
DSC 140B - Homework 05

Due: Wednesday, February 11

Instructions:

- Write your solutions to the following problems **by hand**, either on another piece of paper that you scan or using a tablet. Typed solutions will not be accepted for credit!
- Unless otherwise noted by the problem's instructions, show your work or provide some justification for your answer.
- Homework problems are graded pass/fail on completeness and effort, not correctness.
- Homeworks are due via Gradescope at 11:59 PM.

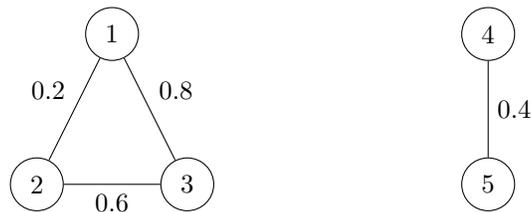
Problem 1. (0.5 credits)

In lecture, we saw that one minimizer of the cost function

$$\text{Cost}(\vec{f}) = \frac{1}{2} \sum_i \sum_j w_{ij} (f_i - f_j)^2$$

is the vector $\vec{f} = \frac{1}{\sqrt{n}}(1, 1, \dots, 1)^T$ which embeds all of the nodes of the graph to exactly the same number. The cost of this trivial embedding is zero, and we typically ignore it in favor of the eigenvector of the graph Laplacian with the *next* smallest positive eigenvalue for the embedding.

In the case where the graph has multiple connected components, however, there may be additional embeddings which have a cost of zero. For example, consider the similarity graph G shown below:



The weight of each edge is shown; the weight of non-existing edges is zero.

Find a normalized embedding vector \vec{g} for the above graph such that $\text{Cost}(\vec{g}) = 0$ and $\vec{g} \perp \frac{1}{\sqrt{n}}(1, 1, \dots, 1)^T$; that is, \vec{g} is orthogonal to the vector of all ones. Check that this is indeed a zero-cost embedding and show your work. Explain in words: why is this not a good choice for an embedding of the graph?

Problem 2. (1 credit)

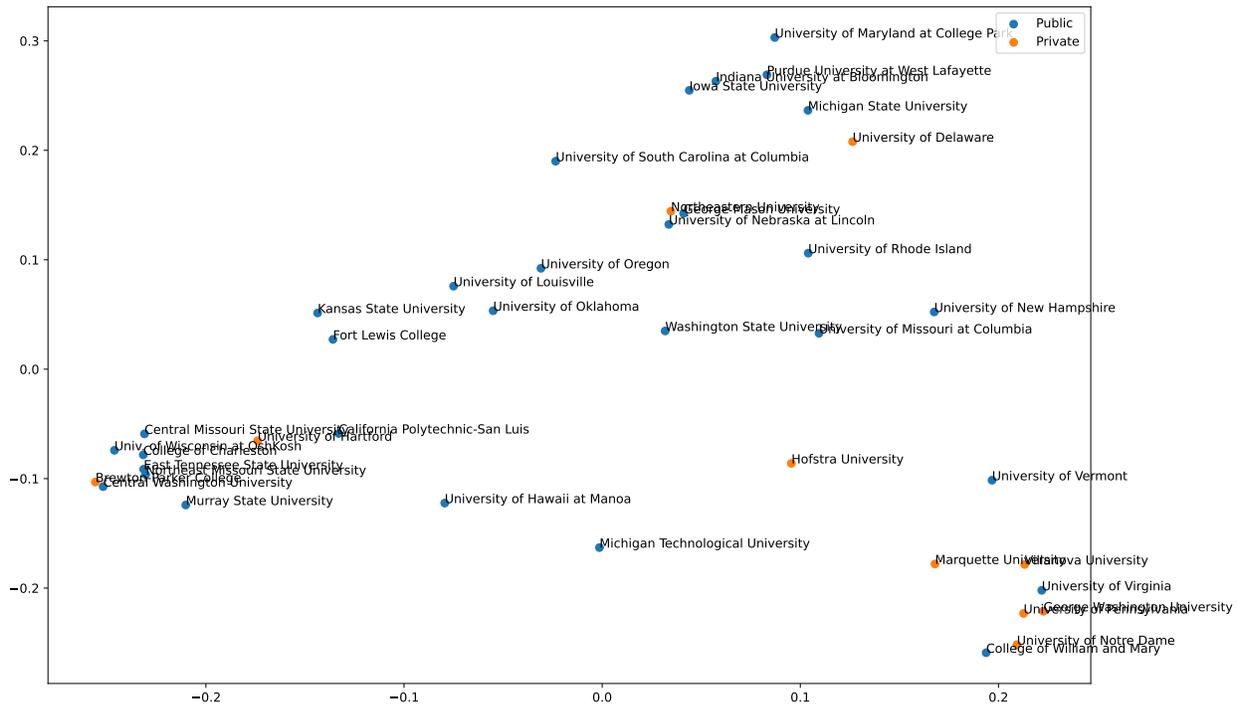
The data at the link below contains information on 40 colleges in the US:

