

Assignment #5

Due: 12:00 noon EDT, Oct 26, 2020

Upload at: <https://www.gradescope.com/courses/135869/assignments/746434>

Remember to append your Colab PDF as explained in the first homework, with all outputs visible.
When you print to PDF it may be helpful to scale at 95% or so to get everything on the page.

Problem 1 (18pts)

We sometimes need to introduce norms for matrices. This comes up when we want to talk about the magnitude of a matrix, or when we need a notion of distance between matrices. One important matrix norm is the Frobenius norm, which can be written as the square root of the sum of squared matrix elements for a matrix $A \in \mathbb{R}^{m \times n}$:

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{i,j}^2}$$

It turns out that there is a particularly compact and elegant way to rewrite this:

$$\|A\|_F = \sqrt{\text{trace}(A^T A)}$$

- (A) Show that the Frobenius norm is also the square root of the sum of the squared singular values.
- (B) Assume that A is square and invertible, but with a very small Frobenius norm. What kind of value would you expect to get for the Frobenius norm of A^{-1} ?

Problem 2 (40pts)

In this problem you will use SVD to model a text corpus, a small subset of New York Times articles. You will use a “bag of words” representation in which documents are represented by the counts of words, usually excluding very common “stop words” like *the* and *and*. Download `nyt.pkl.gz` and upload it to your Google drive so you don’t have to upload it every time you open the Colab; it’s a fairly big file. Here’s some code to get you going:

```
import pickle as pkl
import numpy as np
import gzip

filename = 'drive/My Drive/COS 302/nyt.pkl.gz'
with gzip.open(filename, 'rb') as fh:
    nyt = pkl.load(fh)
documents = nyt['docs']
vocab = nyt['vocab']

# Create reverse lookup table.
vocab_indices = dict([(w, i) for (i, w) in enumerate(vocab)])

M = len(documents)
N = len(vocab)
print('%d documents, %d words' % (M,N))

count_mat = np.zeros((M,N))
for mm, doc in enumerate(documents):
    for word, count in doc.items():
        count_mat[mm, vocab_indices[word]] = count
```

- (A) Typically, raw counts don’t lead to discovery of interesting structure. Instead, it is common to use something like **TF-IDF**, which stands for *term frequency-inverse document frequency*. Term frequency is the number of times word n appeared in document m , divided by the total number of words in document m .

$$tf_{m,n} = \frac{\text{\# times word } n \text{ appears in doc } m}{\text{total \# of words in doc } m}$$

Transform `count_mat` from the code above into a term frequency matrix.

- (B) Inverse document frequency is typically the natural log of the number of documents, divided by the number of documents in which word m appears. (The plus-one ensures that you’re not dividing by zero.)

$$idf_n = \log \frac{\text{total \# of documents}}{1 + \text{\# of documents with word } n}$$

Compute the idf vector from `count_mat`.

- (C) Now compute the TF-IDF matrix by multiplying (broadcasting) $tf_{m,n}$ and idf_n . Use `numpy.linalg.svd` to take the SVD of this TF-IDF matrix; it may take a minute since the matrix is relatively large. Plot the singular values in decreasing order.
- (D) The right singular vectors (columns of the V matrix) will ideally represent interesting topical dimensions. For each of the top 20 right singular vectors: identify the words in the vocabulary that have the largest entries in the vector and print them. That will probably mean looping over the first 20 right singular vectors, doing an appropriate `argsort` and then finding that entry in the `vocab` variable. Do you see any interesting qualitative structure in these groups of words?

Problem 3 (40pts)

In various data analysis problems, it is often easier to reason about the pairwise distances between data, rather than features of the data directly. This comes up particularly when dealing with discrete data where there aren't vectors, but perhaps there is a sensible concept of distance. A prime example is strings, in which there are various sensible **edit distances**. We previously examined principal component analysis (PCA) as a way to find low-dimensional representations of data, but if you only have distances, vanilla PCA doesn't apply. Instead, the classic approach is to use principal *coordinates* analysis, also called **multidimensional scaling** (MDS) to map the data to \mathbb{R}^d in such a way that the pairwise distances are approximately preserved. That is, you are given a matrix of squared distances $D \in \mathbb{R}^{n \times n}$, and your goal is to discover reasonable locations $\{x_i\}_{i=1}^n$ for, say, $x \in \mathbb{R}^2$, so that $D_{ij} \approx \|x_i - x_j\|_2^2$. The details of MDS are beyond the scope of this course, but the steps are straightforward applications of tools that you've been learning to use in COS 302.

- (1) Compute a squared distance matrix D .
- (2) Subtract off the column-wise means and then the row-wise means of D , i.e., perform "double centering". Be sure to compute the row-wise means *after* you've subtracted the column-wise means.
- (3) Compute the largest eigenvalues $\lambda_1, \lambda_2, \dots$ and associated eigenvectors v_1, v_2, \dots of $-\frac{1}{2}D$.
- (4) The j th entry of the vector $\sqrt{\lambda_i}v_i$ can now be used as the i th coordinate of x_j . That is, for mapping data into \mathbb{R}^2 , you would form a $n \times 2$ matrix $X = [\sqrt{\lambda_1}v_1 \quad \sqrt{\lambda_2}v_2]$ whose rows are locations to embed each datum.

In this problem, you're going to take a list of strings, find edit distances between them all, and then make a visualization of them in \mathbb{R}^2 . By default, you can use `dog_names1000.txt`, a list of dog names taken from a random subset of **registered dogs in Anchorage, Alaska**; you can use any list of strings you want, however, as long as it has at least 500 or so entries: lists of cities, street names, **metal bands**, etc. You'll need to figure out how to load the string from file into a list, but then the function below will compute a squared distance matrix for you.

```
import numpy as np
import editdistance
def sq_distances(names):
    '''Takes a list of strings. Returns squared edit distances.'''
    N = len(names)
    sq_dists = np.zeros((N,N))
    for ii, name1 in enumerate(names):
        for jj, name2 in enumerate(names[:ii]):
            sq_dists[ii, jj] = (editdistance.eval(name1, name2) \
                               / np.maximum(len(name1), len(name2)))**2
            sq_dists[jj, ii] = sq_dists[ii, jj]
    return sq_dists
```

Implement the MDS procedure above to put your strings into locations in \mathbb{R}^2 and then plot them. Assuming you have a $n \times 2$ matrix X , the code below should get you started making the figure.

```
import matplotlib.pyplot as plt
plt.figure(figsize=(20,20))
plt.plot(X[:,0], X[:,1], '.')
for ii in range(n):
    plt.text(X[ii,0], X[ii,1], names[ii], fontsize=10)
plt.show()
```

Problem 4 (2pts)

Approximately how many hours did this assignment take you to complete?

My notebook URL: <https://colab.research.google.com/XXXXXXXXXXXXXXXXXXXXXXXXXX>

Changelog

- 21 October 2020 – Fixed sample code and clarified double centering in Problem 3.
- 9 October 2020 – Updated for Fall, 2020.
- 5 March 2020 – Initial version.