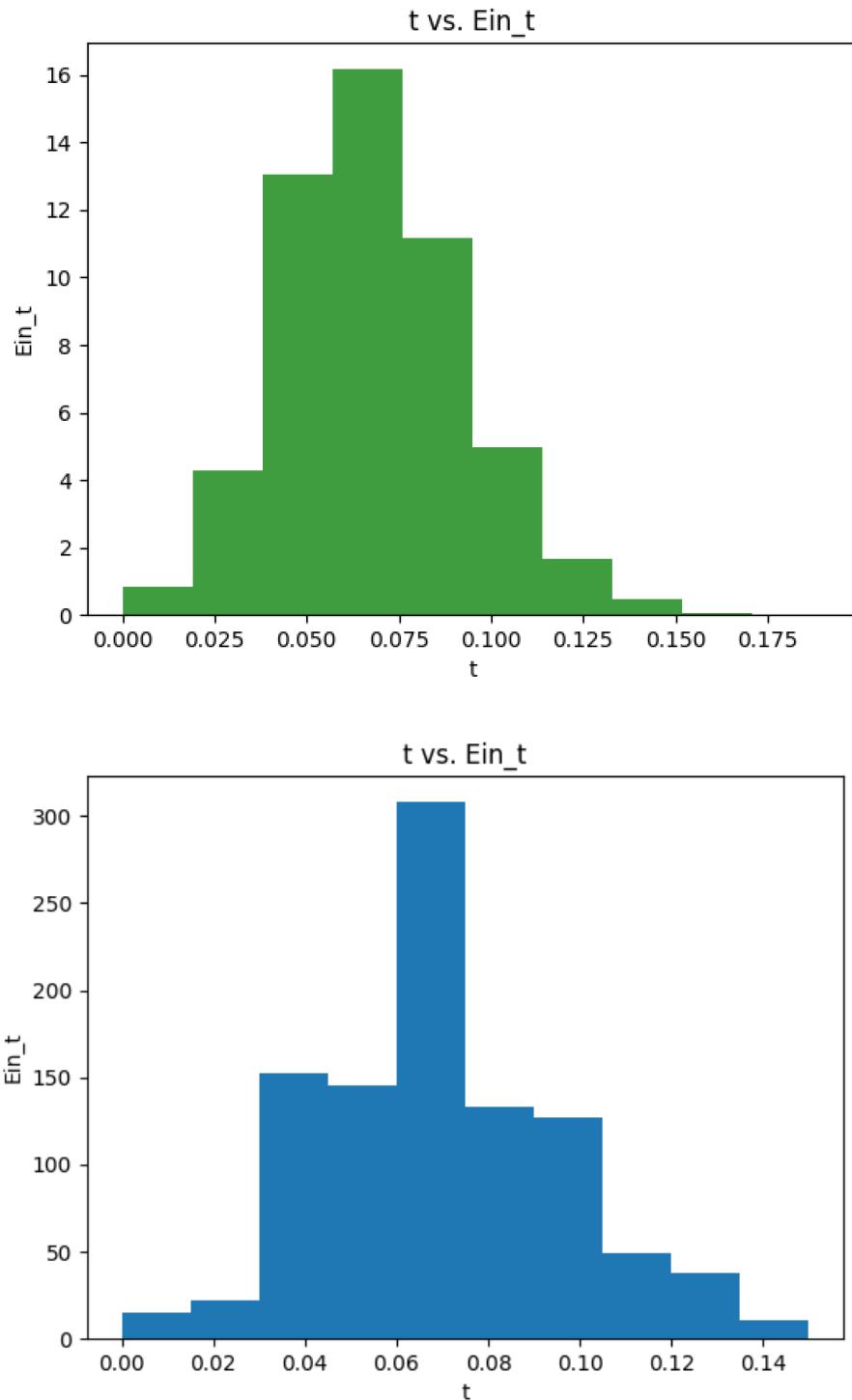
Problem 11,12:

Problem 13:

當我們增加tree的高度的時候，確實可以發現我們的正確率提高不少，而在 testing set 上頭做的效果比training data要來的錯誤率高一些。
總體樹高為5。