<https://f000.backblazeb2.com/file/jeldridge-data/006-colleges/colleges-sample.csv>

The data contains 17 numerical features measuring the size of the student body of each college, the graduation rate, etc., and one Boolean feature (“Private”) saying whether the school is private or public. If we consider only the numerical features, each college is a point in \mathbb{R}^{17} .

Using Laplacian Eigenmaps, embed each college as a point in two dimensions and plot the result to obtain a similar figure to that shown below. You should construct a *symmetric* similarity matrix by building a k -neighbors graph (you will need to choose k). You should also consider whether the data needs to be

standardized. You should drop the “Private” column for the purpose of computing the embedding – use only the numerical data.



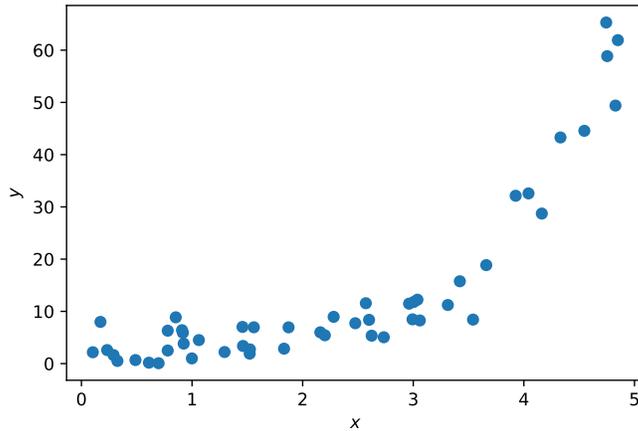
Like the plot above, each point in your plot should be annotated with the name of the college, and the color of the point should show whether the college is public or private. However, your plot may look significantly different from the plot above due to the sign of the eigenvectors your code finds (it is arbitrary), as well as your choice of k in constructing the k -neighbors graph, etc. We will look to see if it shows the same general *patterns* as the plot above. For example, notice how there is a group of large midwestern state schools (including “Michigan State University”, etc.), another group of large private universities (including “Notre Dame”), etc. Your plot should show the same.

You may use packages like `sklearn` to compute the k -neighbors graph, but not to do the actual embedding itself (e.g., do not use anything from `sklear.manifold`). Hint: the output of some `sklearn` functions is a sparse matrix; it is OK in this case to convert it to a dense `numpy` array. Also note that `sklearn`’s function for building a k -neighbors graph returns a *directed* graph, not an undirected graph, and therefore produces an asymmetric weight matrix. You will need to think about how to “symmetrize” it.

Note: you do not need to handwrite your solution to this problem, but you should include a screenshot of your code as part of your submission.

Problem 3. (1.5 credits)

Consider fitting a prediction function to the data shown below:



The pattern in this data is not linear, but we can still fit it using a linear prediction function and an appropriate choice of basis functions.

The data for this problem can be found at:

https://f000.backblazeb2.com/file/jeldridge-data/007-noisy_exponential/data.csv

Note: for this problem, you do not need to handwrite any code; you can include screenshots of your code as part of your submission.

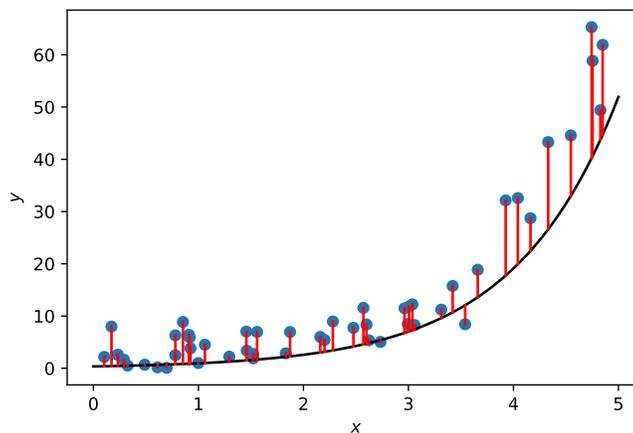
- a) For this and the next few subproblems, consider fitting a function of the form $H_1(x; w) = we^x$. Note that this function has just one parameter to be learned: w .

