

$$1. \quad g_{S, i, \theta}(x) = S_i \cdot \text{sign}(x_i - \theta) \Rightarrow$$

$$S = \{+1, -1\}$$

$$i = \{1, \dots, d\}$$

$$\theta = \{L + 0.5, L + 1.5, \dots, R - 0.5\} \quad L, R \in \mathbb{R}$$

$$K_{ds}(x, x') = (g_{+1, 1, \theta_1}(x), g_{+1, 1, \theta_2}(x), \dots, g_{+1, 1, \theta_k}(x), g_{-1, d, \theta_k}(x)) \cdot$$

$$(g_{+1, 1, \theta_1}(x'), g_{+1, 1, \theta_2}(x'), \dots, g_{+1, 1, \theta_k}(x'), g_{-1, d, \theta_k}(x'))$$

$$= \sum_{S=\{+1, -1\}} \sum_{T=1}^d \sum_{\theta=L+0.5}^{R-0.5} [g_{S, T, \theta}(x) \ g_{S, T, \theta}(x')]$$

$$= \sum_{T=1}^d \sum_{\theta=L+0.5}^{R-0.5} [(+1) \text{sign}(x_T - \theta) (+1) \text{sign}(x'_T - \theta) + (-1) \text{sign}(x_T - \theta) (-1) \text{sign}(x'_T - \theta)]$$

$$= 2 \sum_{T=1}^d \sum_{\theta=L+0.5}^{R-0.5} [\text{sign}(x_T - \theta) \text{sign}(x'_T - \theta)]$$

↓

$$\left\{ \begin{array}{l} x_i, x'_i > \theta \\ x_i, x'_i < \theta \end{array} \right\} \rightarrow \text{sign}(x_i - \theta) \text{sign}(x'_i - \theta) = 1$$

$$\left\{ \begin{array}{l} x_i > \theta \quad \& \quad x'_i < \theta \\ x_i < \theta \quad \& \quad x'_i > \theta \end{array} \right\} \rightarrow \text{sign}(x_i - \theta) \text{sign}(x'_i - \theta) = -1$$

$$\left\{ \begin{array}{l} x_i = 0 \quad \text{or} \quad x'_i = 0 \end{array} \right\} \rightarrow \text{sign}(x_i - \theta) \text{sign}(x'_i - \theta) = 0$$

ex.

$$\theta \text{ is in } [x_i, x'_i] \quad \begin{array}{ccccccc} & & & \theta & & & \\ & | & | & | & | & | & \rightarrow \\ L & x_i & x'_i & & & R & \Rightarrow -1 \end{array}$$

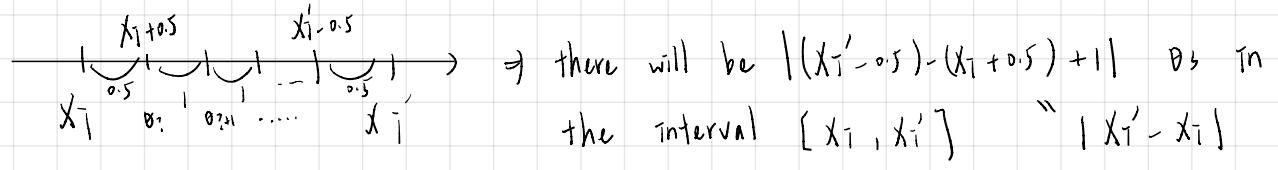
$$\theta \text{ is out side of } [x_i, x'_i] \quad \begin{array}{ccccccc} & & & \theta & & & \\ & | & | & | & | & | & \rightarrow \\ L & x_i & x'_i & & & R & \Rightarrow 1 \end{array}$$

