

Homework 1

Instructions:

- You may discuss problems with your study group, but ultimately all your work (mathematical problems, code, experimental details) must be individual.
- Your solutions must be **typeset** and uploaded to Gradescope by 11.59PM on Thursday April 10. No late homeworks will be accepted under any circumstances, so you are encouraged to upload early.
- A subset of the problems will be graded.

Conceptual and mathematical problems

1. *Casting an image into vector form.* We have a 15×15 image that we would like to represent as a vector in \mathbb{R}^d .
 - (a) If the image is greyscale, we can use one coordinate of the vector for each pixel. In this case, what is d ?
 - (b) On the other hand, for a color image we would need three coordinates for each pixel (to represent R,G,B values). In this case, what is d ?
2. *Euclidean distance.* What is the Euclidean distance between the following two points in \mathbb{R}^3 ?

$$\begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}, \begin{bmatrix} 3 \\ 2 \\ 1 \end{bmatrix}$$

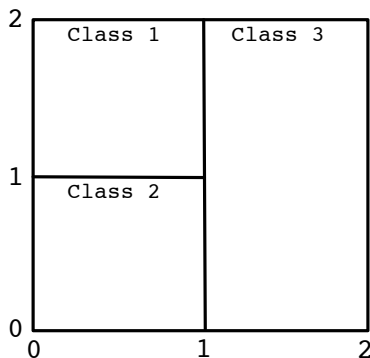
3. *Accuracy of a random classifier.* A particular data set has 4 possible labels, with the following frequencies:

Label	Frequency
A	50%
B	20%
C	20%
D	10%

- (a) What is the error rate of a classifier that picks a label (A, B, C, D) at random, each with probability $1/4$?
 - (b) One very simple type of classifier just returns the same label, always.
 - What label should it return?
 - What will its error rate be?
4. *Decision boundary of the nearest neighbor classifier.* In this problem,

- The data space is $\mathcal{X} = [0, 2]^2$: each point has two coordinates, and they lie between 0 and 2.
- The labels are $\mathcal{Y} = \{1, 2, 3\}$.
- The distance function is ℓ_2 (Euclidean distance).

The correct labels in different parts of \mathcal{X} are as shown below.

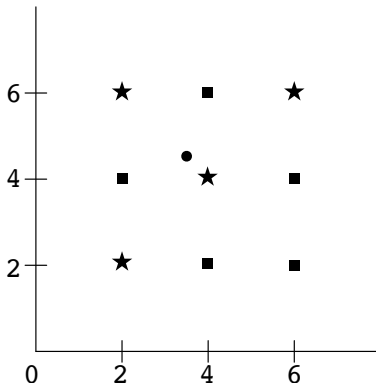


- (a) What is the label of point $\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$?

Now suppose you have a training set consisting of just two points, located at

$$\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}, \begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix}.$$

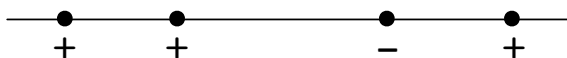
- (b) What label will the nearest neighbor classifier assign to point $\begin{bmatrix} 1.5 \\ 0.5 \end{bmatrix}$?
- (c) What label will the nearest neighbor classifier assign to point $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$?
- (d) Which label will this classifier never predict?
- (e) Now suppose that when the classifier is used, the test points are uniformly distributed over the square \mathcal{X} . What is the error rate of the 1-NN classifier?
5. In the picture below, there are nine training points, each with label either **square** or **star**. These will be used to guess the label of a query point at $\begin{bmatrix} 3.5 \\ 4.5 \end{bmatrix}$, indicated by a circle.



Suppose Euclidean distance is used.

