

1. (20 points) When talking about non-uniform voting in aggregation, we mentioned that α can be viewed as a weight vector learned from any linear algorithm coupled with the following transform:

$$\phi(\mathbf{x}) = (g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_T(\mathbf{x})).$$

When studying kernel methods, we mentioned that the kernel is simply a computational short-cut for the inner product $(\phi(\mathbf{x}))^T(\phi(\mathbf{x}'))$. In this problem, we mix the two topics together using the decision stumps as our $g_i(\mathbf{x})$.

Assume that the input vectors contain only integers between (including) L and R , where $L < R$.

Consider the decision stumps $g_{s,i,\theta}(\mathbf{x}) = s \cdot \text{sign}(x_i - \theta)$, where

- $i \in \{1, 2, \dots, d\}$,
- d is the finite dimensionality of the input space,
- $s \in \{-1, +1\}$,
- $\theta \in \{\theta_1 = L + 0.5, \theta_2 = L + 1.5, \dots, \theta_k = R - 0.5\}$

Define $\phi_{ds}(\mathbf{x}) = \left(g_{+1,1,\theta_1}(\mathbf{x}), g_{+1,1,\theta_2}(\mathbf{x}), \dots, g_{+1,1,\theta_k}(\mathbf{x}), \dots, g_{-1,d,\theta_k}(\mathbf{x}) \right)$. What is $K_{ds}(\mathbf{x}, \mathbf{x}') = (\phi_{ds}(\mathbf{x}))^T (\phi_{ds}(\mathbf{x}'))$? Prove your answer.

(Hint: This result shows that aggregation learning with SVMs is possible. Those who are interested in knowing how perceptrons and decision trees can be used instead of decision stumps during kernel construction can read this early work [1].)

$$q(x) \cdot (g_1(x) \cdot g_2(x) \cdot g_3(x) \cdots g_T(x))$$

$$L \xrightarrow{\theta_1} \theta_2 \xrightarrow{\dots} \theta_R \xrightarrow{R}$$

$$g_{s,\bar{n},\theta}(1) = s \cdot \text{sign}(\bar{n} - \theta)$$

$$k_{ds}(x, x') = \phi_{ds}(x)^T \phi_{ds}(x')$$

$$= \sum_{\tilde{\lambda}=1}^d \sum_{n=1}^k \frac{g_{+,\tilde{\lambda},n}(x) g_{+,\tilde{\lambda},n}(y)}{} + \lambda$$

$$= 2 \cdot \sum_{i=1}^d \sum_{n=1}^k \text{sign}(x_i - \theta_n) \text{sign}(x_i' - \theta_n)$$

$$\begin{aligned} & \xrightarrow{\text{Sigmoid}} [s.\text{sign}(x_n' - \theta_n)] \\ &= \text{sign}(x_n - \theta_n) \cdot \text{sign}(x_n' - \theta_n) \end{aligned}$$

$$= \text{sign}(x_i - \theta_n) \cdot \text{sign}(x_n - \theta_n)$$

$$= 2 \sum_{\bar{X}=1}^d ((R-L) - 2(X_{\bar{X}} - X'_{\bar{X}})) = 2d(R-L) - 4 \sum_{\bar{X}=1}^d (X_{\bar{X}} - X'_{\bar{X}}) \#$$

2. (20 points) For a given valid kernel K , consider a new kernel $\tilde{K}(\mathbf{x}, \mathbf{x}') = uK(\mathbf{x}, \mathbf{x}') + v$ for some $u > 0$. Note that v can be any real value. When $v \geq 0$, it is easy to prove that \tilde{K} is still a valid kernel; when $v < 0$, however, \tilde{K} may not always be a valid kernel. Prove that for the dual of soft-margin support vector machine, using \tilde{K} along with a new $\tilde{C} = \frac{C}{u}$ instead of C with the original C leads to an equivalent g_{SVM} classifier. That is, the optimal (α, b) obtained leads to the same decision boundary.

