

## CS 184A/284A: Artificial Intelligence in Biology and Medicine

**Homework 1**Due date: **Oct. 17, 2025**

Instructor: Xiaohui Xie

This homework (and many subsequent ones) will involve data analysis and reporting on methods and results using Python code. We have provided a **Jupyter/iPython notebook** which you need to write both your code and reasoning for each question on the cell specific to that question.

You have to submit **a single PDF file** that contains everything to Gradescope, and associated each page of the PDF to each problem. This includes any text you wish to include to describe your results, the complete code snippets of how you attempted each problem, and any figures that were generated. 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.

Jupyter/iPython notebooks provide an easy way to **export your results to a PDF** file<sup>1</sup>. 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 iPython notebook<sup>2</sup>. If you have any questions/concerns about using iPython, ask us on Ed Discussion.

**Summary so far:** (1) submit a single, standalone PDF report, with all code; (2) Use provided Jupyter notebooks.

**Points:** This homework adds up to a total of **55 points**, as follows:

Problem 0: Get Connected	5 points
Problem 1: Python & Data Exploration	20 points
Problem 2: kNN Predictions	25 points
Statement of Collaboration	5 points

## Problem 0: Get Connected (5 points)

Please visit our class forum on Ed Discussion. Ed Discussion will be the place to post your questions and discussions, rather than by email to me or the TAs, since chances are that other students have the same or similar questions, and will be helped by seeing the discussion. Remember, your Ed Discussion participation will be taken into account for the participation grade as well. You do not need to mention anything regarding this in the report, we will be able to check whether you have visited Ed Discussion or not.

## Problem 1: Python & Data Exploration (20 points)

In this problem, we will explore some basic statistics and visualizations of an example data set. First, install the required libraries (**NumPy**, **matplotlib**, and **Scikit-Learn**), then download the zip file for Homework 1, which contains the “Fisher iris” data set, and load the latter into Python:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 iris = np.genfromtxt("data/iris.txt", delimiter=None) # load the text file
5 Y = iris[:, -1] # target value is the last column
6 X = iris[:, 0:-1] # features are the other columns
```

<sup>1</sup>For example, by doing a **Print Preview** in Chrome and **printing** it to a PDF.

<sup>2</sup>Tips from Gradescope: [http://gradescope-static-assets.s3-us-west-2.amazonaws.com/help/submitting\\_hw\\_guide.pdf](http://gradescope-static-assets.s3-us-west-2.amazonaws.com/help/submitting_hw_guide.pdf)

The Iris data consist of four real-valued features used to predict which of three types of iris flower was measured (a three-class classification problem).

1. Use `X.shape` to get the number of features and the data points. Report both numbers, mentioning which number is which. (5 points)
2. For each feature, plot a histogram (`plt.hist`) of the data values (5 points)
3. Compute the mean & standard deviation of the data points for each feature (`np.mean`, `np.std`) (5 points)
4. For each pair of features (1,2), (1,3), and (1,4), plot a scatterplot (see `plt.plot` or `plt.scatter`) of the feature values, colored according to their target value (class). (For example, plot all data points with  $y = 0$  as blue,  $y = 1$  as green, etc.) (5 points)

## Problem 2: kNN predictions (25 points)

In this problem, you will continue to use the Iris data and explore a KNN classifier using `KNeighborsClassifier` class from scikit-learn. While doing the problem, please explore the implementation to become familiar with how it works.

First, we will split the data into training and validation subsets:

```
1 iris = np.genfromtxt("data/iris.txt", delimiter=None) # load the data
2 Y = iris[:, -1]
3 X = iris[:, 0:-1]
4
5 # Note: indexing with ":" indicates all values (in this case, all rows);
6 # indexing with a value ("0", "1", "-1", etc.) extracts only that value (here, columns);
7 # indexing rows/columns with a range ("1:-1") extracts any row/column in that range.
8
9 from sklearn.model_selection import train_test_split
10
11 Xtr, Xva, Ytr, Yva = train_test_split(
12     X, Y, test_size=0.25, random_state=0, stratify=Y, shuffle=True
13 ) # split data into 75/25 train/validation
```

You may also find it useful to set the random number seed at the beginning (in general, for every assignment), e.g., `numpy.random.seed(0)`, to ensure consistent behavior each time.

You can build now and `train` a kNN classifier on `Xtr, Ytr` and make predictions on some data `Xva` with it:

