**Instructions**: Present your work in a neat and organized manner. Please **use** either the $8.5 \times 11$ size paper or the filler paper with pre-punched holes. Please do **not use** paper which has been torn from a spiral notebook. Please secure all your papers by using either a staple or a paper clip, but **not** by folding its (upper left) corner.

You **must** show **all of the essential details** of your work to get full credit. If you used *Mathematica* for some of your calculations, attach a printout showing your commands and their output. If I am forced to fill in gaps in your solution by using notrivial (at my discretion) steps, I will also be forced to reduce your score.

Please refer to the syllabus for the instructions on working on homework assignments with other students and on submitting **your own** work.

# Homework Assignment # 12

1. This problem is worth **1.5 points**.

   (a) Verify[1] that vector $\mathbf{v}_1$ defined after Eq. (1) of the Notes is the *only* eigenvector of matrix $A$ of Eq. (1).

   (b) For that $A$ and vector $\mathbf{u}$ defined before Eq. (2a), find $A^{-1}\mathbf{u}$.
   *Hint*: $A^{-1}\mathbf{u} = \mathbf{w}$ for some $\mathbf{w}$. Now, since $\mathbf{u}$ and $\mathbf{v}_1$ have been chosen to form a basis, then $\mathbf{w} = p\mathbf{v}_1 + q\mathbf{u}$ for some $p$ and $q$. Thus,
   $$A^{-1}\mathbf{u} = p\mathbf{v}_1 + q\mathbf{u}. \tag{HW12.1}$$

   To find $p$ and $q$, multiply both sides of (HW12.1) by $A$. In subsequent calculations, use the fact that $\mathbf{v}_1$ and $\mathbf{u}$ are linearly independent (since they form a basis). If you forgot what mathematical formula expresses the condition of linear independence of vectors, consult a Linear Algebra textbook.
   *Note*: **No credit** will be given if you find $A^{-1}u$ by a direct calculation involving the explicit form of *either* $A^{-1}$ *or* $A$. *Moreover*, you are also *not* allowed to use the explicit form of $\mathbf{v}_1$ and $\mathbf{u}$. The only pieces of information you are allowed to use are: (i) How $A$ acts on its eigenvector $\mathbf{v}_1$, and (ii) How $A$ acts on $\mathbf{u}$ (see Eq. (2b) in the Notes).

   (c) Using the result of part (b) and an approach similar to the one used to prove Eq. (14), show that Eq. (4) also holds for integer $n < 0$.

   (d) Verify[2] that vector $\mathbf{v}_1$ defined after Eq. (5) is the *only* eigenvector of matrix $A$ of Eq. (5).

   (e) For that $A$ and vector $\mathbf{u}_2$ defined in (6), find $A^{-1}\mathbf{u}_2$.
   *Note*: Follow the method of part (b). The Note for part (b) also applies here.

2. This problem is worth **2.5 points**.

   As announced in the preamble to Lecture 12, properties of a *diagonalizable* matrix $C$ that is in some sense close to a non-diagonalizable one are close to the properties of the latter matrix. In particular:

   - The eigenvalues of such a "diagonalizable, but almost non-diagonalizable" matrix do *not* contain adequate information about the large-$n$ behavior of $C^n\underline{x}$.
     In layman terms, one can put it as:

   - One cannot trust the information provided by eigenvalues of a matrix that is diagonalizable but is close to being non-diagonalizable.

---

[1]doing the calculations by hand
[2]doing the calculations by hand

In parts (a)–(c) of this problem, you will verify this statement for a diagonalizable matrix $C$ that is close to a non-diagonalizable matrix $A$ given by Eq. (16) of the Notes.

Now, if eigenvalues of an "almost nondiagonalizable" matrix cannot be trusted to predict its action, then is there a criterion by which one can tell that such a matrix is close to some nondiagonalizable matrix? You will explore this question in part (d).

Finally, in part (e), you will be guided to uncover a reason (in general terms) behind the phenomenon on which you (should) have reported in part (c). Then later, in Problem 3, you will further explore how that "general-terms" reason is manifested in a specific example.

(a) Download the file `demo_for_lecture12.m` (it is posted next to Lecture 12). Set matrix $A$ defined by this code to be non-normal. Run the code and **do not** close figures or clear Matlab's workspace.

What is the maximum size, $b_{max}$, of the entries of $B$? You can look inside the code or use the command `max`. In the latter case, type `help max` first; also, round your answer to one significant figure. (Note that since the maximum entry of $A$ is 1, then this $b_{max}$ also gives the *relative* change of entries of $(A + B)$ relative to entries of $A$.)

What is the change of the eigenvalues, $|\lambda_{(A+B)} - \lambda_A|$, that adding $B$ to $A$ has caused? Is your answer consistent with formula (27)?