(Note: This result can be used to shift and scale the kernel function derived in the previous problem to a simpler form (even not as a valid kernel) when coupling it with the soft-margin SVM.)

$$\tilde{C} = \frac{C}{u} \quad . \quad \tilde{K} = uK(x_i x') + v$$

$$b = y_s - \sum_{sv} \alpha_n y_n K(x_n, x_s), \quad 0 < \alpha_n < C$$

$$\sum_{n=1}^N \alpha_n y_n = 0$$

$$K' = \min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m (uK(x_n, x_m) + v) - \sum_{n=1}^N \alpha_n, \quad \text{subject to } 0 \leq \alpha_n < \frac{C}{u}$$

$$\begin{aligned} g_{\text{SVM}}(x) &= \text{sign} \left(\sum_{sv} \frac{\alpha_n}{u} y_n \tilde{K}(x_n, x) + b \right) \\ &= \text{sign} \left(\sum_{sv} \frac{\alpha_n}{u} y_n (u \times K(x_n, x) + v + b) \right) \\ &= \text{sign} \left(\sum_{sv} \alpha_n y_n K(x_n, x) + y_s - \sum_{sv} \alpha_n y_n K(x_n, x_s) \right) \\ &= \text{sign} \left[\sum_{sv} \alpha_n y_n K(x_n, x) + b \right] \rightarrow \text{equivalent } g_{\text{SVM}}. \end{aligned}$$

3. (20 points) Consider an aggregated binary classifier G that is constructed by uniform blending on 17 binary classifiers $\{g_t\}_{t=1}^{17}$. That is,

$$G(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^{17} g_t(\mathbf{x}) \right)$$

Assume that each g_t is of test 0/1 error $E_{\text{out}}(g_t) = e_t > 0$. Define the total error $E = \sum_{t=1}^{17} e_t$. What is the largest value of $\frac{E_{\text{out}}(G)}{E}$? Prove your result.

(Note: The ratio roughly shows how G reduces the total error made by the hypotheses.)

$$g_t(x) = \text{sign} \left(\sum_{t=1}^n g_t(x) \right)$$

for $y=+1$, there must be more g_t that gives wrong prediction than g_t that have right prediction at least q g_t classify as + to $E_{\text{out}}(g_t)=1$

$$E = \sum_{t=1}^{17} E_{\text{out}}(g_t) \geq q , \quad \max \left(\frac{E_{\text{out}}(g_t)}{E} \right) = \frac{1}{q} *$$

4. (20 points) If bootstrapping is used to sample exactly $\frac{3}{4}N$ examples out of N , what is the probability that an example is *not* sampled when N is very large? List your derivation steps.

(Note: This is just an "easy" exercise of letting you think about the amount of OOB data in bagging/random forest.)

$$\begin{aligned}
 \lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^{\frac{3}{4}N} &= \lim_{N \rightarrow \infty} \left(\frac{N-1}{N}\right)^{\frac{3}{4}N} \\
 &= \lim_{N \rightarrow \infty} \frac{1}{\left(\frac{N}{N-1}\right)^{\frac{3}{4}N}} \\
 &= \lim_{N \rightarrow \infty} \frac{1}{\left(1 + \frac{1}{N-1}\right)^{\frac{3}{4}N}} \quad \lim_{N \rightarrow \infty} \left(1 + \frac{1}{N}\right)^N = e. \\
 &= \lim_{N \rightarrow \infty} \left(\frac{1}{\left(1 + \frac{1}{N-1}\right)^{\frac{3}{4}}}^N\right) \quad , \quad \lim_{N \rightarrow \infty} \left(\frac{1}{1 + \frac{1}{N}}\right)^N = \frac{1}{e} \\
 &\stackrel{\text{as } N \rightarrow \infty}{\longrightarrow} \left(\frac{1}{e}\right)^{\frac{3}{4}} \#
 \end{aligned}$$

5. (20 points) Consider applying the AdaBoost algorithm on a binary classification data set where 98% of the examples are positive. Because there are so many positive examples, the base algorithm within AdaBoost returns a constant classifier $g_1(\mathbf{x}) = +1$ in the first iteration. Let $u_n^{(2)}$ be the individual example weight of each example in the second iteration. What is

$$\frac{\sum_{n: y_n > 0} u_n^{(2)}}{\sum_{n: y_n < 0} u_n^{(2)}}?$$

Prove your answer.

