

cs273-hw5-Abdul_Kalam_Syed

November 23, 2025

1 CS273A Homework 5

1.1 Due: Wednesday Nov 19 2025 (11:59pm)

1.2 Instructions

This homework (and subsequent ones) will involve data analysis and reporting on methods and results using Python code. You will submit a **single PDF file** that contains everything to Gradescope. This includes any text you wish to include to describe your results, the complete code snippets of how you attempted each problem, any figures that were generated, and scans of any work on paper that you wish to include. It is important that you include enough detail that we know how you solved the problem, since otherwise we will be unable to grade it.

Your homeworks will be given to you as Jupyter notebooks containing the problem descriptions and some template code that will help you get started. You are encouraged to use these starter Jupyter notebooks to complete your assignment and to write your report. This will help you not only ensure that all of the code for the solutions is included, but also will provide an easy way to export your results to a PDF file (for example, doing *print preview* and *printing to pdf*). I recommend liberal use of Markdown cells to create headers for each problem and sub-problem, explaining your implementation/answers, and including any mathematical equations. For parts of the homework you do on paper, scan it in such that it is legible (there are a number of free Android/iOS scanning apps, if you do not have access to a scanner), and include it as an image in the Jupyter notebook.

Double check that all of your answers are legible on Gradescope, e.g. make sure any text you have written does not get cut off.

If you have any questions/concerns about using Jupyter notebooks, ask us on EdD. If you decide not to use Jupyter notebooks, but go with Microsoft Word or LaTeX to create your PDF file, make sure that all of the answers can be generated from the code snippets included in the document.

Summary of Assignment: 100 total points - Problem 1: Decision Trees by Hand (25 points)
- Problem 1.1: Shannon Entropy (5 points) - Problem 1.2: Information Gain (10 points) - Problem 1.3: Full Tree (10 points) - Problem 2: Decision Trees in Python (34 points) - Problem 2.1: Feature Statistics (8 points) - Problem 2.2: Initial Tree (6 points) - Problem 2.3: Exploring Depth Control (10 points) - Problem 2.4: Exploring Leaf Size (10 points) - Problem 3: Random Forests (20 points) - Problem 3.1: Training Members (10 points) - Problem 3.2: Ensemble Prediction (10 points) - Problem 4: VC Dimension (16 points) - Problem 4.1: Model A (4 points) - Problem 4.2: Model B (4 points) - Problem 4.3: Model C (4 points) - Problem 4.4: Model D (4 points) - Statement of Collaboration (5 points)

Before we get started, let's import some libraries that you will make use of in this assignment. Make sure that you run the code cell below in order to import these libraries.

Important: In the code block below, we set `seed=1234`. This is to ensure your code has reproducible results and is important for grading. Do not change this. If you are not using the provided Jupyter notebook, make sure to also set the random seed as below.

Important: Do not change any codes we give you below, except for those waiting for you to complete. This is to ensure your code has reproducible results and is important for grading.

```
[1]: import numpy as np
      import matplotlib.pyplot as plt

      import requests                                     # reading data
      from io import StringIO

      from sklearn.datasets import fetch_openml          # common data set access
      from sklearn.preprocessing import StandardScaler    # scaling transform
      from sklearn.model_selection import train_test_split # validation tools
      from sklearn.metrics import zero_one_loss as J01

      import sklearn.tree as tree

      # Fix the random seed for reproducibility
      # !! Important !! : do not change this
      seed = 1234
      np.random.seed(seed)
```

1.3 Problem 1: Decision Trees for Spam

In order to reduce my email load, I decide to implement a machine learning algorithm to decide whether or not I should read an email, or simply file it away instead. To train my model, I obtain the following data set of binary-valued features about each email, including whether I know the author or not, whether the email is long or short, and whether it has any of several key words, along with my final decision about whether to read it ($y = +1$ for “read”, $y = -1$ for “discard”).

