5.

$$G(x) = \text{sign}\left(\sum_{t=1}^{2M+1} g_t(x)\right)$$

$$E_{out}(G) = [\![ y \neq \text{sign}\left(\sum_{t=1}^{2M+1} g_t(x)\right) ]\!]$$

at least $M+1$ classifiers $y \neq \text{sign}(g_t(x))$   s.t. $E(G)=1$

$$\Rightarrow \sum_{t=1}^{2M+1} e_t \geq M+1$$

$$\Rightarrow 1 \leq \frac{1}{M+1} \sum_{t=1}^{2M+1} e_t$$

$$\Rightarrow E(G) \leq \frac{1}{M+1} \sum_{t=1}^{2M+1} e_t$$

collaborator :
B11901040 項蓮均

6.

collaborator:
B11901040 項蓮均

Incorrect: $u_n^{(t+1)} \leftarrow u_n^{(t)} \times \sqrt{\dfrac{1-\epsilon_t}{\epsilon_t}}$

Correct: $u_n^{(t+1)} \leftarrow u_n^{(t)} \Big/ \sqrt{\dfrac{1-\epsilon_t}{\epsilon_t}}$

$$U_{t+1} = \sum_{n=1}^{N} u_n^{(t+1)}$$

$$= \sum_{\text{incorrect } n} \left( u_n^{(t)} \times \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} + \sum_{\text{correct } n} u_n^{(t)} \Big/ \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} \right)$$

$$= \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} \sum_{\text{incorrect } n} u_n^{(t)} + \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} \sum_{\text{correct } n} u_n^{(t)}$$

$$= \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} \, U_t \times \epsilon_t + \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} \, U_t (1-\epsilon_t)$$

$$= U_t \left( \sqrt{\epsilon_t(1-\epsilon_t)} + \sqrt{\epsilon_t(1-\epsilon_t)} \right)$$

$$= 2\sqrt{\epsilon_t(1-\epsilon_t)} \, U_t$$

$$\therefore \quad \frac{U_{t+1}}{U_t} = 2\sqrt{\epsilon_t(1-\epsilon_t)}$$

Q.E.D.

**7.**

$$\min_{\alpha} \frac{1}{N} \sum_{n=1}^{N} \left( (y_n - s_n) - \alpha \, g_t(\underline{x_n}) \right)^2$$

$$\frac{1}{N} \sum_{n=1}^{N} \left( (y_n - s_n) - \alpha \, g_t(\underline{x_n}) \right)^2$$

$$= \frac{1}{N} \sum_{n=1}^{N} \left( (y_n - s_n)^2 - 2\alpha (y_n - s_n) g_t(\underline{x_n}) + \alpha^2 g_t^2(\underline{x_n}) \right)$$

$$\frac{\partial}{\partial \alpha} \frac{1}{N} \sum_{n=1}^{N} \left( (y_n - s_n)^2 - 2\alpha (y_n - s_n) g_t(\underline{x_n}) + \alpha^2 g_t^2(\underline{x_n}) \right) = 0$$

$$\Rightarrow \sum_{n=1}^{N} \left( -2(y_n - s_n) g_t(\underline{x_n}) + 2\alpha \, g_t^2(\underline{x_n}) \right) = 0$$

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

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

From the solution of linear regression, $((y_n - s_n) - g_t(\underline{x_n}))$ are orthogonal to $g_t(\underline{x_n})$

$$\Rightarrow \sum_{n=1}^{N} \left( (y_n - s_n) - g_t(\underline{x_n}) \right) \cdot g_t(\underline{x_n}) = 0$$

$$\Rightarrow \sum_{n=1}^{N} (y_n - s_n) g_t(\underline{x_n}) - \sum_{n=1}^{N} g_t^2(\underline{x_n}) = 0$$

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

$$\therefore \alpha_1 = \frac{\sum_{n=1}^{N} (y_n - s_n) g_t(\underline{x_n})}{\sum_{n=1}^{N} g_t^2(\underline{x_n})} = 1 \qquad Q.E.D.$$

8.

collaborator :
B11901040 項達均

perform linear regression on $\{(g_t(\underline{x}_n), y_n - s_n)\}_{n=1}^{N}$

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

$$\xrightarrow{\frac{\partial}{\partial \eta}} \frac{1}{N} \sum_{n=1}^{N} 2(-g_t(\underline{x}_n))((y_n - s_n) - \eta g_t(\underline{x}_n)) = 0$$

$$\Rightarrow \sum_{n=1}^{N} g_t(\underline{x}_n)((y_n - s_n) - \eta g_t(\underline{x}_n)) = 0$$

$\eta = \alpha_t$

update  $s_n \longleftarrow s_n + \alpha_t g_t(\underline{x}_n)$

$$\therefore \sum_{n=1}^{N} g_t(\underline{x}_n)(y_n - (s_n + \alpha_t g_t(\underline{x}_n))) = 0$$

$$\Rightarrow \sum_{n=1}^{N} g_t(\underline{x}_n)(y_n - s_n) = 0 \qquad Q.E.D.$$

9.