(Note: This is designed to help you understand how AdaBoost can deal with "imbalanced" data immediately after the first iteration.)

prediction :

If all correct $\rightarrow y_n > 0$

If all wrong $\rightarrow y_n < 0$

$\rightarrow y_n > 0$, weighted correct rate 2%

$$\frac{\sum_n y_n > 0 \ u_n^{(2)}}{\sum_n y_n < 0 \ u_n^{(2)}} = \frac{\sum_n y > 0 \ u_n^{(1)} \times 0.02}{\sum_n y < 0 \ u_n^{(1)} \times 0.98} = | *$$

$\downarrow y_n < 0$, weighted incorrect rate 98%

6. (20 points) For the AdaBoost algorithm introduced in Lecture 212, let $U_t = \sum_{n=1}^N u_n^{(t)}$. That is, $U_1 = 1$ (you are very welcome! ;)). Assume that $0 < \epsilon_t < \frac{1}{2}$ for each hypothesis g_t . Express $\frac{U_{T+1}}{U_1}$ in terms of $\epsilon_1, \epsilon_2, \dots, \epsilon_T$. Prove your answer.

(Hint: Consider checking $\frac{U_{t+1}}{U_t}$ first. The result is the backbone of proving that AdaBoost will converge within $O(\log N)$ iterations.)

$$U_t \cdot \varepsilon_t = \sum U_n^{(t)} \mathbb{I}[y_n \neq g_t(x_n)]$$

$$U_t (1 - \varepsilon_t) = \sum U_n^{(t)} \mathbb{I}[y_n = g_t(x_n)]$$

$$\begin{aligned} U_{t+1} &= U_t \cdot \varepsilon_t \sqrt{\frac{1-\varepsilon_t}{\varepsilon_t}} + (U_t - U_t \cdot \varepsilon_t) \sqrt{\frac{\varepsilon_t}{1-\varepsilon_t}} \\ &= U_t \left(\varepsilon_t \sqrt{\frac{1-\varepsilon_t}{\varepsilon_t}} + (U_t - U_t \cdot \varepsilon_t) \sqrt{\frac{\varepsilon_t}{1-\varepsilon_t}} \right) \\ &= U_t \left(\varepsilon_t \sqrt{\frac{1-\varepsilon_t}{\varepsilon_t}} + (1 - \varepsilon_t) \sqrt{\frac{\varepsilon_t}{1-\varepsilon_t}} \right) = 2 \sqrt{\varepsilon_t(1-\varepsilon_t)} U_t \end{aligned}$$

since $U_{t+1} = 2 \sqrt{\varepsilon_t(1-\varepsilon_t)} U_t$

$$U_{t+2} = 2 \sqrt{\varepsilon_t(1-\varepsilon_t)} U_{t+1} = (2 \sqrt{\varepsilon_t(1-\varepsilon_t)})^2 U_{t+1}$$

$$U_{t+3} = (2 \sqrt{\varepsilon_t(1-\varepsilon_t)})^3 U_t$$

\vdots

$$U_{t+T} = (2 \sqrt{\varepsilon_t(1-\varepsilon_t)})^T U_t$$

$$A_t = 1, \quad \frac{U_{T+1}}{U_1} = 2^T \prod_{t=1}^T \sqrt{\varepsilon_t(1-\varepsilon_t)} \#$$

7. (20 points) For the gradient boosted decision tree, after updating all s_n in iteration t using the steepest η as α_t , what is the value of $\sum_{n=1}^N s_n g_t(\mathbf{x}_n)$? Prove your answer.

(Note: This special value may tell us some important physical property on the relationship between the vectors $[s_1, s_2, \dots, s_N]$ and $[g_t(\mathbf{x}_1), g_t(\mathbf{x}_2), \dots, g_t(\mathbf{x}_N)]$.