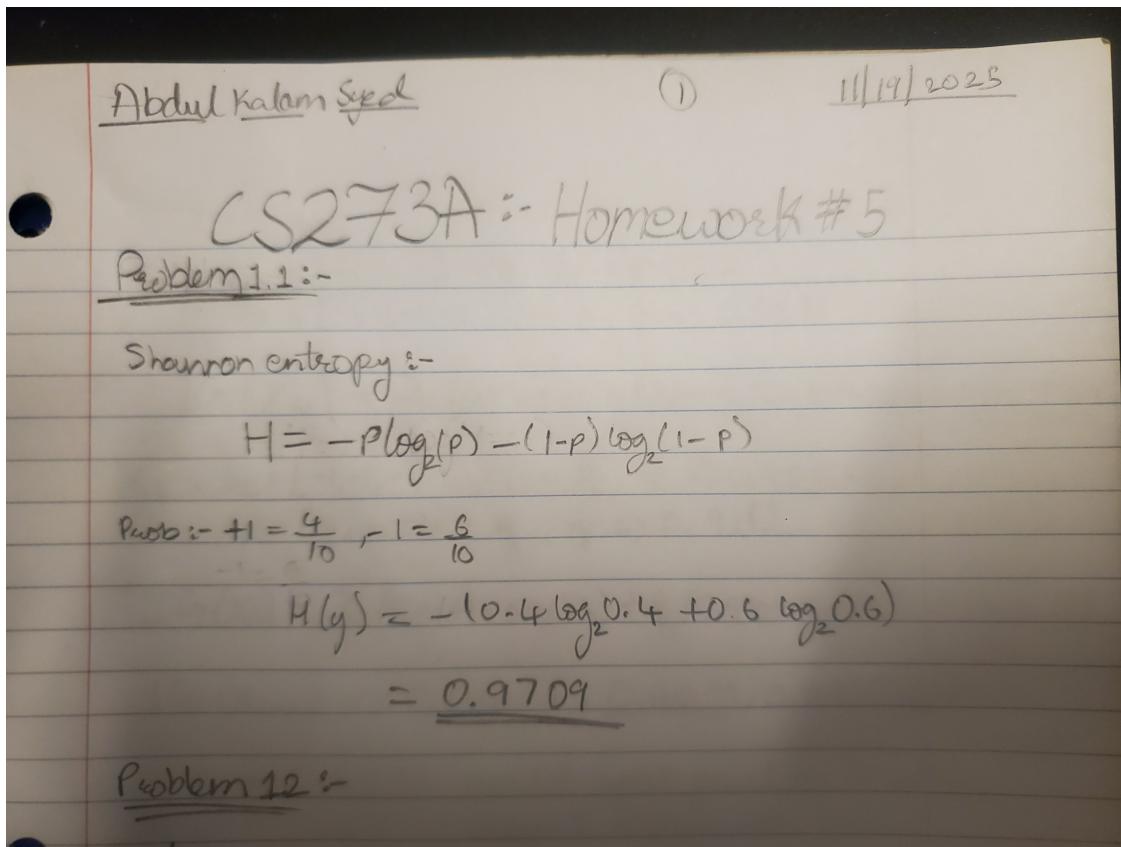
X1	X2	X3	X4	X5	y
(know author?)	(is long?)	(has ‘research’?)	(has ‘grade’?)	(has ‘lottery’?)	(read?)
0	0	1	1	0	-1
1	1	0	1	0	-1
0	1	1	1	1	-1
1	1	1	1	0	-1
0	1	0	0	0	-1
1	0	1	1	1	1

X1	X2	X3	X4	X5	y
0	0	1	0	0	1
1	0	0	0	0	1
1	0	1	1	0	1
1	1	1	1	1	-1

In the case of any ties where both classes have equal probability, we will prefer to predict class +1.
Solve the following problems “by hand” (you can use python for logarithms, etc.)

1.3.1 Problem 1.1 (5 points)

Calculate the Shannon entropy $H(y)$ of the binary class variable y , in bits. **Hint:** Your answer should be a number between 0 and 1.



1.3.2 Problem 1.2 (5 points)

Calculate the information gain for each feature x_i . Which feature should be split first?