- (a) How will the point be classified by 1-NN? The options are **square**, **star**, or **ambiguous**.
 - (b) By 3-NN?
 - (c) By 5-NN?
6. We decide to use 4-fold cross-validation to figure out the right value of k to choose when running k -nearest neighbor on a data set of size 10,000. When checking a particular value of k , we look at four different training sets. What is the size of each of these training sets?
 7. An extremal type of cross-validation is *n-fold cross-validation* on a training set of size n . If we want to estimate the error of k -NN, this amounts to classifying each training point by running k -NN on the remaining $n - 1$ points, and then looking at the fraction of mistakes made. It is commonly called *leave-one-out cross-validation* (LOOCV).

Consider the following simple data set of just four points:



What is the LOOCV error for 1-NN? For 3-NN?

8. We build a nearest neighbor classifier using a data set of n points in \mathbb{R}^d .
 - If we use Euclidean distance, then computing the distance between two points takes $O(d)$ time.
 - Therefore, computing the nearest neighbor of a test point takes time $O(nd)$ if done using a naive brute-force search.

How long does it take to estimate the error of this classifier using leave-one-out cross-validation (assuming naive search)? Your answer should be a function of n and d .

Programming problems

Before attempting these problems, make sure that Python 3 and Jupyter are installed on your computer. Supporting files are in the archive `hw1.zip`, available from the course website.

9. *Nearest neighbor on MNIST*. The Jupyter notebook `nn-mnist.ipynb` (in the archive file mentioned above) implements a basic 1-NN classifier for a subset of the MNIST data set. It uses a separate training and test set. Begin by going through this notebook, running each segment and taking care to understand exactly what each line is doing.

Now do the following.

- (a) For test point #100 (which is actually the 101st point in the test set due to Python's zero-indexing), print its image as well as the image of its nearest neighbor in the training set. Put these images in your writeup. Is this test point classified correctly?
- (b) The *confusion matrix* for the classifier is a 10×10 matrix N_{ij} with $0 \leq i, j \leq 9$, where N_{ij} is the number of test points whose true label is i but which are classified as j . Thus, if all test points are correctly classified, the off-diagonal entries of the matrix will be zero.
 - Compute the matrix N for the 1-NN classifier and print it out.

- Which digit is misclassified most often? Least often?
- (c) For each digit $0 \leq i \leq 9$: look at all training instances of image i , and compute their mean. This average is a 784-dimensional vector. Use the `show_digit` routine to print out these 10 average-digits.
10. *Classifying back injuries.* In this problem, you will use nearest neighbor to classify patients' back injuries based on measurements of the shape and orientation of their pelvis and spine.

The data set contains information from 310 patients. For each patient, there are: six numeric features (the x) and a label (the y): 'NO' (normal), 'DH' (herniated disk), or 'SL' (spondilolysthesis). We will divide this data into a training set with 250 points and a separate test set of 60 points.

- Make sure you have the data set `spine-data.txt`. You can load it into Python using the following.

```
import numpy as np
# Load data set and code labels as 0 = 'NO', 1 = 'DH', 2 = 'SL'
labels = [b'NO', b'DH', b'SL']
data = np.loadtxt('spine-data.txt', converters={6: lambda s: labels.index(s)})
```

This converts the labels in the last column into 0 (for 'NO'), 1 (for 'DH'), and 2 (for 'SL').

- Split the data into a training set, consisting of the *first* 250 points, and a test set, consisting of the remaining 60 points.
- Code up a nearest neighbor classifier based on this training set. Try both ℓ_2 and ℓ_1 distance. Recall that for $x, x' \in \mathbb{R}^d$:

$$\|x - x'\|_2 = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}$$

$$\|x - x'\|_1 = \sum_{i=1}^d |x_i - x'_i|$$

Now do the following exercises, to be turned in.

- (a) What error rates do you get on the test set for each of the two distance functions?
- (b) For each of the two distance functions, give the *confusion matrix* of the NN classifier. This is a 3×3 table of the form:

	NO	DH	SL
NO			
DH			
SL			

The entry at row DH, column SL, for instance, contains the number of test points whose correct label was DH but which were classified as SL.