Hence we need to find how many θ is in and out of

the interval $[x_i, x'_i]$ (Suppose $x_i \leq x'_i$) to get the value of

$$\sum_{\theta=L+0.5}^{R-0.5} \text{sign}(x_i - \theta) \text{sign}(x'_i - \theta)$$

there will be $R-L$ θ_s in the interval $[L+0.5, R-0.5]$



$$= \sum_{i=1}^d (\text{count of } \theta_s \text{ outside of } [x_i, x_i'] - \text{count of } \theta_s \text{ in } [x_i, x_i'])$$

$$= \sum_{i=1}^d \left\{ [(R-L) - |x_i' - x_i|] - |x_i' - x_i| \right\}$$

$$= \sum_{i=1}^d ((R-L) - 2|x_i' - x_i|)$$

$$= d(R-L) - 4||x' - x||$$

21

$$\text{Dual of soft-margin SVM: } \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m \tilde{K}(z_n, z_m) - \sum_{n=1}^N \alpha_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha_n = 0 ;$$

$$0 \leq \alpha_n \leq C, \text{ for } n=1, 2, \dots, N;$$

$$\text{implicitly } w = \sum_{n=1}^N \alpha_n y_n z_n ;$$

$$\beta_n = b - \alpha_n \text{ for } n=1, 2, \dots, N$$

$$g_{\text{SVM}}(x) = \text{sign} \left(\sum_{\substack{n \\ \text{N indicates } n}} \alpha_n y_n \tilde{K}(z_n, x) + b \right)$$

$$\text{expand } K' \rightarrow \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m (u K(z_n, z_m) + v) - \sum_{n=1}^N \alpha_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha_n = 0, \quad 0 < \alpha_n < \frac{C}{u}$$

$$\text{expand } \rightarrow \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m u K(z_n, z_m) + \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m - \sum_{n=1}^N \alpha_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha_n = 0, \quad 0 < \alpha_n < \frac{C}{u}$$

$$\begin{matrix} \times u \\ \parallel n > 0 \end{matrix} \rightarrow \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (u \alpha_n) (u \alpha_m) y_n y_m (K(z_n, z_m) + v) - \sum_{n=1}^N (u \alpha_n)$$

$$\text{subject to } \sum_{n=1}^N u y_n \alpha_n = 0, \quad 0 < (u \alpha_n) < C$$

$$\begin{matrix} \text{take } \alpha' = u \alpha_n \\ \rightarrow \end{matrix} \min_{\alpha'} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha'_n \alpha'_m y_n y_m (K(z_n, z_m) + v) - \sum_{n=1}^N \alpha'_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha'_n = 0, \quad 0 < \alpha'_n < C$$

$$\text{expand } \rightarrow \min_{\alpha'} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha'_n \alpha'_m y_n y_m K(z_n, z_m) + \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha'_n \alpha'_m y_n y_m - \sum_{n=1}^N \alpha'_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha'_n = 0, \quad 0 < \alpha'_n < C$$

$$\downarrow$$

$$= \frac{1}{2} \left(\sum_{n=1}^N \alpha'_n y_n \right) \left(\sum_{m=1}^N \alpha'_m y_m \right) = 0$$

$$\rightarrow \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha'_n \alpha'_m y_n y_m K(z_n, z_m) - \sum_{n=1}^N \alpha'_n$$

$$\text{subject to } \sum_{n=1}^N y_n \alpha'_n = 0, \quad 0 < \alpha'_n < C$$

Hence we can solve soft-margin SVM by the new optimization problem above.

$$\begin{aligned} \text{min}_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m K(x_n, x_m) - \sum_{n=1}^N \alpha_n \\ \text{subject to } \sum_{n=1}^N y_n \alpha_n = 0, \quad 0 < \alpha_n < C \end{aligned}$$

which is exactly the same problem as using kernel $= K(x, x')$
and $C = C$

Hence $\tilde{g}_{\text{SVM}}(x) = g_{\text{SVM}}(x)$

$$3. E_{out}(h) = \text{avg}(E_{out}(g_t)) - \text{avg}(\epsilon(g_t-h)^2)$$

$$= \frac{1}{N} \sum_{t=1}^N \epsilon_t - \text{avg}(\epsilon(g_t-h)^2)$$

$$\Rightarrow \frac{E_{out}(h)}{E} = \frac{1}{N} - \frac{\text{avg}(\epsilon(g_t-h)^2)}{E}$$