Problem 12 :-

$$H(y) = 0.9709$$

$$x_1 \Rightarrow S_L = \{i : x_1^{(i)} = 0\} \Rightarrow [-1, -1, -1, 1] \Rightarrow \frac{3}{4}, \frac{1}{4} H(P_{S_L}) = 0.8113$$

$$S_R = \{i : x_1^{(i)} = 1\} \Rightarrow [-1, -1, 1, 1, -1] \Rightarrow \frac{3}{5}, \frac{2}{5} H(P_{S_R}) = 1 - 0.8113$$

$$IG(x_1) = \frac{4}{10} (0.9709 - 0.8113) + \frac{6}{10} (0.9709 - 1) = \underline{\underline{0.0464}}$$

$$x_2 \Rightarrow S_L = [-1, 1, 1, 1, 1] \Rightarrow \frac{1}{5}, \frac{4}{5} \Rightarrow 0.7219$$

$$S_R = [-1, -1, -1, 1, -1] = 1, 0 \Rightarrow 0$$

$$IG(x_2) = \frac{5}{10} (0.9709 - 0.7219) + \frac{5}{10} (0.9709 - 0) = \underline{\underline{0.6099}}$$

$$x_3 = S_L = [-1, -1, 1] \Rightarrow \frac{2}{3}, \frac{1}{3} \Rightarrow 0.9183$$

$$S_R = [-1, -1, -1, 1, 1, -1] \Rightarrow \frac{4}{7}, \frac{3}{7} \Rightarrow 0.9852$$

$$IG(x_3) = \frac{3}{10} (0.9709 - 0.9183) + \frac{7}{10} (0.9709 - 0.9852) = \underline{\underline{0.0058}}$$

$$x_4 = S_L = [-1, 1, 1] = \frac{1}{3}, \frac{2}{3} \Rightarrow 0.9183$$

$$S_R = [-1, -1, -1, -1, 1, 1, -1] \Rightarrow \frac{5}{7}, \frac{2}{7} \Rightarrow 0.8631$$

$$IG(x_4) = \frac{3}{10} (0.9709 - 0.9183) + \frac{7}{10} (0.9709 - 0.8631) = \underline{\underline{0.09124}}$$

$$x_5 = S_L = [-1, -1, -1, -1, 1, 1, 1] = \frac{4}{7}, \frac{3}{7} \Rightarrow 0.9852$$

$$S_R = [-1, 1, -1] = \frac{2}{3}, \frac{1}{3} \Rightarrow 0.9183$$

$$IG(x_5) = \frac{7}{10} (0.9709 - 0.9852) + \frac{3}{10} (0.9709 - 0.9183) = \underline{\underline{0.0058}}$$

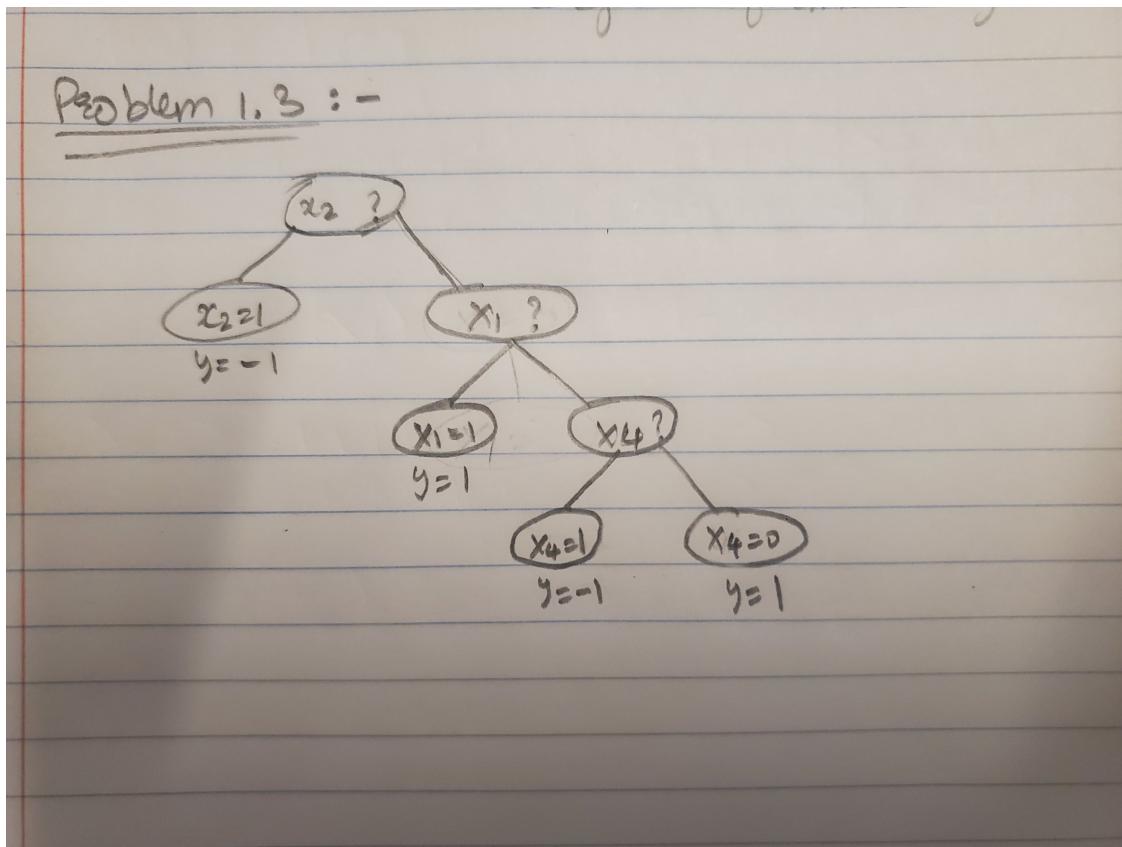
$$IG(x_{1-5}) \Rightarrow 0.0464, \underline{\underline{0.6099}}, 0.0058, 0.09124, \\ 0.0058$$

