Q1. S.7 bootstrapping each sample contains about
2/3 unique examples

Say, 'N' observations in total
— we randomly draw 'N' observations
from the dataset w replacement
— ∴ Prob for each select $= 1/N$

Prob of not selecting ⟹ $(1 - 1/N)$
Prob of 1 obs NEVER getting ⟹ ∴ $N \times \left(1 - \frac{1}{N}\right)$
selected for all 'N'.

∴ Draws are independent

$$P(unique) = 1 - \left(1 - \frac{1}{N}\right)^N \quad = 1 - \left(\frac{N-1}{N}\right)^N$$

Sub w N=3
for simplicity ⟹ $1 - \left(\frac{2}{3}\right)^3 = 1 - \frac{8}{27} = \frac{19}{27} \approx 0.70$

∴ 70% ⟹ unique → $\approx \frac{2}{3}$

Q2) Simple majority among K learned hyp.
Each hyp has error $\varepsilon$
$\hookrightarrow$ maybe not independent of each other

Is the error of the ensemble never worse than $\varepsilon$?

It can be worse than $\varepsilon$.

Proof:
Say we have 5 cases, and 3 hypothesis
Errors are:
i) Hypothesis 1 fails for cases 1, 3, 4 $\left. \right\}$ $\varepsilon = 3/5$
ii) Hypothesis 2 fails for cases 1, 2, 5
iii) Hypothesis 3 fails for cases 2, 3, 4

$\therefore$ The simple majority decides
1, 2, 3, 4 $\longrightarrow$ fails erroreously which means
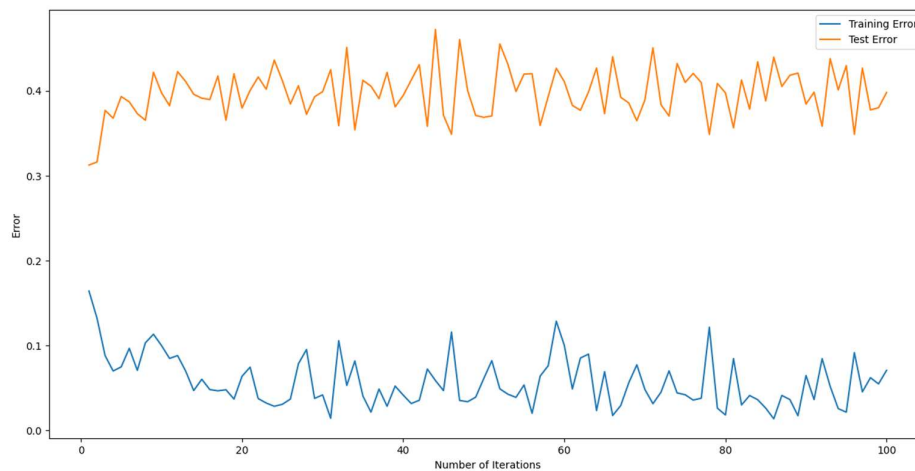Overall error $= 4/5 = \varepsilon'$

$$\boxed{\varepsilon < \varepsilon'}$$

Question 3:



Question 4:

**b) Plot the training error as well as test error and discuss its behavior.**
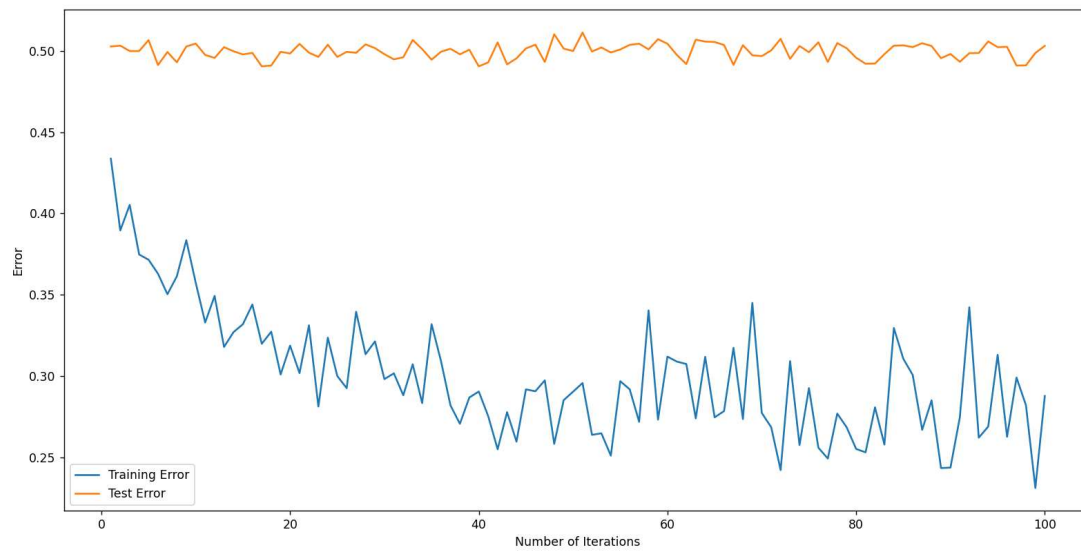


As expected, training error is lower than test error and although there are irregularities throughout, there is a reduction in both errors.

**c) Investigating the number of iterations needed for the test error to start rising**
It starts rising from the first iteration, at least in my code.

**d) Repeat the AdaBoost experiments with new dataset**



The reduction in training error is much more evident and thus the model gets better with more iterations.