$\epsilon(g_t-h)^2 \geq 0$ (it equals to 0 when under any condition, all $g_t = h$, which means $g_t = \text{every other } g_i \text{ in } G$)

$$\therefore \text{avg}(\epsilon(g_t-h)^2) \geq 0$$

$$\text{Hence } \frac{E_{out}(h)}{E} \geq \frac{1}{N} *$$

4. Probability for an example not to be sampled, if we only sample 1 example: $\frac{N-1}{N}$

Probability for an example not to be sampled, if we sample $\frac{3}{4}N$ examples: $(\frac{N-1}{N})^{\frac{3}{4}N}$

$$\lim_{N \rightarrow \infty} \left(\frac{N-1}{N} \right)^{\frac{3}{4}N} = \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N} \right)^{\frac{3}{4}N} = \lim_{N \rightarrow \infty} \left[\left(1 - \frac{1}{N} \right)^N \right]^{\frac{3}{4}} = \left[\lim_{N \rightarrow \infty} \left(1 - \frac{1}{N} \right)^{\frac{1}{N}} \right]^{\frac{3}{4}} = e^{-0.75} \approx 0.47237$$

5. 98% data + , 2% data -

Suppose $N=100$, we'll have $y_n > 0$, $\neg y_n < 0$

$$\text{then } \frac{\sum_{n:y_n>0} u_n^{(t)}}{\sum_{n:y_n<0} u_n^{(t)}} = \frac{\sum_{n:y_n>0} u_n^{(t)} / \phi_t}{\sum_{n:y_n<0} u_n^{(t)} \times \phi_t}, \text{ where } \phi_t = \sqrt{\frac{1-t_t}{t_t}} = \sqrt{\frac{98\%}{2\%}} = 1$$

$$= \frac{\frac{1}{100} \times 98 \times \frac{1}{1}}{\frac{1}{100} \times 2 \times 1} = 1 \quad \times$$

$$b. V_t = \sum_{n=1}^N u_n^{(t)}$$

$$t_t = \frac{\sum_{n=1}^N u_n^{(t)} [y_n \neq g_t(x)]}{\sum_{n=1}^N u_n^{(t)}}$$

$$\Rightarrow V_t \cdot t_t = \sum_{n=1}^N u_n^{(t)} [y_n \neq g_t(x)], V_t \cdot (1-t_t) = \sum_{n=1}^N u_n^{(t)} [y_n = g_t(x)]$$

$$V_{t+1} = \sum_{n=1}^N u_n^{(t+1)} = \sum_{n=1}^N u_n^{(t)} [y_n \neq g_t(x)] \times \phi_t \\ + \sum_{n=1}^N u_n^{(t)} [y_n = g_t(x)] / \phi_t$$

$$= V_t \cdot t_t \times \phi_t + V_t \cdot (1-t_t) / \phi_t$$

$$= V_t \cdot \sqrt{\frac{1-t_t}{t_t}} \times t_t + V_t \cdot (1-t_t) \cdot \sqrt{\frac{t_t}{1-t_t}} \\ = \sqrt{V_t \cdot t_t (1-t_t)}$$

$$\Rightarrow \frac{V_{t+1}}{V_t} = \sqrt{t_t (1-t_t)}$$

$$\text{Then } \frac{V_{T+1}}{V_T} \times \frac{V_T}{V_{T-1}} \times \cdots \times \frac{V_2}{V_1} = \sqrt{t_1 (1-t_1)} \cdot \sqrt{t_2 (1-t_2)} \cdots \sqrt{t_T (1-t_T)} \\ = \sqrt{t_1 (1-t_1)} \sqrt{t_2 (1-t_2)} \cdots \sqrt{t_T (1-t_T)} \quad \times$$

