

## Homework 3: I Heard You Like Math

---

DUE: Tuesday, October 10 by 11:59:59pm

Out September 28, 2023

### Questions

This homework assignment focuses on machine learning theory, linear and ensemble models, and graph theoretic methods.

#### 1 REGULARIZATION IN LINEAR REGRESSION [35PTS]

Recall our high-dimensional linear regression equation from the previous homework assignment, where we needed to find the  $\beta$  that minimized the squared-error loss function:

$$J(\beta) = \sum_{i=1}^n (y_i - \vec{x}_i^T \beta)^2,$$

or more simply in matrix form:

$$J(\beta) = (X\beta - Y)^T(X\beta - Y), \tag{1}$$

When the number of features  $m$  is much larger than the number of training examples  $n$ , or very few of the features are non-zero (as we saw in Homework 1), the matrix  $X^T X$  is not full rank, and therefore cannot be inverted. This wasn't a problem for logistic regression which didn't have a closed-form solution anyway; for "vanilla" linear regression, however, this is a show-stopper.

Instead of minimizing our original loss function  $J(\beta)$ , we minimize a new loss function  $J_R(\beta)$  (where the  $R$  is for "regularized" linear regression):

$$J_R(\beta) = \sum_{i=1}^n (Y_i - X_i^T \beta)^2 + \lambda \sum_{j=1}^m \beta_j^2,$$

which can be rewritten as:

$$J_R(\beta) = (X\beta - Y)^T (X\beta - Y) + \lambda \|\beta\|^2 \quad (2)$$

**[5pts]** Explain what happens as  $\lambda \rightarrow 0$  and  $\lambda \rightarrow \infty$  in terms of  $J$ ,  $J_R$ , and  $\beta$ .

**[10pts]** Rather than viewing  $\beta$  as a fixed but unknown parameter (i.e. something we need to solve for), we can consider  $\beta$  as a random variable. In this setting, we can specify a prior distribution  $P(\beta)$  on  $\beta$  that expresses our prior beliefs about the types of values  $\beta$  should take. Then, we can estimate  $\beta$  as:

$$\beta_{\text{MAP}} = \operatorname{argmax}_{\beta} \prod_{i=1}^n P(Y_i | X_i; \beta) P(\beta), \quad (3)$$

where MAP is the *maximum a posteriori* estimate.

(aside: this is different from the MLE, which is the frequentist strategy for solving for a parameter. think of MAP as the Bayesian version.)

Show that maximizing Equation 3 can be expressed as *minimizing* Equation 2 with the assumption of a Gaussian prior on  $\beta$ , i.e.  $P(\beta) \sim \mathcal{N}(0, I\sigma^2/\lambda)$ . In other words, show that the  $L_2$ -norm regularization term in Equation 2 is effectively imposing a Gaussian prior assumption on the parameter  $\beta$ .

*Hint #1:* Start by writing out Equation 3 and filling in the probability terms.

*Hint #2:* Logarithms nuke pesky terms with exponents without changing linear relationships.

*Hint #3:* Multiplying an equation by -1 will switch from “argmin” to “argmax” and vice versa.

**[10pts]** What is the probabilistic interpretation of  $\lambda \rightarrow 0$  under this model? What about  $\lambda \rightarrow \infty$ ? Take note: this is asking a related but *different* question than the first part of this problem!

*Hint:* Consider how the prior  $P(\beta)$  is affected by changing  $\lambda$ .

**[10pts]** We have two data points in  $\mathbb{R}^3$ :

$$\begin{aligned}\vec{x}_1 &= [2, 1]^T, y_1 = 7 \\ \vec{x}_2 &= [1, 2]^T, y_2 = 5\end{aligned}$$

We know that for linear regression with a bias/intercept term and mean-squared objective function, there are *infinite* solutions with these two points (i.e., any line in  $\mathbb{R}^3$  can be made to cross through these two points).

Give a specific third point  $\langle \vec{x}_3, y_3 \rangle$  such that, when included with the first two above, will cause linear regression to *still have infinite solutions*. Your  $\vec{x}_3$  should not equal  $\vec{x}_1$  nor  $\vec{x}_2$ , nor should your  $y_3$  equal either  $y_1$  or  $y_2$ .

## 2 SPECTRAL CLUSTERING **[35PTS]**

The general idea behind spectral clustering is to construct a mapping of data points to an eigenspace of a graph-induced affinity matrix  $A$ , with the hope that the points are well-separated in the eigenspace to the point where something simple like k-means will work well on the embedded data.