What property of matrix $A$ is responsible for the relation between $|\lambda_{(A+B)} - \lambda_A|$ and $b_{max}$?

Summarize the result of this part by copying to your paper the following sentence where you need to fill in the blanks, and $\boxed{\text{put it in a box}}$ :

_____ of a matrix that is close to _____ are very _____ to small perturbations of the matrix.

(b) Define a column vector $\underline{u}$ as in Eq. (19) of the Notes.[3] Compute $\underline{p} = A^{50}\underline{u}$. [4] Plot it versus its index (i.e., simply type `plot(p)`). Discuss how your plot agrees (within an order of magnitude or so) with formula (20). [5]

(c) This part addresses the first **main point** of this exercise. (Reread the preamble if you have forgotten what it is supposed to be.)

Compute vector $\underline{q} = (A + B)^{50}\underline{u}$ and plot it. Give your answers to the questions below.

- Note that, from part (a), all eigenvalues of $(A + B)$ are different.
  (1) Quote a theorem from your Linear Algebra course stating whether a matrix with all distinct eigenvalues can be diagonalizable or not.
  (2) Further, state whether eigenvectors of such a matrix form a basis in $\mathbb{R}^M$, where $M \times M$ is the size of $(A + B)$.
  Given *your answers to these two questions*, and also from what we observed in *Lecture 5*, what should the behavior of $(A + B)^n\underline{x}$ for any vector $\underline{x}$ and a sufficiently large $n$ be determined by?

- What is the largest absolute value[6] of the eigenvalues of $(A + B)$?

- Does this largest eigenvalue explain the magnitude of the entries of $\underline{q}$? Justify your answer with a short calculation and a brief explanation.

---

[3] You can define a zero column vector using the command `zeros` and then overwrite its last entry by the command `u(end)=1`.

[4] For the correct syntax of raising a matrix to a power and of matrix-vector multiplication, either see the Matlab Primer posted on the course webpage, or find an example in `hw11_p1.m`.

[5] If you would like to obtain a closer agreement, note that the largest coefficient in (20) is actually the binomial coefficient "n-choose-(M-1)", i.e. $n!/\big((n - (M - 1))!\,(M - 1)!\big)$. Matlab has a command `nchoosek` to compute it.

[6] Matlab's command for the absolute value is `abs`.

In your paper, (re)state the main point that you have demonstrated by doing this part of the Exercise. Put this statement in a box.

(d) Here you will address the question:

*Is there a quantitative criterion* by which the large-$n$ behavior of $(A + B)^n$ is similar to that of $A^n$, where $A$ and $B$ are the matrices defined in part (a). You already know, from part (a), that eigenvalues is *not* such a criterion.

Given the vectors $p$ and $q$ computed in parts (b) and (c), plot, in a new figure, their difference $(q - p)$.

The following statements would have been true if $A$ had been (close to) normal.

- Since all eigenvalues of $A$ equal 1, then $\|A^n \underline{u}\|$, where $\underline{u}$ is defined in part (b), would be $1^n \|\underline{u}\|$. Here $\| \ldots \|$ denotes the length of a vector.

- Eigenvalues $\lambda_{A+B}$ would be on the order of $\lambda_A + b_{max}$, where $b_{max}$ is the maximum (in magnitude) entry of $B$.

- One would then expect that $\|(A + B)^n \underline{u} - A^n \underline{u}\|$ would be on the order:

$$( (1 + b_{max})^n - 1^n ) \|\underline{u}\|, \qquad\qquad \text{(HW12.2)}$$

You know, of course, that none of the above statements apply to $A$ since $A$ is not normal. In particular, if you look at the length of vector $(p - q)$, you can see that it does not at all follow estimate (HW12.2).

Yet, is there some combination of any of the quantities $p$, $q$, $(p - q)$, or their lengths, that does have the order of magnitude of (HW12.2)? Based on your work in this part, choose (and complete, if needed) one of the two sentences below, copy the result to your paper, and put it in a box.

"None of the quantities $p = A^n \underline{u}$, $q = (A + B)\underline{u}$, $(p - q)$, their lengths, or combinations thereof can be used to quantify how close a given matrix ($A + B$ in this case) is to some nondiagonalizable matrix ($A$ in this case)."

or

"Quantity _____, as defined above, is a quantitative measure of how close a given matrix ($A + B$ in this case) is to some nondiagonalizable matrix ($A$ in this case)."