1. Let $\vec{s}_t = [s_1, s_2, \dots, s_N]$, $\vec{g}_t = [g_t(x_1), g_t(x_2), \dots, g_t(x_N)]$, $\vec{y} = [y_1, y_2, \dots, y_N]$
 y at round t

the optimal η is obtain by solving $\min_{\eta} \frac{1}{N} \sum_{n=1}^N ((y_n - s_n) - \eta g_t(x_n))^2$

$$\text{hence } \eta = \frac{\sum_{n=1}^N g_t(x_n)(y_n - s_n)}{\sum_{n=1}^N g_t^2(x_n)} = \frac{\vec{g}_t (\vec{y} - \vec{s}_{t-1})^\top}{\vec{g}_t \vec{g}_t^\top} = \alpha_t$$

$$\text{we know } \vec{s}_t = \alpha_t \vec{g}_t + \vec{s}_{t-1}$$

$$\begin{aligned} \text{then } \sum_{n=1}^N s_n g_t(x_n) &= \vec{s}_t \vec{g}_t^\top = (\alpha_t \vec{g}_t + \vec{s}_{t-1}) \vec{g}_t^\top = \alpha_t \vec{g}_t \vec{g}_t^\top + \vec{s}_{t-1} \vec{g}_t^\top \\ &= \vec{g}_t (\vec{y} - \vec{s}_{t-1})^\top + \vec{g}_t \vec{s}_{t-1}^\top \\ &= \vec{g}_t \vec{y}^\top \\ &= \sum_{n=1}^N g_t(x_n) y_n \end{aligned}$$

8.

By gradient descent update function : $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)}$

Because $w^{(l)} = 0$ (initialization), we know $w_{ij}^{(l)} = 0$, for any i, j, l .

$$\begin{aligned} \text{Get } \delta_j^{(l)} &= \sum_{k=1}^{d(l+1)} (\delta_k^{(l+1)}) (w_{jk}^{(l+1)}) (\tanh' s_j^{(l)}) \\ &= \sum_{k=1}^{d(l+1)} (\delta_k^{(l+1)}) (w_{jk}^{(l+1)}) (\tanh' \left(\sum_{i=1}^{d(l+1)} w_{ij}^{(l)} x_i^{(l)} \right)) \\ &= 0 \end{aligned}$$

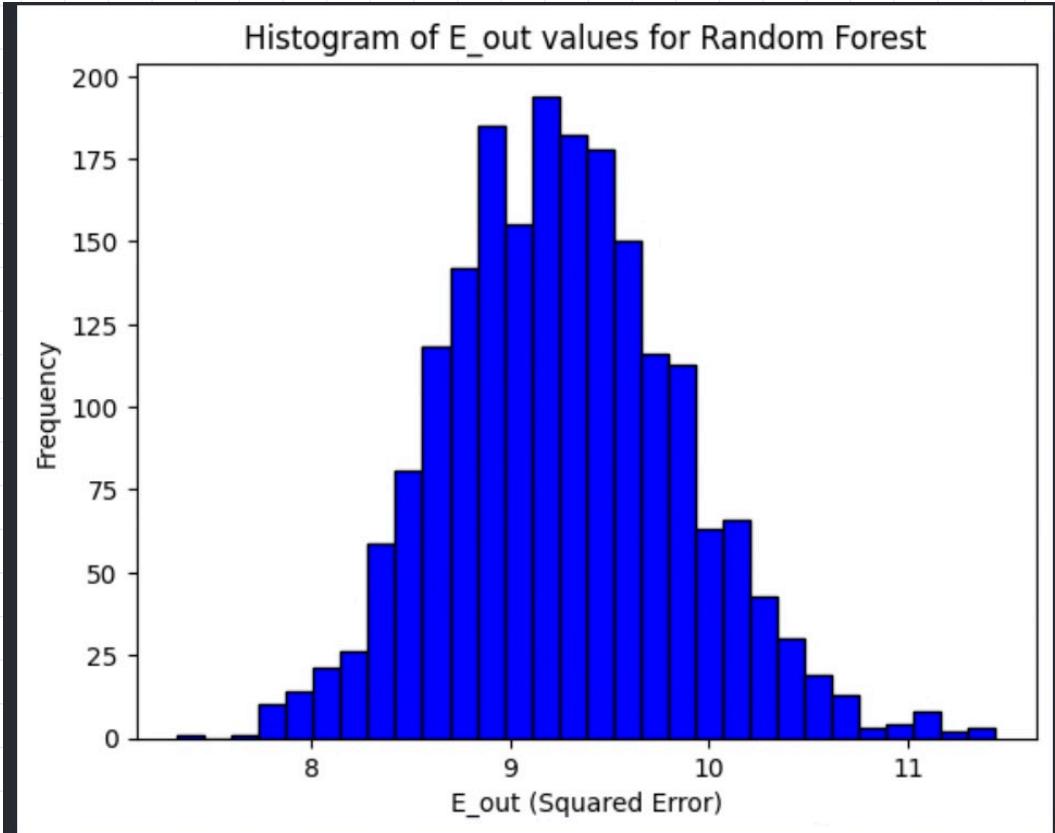
Hence, the update function becomes $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - 0$, for every i, j, l .