A very simple affinity matrix can be constructed as follows:

$$A_{i,j} = A_{j,i} = \begin{cases} 1 & \text{if } d(\vec{x}_i, \vec{x}_j) \leq \Theta \\ 0 & \text{otherwise} \end{cases}$$

where  $d(\vec{x}_i, \vec{x}_j)$  denotes Euclidean distance between points  $\vec{x}_i$  and  $\vec{x}_j$ .

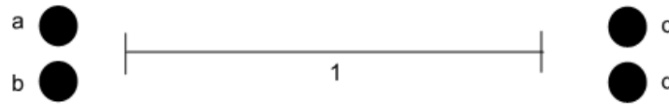


Figure .1: Simple toy dataset.

As an example, consider forming an affinity matrix for the dataset in Figure 1 using the affinity equation above, using  $\Theta = 1$ . Then we get the affinity matrix in Figure 2.

For this particular example, the clusters  $\{a, b\}$  and  $\{c, d\}$  show up as nonzero blocks in the affinity matrix. This is, of course, artificial since we could have constructed the matrix  $A$  using any ordering of  $\{a, b, c, d\}$ . For example, another possible affinity matrix  $A$  could have been as in Figure 2(b).

$$A = \begin{array}{c|cccc} & a & b & c & d \\ \hline a & 1 & 1 & 0 & 0 \\ b & 1 & 1 & 0 & 0 \\ c & 0 & 0 & 1 & 1 \\ d & 0 & 0 & 1 & 1 \end{array} \quad \tilde{A} = \begin{array}{c|cccc} & a & c & b & d \\ \hline a & 1 & 0 & 1 & 0 \\ c & 0 & 1 & 0 & 1 \\ b & 1 & 0 & 1 & 0 \\ d & 0 & 1 & 0 & 1 \end{array}$$

(a) (b)

Figure .2: Affinity matrices of Fig. 1 with  $\Theta = 1$ .

The key insight here is that the eigenvectors of both  $A$  and  $\tilde{A}$  have the same entries, just permuted. The eigenvectors with nonzero eigenvalues of  $A$  are  $\vec{e}_1 = [0.7, 0.7, 0, 0]^T$  and  $\vec{e}_2 = [0, 0, 0.7, 0.7]^T$ . Likewise, the nonzero eigenvectors of  $\tilde{A}$  are  $\vec{e}_1 = [0.7, 0, 0.7, 0]^T$  and  $\vec{e}_2 = [0, 0.7, 0, 0.7]^T$ .

Spectral clustering embeds the original data points in a new space by using the coordinates of these eigenvectors. Specifically, it maps the point  $\vec{x}_i$  to the point  $[e_1(i), e_2(i), \dots, e_k(i)]$ , where  $\vec{e}_1, \dots, \vec{e}_k$  are the top  $k$  eigenvectors of  $A$ . We refer to this mapping as the *spectral embedding*. See Figure 3 for an example.

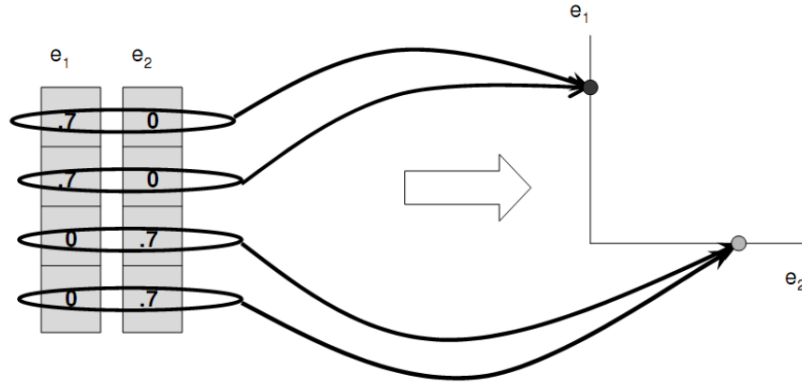


Figure .3: Using the eigenvectors of  $A$  to embed the data points. Notice that the points  $\{a, b, c, d\}$  are tightly clustered in this space.

In this problem, we'll analyze how spectral clustering works on the simple dataset shown in the next figure.

