

1.

$$\text{Gini impurity: } 1 - \sum_{k=1}^K p_k^2$$

$$\max \text{ Gini impurity} \Rightarrow \min \sum_{k=1}^K p_k^2$$

設一向量 $v = \sum_{k=1}^K p_k \hat{i}_k$ 其中 \hat{i}_k 是第 k 维度的单位向量

$$u = \sum_{k=1}^K \hat{i}_k$$

$$u \cdot v = |u||v| \cos \theta \leq |u||v| \Rightarrow \sum_{k=1}^K p_k (\hat{i}_k \cdot u) = \sum_{k=1}^K p_k \leq \sqrt{\left(\sum_{k=1}^K p_k^2 \hat{i}_k \cdot \hat{i}_k\right) \cdot \left(\sum_{k=1}^K \hat{i}_k \cdot \hat{i}_k\right)}$$

$$\sum_{k=1}^K p_k \leq \sqrt{\sum_{k=1}^K p_k^2 \cdot K} \Rightarrow 1 \leq \sqrt{\sum_{k=1}^K p_k^2 \cdot K} \Rightarrow \frac{1}{K} \leq \sum_{k=1}^K p_k^2$$

$$\text{Gini impurity} = 1 - \sum_{k=1}^K p_k^2 \geq 1 - \frac{1}{K} = \frac{K-1}{K}$$

2.

$$\begin{aligned} M_+ (1 - M_+ - M_-)^2 + M_- (1 - (M_+ - M_-))^2 &= M_+ (2 - 2M_+)^2 + (1 - M_+) (-2M_+)^2 \\ &= 4M_+ (1 - M_+) [(1 - M_+) + M_+] \\ &= 4M_+ (1 - M_+) = 4M_+ - 4M_+^2 \end{aligned}$$

$$1 - M_+^2 - M_-^2 = 1 - M_+^2 - (1 - M_+)^2 = 1 - M_+^2 - (M_+^2 - 2M_+ + 1) = 2M_+ - 2M_+^2$$

$$2 \times \text{Gini impurity} = 4M_+ - 4M_+^2 = \text{SRE when using binary classification.}$$

3.

$$\text{一個數據沒被選到的機率: } \left(1 - \frac{1}{N}\right)^N = \left(1 - \frac{1}{N}\right)^N = \left[\left(1 - \frac{1}{N}\right)^N\right]^P$$

$$\text{for } N \text{ large, } \left[\left(1 - \frac{1}{N}\right)^N\right]^P \approx e^{-P}$$

共有約 $N \times e^{-P}$ 個數據不會被選到。

4. $G = \text{Uniform}(g_t)$, 对一类 x 与 y , 若分错错误, $\sum_{k=1}^K [g_k(x) \neq y] \geq \frac{K+1}{2}$, 要有超过一半 g_k 分错错误
 设有 N 个类, 全部的 $tree$ 共分错 $\sum_{k=1}^K e_k \cdot N$ 次 $\Rightarrow G$ 最多错 $N \cdot \frac{\sum_{k=1}^K e_k}{(K+1)}$ 个类

$$Error(G) = \frac{\sum_{n=1}^N [G(x_n) \neq y_n]}{N} \leq N \cdot \frac{\sum_{k=1}^K e_k}{(K+1)} \cdot N = \frac{2}{K+1} \sum_{k=1}^K e_k$$

5.

$$\alpha_1 = \min_{\eta} \frac{1}{N} \sum_{n=1}^N ((y_n - s_n) - \eta g_1(x_n))^2$$

\Rightarrow 对 η 做偏微

$$= -\frac{2}{N} \sum_{n=1}^N ((y_n - s_n) - \eta g_1(x_n)) \cdot g_1(x_n)$$

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

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

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

6.