$\because$ one hidden layer including the output neuron

$\therefore L = 2$

$$\begin{cases} s_j^{(2)} = \sum_{i=0}^{d^{(1)}} w_{ij}^{(2)} x_i^{(1)} \\ s_j^{(1)} = \sum_{i=0}^{d^{(0)}} w_{ij}^{(1)} x_i^{(0)} \quad , \quad s_{j+1}^{(1)} = \sum_{i=0}^{d^{(0)}} w_{i\,j+1}^{(1)} x_i^{(0)} \\ \delta_j^{(2)} = -2(y_n - s_j^{(2)}) \end{cases}$$

$w_{ij}^{(1)} = 0.5 \quad \forall\, i \text{ and } 1 \le j < d^{(1)}$

$$\begin{aligned} \delta_j^{(1)} &= \sum_k \delta_k^{(2)} w_{jk}^{(2)} \tanh'(s_j^{(1)}) \\ &= \sum_k \delta_k^{(2)} w_{j+1\,k}^{(2)} \tanh'(s_{j+1}^{(1)}) \\ &= \delta_{j+1}^{(1)} \end{aligned}$$
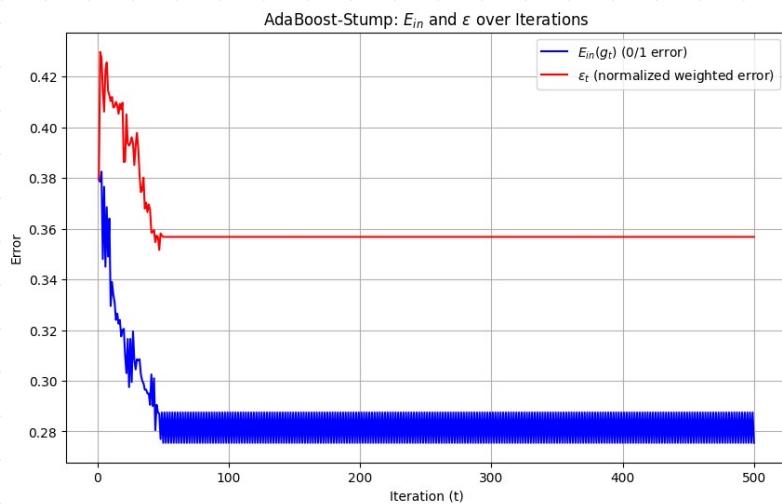
for the backprop algorithm,

$$w_{ij}^{(l)} = w_{ij}^{(l)} - \eta\, x_i^{(l-1)} \delta_j^{(l)}$$

$$\begin{aligned} w_{ij}^{(1)} &= w_{ij}^{(1)} - \eta\, x_i^{(0)} \delta_j^{(1)} \\ &= w_{i\,j+1}^{(1)} - \eta\, x_i^{(0)} \delta_{j+1}^{(1)} \\ &= w_{i\,j+1}^{(1)} \end{aligned}$$

Q.E.D.

10. Collaborators: B11901073 林禹融

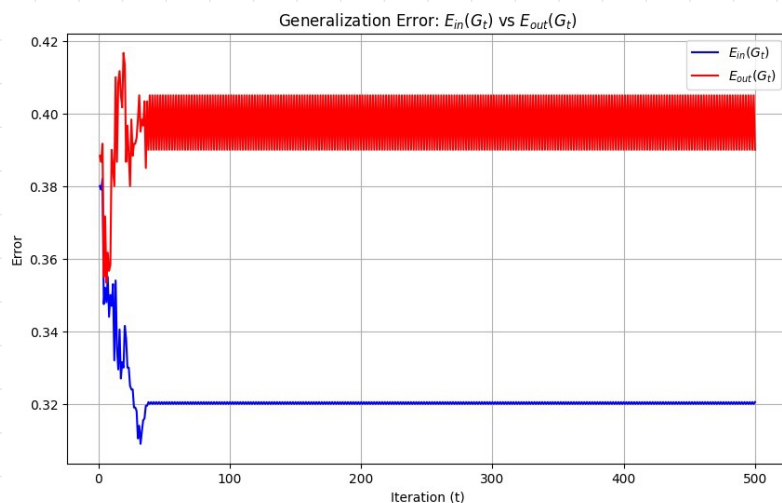AdaBoost-Stump: $E_{in}$ and $\varepsilon$ over Iterations



