



DDA 3005 — Numerical Methods

Solutions 4

Problem 1 (Power Iteration with Shift):

(approx. 20 points)

Let the matrix $\mathbf{A} \in \mathbb{R}^{4 \times 4}$ and the initial point $\mathbf{x}^0 \in \mathbb{R}^4$ be given via

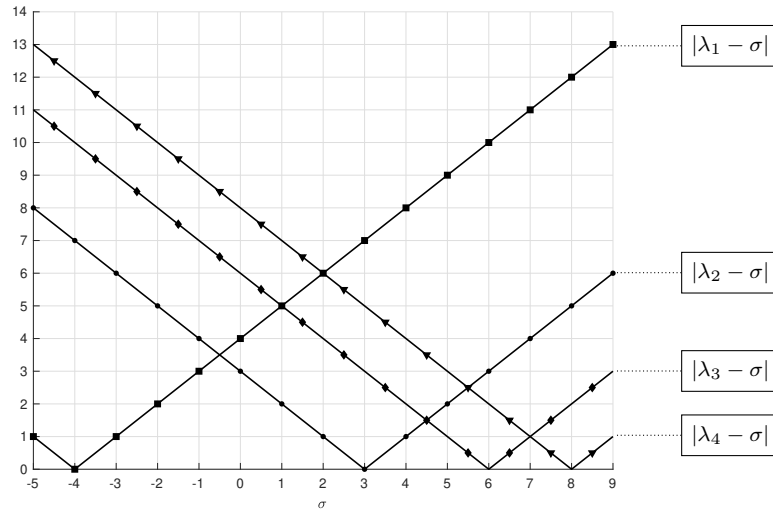
$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 5 & 0 \\ 0 & 3 & 0 & 0 \\ 5 & 0 & 1 & 0 \\ 0 & 0 & 0 & 8 \end{bmatrix}, \quad \tilde{\mathbf{x}}^0 = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}, \quad \mathbf{x}^0 = \frac{\tilde{\mathbf{x}}^0}{\|\tilde{\mathbf{x}}^0\|}. \quad (1)$$

In this problem, we want to apply the (normalized) power iteration with shift $\sigma \in \mathbb{R}$ to \mathbf{A} :

$$\tilde{\mathbf{x}}^k = (\mathbf{A} - \sigma \mathbf{I})\mathbf{x}^k, \quad \mathbf{x}^k = \tilde{\mathbf{x}}^k / \|\tilde{\mathbf{x}}^k\|, \quad \sigma_k = (\mathbf{x}^k)^\top \mathbf{A} \mathbf{x}^k, \quad k = 1, 2, \dots$$

Let $(\lambda_i, \mathbf{v}_i)$, $i = 1, \dots, 4$, further denote the corresponding eigenpairs of \mathbf{A} . It can be shown that $(\mathbf{x}^0)^\top \mathbf{v}_i \neq 0$ for all $i = 1, \dots, 4$.

The plot below depicts the mappings $\sigma \mapsto |\lambda_i - \sigma|$, $i = 1, 2, 3, 4$ for the different eigenvalues of \mathbf{A} and choices of σ .



- State the eigenvalues $\lambda_1, \dots, \lambda_4$ of \mathbf{A} .
- Discuss which eigenpairs $(\lambda_i, \mathbf{v}_i)$ of \mathbf{A} can be recovered by the (normalized) power iteration using suitable choices of the shift σ . Can all eigenpairs of \mathbf{A} be recovered? Provide detailed explanations!
- Derive the optimal choice of the shift σ for which the power iteration (1) converges with the fastest possible (optimal) convergence rate.