We know that every gradient component in this NN is 0.

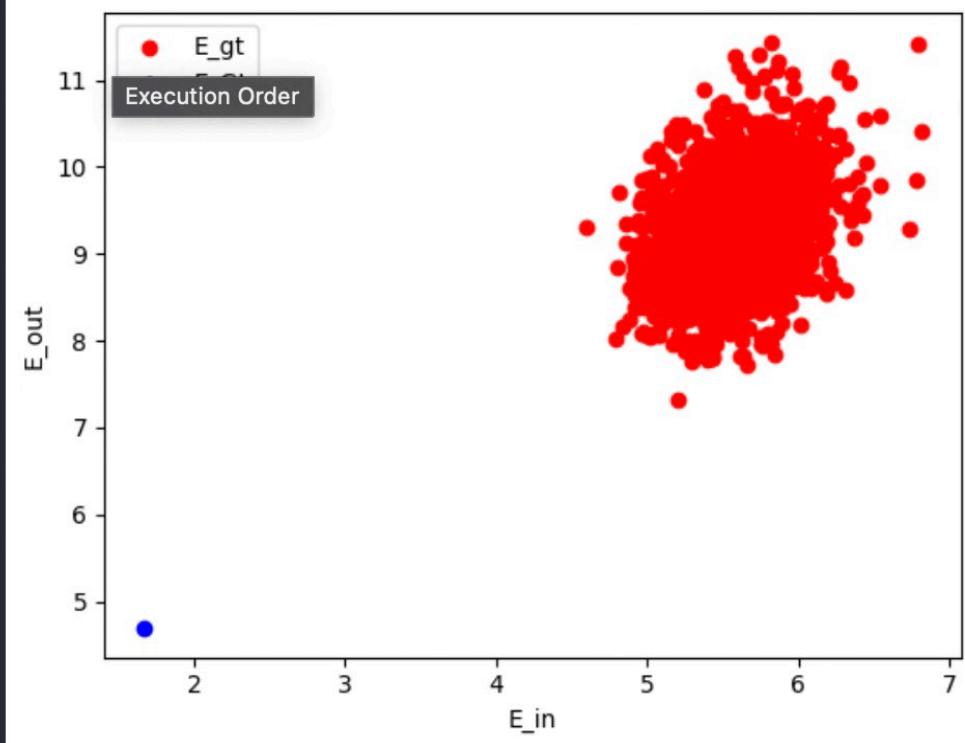
9.

8.79324462640737

10.