根据上题: $\alpha_t = \min_{\eta} \frac{1}{N} \sum_{n=1}^N (y_n - s_n - \eta g_t(x_n))^2$

\Rightarrow 对 η 做偏微

$$= -\frac{2}{N} \sum_{n=1}^N (y_n - s_n - \eta g_t(x_n)) \cdot g_t(x_n) = 0$$

$$\Rightarrow \sum_{n=1}^N (y_n - s_n) \cdot g_t(x_n) = \eta \sum_{n=1}^N g_t^2(x_n) \Rightarrow \alpha_t = \frac{\sum_{n=1}^N g_t(x_n) \cdot (y_n - s_n)}{\sum_{n=1}^N g_t^2(x_n)}$$

$$\alpha_t \sum_{n=1}^N g_t^2(x_n) = \sum_{n=1}^N g_t(x_n) \cdot y_n - \sum_{n=1}^N g_t(x_n) \cdot s_n \Rightarrow \sum_{n=1}^N g_t(x_n) \cdot s_n = \sum_{n=1}^N g_t(x_n) \cdot y_n - \alpha_t \sum_{n=1}^N g_t^2(x_n)$$

但此时 s_n 还未更新, 新的 s_n 为 $s_n + \alpha_t g_t(x_n)$

$$\sum_{n=1}^N s_n g_t(x_n) = \sum_{n=1}^N (s_n + \alpha_t g_t(x_n)) \cdot g_t(x_n) = \sum_{n=1}^N s_n g_t(x_n) + \alpha_t \sum_{n=1}^N g_t^2(x_n)$$

$$= \sum_{n=1}^N g_t(x_n) \cdot y_n - \alpha_t \sum_{n=1}^N g_t^2(x_n) + \alpha_t \sum_{n=1}^N g_t^2(x_n) = \sum_{n=1}^N g_t(x_n) \cdot y_n$$

7. squared error polynomial regression: $\beta = (\lambda I + K)^{-1} y$.
 \therefore the regression is without regularization $\Rightarrow \lambda = 0 \Rightarrow \beta = K^{-1} y \Rightarrow K\beta = y$
 根據 problem 5. $\alpha_i = \min_{\eta} \frac{1}{N} \sum_{n=1}^N ((y_n - \eta g(x_n)) - \eta g(x_n))^2$
 $\Rightarrow -\frac{2}{N} \sum_{n=1}^N (y_n - \eta g(x_n)) \cdot g(x_n)$
 對 η 偏微 \Rightarrow 對於所有 $x_n, n \in N, g(x_n) = \beta^T z_n = \beta \cdot K_n \Rightarrow$ 已知 $K\beta = y \Rightarrow y_n = K_n \beta = \beta \cdot K_n$
 \Rightarrow 對於所有 $x_n, g(x_n) = y_n$
 $-\frac{2}{N} \sum_{n=1}^N (y_n - \alpha_i g(x_n)) \cdot g(x_n) = -\frac{1}{N} \sum_{n=1}^N (y_n - \alpha_i y_n) \cdot y_n = 0 \Rightarrow \alpha_i = 1 \neq$

8.
 $OR(x_1, x_2, \dots, x_d) = \begin{cases} \text{FALSE} & \text{iff } x_1 = x_2 = \dots = x_d = \text{FALSE} \\ \text{TRUE} & \text{else.} \end{cases}$

set $w_0 = d - 0.5, w_1 = w_2 = \dots = w_d = 1$
 如此一來 $\sum_{i=0}^d w_i x_i < 0$ when $x_1 = x_2 = \dots = x_d = -1$
 $\sum_{i=0}^d w_i x_i > 0$ for other cases.

9.
 initial $s_j^{(l)} = \sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_j^{(l-1)} = 0$ for $l = 1, 2, \dots, L$.
 $s_i^{(1)} = -2(y_n - s_i^{(1)}) = -2 \cdot y_n \neq 0$
 $s_j^{(l)} = \frac{\partial \mathcal{L}_n}{\partial s_j^{(l)}} = \sum_k (s_k^{(l+1)})(w_{jk}^{(l+1)})(\tanh'(s_j^{(l)})) = 0$ for $l = 1, \dots, L-1$. #

10. $\frac{\partial e}{\partial S_k^{(u)}}$ 有 2 种情形: ① $y = k \neq k \Rightarrow V_k = 0$ ② $y = k \Rightarrow V_k = 1$