**[5pts]** For the dataset in Figure 4, assume that the first cluster has  $m_1$  points in it, and the second one has  $m_2$  points. If we use the affinity equation from before to compute the

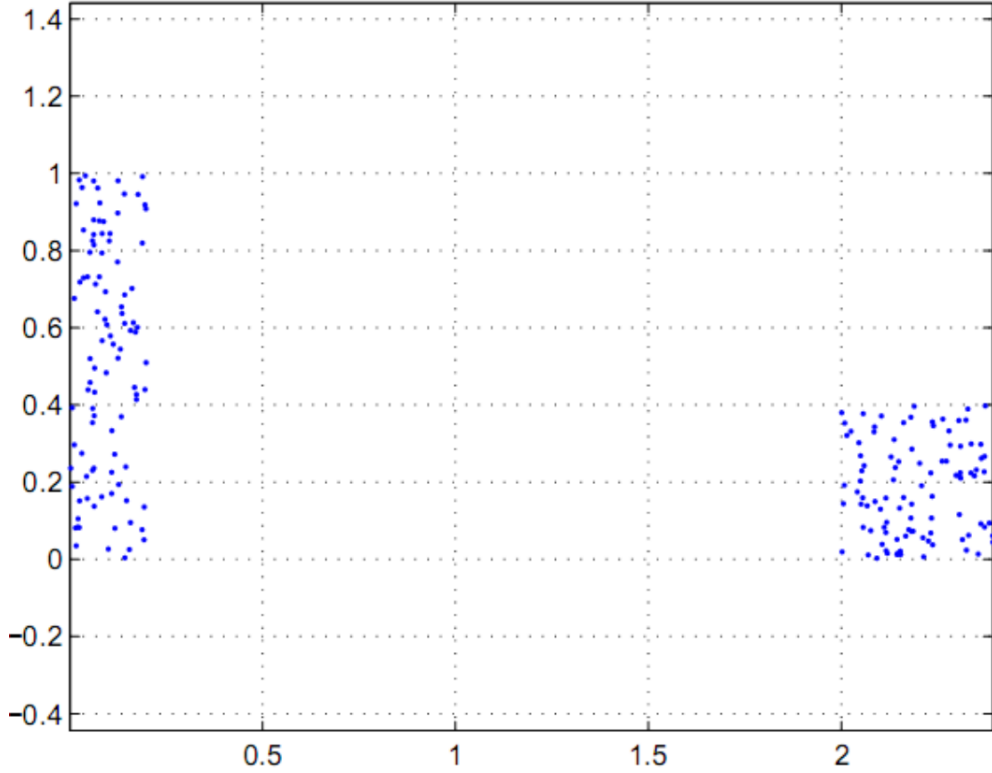


Figure .4: Dataset with rectangles.

affinity matrix  $A$ , what  $\Theta$  value would you choose and why?

**[10pts]** The second step is to compute the first  $k$  dominant eigenvectors of the affinity matrix, where  $k$  is the number of clusters we want to have. For the dataset in the above figure, and the affinity matrix defined by the previous equation, is there a value of  $\Theta$  for which you can analytically compute the first two eigenvalues and eigenvectors? If not, explain why not. If yes, compute and record these eigenvalues and eigenvectors. What are the other  $((m_1 + m_2) - k)$  eigenvalues? Explain briefly.

Spectral clustering algorithms often make use a graph Laplacian matrix,  $L$ . A favorite variant is the *normalized* graph Laplacian,  $L = D^{-1/2}AD^{-1/2}$ , as this formulation has many convenient properties ( $D$  is a diagonal matrix whose  $i^{th}$  diagonal element,  $d_{ii}$ , is the sum of the  $i^{th}$  row of  $A$ ).

**[10pts]** Show that a vector  $\vec{v} = [\sqrt{d_{11}}, \sqrt{d_{22}}, \dots, \sqrt{d_{nn}}]^T$  is an eigenvector of  $L$  with corresponding eigenvalue  $\lambda = 1$ .

One of the convenient properties of normalized graph Laplacians is the eigenvalue  $\lambda_1$  of the leading eigenvector is, at most, 1; all other eigenvalues  $\lambda_2, \dots, \lambda_n$  have values strictly smaller than 1.

Consider a matrix  $P$ , where  $P = D^{-1}A$ , where  $A$  is our affinity matrix and  $D$  is the diagonal matrix. Each  $p_{ij} = a_{ij}/d_{ii}$ . Note the intuition of this operation: we are normalizing each edge by the total degree of the incoming vertex, essentially creating a “transition probability”  $p_{ij}$  of transitioning from vertex  $i$  to vertex  $j$ . In other words, each row of  $P$  sums to 1, so it is therefore a valid probability transition matrix. Hence,  $P^t$  is a matrix whose  $\{i, j\}^{th}$  element shows the probability of being at vertex  $j$  after  $t$  number of steps, if one started at vertex  $i$ .