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

$$\frac{\partial E}{\partial \eta} = \frac{1}{N} \sum_{n=1}^N 2((y_n - s_n) - \eta g_t(x_n))(-g_t(x_n)) \Rightarrow \\ \Rightarrow \sum_{n=1}^N (y_n - s_n) g_t(x_n) - \eta \sum_{n=1}^N (g_t(x_n))^2 = 0 \quad , \quad s_n \leftarrow s_n + \eta g_t(x_n)$$

$$\begin{aligned} \sum_{n=1}^N (y_n - s_n) g_t(x_n) &= \sum_{n=1}^N (y_n - s_n - \eta g_t(x_n))(g_t(x_n)) \\ &= \sum_{n=1}^N (y_n - s_n) g_t(x_n) - \eta \sum_{n=1}^N (g_t(x_n))^2 = 0 \quad \# \end{aligned}$$

Neural Networks

8. (20 points) For a Neural Network with at least one hidden layer and $\tanh(s)$ as the transformation functions on all neurons (including the output neuron), when all the initial weights $w_{ij}^{(\ell)}$ are set to 0, what gradient components are also 0? Prove your answer.

(Note: This result tells you that all-0 initialization may make it impossible for back-propagation to update some weights.)

$$x_j^{(1)} = \tanh\left(\sum_i w_{ij}^{(1)} x_i^{(0)}\right) = \tanh(0) = 0$$

$$x_k^{(2)} = \tanh\left(\sum_j w_{jk}^{(2)} x_j^{(1)}\right) = \tanh(0) = 0$$

⋮

$$g_{nn}(x) = \dots \tanh\left(\sum_j w_{jk}^{(L)} \tanh\left(\sum_i w_{ij}^{(L)} \cdot x_i^{(0)}\right)\right) = 0$$

The output will be 0.

$$\frac{\partial e_n}{\partial w_{ij}^{(L)}} = \frac{\partial e_n}{\partial s_j^{(L)}} \cdot \frac{\partial s_j^{(L)}}{\partial w_{ij}^{(L)}} = s_j^{(L)} \cdot \underline{w_{ij}^{(L)}} = 0, \quad w_{ij}^{(L)} = 0$$

$$w_{ij}^{(L)} \leftarrow w_{ij}^{(L)} - \gamma x_i^{(L-1)} \delta_j^{(L)}, \quad x_i^{(L-1)} = 0$$

$$\Rightarrow w_{ij}^{(L)} \leftarrow w_{ij}^{(L)} \rightarrow \text{weight always } 0 \#$$

9. (20 points, *) First, let's implement a simple C&RT algorithm without pruning using the squared error as the impurity measure, as introduced in the class. For the decision stump used in branching, if you are branching with feature i , please sort all the $x_{n,i}$ values to form (at most) $N + 1$ segments of equivalent θ , and then pick θ within the median of the segment. If multiple (i, θ) produce the best split, pick the one with the smallest i (and if there is a tie again, pick the one with the smallest θ).

Please run the algorithm on the following set for training:

http://www.csie.ntu.edu.tw/~htlin/course/ml23fall/hw6/hw6_train.dat

and the following file as our test data set for evaluating E_{out} :

http://www.csie.ntu.edu.tw/~htlin/course/ml23fall/hw6/hw6_test.dat

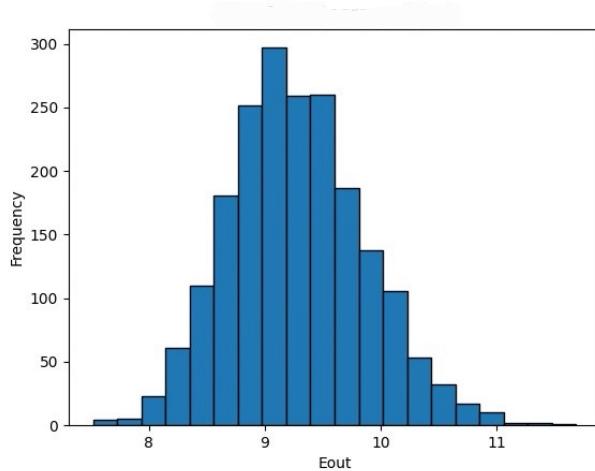
The datasets are in LIBSVM format and are processed from