① $y = k \neq k$:

$$e = -\sum_{k=1}^K V_k \ln q_k = -\ln q_k$$

$$\frac{\partial e}{\partial S_k^{(u)}} = \frac{\partial}{\partial S_k^{(u)}} \left(-\ln \left(\frac{\exp(S_k^{(u)})}{\sum_{t=1}^K \exp(S_t^{(u)})} \right) \right) = \frac{\partial}{\partial S_k^{(u)}} \left(-\ln \exp(S_k^{(u)}) + \ln \left(\sum_{t=1}^K \exp(S_t^{(u)}) \right) \right)$$

$$= 0 + \frac{1}{\sum_{t=1}^K \exp(S_t^{(u)})} \cdot \frac{\partial}{\partial S_k^{(u)}} \left(\sum_{t=1}^K \exp(S_t^{(u)}) \right) = \frac{\exp(S_k^{(u)})}{\sum_{t=1}^K \exp(S_t^{(u)})} = q_k - V_k$$

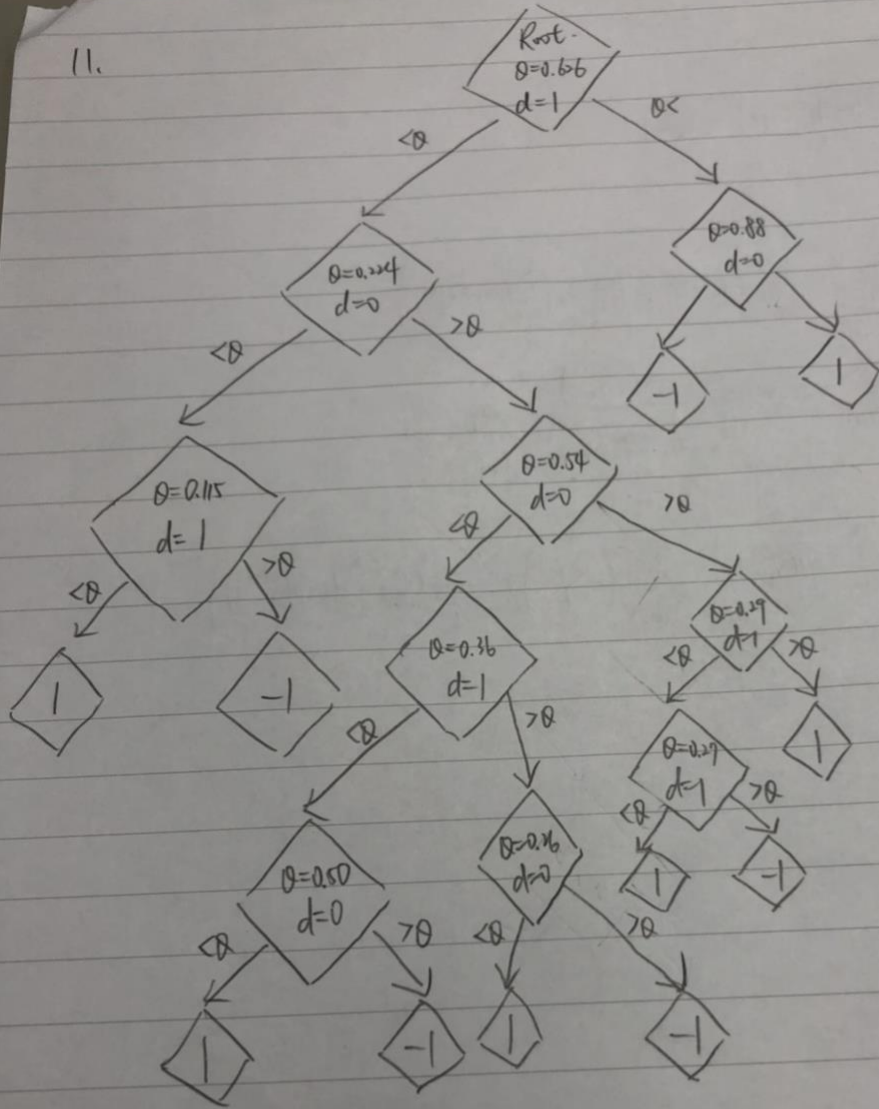
② $y = k$:

$$\frac{\partial e}{\partial S_k^{(u)}} = \frac{\partial}{\partial S_k^{(u)}} \left(-\ln \exp(S_k^{(u)}) + \ln \left(\sum_{t=1}^K \exp(S_t^{(u)}) \right) \right) = \frac{\partial}{\partial S_k^{(u)}} \left(-S_k^{(u)} \right) + \frac{\partial}{\partial S_k^{(u)}} \left(\ln \left(\sum_{t=1}^K \exp(S_t^{(u)}) \right) \right)$$