**[10pts]** Show that  $P^\infty = D^{-1/2} \vec{v}_1 \vec{v}_1^T D^{1/2}$ . This property shows that if points are viewed as vertices according to a transition probability matrix, then  $\vec{v}_1$  is the only eigenvector needed to compute the probability distribution over  $P^\infty$ .

### 3 CODING [40PTS]

In this question, you’ll be implementing a slightly simplified version of the MultiRankWalk (MRW) semi-supervised learning algorithm discussed in lecture. The paper is here: <https://lti.cs.cmu.edu/sites/default/files/research/reports/2009/cmuli09017.pdf>

The basic procedure of MRW is similar to other graph-based random walk algorithms such as PageRank. For a graph  $G$  defined by the set of vertices  $V$  and edges  $E$ , the MRW procedure is as follows:

$$\vec{r} = (1 - d)\vec{u} + dW\vec{r}$$

where  $W$  is the weighted transition matrix of graph  $G$  from vertex  $i$  to  $j$  is given by  $W_{ij} = A_{ij}/d_{ii}$ , where  $d_{ii}$  is the degree of the  $i^{th}$  vertex.  $\vec{u}$  is the normalized teleportation vector, where  $|\vec{u}| = |V|$  and  $||\vec{u}||_1 = 1$ .  $d$  is a constant damping factor, controlling how often random jumps are made.

The value  $A_{ij}$  comes from our use of an affinity matrix in representing the graph. **This is a deviation from the MRW paper**, which assumes a simple adjacency matrix. The affinity matrix  $A$  will be determined using the radial-basis function kernel, also known as the Gaussian kernel or heat kernel. It has the form  $A_{ij} = A_{ji} = e^{-\gamma||\vec{x}_i - \vec{x}_j||^2}$ , and is implemented in scikit-learn’s `sklearn.metrics.pairwise` module as `rbf_kernel()`. Once you have the affinity matrix  $A$ , the diagonal (degree) matrix  $D$  can be found by summing the rows of  $A$ , i.e.  $D_{ii} = \sum_j A_{ij}$ . Finally, the weighted transition probability matrix  $W$  can be found using  $A$  and  $D$  and the above formulation.

Your task is to solve for the ranking vector  $\vec{r}$  by iteratively substituting  $\vec{r}^{t-1}$  with  $\vec{r}^t$  until convergence or a set number of iterations.

In this implementation, the  $\vec{u}$  vector actually functions as a *seed vector*: this identifies vertices that are labeled and function as seeds for the subsequent label-spreading. “Seeds” are labeled data points used to initiate the label-spreading of the MRW algorithm and predict classes for unlabeled data. The original MRW paper cites several methods, including using PageRank to initially rank labeled vertices in terms of preference as seed vertices to MRW. Your code will need to implement both random seed selection, and degree-based seed selection. In the former, you’ll randomly pick  $k$  labeled data points from each class and use them as seeds. In the latter, you’ll rank the labeled vertices of each class by their degree (i.e. sums of the rows of  $A$ ) and select the top  $k$  in each class.

Critically, you will need to perform MRW for **each distinct class  $c$  in the data**. Specifically, when initializing the labeled seeds in  $\vec{u}$ , you need to set each corresponding element  $\vec{u}_i = 1$  such that  $\vec{y}_i = c$ . All other entries of  $\vec{u}$  should be 0. Once this step is completed, you will need to normalize  $\vec{u}$  such that  $\|\vec{u}\|_1 = 1$ . Next, you can proceed with MRW. Finally, you will repeat this process again for all unique labels  $c$  in your dataset, so that at the end you’ll have a set of ranking vectors  $\vec{r}_1, \vec{r}_2, \dots, \vec{r}_c$  for each class.

Once you have generated a ranking vector  $\vec{r}$  for each class, you’ll then assign labels to all your unlabeled data. For the  $i^{th}$  vertex, whichever ranking vector  $\vec{r}$ ’s  $i^{th}$  element is largest, assign the corresponding class label represented by that ranking vector to the unlabeled data point. Continue for all unlabeled data.

Your code should be able to process: an input file containing the  $n$   $m$ -dimensional data points, the number of labeled data points  $k$  to use from each class as seeds, whether to choose seeds randomly or by vertex degree, the damping factor  $d$ , and an output file to write the predicted classes for all data.