Solution :

- a) We can directly read the eigenvalues of \mathbf{A} from the plot: $\lambda_1 = -4$, $\lambda_2 = 3$, $\lambda_3 = 6$, and $\lambda_4 = 8$.
- b) The power iteration with shift σ “converges” to the eigenpair of $\mathbf{A} - \sigma \mathbf{I}$ corresponding to the largest eigenvalue of $\mathbf{A} - \sigma \mathbf{I}$ (in magnitude). As the eigenpairs of $\mathbf{A} - \sigma \mathbf{I}$ are given by $(\lambda_i - \sigma, \mathbf{v}_i)$, we can use the depicted plot to determine which eigenvalue has the largest magnitude. In particular, we are interested in finding

$$\max_{i=1,2,3,4} |\lambda_i - \sigma|$$

Clearly, for $\sigma < 2$, $\lambda_4 - \sigma = 8 - \sigma$ has the largest magnitude; for $\sigma > 2$, $\lambda_1 - \sigma = -4 - \sigma$ has the largest magnitude. For $\sigma = 2$, it holds that $|\lambda_4 - 2| = |\lambda_1 - 2|$, i.e., there is no gap between the first and second largest eigenvalue and the power method does not converge.

As $(\mathbf{x}^0)^\top \mathbf{v}_1 \neq 0$ and $(\mathbf{x}^0)^\top \mathbf{v}_4 \neq 0$, we can conclude that the power method with shift only allows to recover the eigenpairs $(-4, \mathbf{v}_1)$ and $(8, \mathbf{v}_4)$.

- c) The rate of convergence depends on the ratio between the second and first largest eigenvalue of $\mathbf{A} - \sigma \mathbf{I}$ (in magnitude). Let us denote those two eigenvalues by $\tilde{\lambda}_1 - \sigma$ and $\tilde{\lambda}_2 - \sigma$ ($\tilde{\lambda}_1$ and $\tilde{\lambda}_2$ can change with the choice of σ). Then, the rate is given by:

$$\frac{|\tilde{\lambda}_2 - \sigma|}{|\tilde{\lambda}_1 - \sigma|} = 1 - \frac{|\tilde{\lambda}_1 - \sigma| - |\tilde{\lambda}_2 - \sigma|}{|\tilde{\lambda}_1 - \sigma|}$$

We see that the best rate can be obtained when $|\tilde{\lambda}_1 - \sigma| - |\tilde{\lambda}_2 - \sigma|$ is as large as possible and $|\tilde{\lambda}_1 - \sigma|$ is as small as possible. Using the plot, we can easily find $\tilde{\lambda}_1$ and $\tilde{\lambda}_2$. For $\sigma \leq 1$, we have

$$\frac{|\tilde{\lambda}_1 - \sigma| - |\tilde{\lambda}_2 - \sigma|}{|\tilde{\lambda}_1 - \sigma|} = \frac{2}{8 - \sigma} \implies \text{best value at } \sigma = 1: \frac{2}{7}.$$

For $\sigma \geq 5.5$, we have

$$\frac{|\tilde{\lambda}_1 - \sigma| - |\tilde{\lambda}_2 - \sigma|}{|\tilde{\lambda}_1 - \sigma|} = \frac{7}{4 + \sigma} \implies \text{best value at } \sigma = 5.5: \frac{14}{19}.$$

For $\sigma \in (1, 2]$, it holds that $\frac{|\tilde{\lambda}_1 - \sigma| - |\tilde{\lambda}_2 - \sigma|}{|\tilde{\lambda}_1 - \sigma|} = \frac{8 - \sigma - 4 - \sigma}{8 - \sigma} = \frac{4 - 2\sigma}{8 - \sigma}$ – this is largest for $\sigma = 1$ (again yielding $\frac{2}{7}$).

For $\sigma \in [2, 5.5)$, it holds that $\frac{|\tilde{\lambda}_1 - \sigma| - |\tilde{\lambda}_2 - \sigma|}{|\tilde{\lambda}_1 - \sigma|} = \frac{4 + \sigma - 8 + \sigma}{4 + \sigma} = \frac{2\sigma - 4}{4 + \sigma}$ – which is largest for $\sigma = 5.5$.

Overall, we can conclude that the optimal rate is attained for the shift $\sigma = 5.5$.

Problem 2 (Power Methods):

(approx. 30 points)

In this exercise, we want to calculate eigen-pairs of the 5×5 matrix using power methods:

$$\mathbf{A} = \begin{bmatrix} 46 & 99 & 45 & 24 & -27 \\ -6 & -32 & -6 & 4 & 18 \\ 18 & 18 & 19 & 6 & 0 \\ -9 & -90 & -9 & -47 & 27 \\ 12 & 39 & 12 & 28 & 19 \end{bmatrix}.$$

- a) Implement the power method and compute the largest eigenvalue and the corresponding eigenvector of \mathbf{A} , e.g., using the initial point $\mathbf{x}^0 = \mathbf{e}_1 = [1, 0, 0, 0, 0]^\top$.