$$= -1 + \frac{\exp(S_k^{(u)})}{\sum_{t=1}^K \exp(S_t^{(u)})} = q_k - V_k$$

$$\Rightarrow \frac{\partial e}{\partial S_k^{(u)}} = q_k - V_k \neq$$

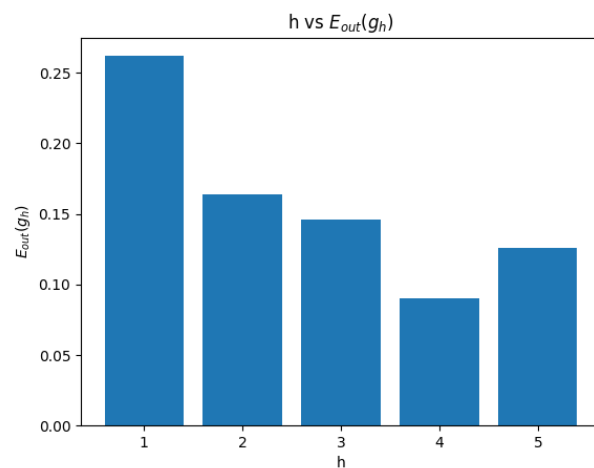
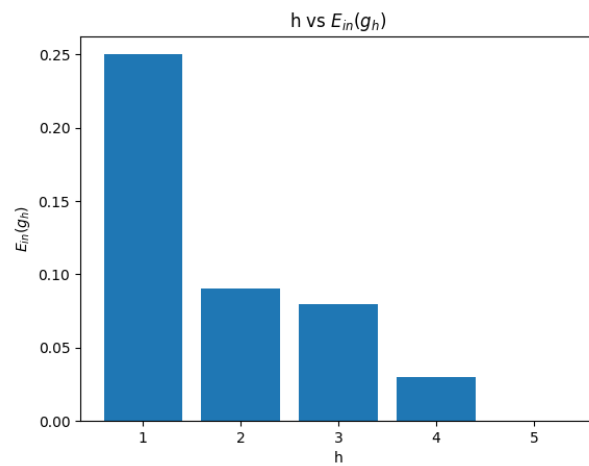
11.



12.

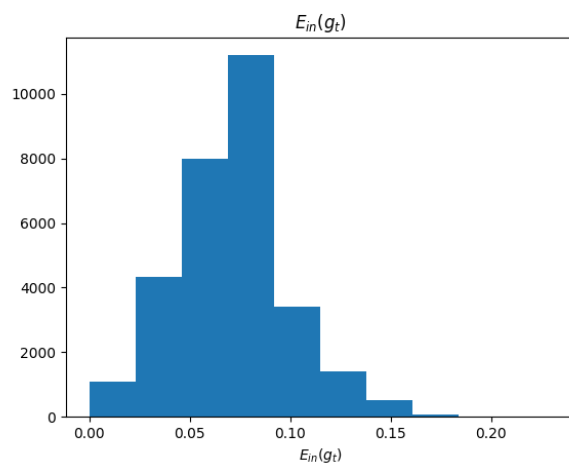
$E_{in} = 0.0$, $E_{out} = 0.126$

13.

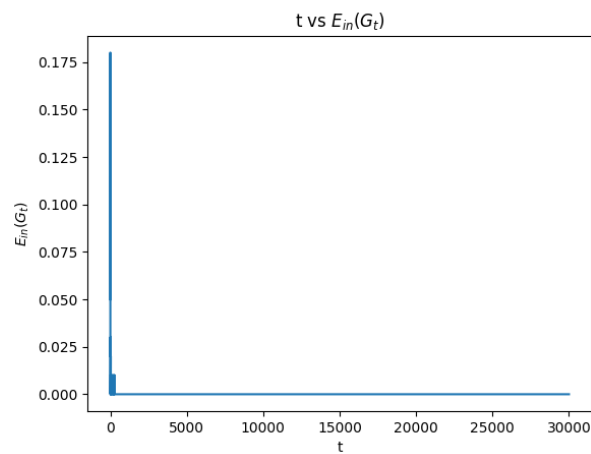


可以看出 $E_{in}(g_h)$ 隨 h 嚴格遞減，因為 tree 越深只會分越細， $E_{in}(g_h)$ 必然越來越小。而 $E_{out}(g_h)$ 雖然有遞減趨勢，但卻沒有嚴格遞減。可以看出整體來說 h 越大， $E_{in}(g_h)$ 和 $E_{out}(g_h)$ 都會越來越小，直到樹完全長完為止(此例 h 最高為 5)。