(e) Finally, let us understand the reason behind your observation in part (c). Recall that eigenvalues of a non-diagonalizable matrix $C$ do not describe the large-$n$ behavior of $C^n \underline{x}$ simply because $C$ does not have enough eigenvectors for a basis. However, as you have, hopefully, explained it in part (c), matrix $(A + B)$ *is* diagonalizable and hence its eigenvectors *do* form a basis. But,

the **the Key question** is: *How good is that basis?*

To answer it, find the eigenvectors of $(A + B)$ using the command `eig` (read `help` for it first). Matlab gives you those eigenvectors organized into the columns of a square matrix; let us call it $S$. Find `cond(S)`. Using this number, answer the key question highlighted above in *italic* following these logical steps:

(1) Based on the condition number of $S$, say whether it is singular, close to singular, not so close to singular, or very far from singular. (Review Sec. 10.5 about the condition number as needed.)

(2a) What can one say of columns of a singular matrix? (If you forgot, find the answer in a Linear Algebra textbook or elsewhere.)

(2b) Do column of a singular matrix form a basis?

(3a) Given your answer to question (2a), what could you then reasonably say about columns of a near-singular matrix?

3

(3b)  Based on your answer to (3a), do columns of a near-singular matrix form a basis?

(4)  Based on your answers to steps (1) – (3), answer the Key question highlighted above in italic.

(5)  The answer you gave in (4) is the reason behind your observations in parts (c) and (d). In your paper, state this in the form:

Since ⟨answer to Key question⟩,  then  ⟨your observation in part (c)⟩

and similarly for part (d); here you need to provide details for the phrases inside the angle brackets. Then, | put each of these two statements in a box |.


3. This problem is worth **2.5 points**.

This problem expands on Problem 2 in two ways. First, you will see that a matrix may not *look* like being close to a non-diagonalizable matrix, and yet it will behave as such when acting on certain vectors. Second, you will visualize what a poor basis looks like, and *what is poor* about it.

Consider matrix

$$A = \begin{pmatrix} 0.7 & 1 \\ 0 & 0.9 \end{pmatrix}.$$

Its eigenvalues, $0.7$ and $0.9$, are different. So, it does not appear to be close to a non-diagonalizable matrix. Yet, it *behaves* as one, as you will see below. The features that are responsible for this are:
(i)  $A$ is a triangular matrix (like the non-diagonalizable matrices considered in Lecture 12),  and
(ii)  the above-diagonal entry, $1$, is significantly greater than the difference of the eigenvalues, $0.2$.

(a)  Let  $\underline{v}_1 = (1\ 0)^T$. Make a table of the lengths of vectors  $A^n \underline{v}_1$ for  $n = 1, \ldots, 10$. Round your answers to two significant digits. Is the behavior you have observed explained by the eigenvalue(s) of $A$?

*Note 1* :   The length of a vector $(x_1\ x_2)^T$ is $\sqrt{x_1^2 + x_2^2}$. It is a simple measure that allows one to tell whether a repeated action of $A$ reduces or magnifies a given vector.

*Note 2*:  This is not the final version of the table that you will be asked to submit. In part (e) you will be asked to expand it, and a similar table in part (b), in a certain way.

(b)  Let  $\underline{v}_2 = (0\ 1)^T$. Augment your table with the lengths of vectors  $A^n \underline{v}_2$ for  $n = 1, \ldots, 10$. Is the behavior you have observed explained by the eigenvalue(s) of $A$?

(c)  Find the eigenvectors $\underline{s}_1$ and $\underline{s}_2$ of matrix $A$. You can do so either by hand or on a computer, but make sure that both eigenvectors *have length one* (Matlab does this by default).

Draw $\underline{s}_1$ and $\underline{s}_2$ to scale.

Find the condition number of a matrix  $S \equiv [\underline{s}_1,\ \underline{s}_2]$. (The number you will find is many times smaller than that in Problem 2(e), but keep in mind that here you have a small, $2 \times 2$, matrix, while in Problem 2(e) the matrices were $20 \times 20$.)

Describe how your drawing of $\underline{s}_1$ and $\underline{s}_2$, the condition number, and your conclusion for Problem 2 all together give a reason for the behavior observed in either part (a) or part (b) of this problem.

(d)  Let us now explain what property of a basis makes it "poor". This property is *inefficiency*. To explain this, let us first give a
*Layman Example of Inefficiency*:  Suppose that you are assigned to dig a 2-ft hole in your backyards in 2 days.  The most efficient way to do so would be to dig about a foot on each of the days. Digging the entire hole in just one day is still reasonably efficient.   An inefficient way would be to dig 5 ft on day 1 and then put back 3 ft of dirt on the second day.