Choose a proper stopping criterion for the power method and state how many iterations your algorithm requires to recover the largest eigenvalue (using modest or high accuracy). Explain the choice of your termination criterion.

- b) Implement the inverse iteration (with shift σ) to calculate the remaining eigen-pairs of the matrix \mathbf{A} (again using $\mathbf{x}^0 = \mathbf{e}_1$ as initial point). Possible choices for the shift σ are $\sigma \in \{0, 46, -32, 19, -47\}$.
- c) Implement the Rayleigh Quotient iteration using the five initial points $\mathbf{x}^0 \in \{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4, \mathbf{e}_5\}$ and state the respective obtained eigen-pairs. Compare the performance of the Rayleigh Quotient iteration and the power iterations tested in part a) and b).

Solution :

- a) An implementation of the power method is shown below:

```

1  % problem 2a: power method
2
3  A = [46,99,45,24,-27;-6,-32,-6,4,18;18,18,19,6,0;-9,-90,-9,-47,27;12,39,12,28,19];
4  maxit = 10000;
5
6  x = [1;0;0;0;0];
7  nx = norm(x);
8
9  for i = 1:maxit
10     x      = A*x;
11     nx      = norm(x);
12     x      = x/nx;
13     sig     = x'*(A*x);
14     tol     = norm(A*x-sig*x);
15
16     if tol <= 1e-12
17         break
18     end
19
20     fprintf(1,['%i ; %1.4e ; %1.10f\n',i,tol,sig]);
21 end

```

The method converges to the eigenvalue 55 within 848 iterations. We choose $\|\mathbf{Ax} - \sigma\mathbf{x}\| \leq \text{tol}$ as stopping criterion. In the case $\text{tol} = 0$, this implies that σ is an eigenvalue of \mathbf{A} with eigenvector \mathbf{x} .

- b) An implementation of the inverse iteration is shown below:

```

1  % problem 2b: inverse iteration
2
3  A = [46,99,45,24,-27;-6,-32,-6,4,18;18,18,19,6,0;-9,-90,-9,-47,27;12,39,12,28,19];
4  maxit = 10000;
5
6  for mu = [0,46,-32,19,-47]
7     x = [1;0;0;0;0];
8     nx = norm(x);
9

```

```

10     fprintf(1, '\nInverse iteration with shift mu = %g\n', mu);
11
12     for i = 1:maxit
13         x = (A-mu*eye(5))\x;
14         nx = norm(x);
15         x = x/nx;
16         sig = x'*(A*x);
17         tol = norm(A*x-sig*x);
18
19         fprintf(1, '[%i] ; %1.4e ; %1.10f\n', i, tol, sig);
20         if tol <= 1e-13
21             break
22         end
23     end
24 end

```

As suggested in the problem description, we use the shifts $\sigma \in \{0, 46, -32, 19, -47\}$. The returned results are as follows:

σ	iterations	eigenvalue λ	eigenvector \mathbf{v}
0	11	1.0000	$\mathbf{v}^T = [0.7071, -0.0000, -0.7071, 0.0000, -0.0000]$
46	47	55.0000	$\mathbf{v}^T = [0.8165, 0.0000, 0.4082, -0.0000, 0.4082]$
-32	26	-26.0000	$\mathbf{v}^T = [-0.0000, 0.3015, -0.0000, -0.9045, 0.3015]$
19	49	28.0000	$\mathbf{v}^T = [0.3015, -0.3015, 0.0000, -0.0000, -0.9045]$
-47	26	-53.0000	$\mathbf{v}^T = [-0.3015, 0.0000, 0.0000, 0.9045, -0.3015]$