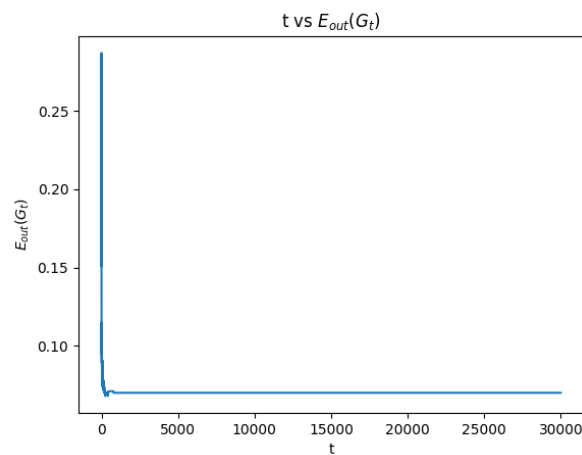
14.



15.



16.



15, 16 題的兩張圖非常相似，兩個都在 tree 數量小於 1000 時就趨於穩定。第一張圖的 $E_{in}(G_t)$ 最後穩定為 0，而 $E_{out}(G_t)$ 則在 0.05 左右。所以其實樹的數量並不需要太多，畢竟太多只會增加 model 的計算量，對降低 error rate 並沒有實質幫助。最好是能夠讓 tree 的 random 因素多一點增加 diversity 才會讓 forest 更有效。

17.

已知 $XOR(x_1, x_2, \dots, x_d)$ 即为计算当中为 TRUE 的变量个数，

奇数则 return TRUE, 偶数则 return FALSE

对于 x_1, x_2, \dots, x_d , construct 一个 $d-d-1$ feed forward network

such that

$$\begin{cases} w_{0k}^{(1)} = d+0.5-2k & \text{for } k \leq d \\ w_{ij}^{(1)} = 1 & \text{for } i, j \leq d, i \neq 0 \end{cases}$$

→ 如此一来, 若 $x_i^{(1)} = \text{TRUE}$, 则代表 $x_1 \sim x_d$ 中至少有 i 个 1

$$\begin{cases} w_{01}^{(2)} = -0.5 \\ w_{i1}^{(2)} = -i+1 & \text{for } i \leq d, i \neq 0 \end{cases} \Rightarrow w_{i1}^{(2)} = (-0.5, 1, -1, 1, -1, \dots)$$

当 x_1, \dots, x_d 有 k 个 TRUE, k 为奇数: 前 k 项 $x_i^{(1)}$ 为 1 后 $d-k$ 项 $x_i^{(1)}$ 为 -1

$$\begin{aligned} f_d(x_1, \dots, x_d) &= -0.5 + (1+(-1)+1+(-1)+\dots+1)x + ((-1)+1+(-1)+1+\dots)x(-1) \\ &= -0.5 + 1 + (-1, 0)x - 1 = 0.5 + (1, 0)x > 0 \end{aligned}$$

→ 取决于 d 为奇或偶

$$XOR(x_1, x_2, \dots, x_d) = \text{sign}(f_d(x_1, \dots, x_d)) = 1$$

当 x_1, \dots, x_d 有 k 个 TRUE, k 为偶数:

$$\begin{aligned} f_d(x_1, \dots, x_d) &= -0.5 + (1+(-1)+1+(-1)+\dots+1+(-1))x + (1+(-1)+1+(-1)+\dots)x(-1) \\ &= -0.5 + 0 + (-1, 0)x - 1 = (-1.5, -0.5)x < 0 \end{aligned}$$

→ 取决于 d 为奇或偶

$$XOR(x_1, x_2, \dots, x_d) = \text{sign}(f_d(x_1, \dots, x_d)) = -1$$

⇒ 此 $d-d-1$ feed-forward NN 即可 implement $XOR(x_1, x_2, \dots, x_d)$