You’ll also be provided the boilerplate to read in the necessary command-line parameters:

1. **-i**: a file path to a text file containing the data
2. **-d**: the damping factor (float between 0 and 1)
3. **-k**: number of data points per class to use as seeds
4. **-t**: type of seed selection to use, “random” or “degree”
5. **-e**: the epsilon threshold, or squared difference of  $\vec{r}^t$  and  $\vec{r}^{t+1}$  to determine convergence
6. **-g**: value of gamma for the pairwise RBF affinity kernel
7. **-o**: a file path to an output file, where the predicted labels will be written

The format of the input file will be tab-delimited, where a single data point will be on one line. The first column will be the labels: any unlabeled data will have a label of -1. Functions are already written in the `homework3.py-TEMPLATE` file that will handle reading in data and parsing command-line arguments.

The format of the output file should be one label prediction per line; therefore, the number of lines in the input file and the output file should match exactly (so for the labeled data, you can either use the labels you read in from the file or the labels that are predicted from your ranking vectors, though in theory they should be the same). Essentially, fill in the -1 values in your initial label vector, then just write the vector to a text file, such that each element of the vector is on its own line. For your convenience, the ground-truth label files `y_easy.txt` and `y_hard.txt` for the full datasets are provided; you can use these to check how well your code is predicting the -1 labels.

**HINT 1:** The value of gamma can substantially affect the accuracy of your method. Larger values shrink the neighborhoods and isolate points from each other; smaller values expand the neighborhoods and make everything look the same distance. If in doubt, plot the affinity matrix using `matplotlib.pyplot.imshow`, and you should see a block-diagonal-ish structure. For the easy dataset, try values around 0.5. For the harder dataset, try values in the 10-50 range.

**HINT 2:** At the same time, adding more seeds per class can help immensely. The default value in the template script is only 1 seeded value per class; while you can still attain high-90s accuracy with proper values of gamma on the hard dataset, it's almost impossible to hit perfect accuracy without increasing the number of seeds.

**HINT 3:** The two test datasets provided should not require any more than 100 iterations to converge using the default epsilon.

## Administration

### 1 SUBMITTING

All submissions will go to **AutoLab**. You can access AutoLab at:

- <https://autolab.cs.uga.edu>

You can submit deliverables to the **Homework 3** assessment that is open. When you do, you'll submit two files:

1. `homework3.py`: the Python script that implements your algorithms, and
2. `homework3.pdf`: the PDF write-up with any questions that were asked



These should be packaged together in a tarball; the archive can be named whatever you want when you upload it to AutoLab, but the files in the archive should be named **exactly** what is above. Deviating from this convention could result in the autograder failing!

To create the tarball archive to submit, run the following command (on a \*nix machine):

```
> tar cvf homework3.tar homework3.py homework3.pdf
```

This will create a new file, **homework3.tar**, which is basically a zip file containing your Python script and PDF write-up. Upload the archive to AutoLab. There's no penalty for submitting as many times as you need to, but keep in mind that swamping the server at the last minute may result in your submission being missed; AutoLab is programmed to close submissions *promptly* at 11:59pm on October 10, so give yourself plenty of time! A late submission because the server got hammered at the deadline will *not* be acceptable (there is a *small* grace period to account for unusually high load at deadline, but I strongly recommend you avoid the problem altogether and start early).

Also, to save time while you're working on the coding portion, you are welcome to create a tarball archive of just the Python script and upload that to AutoLab. Once you get the autograder score you're looking for, you can then include the PDF in the folder, tarball everything, and upload it. AutoLab stores the entire submission history of every student on every assignment, so your autograder (code) score will be maintained and I can just use your most recent submission to get the PDF.

## 2 REMINDERS

- If you run into problems, ping the **#questions** room of the Discord server. If you still run into problems, ask me. But please please please, **do NOT** ask Google to give you the code you seek! I will be on the lookout for this (and already know some of the most popular venues that might have solutions or partial solutions to the questions here).
- Prefabricated solutions (e.g. **scikit-learn**, OpenCV) are NOT allowed! You have to do the coding yourself! But you **can** use the pairwise metrics in scikit-learn, as well as the vector norm in SciPy.
- If you collaborate with anyone or anybot, just mention their names in a code comment and/or at the top of your homework writeup.
- Cite any external and/or non-course materials you referenced in working on this assignment.