```
1 from sklearn.neighbors import KNeighborsClassifier
2
3 K = 1 # K is an integer, e.g. 1 for nearest neighbor prediction
4 knn = KNeighborsClassifier(n_neighbors=K) # create the knn classifier model
5 knn.fit(Xtr, Ytr) # train the classifier
6 Yva_hat = knn.predict(Xva) # get estimates of y for each data point in Xva
```

If your data are 2D, you can visualize a data set and a classifier's decision regions using the following function:

```
1
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from matplotlib.colors import ListedColormap
5
6 def plotClassify2D(model, X2d, Y, grid_step=0.02, padding=0.5):
7     """
8     Plot 2D decision regions for a trained classifier with consistent colors.
9
10    Parameters
```

```

11 -----
12 model : fitted classifier with .predict()
13 X2d : (n_samples, 2)
14 Y : (n_samples,)
15 grid_step : mesh resolution
16 padding : extra margin around data
17 """
18 # Define consistent color maps
19 unique_classes = np.unique(Y)
20 n_classes = len(unique_classes)
21
22 # Use a subset of matplotlib's default tab10 colormap
23 cmap = plt.cm.get_cmap('tab10', n_classes)
24 colors = [cmap(i) for i in range(n_classes)]
25 cmap_background = ListedColormap(colors)
26
27 # Create grid
28 x_min, x_max = X2d[:, 0].min() - padding, X2d[:, 0].max() + padding
29 y_min, y_max = X2d[:, 1].min() - padding, X2d[:, 1].max() + padding
30 xx, yy = np.meshgrid(
31     np.arange(x_min, x_max, grid_step),
32     np.arange(y_min, y_max, grid_step)
33 )
34
35 # Predict over grid
36 Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
37 Z = Z.reshape(xx.shape)
38
39 # Plot decision regions
40 plt.contourf(xx, yy, Z, alpha=0.3, cmap=cmap_background)
41
42 # Plot points with matching colors
43 for i, c in zip(unique_classes, colors):
44     plt.scatter(X2d[Y == i, 0], X2d[Y == i, 1],
45                 c=[c], edgecolor='k', s=20, label=f"class {i}")
46
47 plt.xlim(xx.min(), xx.max())
48 plt.ylim(yy.min(), yy.max())
49 plt.legend(frameon=False)
50
51 plotClassify2D(knn, Xtr, Ytr) # make 2D classification plot with data (Xtr,Ytr)

```

This function plots the training data and colored points as per their labels, then calls `knn`'s `predict` function on a densely spaced grid of points in the 2D space, and uses this to produce the background color.

1. Modify the code listed above to use only the first two features of  $X$  (e.g., let  $X$  be only the first two columns of `iris`, instead of the first four), and visualize (plot) the classification boundary for varying values of  $K = [1, 5, 10, 50]$  using `plotClassify2D`. (10 points)
2. Again using only the first two features, compute the error rate (number of misclassifications) on both the training and validation data as a function of  $K = [1, 2, 5, 10, 50, 100, 200]$ . You can do this most easily with a for-loop:

```

1 K=[1,2,5,10,50,100];
2 for i,k in enumerate(K):
3     learner = KNeighborsClassifier(...) # TODO: complete code to train model
4     Yhat = learner.predict(...) # TODO: predict results on training data
5     errTrain[i] = ... # TODO: count what fraction of predictions are wrong
6     #TODO: repeat prediction / error evaluation for validation data
7

```

```
8 plt.semilogx(... #TODO: average and plot results on semi-log scale
9
```

Plot the resulting error rate functions using a semi-log plot (`semilogx`), with training error in red and validation error in green. Based on these plots, what value of  $K$  would you recommend? (10 points)

3. Provide the same plots as the previous, but with all the features in the dataset. Are the plots very different? Is your recommendation different? (5 points)

## 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 Ed Discussion) 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 Ed Discussion, etc.). Especially **after** you have started working on the assignment, try to restrict the discussion to Ed Discussion as much as possible, so that there is no doubt as to the extent of your collaboration.