$x_2 = 0.6099$ should split first as it is the largest information gain.

Problem 1.3 :-

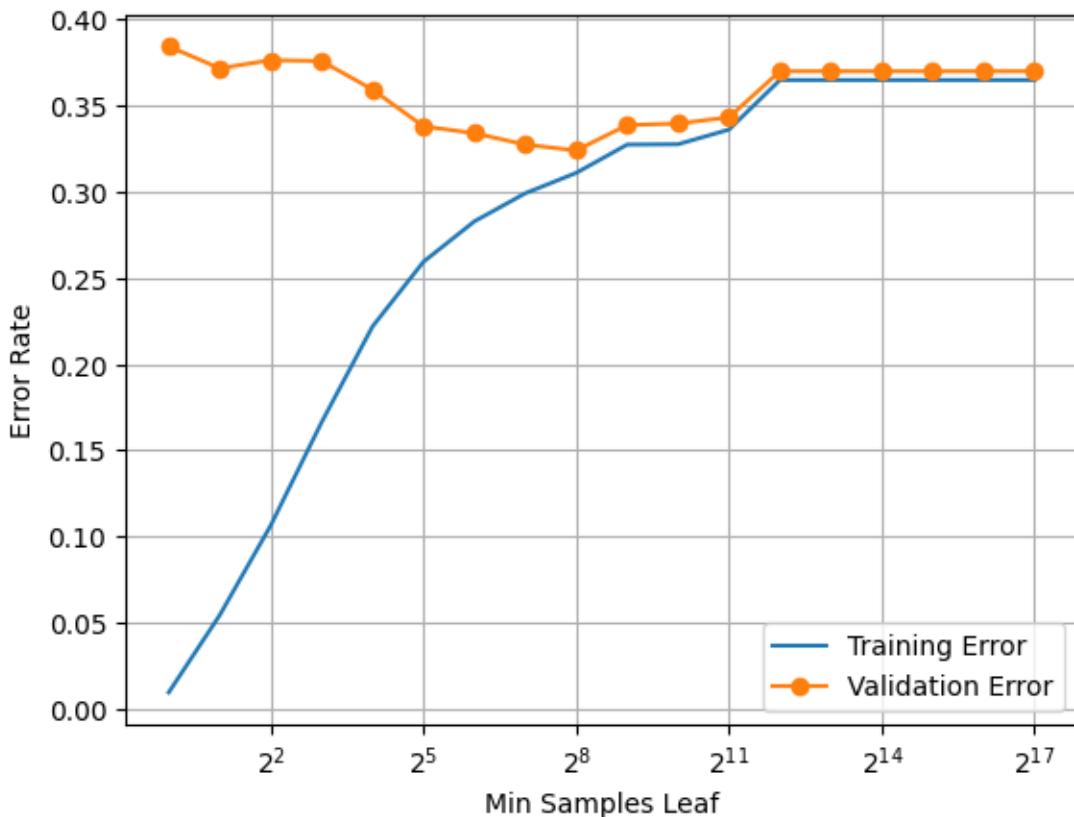
1.3.3 Problem 1.3 (5 points)

Draw (or otherwise illustrate) the complete decision tree that will be learned from these data.