The AdaBoost algorithm reduces $E_{in}(g_t)$ and $E_t$ quickly and remain same value respectively.

```python
  9  class DecisionStump:
 10      def __init__(self):
 11          self.feature_index = None
 12          self.threshold = None
 13          self.sign = None
 14
 15      def train(self, X, y, weights):
 16          _, n = X.shape
 17          best_error = float('inf')
 18
 19          # Search for the best dimension, threshold, and sign
 20          for feature_index in range(n):
 21              thresholds = np.unique(X[:, feature_index])
 22              for threshold in thresholds:
 23                  for sign in [-1, 1]:
 24                      predictions = sign * np.sign(X[:, feature_index] - threshold)
 25                      error = np.sum(weights * (predictions != y))
 26                      if error < best_error:
 27                          best_error = error
 28                          self.feature_index = feature_index
 29                          self.threshold = threshold
 30                          self.sign = sign
 31
 32          return best_error
 33
 34      def predict(self, X):
 35          predictions = self.sign * np.sign(X[:, self.feature_index] - self.threshold)
 36          return predictions
 37
 38  # AdaBoost Algorithm
 39  def adaBoost(X_train, y_train, X_test, y_test, T=500):
 40      m = X_train.shape[0]
 41      weights = np.ones(m) / m
 42      classifiers = []
 43      alphas = []
 44
 45      # Metrics for plotting
 46      E_in = []
 47      epsilons = []
 48
 49      for t in range(T):
 50          stump = DecisionStump()
 51          error = stump.train(X_train, y_train, weights)
 52
 53          # Avoid divide-by-zero
 54          if error == 0:
 55              error = 1e-10
 56
 57          # Compute alpha
 58          alpha = 0.5 * np.log((1 - error) / error)
 59          alphas.append(alpha)
 60          classifiers.append(stump)
 61
 62          # Update weights
 63          predictions = stump.predict(X_train)
 64          weights *= np.exp(-alpha * y_train * predictions)
 65
 66          # Compute 0/1 error (E_in) and epsilon
 67          E_in.append(zero_one_loss(y_train, np.sign(sum(alpha * clf.predict(X_train) for clf, alpha in zip(classifiers, alphas)))))
 68          epsilons.append(error)
 69
 70          print(f"Iteration {t+1}: Error = {error:.4f}, Alpha = {alpha:.4f}, E_in = {E_in[-1]:.4f}")
 71
 72      # Final testing error
 73      final_predictions = np.sign(sum(alpha * clf.predict(X_test) for clf, alpha in zip(classifiers, alphas)))
 74      E_out = zero_one_loss(y_test, final_predictions)
 75
 76      print(f"Final Test Error (E_out): {E_out:.4f}")
 77
 78      return E_in, epsilons
 79
```

11.

Collaborators: B11901073 林禹融

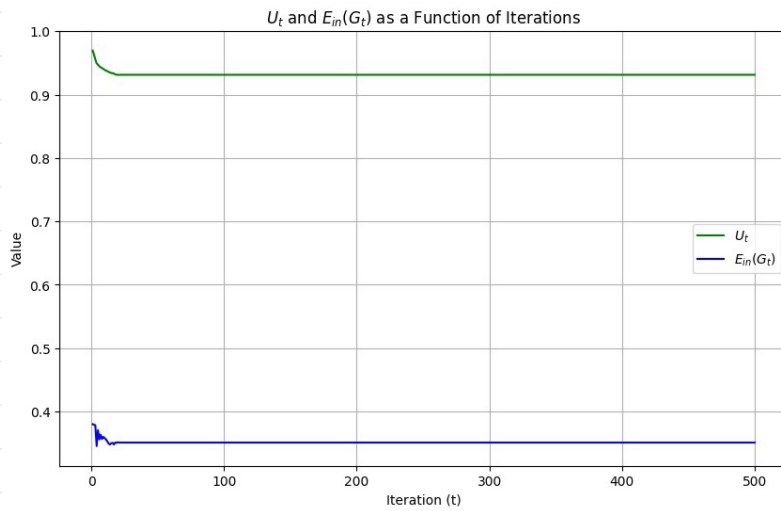Generalization Error: $E_{in}(G_t)$ vs $E_{out}(G_t)$



① Although $E_{in}(G_t)$ is relatively high initially, it decreases fast and remains around 0.32 after about 50 iterations.

② $E_{out}(G_t)$ oscillates between about 0.39 and 0.41 after about 50 iterations. Each new weak classifier added to $G_t$ contributes a small correction based on its $\alpha_t$. If the weak classifier's corrections are small and do not generalize well to the test data, they may slightly shift the decision boundary of $G_t$ in an inconsistent manner. However, the mean value of the oscillation remains about 0.395 after about 50 iterations.

① & ② indicate that after a certain number of iterations, the model reaches its generalization capacity.

# 12. no collaborators

**$U_t$ and $E_{in}(G_t)$ as a Function of Iterations**



$U_t$ decreases steeply in the first few iterations and remains at about 0.93, which suggest that AdaBoost algorithm quickly adjusts the sample weights and coverges efficiently. However, it doesn't reach perfect accuracy.