c) An implementation of the Rayleigh Quotient iteration is shown below:

```

1  % problem 2c: Rayleight Quotient iteration
2
3  A = [46,99,45,24,-27;-6,-32,-6,4,18;18,18,19,6,0;-9,-90,-9,-47,27;12,39,12,28,19];
4  maxit = 10000;
5
6  for i = 1:5
7      x = zeros(5,1);
8      x(i) = 1;
9
10     fprintf(1, '\nRayleigh Quotient iteration with x0 = e%i\n', i);
11
12     sig = x'*(A*x)/norm(x);
13
14     for j = 1:maxit
15         x = (A-sig*eye(5))\x;
16         nx = norm(x);
17         x = x/nx;
18         sig = x'*(A*x);
19         tol = norm(A*x-sig*x);
20
21         fprintf(1, '[%i] ; %1.4e ; %1.10f\n', j, tol, sig);
22         if tol <= 1e-13
23             break
24         end
25     end
26 end

```

Running this code, we can obtain the following results:

\mathbf{x}^0	iterations	eigenvalue λ	eigenvector \mathbf{v}
\mathbf{e}_1	5	55.0000	$\mathbf{v}^\top = [-0.8165, 0.0000, -0.4082, 0.0000, -0.4082]$
\mathbf{e}_2	6	-26.0000	$\mathbf{v}^\top = [0, 0.3015, 0.0000, -0.9045, 0.3015]$
\mathbf{e}_3	8	1.0000	$\mathbf{v}^\top = [-0.7071, -0.0000, 0.7071, 0.0000, -0.0000]$
\mathbf{e}_4	5	-53.0000	$\mathbf{v}^\top = [0.3015, 0.0000, -0.0000, -0.9045, 0.3015]$
\mathbf{e}_5	5	28.0000	$\mathbf{v}^\top = [0.3015, -0.3015, -0.0000, -0.0000, -0.9045]$

Hence, in this example, the Rayleigh Quotient iteration also manages to find all eigenvalues and eigenvectors of \mathbf{A} . The method converges significantly faster than the inverse iteration even though \mathbf{A} is not a symmetric matrix. However, we cannot necessarily ensure and observe fast cubic convergence in this case!

Problem 3 (Computing Roots of Polynomials):

(approx. 50 points)

For given degree $n \in \mathbb{N}$ and coefficients $a_i \in \mathbb{R}$, $i = 0, \dots, n-1$, let us consider the univariate polynomial

$$p(x) := x^n + a_{n-1}x^{n-1} + \dots + a_2x^2 + a_1x + a_0.$$

In this exercise, we want to develop an algorithm that allows computing all roots of the polynomial p using eigenvalues and QR iterations.

- a) The so-called companion matrix associated with the polynomial p is given by

$$\mathbf{C}_p := \begin{bmatrix} 0 & 0 & \dots & 0 & -a_0 \\ 1 & 0 & \dots & 0 & -a_1 \\ 0 & 1 & \dots & 0 & -a_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 1 & -a_{n-1} \end{bmatrix}.$$

Show that $\det(\mathbf{C}_p - \lambda \mathbf{I}) = (-1)^n p(\lambda)$ for all λ .

- b) Implement the QR algorithm (with shift) discussed in the lecture to calculate the eigenvalues of the companion matrix \mathbf{C}_p .

In order to improve the performance of your algorithm, you can utilize the following deflation strategy: let $\mathbf{X}^k \in \mathbb{R}^{n \times n}$ be the current iterate of the QR iteration. We can terminate the iteration, if $\|\mathbf{X}^k(\mathbf{n}, 1:\mathbf{n}-1)\| \leq \text{tol}$, i.e., if the norm of the first $n-1$ entries of the last row of \mathbf{X}^k is small. This ensures that the last row of \mathbf{C}_p has been reduced to upper triangular form and $\mathbf{X}^k(\mathbf{n}, \mathbf{n}) = \mathbf{X}_{nn}^k$ should be a good eigenvalue approximation. We can then continue the QR algorithm with the smaller (deflated) matrix $\mathbf{B} = \mathbf{X}^k(1:\mathbf{n}-1, 1:\mathbf{n}-1) \in \mathbb{R}^{(n-1) \times (n-1)}$ (or $\mathbb{C}^{(n-1) \times (n-1)}$) as new initial point. This process continues until we have recovered all eigenvalues of \mathbf{C}_p . Suitable choices for tol are 10^{-12} or 10^{-13} . You can also use a shifted QR iteration to further enhance performance (e.g., the Rayleigh quotient or Wilkinson shift can be a good strategy).