 2^8 = 256.")
print("This model is worse than the depth controlled model.")

```



Models with higher `min_samples_leaf` have lower complexity because they require more samples to split a node, meaning they create fewer splits.

The best choice of `min_samples_leaf` is the one that gives the lowest validation error. And based on the plot, it seems to be at `min_samples_leaf = 2^8 = 256`. This model is worse than the depth controlled model.

```
data:image/s3,anthropic-data-us-east-2/u/marker_images/0011/0101/0101/10101011/sfishman-chandramapper-0319194802/1d7d2ece4878e46ec94fe4493adccab5.jpg</antml:image>

(Note that no three of the four points are co-linear.)

For each of the learners listed below, answer **(a)** which of the four data sets **can be shattered** by a learner of that form? Give a brief explanation / justification (1-2 sentences). Then use your results to **(b)** guess the VC dimension of the classifier (you do not have to give a formal proof, just your reasoning).

In each learner, the parameters a, b, c, \dots are real-valued scalars, and $T[z]$ is the sign threshold function, i.e., $T[z] = +1$ for $z \geq 0$ and $T[z] = -1$ for $z < 0$.

- $T(a + b x_1)$

```
[ ]: # (a) which data sets & why?
# the first dataset can be shattered because a single point can be labeled as +1
# or -1
# the second dataset can be shattered because two points can also be separated
# the third and fourth data sets cannot be shattered because there is no way to
# separate the points

# (b) VC dimension guess?
# The VC dimension is 2, because it can shatter a maximum of 2 points.
```

- $T((a * b) x_1 + (c/a) x_2)$

```
[ ]: # (a) which data sets & why?
# the first dataset can be shattered because a single point can be separated
# the second dataset can be shattered because two points can also be separated
# the third dataset can be shattered because three points can also be separated
# as any three points are not collinear.
# the fourth dataset cannot be shattered because there is no way to separate
# the points

# (b) VC dimension guess?
# The VC dimension is 3, because it can shatter a maximum of 3 points.
```

- $T((x_1 - a)^2 + (x_2 - b)^2 - c)$

```
[ ]: # (a) which data sets & why?
# the first dataset can be shattered because a single point can be labeled as
↪+1 or -1
# The second dataset can be shattered because a circle can include neither, ↪
↪one, or both points, so all four labelings are possible.
# the third dataset can be shattered because any three points can be assigned ↪
↪all possible labelings using circles.
# the fourth dataset cannot be shattered because there is at least one labeling ↪
↪of four points in general position that no single circle can represent.

# (b) VC dimension guess?
# The VC dimension is 3
```

- $T(a + bx_1 + cx_2) \times T(d + bx_1 + cx_2)$
 - **Hint:** the two linear equations correspond to two parallel lines, since both have the same coefficients for x_1 and x_2 . Then, $T(z) \times T(z') = +1$ if and only if z and z' have the same sign.

```
[ ]: # (a) which data sets & why?
# the first dataset can be shattered because a single point can be labeled as
↪+1 or -1
# the second dataset can be shattered because the strip can include neither, ↪
↪one, or both points, so all four labelings are possible.
# the third dataset can be shattered because the adjustable linear boundary can ↪
↪separate any three points
# the fourth dataset cannot be shattered because there is at least one labeling ↪
↪of four points that cannot be represented

# (b) VC dimension guess?
# The VC dimension is 3
```

<img src="data:image/svg+xml,%3C%3Fxml%20version%3D%221.0%22%20encoding%3D%22UTF-8%22%20standa

1.6.1 Statement of Collaboration (5 points)

It is **mandatory** to include a Statement of Collaboration in each submission, with respect to the guidelines below. Include the names of everyone involved in the discussions (especially in-person ones), and what was discussed.

All students are required to follow the academic honesty guidelines posted on the course website. For programming assignments, in particular, I encourage the students to organize (perhaps using EdD) to discuss the task descriptions, requirements, bugs in my code, and the relevant technical content before they start working on it. However, you should not discuss the specific solutions, and, as a guiding principle, you are not allowed to take anything written or drawn away from these discussions (i.e. no photographs of the blackboard, written notes, referring to EdD, etc.). Especially after you have started working on the assignment, try to restrict the discussion to EdD as much as possible, so that there is no doubt as to the extent of your collaboration.

I have not collaborated with anyone on this assignment.