Let us now return to our basis $\{\underline{s}_1, \underline{s}_2\}$. For each of the vectors $\underline{v}_1$ and $\underline{v}_2$ from parts (a) and (b), find their coordinates $c_1, c_2$ in that basis. That is, solve the vector equation

$$\underline{v} = c_1\underline{s}_1 + c_2\underline{s}_2 \tag{HW12.3}$$

for $c_1$ and $c_2$. See the *Note* below on how to do so.
(Of course, each of $\underline{v}_1$ and $\underline{v}_2$ will have its "own" $c_1$ and $c_2$; i.e., you will have two pairs of $c_1, c_2$.)

Now comes a *critical step*: For each of $\underline{v}_1$ and $\underline{v}_2$, make a *to-scale* sketch illustrating expansion (HW12.3). Be prepared that one of these sketches will be quite large (and you may want to make it in the landscape orientation of the page).

Based on these sketches and the Layman Example of Inefficiency above, *explain* for which of the vectors $\underline{v}_1$ and $\underline{v}_2$, the basis $\{\underline{s}_1, \underline{s}_2\}$ is inefficient.

*Note*: By formula (7) of the document "Background from Linear Algebra", equation (HW12.3) is equivalent to the matrix equation $S\,(c_1, c_2)^T = \underline{v}$ where matrix $S$ is defined in part (c). You can solve this equation in Matlab. Round the answer to one decimal place.

(e) Finally, let us illustrate how this inefficiency led to the results observed in one of the parts (a) or (b).

Recall from Lecture 5 that the result of matrix $A$ acting $n$ times on vector $\underline{v}$ is:

$$A^n\underline{v} = \lambda_1^n c_1\underline{s}_1 + \lambda_2^n c_2\underline{s}_2.$$

The (square of the) length of this vector is found as:

$$\|A^n\underline{v}\|^2 = (\lambda_1^n c_1\underline{s}_1 + \lambda_2^n c_2\underline{s}_2)^T(\lambda_1^n c_1\underline{s}_1 + \lambda_2^n c_2\underline{s}_2). \tag{HW12.4}$$

As you cross-multiply terms, note that $\underline{s}_i^T \underline{s}_i = 1$ ($i = 1, 2$) by Matlab's convention to normalize eigenvectors to have length 1. Further note that since $\underline{s}_1$ and $\underline{s}_2$ are almost parallel, one has

$$\underline{s}_1^T \underline{s}_2 \approx 1. \tag{HW12.5}$$

Verify this with Matlab and give the actual value on the right-hand side of that equation.

Using equation (HW12.5) where you replace "$\approx 1$" with "$= 1$", show that (HW12.4) reduces to

$$\|A^n\underline{v}\|^2 = (\lambda_1^n c_1 + \lambda_2^n c_2)^2. \tag{HW12.6}$$

Now, between the tables in part (a) and (b), focus on the one corresponding to that vector $v$ for which you found expansion (HW12.3) to be inefficient. For that table only, augment it with three columns showing: $\lambda_1^n c_1$, $\lambda_2^n c_2$, and $|\lambda_1^n c_1 + \lambda_2^n c_2|$. Note that the last column in each part should be close[7] to that part's $\|A^n\underline{v}\|$.

FYI − 1: The inefficiency of the basis $\{\underline{s}_1, \underline{s}_2\}$ leads to *cancellation* of two relatively large terms, resulting in a smaller final term. A similar, although more dramatic, cancellation (called catastrophic cancellation) leads to loss of computational accuracy in certain computer calculations.

FYI − 2: You could have noticed that the percentage difference between $\|A^n\underline{v}\|$ and $|\lambda_1^n c_1 + \lambda_2^n c_2|$ (over 10%) is much greater than that between $\underline{s}_1^T \underline{s}_2$ and 1 (2%). The reason is the aforementioned cancellation. Indeed, the percentage difference

$$\frac{\|A^n\underline{v}\|}{|\lambda_1^n c_1 + \lambda_2^n c_2|} - 1$$

---

[7]See the FYI–2 as to why it is not as close as one might expect.

should be compared not with $(\underline{s}_1^T \underline{s}_2 - 1)$, but with

$$(\underline{s}_1^T \underline{s}_2 - 1) \cdot \frac{|2\lambda_1^n \lambda_2^n c_1 c_2|}{(\lambda_1^n c_1 + \lambda_2^n c_2)^2} \, .$$

Given that $c_2 \approx -c_1$, then for $n$ not too large, the two terms in the denominator partially cancel one another, thus making the denominator small compared to the numerator.

**Bonus** (worth **0.5 pt**; credit will be given only if the solution is mostly correct)

In Problem 2 you have explained the size of vector $\underline{p}$. Now explain its shape (that is, *both* the reason for the oscillations and the envelope of the oscillations). You may need to extrapolate results of Sec. 12.3 from the $2 \times 2$ and $3 \times 3$ cases to the $M \times M$ case.