Test your implementation and find all roots of the polynomials:

$$- p(x) = x^4 - 324x^3 + 7175x^2 + 8100x - 180000.$$

$$- q(x) = x^{10} - 167x^9 + 10081x^8 - 251447x^7 + 1676815x^6 + 17367175x^5 - 66421125x^4 - 352378125x^3 + 454612500x^2 + 1949062500x.$$

Report all found roots of p and q (you can order the roots from smallest to largest).

- c) What happens if the coefficients a_2 , a_1 , and a_0 in the polynomial p are changed to $a_2 = 7225$, $a_1 = -8100$, and $a_0 = 180000$? Can you explain your observations?

Hint: Use a Wilkinson shift or plot the adjusted polynomial p .

Acknowledgements: This question is motivated by a conversation and discussion with Xiuyuan Wang.

Solution :

- a) We prove this result via induction. For $n = 1$, the polynomial p reduces to $p(x) = x + a_0$ and the associated companion matrix is given by

$$\mathbf{C}_p^1 = \mathbf{C}_p = -a_0 \implies \det(\mathbf{C}_p^1 - \lambda \mathbf{I}) = -(a_0 + \lambda) = -p(\lambda).$$

Hence, the base case is obviously true. We now assume that the induction hypothesis holds for \mathbf{C}_p^n and we consider $\mathbf{C}_p^{n+1} \in \mathbb{R}^{(n+1) \times (n+1)}$. Let us consider the specific form of $\mathbf{C}_p^{n+1} - \lambda \mathbf{I}$:

$$\mathbf{C}_p^{n+1} - \lambda \mathbf{I} = \begin{bmatrix} -\lambda & 0 & 0 & \cdots & 0 & -a_0 \\ 1 & -\lambda & 0 & \cdots & 0 & -a_1 \\ 0 & 1 & -\lambda & \cdots & 0 & -a_2 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -\lambda & -a_{n-1} \\ 0 & 0 & \cdots & 0 & 1 & -a_n - \lambda \end{bmatrix}.$$

Using the Laplace expansion of the determinant along the 1-st row, it follows:

$$\begin{aligned} \det(\mathbf{C}_p^{n+1} - \lambda \mathbf{I}) &= (-1)^2 \cdot (-\lambda) \cdot \det([\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{1,1}) \\ &\quad + (-1)^{n+2} \cdot (-a_0) \cdot \det([\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{1,n+1}). \end{aligned}$$

Here, $[\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{i,j}$ is the $n \times n$ submatrix of $\mathbf{C}_p^{n+1} - \lambda \mathbf{I}$ resulting from deleting the i -th row and j -th column. As $[\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{1,1}$ is a companion matrix with coefficients a_1, \dots, a_{n-1}, a_n , we can apply the induction hypothesis to obtain:

$$\det([\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{1,1}) = (-1)^n (\lambda^n + a_n \lambda^{n-1} + a_{n-1} \lambda^{n-2} + \cdots + a_2 \lambda + a_1).$$

Furthermore, it holds that

$$[\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{1,n+1} = \begin{bmatrix} 1 & -\lambda & 0 & \cdots & 0 \\ 0 & 1 & -\lambda & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 & -\lambda \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix}.$$

This is an upper triangular matrix and we easily see $\det([\mathbf{C}_p^{n+1} - \lambda \mathbf{I}]_{1,n+1}) = 1$. Combining the results, we obtain

$$\det(\mathbf{C}_p^{n+1} - \lambda \mathbf{I}) = (-1)^{n+1} (\lambda^{n+1} + a_n \lambda^n + \cdots + a_1 \lambda) + (-1)^{n+3} a_0 = (-1)^{n+1} p(\lambda).$$

This finishes our proof.