Write a Python function that takes in one argument, w , and computes the average square loss (that is, the mean squared error) of H_1 on the data set with respect to w . Use your function to compute the average square loss of $w = 0.2$.

Hint: the correct output of your function when called on 0.2 is a number between 100 and 200.

- b) The parameter w controls the value of we^x when x is close to zero. In that sense, it is a scale parameter. Different values of w will change the shape of H_1 , and will incur different amounts of error.

For example, the plot below shows H_1 for a particular choice of w that happens to be too small; H_1 is not a good fit to the data. The red lines show the residuals. Your Python function from the last part computes the average squared length of these red lines.



One way to find the optimal value of w (that is, the value resulting in the smallest average squared residual) is to use a numerical optimizer, such as `scipy.optimize.minimize`.

Use `scipy.optimize.minimize` to find the value of w which minimizes the average square loss (mean squared error).

- c) Another way to find the optimal value of w is to map the data to “feature space” and to fit a straight line in that space. In this case, we have a single basis function $\phi(x) = e^x$. Therefore, $H_1(x) = we^x = w\phi(x)$.

Plot the data in feature space; that is, create a “new” data set where each original point (x, y) becomes $(\phi(x), y)$.

Hint: you should see a linear trend.

- d) Perform linear least squares in feature space. That is, fit a straight line to the “new” data

$$(\phi(x^{(1)}), y_1), \dots, (\phi(x^{(n)}), y_n)$$

and report its slope, w .

Note that the discussion shows you how to perform least squares regression using `np.linalg.lstsq`. In this case, *do not* augment the data, since we are fitting a model we^x *without* a bias term.

You should have found the same value of w in the last part as you did when using `scipy.optimize.minimize`. When we perform least squares regression in feature space, we are finding the value of w which minimizes the average squared residual between the “new” data set and a straight line with slope w . When we used `scipy.optimize.minimize`, we were finding the value of w to minimize the average squared residual between the original data and the curved line, we^x . The fact is – these are the same optimization problem! Whether you fit a straight line in feature space, or a curved line in original space, you will find the same optimal w .

This happened because the function $H_1(x) = we^x$ is a linear function of the parameter, w . Now let’s try the same experiment, but with a *different* prediction function, $H_2(x) = e^{wx}$. This is also a function of one parameter, but the parameter is now in the exponential. This is *not* a linear function of w .

- e) As before, write a Python function that takes in one argument, w , and computes the average square loss (that is, mean squared error) of H_2 on the data set with respect to w . Use your function to compute the average square loss of $w = 0.5$.

Hint: the correct output of your function when called on 0.5 is a number between 200 and 300.

- f) Using `scipy.optimize.minimize`, find the value of w which minimizes the mean squared error of H_2 .

- g) The prediction function of the form $H_2(x; w) = e^{wx}$ does not involve a feature map in the sense that we’ve seen in lecture. Because of that, we cannot map the data to feature space and fit a linear predictor there. However, there is a trick: we can “linearize” the data in another way. Note that if $y \approx e^{wx}$, then $\ln y \approx \ln e^{wx} = wx$.

Plot the “transformed” data set in which the point (x, y) becomes $(x, \ln y)$.

Hint: you should see a linear trend.

- h) Fit a straight line to the “transformed” data set by minimizing mean squared error, and report the slope of this line, w .

- i) You should observe that the w you found in the last step is slightly different than the w that was found “directly” using `scipy.optimize.minimize`. Using your Python function that computes the mean squared error of H_2 , given a choice of w , calculate the mean squared error of both values of w .

You should have found that the w found by “linearizing” the data has a larger mean squared error – it doesn’t “fit” the data quite as well as the true minimizer. That is, linearizing the data and performing least squares *does not* minimize the mean squared error of H_2 . This is because H_2 is not a linear function of the parameter w .

This isn’t to say that it is a *bad* way of fitting the model. The plot below shows H_2 with both choices of w : the “optimal” choice found by the numerical solver, and the choice found by least squares on the “transformed” data. They are not too different!