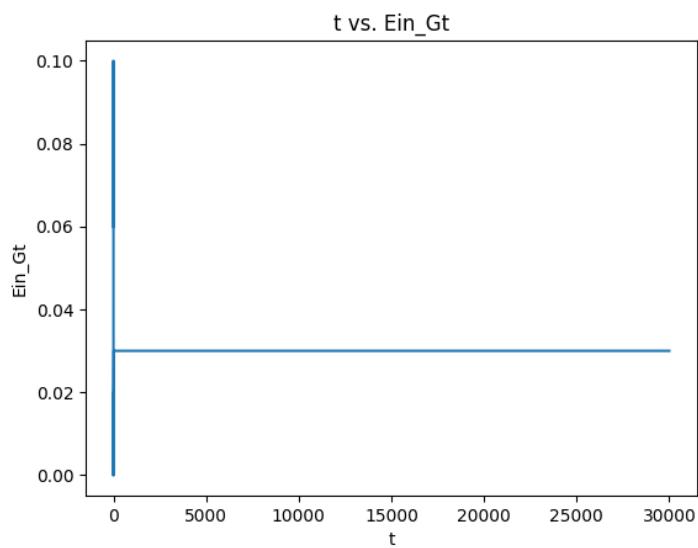
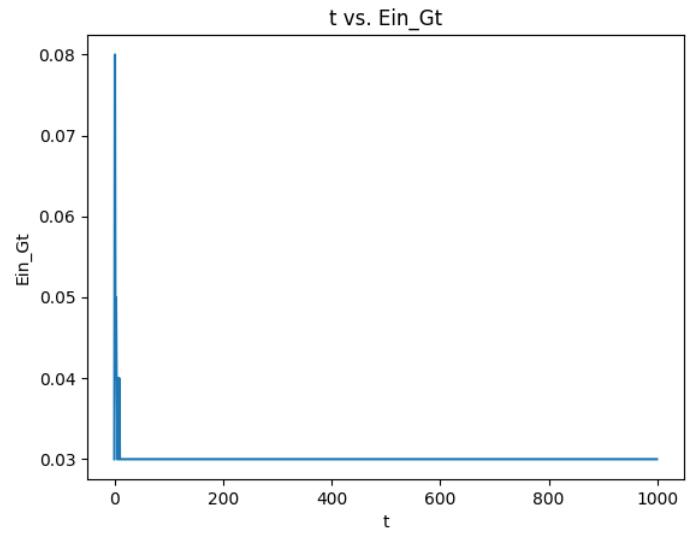
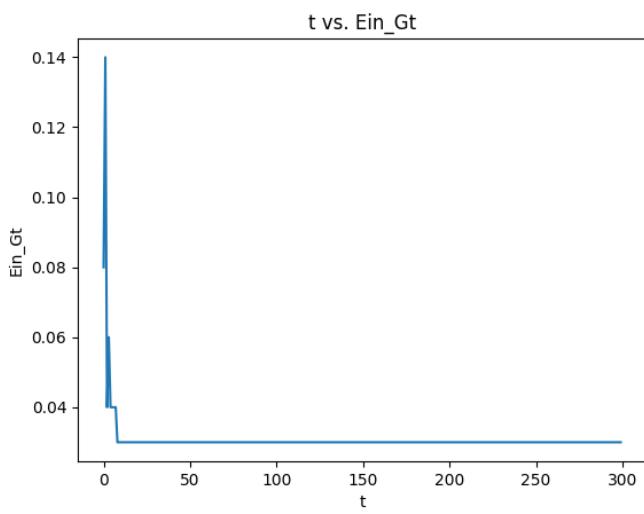
Problem 14:

我們觀察hist上頭的變化，發現最大值大約在0.06左右的錯誤率，當T增加基本上分佈沒有太大改變，主要分布於0.06附近。



Problem 15:

以下是我們觀察 $T=300, T=1000, T=30000$ 的差別，基本上後面變化已經不大了，到達極限值。



Problem 16:

我們觀察在16題所做出來的圖形發現，經過bagging的random forest雖然因為隨機取樣的關係， E_{in} 相較於取全部data來煩而稍微上升一些，但是對於testing的部分 E_{out} 得以下降至最低點大約0.05左右，與沒有做bagging直接使用的DT $E_{out} = 0.126$ 要下降一些。

1000次iteration之後起伏不大，因此我們觀察較明顯起伏的前1000次觀察最低點。