<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/regression/abalone>

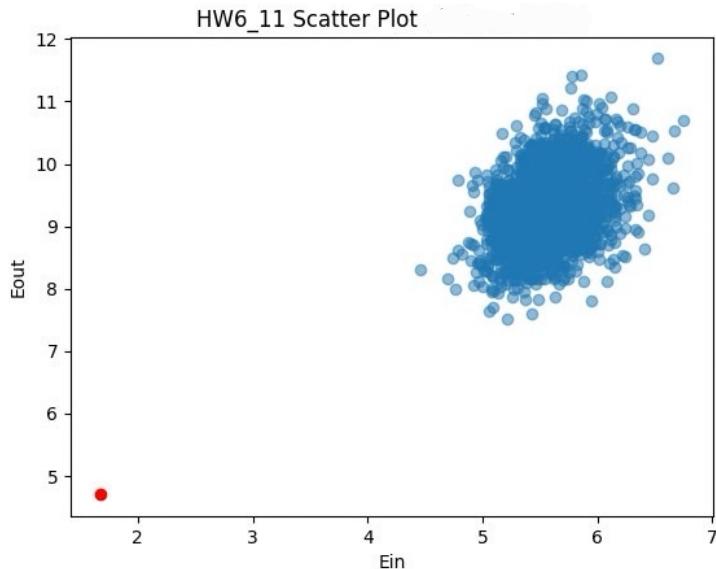
What is the $E_{\text{out}}(g)$, where g is the unpruned decision tree returned from your C&RT algorithm and E_{out} is evaluated using the squared error?

```
Eout(g): 8.79324462640737
```

10. (20 points, *) Next, we implement the random forest algorithm by coupling bagging (by sampling with replacement) with $N' = 0.5N$ with your unpruned decision tree in the previous problem. You do not need to do random feature selection or expansion. Produce $T = 2000$ trees with bagging. Let $g_1, g_2, \dots, g_{2000}$ denote the 2000 trees generated. Plot a histogram of $E_{\text{out}}(g_t)$, where E_{out} is evaluated using the squared error.



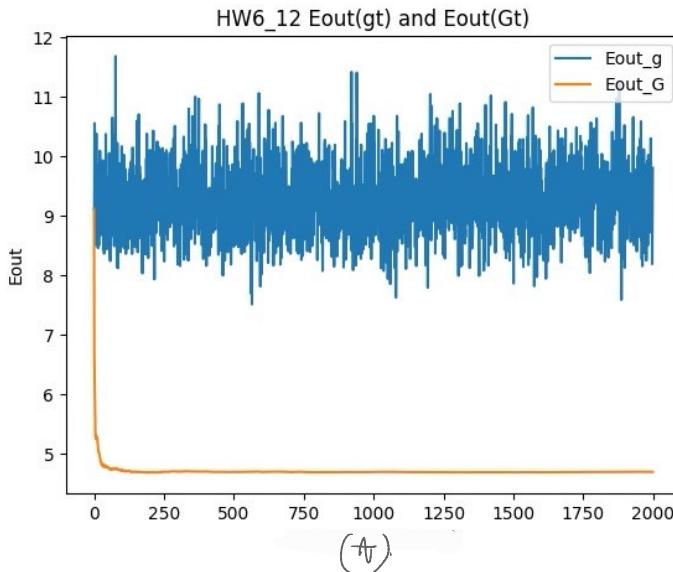
11. (20 points, *) Let $G(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T g_t(\mathbf{x})$ be the random forest formed by the trees above. Plot a scatter plot of $(E_{\text{in}}(g_t), E_{\text{out}}(g_t))$ along with a clear mark of where $(E_{\text{in}}(G), E_{\text{out}}(G))$ is. Describe your findings.



Through this plot, we can clearly find out that $\text{Ein}(G) \gg \text{Ein}(g)$, indicating that forest really reduce training data error . In concern, we should make sure that this algorithm may overfit.

However, $\text{Eout}(G)$ is extremely smaller than $\text{Eout}(g)$, which means that it will also act better on test data.

12. (20 points, *) Let $G_t(\mathbf{x}) = \frac{1}{t} \sum_{\tau=1}^t g_\tau(\mathbf{x})$ be the random forest formed by the first t trees. Plot the $E_{\text{out}}(g_t)$ as a function of t , and plot $E_{\text{out}}(G_t)$ as a function of t on the same figure. Describe your findings.



In the beginning, G act similiarly as g, they both make errors in same situation.
 However;after adding more g into forest, Gout becoming more stable and produce less errors on testing data.
 As we can see, as time go longer, the decay on errors slower than beginning.
 At $t=2000$, stable errors probably stop at $2\sim 3$, comparing beginning error ~ 10 , verified that forest improve the whole algorithm.