

M3/4/5N9 – Project 2

Due 20 Dec 2016

Use the MATLAB format `shorte` when reporting any numerical values.

1. PageRank Algorithm

Google PageRank developed by Larry Page and Sergey Brin revolutionized web search by generating a ranked list of web pages based on the underlying connectivity of the web. The PageRank algorithm is based on an ideal random web surfer who, when reaching a page, goes to the next page by clicking on a link. The surfer has equal probability of clicking any link on the page and, when reaching a page with no links, has equal probability of moving to any other page by typing in its URL. In addition, the surfer may occasionally choose to type in a random URL instead of following the links on a page. The PageRank is the ranked order of the pages from the most to the least probable page the surfer will be viewing.

Suppose the web is composed of N sites. The probability of the surfer being on each page at time k is given by the state vector $\mathbf{x}^{(k)}$, i.e. the probability of being on site j after k steps is the j th component of $\mathbf{x}^{(k)}$. Since $\mathbf{x}^{(k)}$ contains the probabilities of being on a site for all sites, $\sum_{j=1}^N x_j^{(k)} = 1$. With every change of page, the state vector is multiplied by the transition matrix,

$$\mathbf{P} = \alpha \mathbf{G} + \alpha \mathbf{a} \mathbf{1}^T / N + (1 - \alpha) \mathbf{1} \mathbf{1}^T / N \quad (1)$$

to give the new probability, $\mathbf{x}^{(k+1)} = \mathbf{P}^T \mathbf{x}^{(k)}$ the surfer is on any given page. In the equation, \mathbf{G} is the matrix containing the transition probabilities to move to the next page via links, and a site i has links if the $\sum_{j=1}^N G_{ij} = 1$. The entries of the vector \mathbf{a} are 1 if the page has no links and 0 otherwise, and the term containing \mathbf{a} gives the transition probabilities to move to the next page for pages that have no links. Finally, the matrix $\mathbf{1} \mathbf{1}^T / N$, where $\mathbf{1}$ is the vector with each component equal to 1, are the transition probabilities to move to the next page using solely the URL. The value $\alpha \in [0, 1)$ controls the probability of using one transition approach over the other. Thus, $\sum_{j=1}^N P_{ij} = 1$ for each i . PageRank continues to update the state vector until $\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2$ is within a prescribed tolerance. Since Google regularly ranks 4.3 billion pages (thus \mathbf{P} is a $4.3\text{e}9 \times 4.3\text{e}9$ matrix) this computation must be performed as efficiently as possible.

- Based on properties of \mathbf{P}^T , it can be shown that $\mathbf{x}^{(k)}$ converges to a unique \mathbf{x} as $k \rightarrow \infty$. Explain what algorithm PageRank is actually using and explain why the final state vector \mathbf{x} is an eigenvector of \mathbf{P}^T . What is the corresponding eigenvalue?
- The accompanying file `G.mat` contains the matrix \mathbf{G} for a subset of the web (9964 pages) a number of years ago. Construct the matrix \mathbf{P}^T and implement the PageRank algorithm to find the PageRank state if $\alpha = 0.85$, the tolerance set to 10^{-8} , and initial state vector $\mathbf{x}^{(0)} = \mathbf{e}_1$. Generate the list the top 50 sites by their index (you may use `sort` to help generate the list).
- The parameter α and the URL transition matrix must be introduced to ensure convergence to a unique \mathbf{x} . In their original paper, Page and Brin settled on the value of $\alpha = 0.85$, but it's not immediately clear why they chose this value and not another. Perform numerical experiments to study how PageRank performs for different values of α . In particular, note how the final state vector \mathbf{x} and the number of iterations needed to converge vary with α . Discuss your findings in your report and give some indication why the value of $\alpha = 0.85$ was originally chosen. Based on your results, what can you say about how the eigenvalues of \mathbf{P}^T depend on α ?
- Since many pages don't link to any other pages, and those that do have links only link to a handful of pages, the matrix \mathbf{G} , is sparse (you can visualize the sparsity using `spy`) and matrix-vector multiplication involving \mathbf{G}^T can be done rapidly. MATLAB will do this rapid matrix-vector multiplication

automatically if G^T is a sparse data structure. This can be achieved by setting `GT = sparse(G')`. Explain how to take advantage of this in the PageRank algorithm and implement it in MATLAB. Demonstrate the decrease in runtime.