11



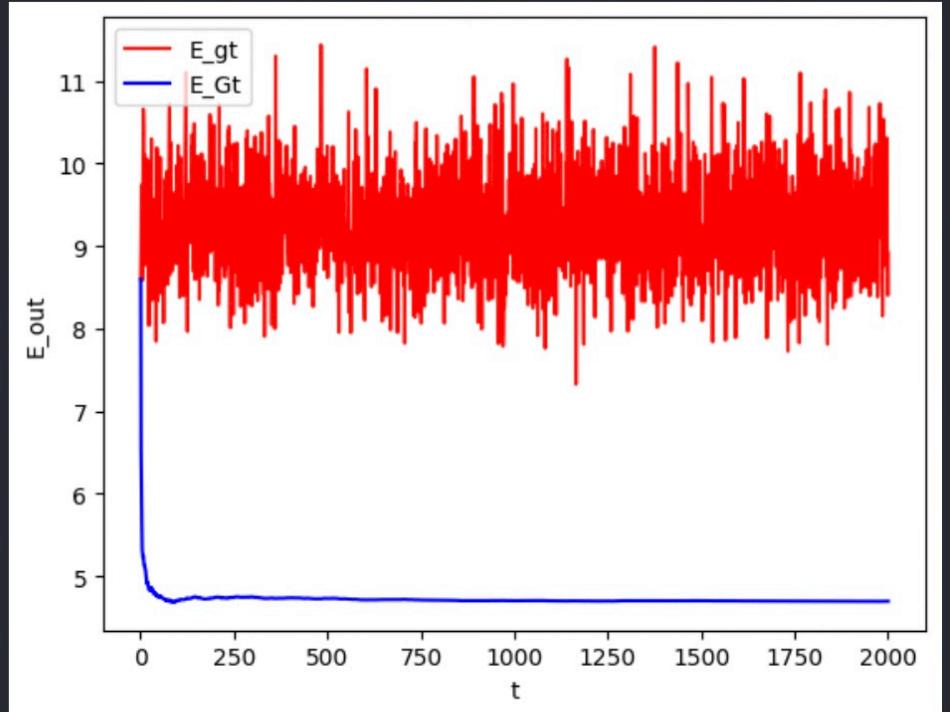
We can see through getting the avg of our random forest, we can significantly let E_{out} and E_{in} to drop.

The E_{in} dropped from 6 to less than 2, and E_{out} dropped from 10 to less than 5.

The E_{out} basicly become the half of those RF.

We can conclude that those RF trained by random sampling are pretty different, so that E_{out} can be reduce that much.

12



As we can see, we don't really need 2000 trees to achieve this amazing performance.

Less than 100 trees would've already do the work.

We can see that the difference for 100 trees are already sufficiently enough.

With more trees, the value only goes towards convergence, which is not a lot of improvement.

If we trained with only 100 trees, we could've save a lot time and achieve almost the same performance.

13.

For $d=2$, $\text{XOR}(g_1, g_2) = \text{or}(\text{AND}(-g_1, g_2), \text{AND}(g_1, -g_2))$

the number of ANDs is the number of $l=1$, which is 2.

So, we can see only a 2-2-1 NN works for $d=2$, because we can't remove any ANDs in the XORL.

For $d=3$, \oplus represents XOR, $+$ represents OR, $g_i g_r$ represents and, g' represents not g

$$\begin{aligned}
 \text{Hence } g_1 \oplus g_2 \oplus g_3 &= (g_1 \oplus g_2) \oplus g_3 \\
 &= (g'_1 g_2 + g_1 g'_2) \oplus g_3 \\
 &= (g'_1 g_2 + g_1 g'_2) g_3' + (g'_1 g_2 + g_1 g'_2)' g_3 \\
 &= g'_1 g_2 g_3' + g_1 g'_2 g_3' + \underline{(g'_1 g_2)' (g_1 g'_2)' g_3} \\
 &\quad \downarrow \\
 &\frac{(g_1 + g_2') (g_1' + g_2)}{g_1 g'_1 g_2 + g_1 g_2 g'_1 + g_1' g_2' g_3 + g_2' g_2 g_3} \\
 &= g_1 g_2' g_3' + g_1' g_2 g_3' + g_1' g_2' g_3 + g_1 g_2 g_3
 \end{aligned}$$

Hence, we can see for $d=3$ to work, it'll yield 4 neurons

meaning the second layer must have size of 4, hence a 3-4-1 NN.

Then we can see the formula for $d=d$,

$$g_1 \oplus g_2 \oplus \dots \oplus g_d = g'_1 g_2 \dots g_d + g_1 g'_2 g_3 \dots g_d + \dots + g_1 g_2 \dots g'_d + g_1 \dots g_d$$

We will need $d+1$ neurons in the middle layer, hence a $d-(d-1)-1$ NN won't do the work.