- b) An implementation of the QR algorithm with deflation and Rayleigh shift is shown below:

```

1  % demo_proots
2
3  %
4  % build companion matrices
5
6  A = [0,0,0,180000;
7       1,0,0,-8100;
8       0,1,0,-7175;
9       0,0,1,+324];
10
11 % A = [0,0,0,0,0,0,0,0,-1949062500 ; ...
12 %      1,0,0,0,0,0,0,0,-454612500; ...
13 %      0,1,0,0,0,0,0,0,352378125; ...
14 %      0,0,1,0,0,0,0,0,66421125; ...
15 %      0,0,0,1,0,0,0,0,-17367175; ...
16 %      0,0,0,0,1,0,0,0,-1676815; ...
17 %      0,0,0,0,0,1,0,0,251447; ...
18 %      0,0,0,0,0,0,1,0,-10081; ...
19 %      0,0,0,0,0,0,0,1,167];
20
21 %
22 % setup QR algorithm
23
24 opts.maxit = 100;
25 opts.tol   = 1e-13;
26 n          = size(A,1);
27
28 lambda     = zeros(n,1);
29 B          = A;
30
31 fprintf(1, '\nQR-Iteration with Deflation and Rayleigh Quotient Shift\n');
32 opts.shift = 'Rayleigh';
33
34 for i = 1:n
35     fprintf(1, 'iteration: %i; matrix size [%i x %i]\n', i, n-i+1, n-i+1);
36     [X, lam] = qr_iteration(B, opts);
37
38     B        = X(1:end-1, 1:end-1);
39     lambda(n-i+1) = lam;
40 end
41
42 %
43 % display
44
45 lambda = sort(lambda, 'ascend');
46
47 fprintf(1, '\n');
48 disp(lambda)
49
50 %
51 % inner function
52
53 function [X, lam] = qr_iteration(B, opts)
54     X      = B;
55     m      = size(B, 1);
56

```

```

57     for i = 1:opts.maxit
58         switch opts.shift
59             case 'pure'
60                 sig = 0;
61             case 'Rayleigh'
62                 sig = X(m,m);
63             case 'Wilkinson'
64                 mu = eig(X(max(m-1,1):m,max(m-1,1):m));
65                 [~,ind] = min(abs(mu-X(m,m)));
66                 sig = mu(ind);
67         end
68
69         [Q,R] = qr(X-sig*eye(m));
70         X = R*Q+sig*eye(m);
71
72         fprintf(1,' — [%i] ; %1.4e ; %1.4f\n',i,norm(X(m,1:m-1)),X(m,m));
73         disp(X(m,m))
74
75         if norm(X(m,1:m-1)) <= opts.tol
76             break
77         end
78     end
79     lam = X(end,end);
80 end

```

The returned roots for p and q are given by:

$$\begin{aligned}
 p: \quad & x_1 = -5, \quad x_2 = 5, \quad x_3 = 24, \quad x_4 = 300, \\
 q: \quad & x_1 = -5, \quad x_2 = -3, \quad x_3 = -3, \quad x_4 = 3, \quad x_5 = 5, \quad x_6 = 30, \quad x_7 = 35, \\
 & x_8 = 50, \quad x_9 = 55.
 \end{aligned}$$

Notice that the polynomial q has the additional root $x_0 = 0$ which was not explicitly included in the companion matrix. Overall, the method converges very quickly — each of the deflation steps only requires 1–10 iterations using a Rayleigh quotient shift. (Only the double root $x_2 = x_3 = -3$ for q needs more steps and convergence requires 23 iterations).

- c) The QR algorithm with Rayleigh shift does not seem to converge in this case. Utilizing the Wilkinson shift, we can recover the roots

$$x_1 = -5i, \quad x_1 = 5i, \quad x_3 = 24, \quad x_4 = 300,$$

i.e., p has complex roots in this case! Hence, the behavior of the QR algorithm with Rayleigh shift is not surprising as this version of the algorithm does not introduce complex numbers (the initial matrix \mathbf{A} is real, the QR factorization produce real matrices \mathbf{Q} and \mathbf{R} , the Rayleigh shift remains real as well). The Wilkinson shift resolves this issue as the eigenvalues of the considered 2×2 submatrix can be complex. (There are also other ways to refine / adjust the QR algorithm to generally handle this case).

The Python code for problem 3 is similar to the presented MATLAB code can be found in the attachments.