2. Principal Component Analysis - Eigenworm

Dimensional reduction is an important part of analysing any large dataset. Consider an $m \times n$ matrix X where the rows contain data measured at n points at different times. One technique to analyse this data is principal component analysis (PCA) which finds the largest eigenvalues and eigenvectors of the covariance matrix $C = X^T X$. Thus, the matrix C is symmetric by construction. In this context, the eigenvectors, v , are referred to as the principal components and the eigenvalues, λ , provide the variance in the direction of the principal component. To reduce the dimension of the data, one can then simply compute the principal components, then project the data onto the p components with the highest variance.

For this problem, we will use data collected from the crawling worm, *C. elegans* (believe it or not, three nobel prizes have been awarded for research involving this organism). PCA is now a regular technique (see attached reference *Dimensionality and Dynamics in the Behavior of C. elegans* by Stephens, Johnson-Kerner, Bailek, and Ryu) used to analyse its motion and often is used to help link behaviour with genetic variations between mutants. The data matrix Θ in the associated file `theta.mat` gives the worm's bending angle at $n = 49$ points along its length measured at $m = 26996$ points in time. PCA will pick out a natural basis (the principal components) for the worm shapes based on those that yield the highest variance in the data. The hope is that very few principal components are needed to reproduce all of the observed shapes allowing for easier interpretation of the experimental data.

- Compute C and modify the code house to turn C into a tridiagonal matrix T via similarity transformations. You will receive *most* marks for using the general algorithm introduced in lecture, but for *full* marks you must take full advantage of symmetry and sparsity in your code. Provide the non-zero entries of T in your report.
- Use the QR algorithm to find the eigenvalues of C and provide the values of the 10 most dominant eigenvalues in your report. Use a tolerance 10^{-14} to determine when an off-diagonal component is zero and an eigenvalue has been found. Full marks will be awarded only if your code uses shifts and takes advantage of the structure of T . Based on the magnitude of the values, do you think that PCA will help us reduce the dimension of the data?
- Find the principal components (eigenvectors) associated with the 10 most dominant eigenvalues. There are different ways of doing this and the choice is up to you, however, for whichever route you choose, you may only use elementary MATLAB operations in you code. Also, if you need to solve a linear system you cannot use backslash, and must solve the system using LU decomposition for full marks. Make a single plot (line plots of the entries vs index) of the principal components associated with the 4 most dominant eigenvalues and provide their numerical values in the report.
- Compute the projection, θ_p , of the first row of Θ onto the space spanned by the first p principal components for $p = 1$ to 10. Provide a table of the projection coefficients for the case $p = 10$. Make a single plot showing the projected data for $p = 2, 4, 6$ and 10, as well as the original data. Make a table the error $\|\theta - \theta_p\|_2$ where θ is the first row of Θ , for $p = 1$ to 10. Discuss your results and provide some insight into how to choose p .

3. Mastery

A study of the various direct methods for solving linear systems via matrix factorization could be viewed as a tour of MATLAB's backslash. In lectures, we only explored in detail the first (QR) and last (LU) stops of backslash as they are the most general methods that apply to virtually all linear systems of interest. The stops backslash takes along the way between these two points explores if the matrix has a particular structure or certain properties allowing either LU to be performed more rapidly or another closely related algorithm to be used instead. These algorithms are covered in Chapter 4 of *Matrix Computations* by Golub and Van Loan (attached).

Take again the linear system $-D_2 u_N = f_N$ from Project 1 and consider the permutation matrix

$$P = [e_1 \ e_N \ e_2 \ e_{N-1} \ e_3 \ \dots] \quad (2)$$

where e_i is the i th column of identity. Since $PP^T = I$, we have that

$$-(P^T D_2 P)(P^T u_N) = (P^T f_N). \quad (3)$$

- (a) For this problem, you are backslash. Using Chapter 4, select the best algorithm to solve Eq. (3) for u_N . In your report, discuss and justify your choice.
- (b) Implement the algorithm, and solve the linear system for u_N for $N = 16, 32, 64$ and 128 , again taking $f_{N,i} = \sin(x_i)$. Explain why this equation can be solved even though we know D_2 is singular. As in project 1, provide the error, E_N , to ensure your code is working properly and use `tic` and `toc` to obtain runtimes.
- (c) Compare your choice of algorithm with those used in Project 1. Is it as fast or even faster than the FFT approach? Discuss your comparison and provide an explanation of the advantages of your chosen algorithm. Also, describe a situation where your chosen algorithm would certainly fail to be the best one!