

Problem 1. For an $n \times n$ matrix A , where

$$A_{1,1} = 2, \quad A_{1,2} = -1, \quad A_{n,n} = 2, \quad A_{n-1,n} = -1$$

and

$$A_{i,i} = 2, \quad A_{i,i+1} = -1, \quad A_{i,i-1} = -1, \quad \forall i \neq 1, n$$

and $A_{i,j} = 0$ otherwise. Compute its LU factorization with MATLAB/Python. Can the LU factorization be obtained faster than $O(n^3)$ complexity? If so, what would the algorithm be?

Solution: The matrix A is the following form:

$$A = \begin{bmatrix} 2 & -1 & 0 & \cdots & 0 & 0 & 0 \\ -1 & 2 & -1 & \ddots & \vdots & \vdots & \vdots \\ 0 & -1 & 2 & \ddots & -1 & 0 & 0 \\ \vdots & \ddots & \ddots & \ddots & 2 & -1 & 0 \\ 0 & \cdots & -1 & 2 & -1 & \ddots & \vdots \\ 0 & \cdots & 0 & -1 & 2 & -1 & \\ 0 & \cdots & 0 & 0 & -1 & 2 & \end{bmatrix}$$

The code is as follows:

```

1  A = generate_tridiagonal_matrix(5);
2  [L_A, U_A] = LU_decomp(A);
3
4
5  function A = generate_tridiagonal_matrix(n)
6      % Generates an n x n symmetric tridiagonal matrix with:
7      % 2 on the diagonal and -1 on the sub- and super-diagonals
8
9      e = ones(n,1);
10     A = 2*diag(e) - diag(e(1:end-1),1) - diag(e(1:end-1),-1);
11 end
12
13 function [L, U] = LU_decomp(A)
14     n_tuple = size(A);
15     assert(n_tuple(1) == n_tuple(2));
16     n = n_tuple(1);
17     U = A;
18     L = eye(n_tuple(1), n_tuple(2));
19     for j = 1:n
20         for i = j+1:n
21             if U(j,j) == 0
22                 error("entry in U is 0");
23             end
24             L(i,j) = U(i,j) / U(j,j);
25             U(i, j:n) = U(i,j:n) - L(i,j) * U(j,j:n);

```

```
26         end
27     end
28 end
```

Note that given the matrix is a tridiagonal matrix, we know that we can exploit the structure of the matrix to get a faster run time. Let m denote the main diagonal, l denote the lower diagonal and u denote the upper diagonal. We can see that it suffices to only use a singular for loop and to construct the LU decomposition as we do not have to iterate through the columns of the matrix, other along the main diagonal. The pseudocode is as follows:

Algorithm 1 LU Decomposition of a Symmetric Tridiagonal Matrix

Require: Main diagonal $a[1..n]$, sub/super-diagonal $b[1..n-1]$

Ensure: Subdiagonal $l[1..n-1]$ of L , diagonal $u[1..n]$ of U

- 1: $u[1] \leftarrow a[1]$
 - 2: **for** $i = 2$ to n **do**
 - 3: $l[i-1] \leftarrow b[i-1]/u[i-1]$
 - 4: $u[i] \leftarrow a[i] - l[i-1] \cdot b[i-1]$
 - 5: **end for**
 - 6: **return** l, u
-

□

Problem 2. Implement a gradient descent method (with fixed step size) to solve

$$\min_{x \in \mathbb{R}^n} \frac{1}{2}(x - x_*)^T A(x - x_*)$$

where $x_* = [1, 2, \dots, n]^T$. What is the complexity per step of gradient descent? (Show the complexity as n grows using a plot). Furthermore, by repeating the experiment for different n , extract the rate of convergence and show its dependency on n using big-O notation.

Solution:

```

1      n_vals = 10:10:400;
2      tolerance = 1.5e-5;
3
4      iters_vec = zeros(length(n_vals), 1);
5      total_flops_vec = zeros(length(n_vals), 1);
6      rel_error_mat = cell(length(n_vals), 1);
7
8      for i = 1:length(n_vals)
9          disp("grad");
10         disp(i)
11         n = n_vals(i);
12         [num_iter, total_flops, rel_errors] = grad_descent(n,
13             tolerance);
14         total_flops_vec(i) = total_flops;
15         iters_vec(i) = num_iter;
16         rel_error_mat{i} = rel_errors;
17     end
18
19     rho_vec = zeros(length(n_vals), 1); % Store convergence
20         rates
21
22     for i = 1:length(n_vals)
23         rel_errors = rel_error_mat{i};
24         k = 0:length(rel_errors)-1;
25         log_err = log(rel_errors(:));
26         coeffs = polyfit(k', log_err, 1);
27         rho_vec(i) = coeffs(1);
28     end
29
30     figure;
31
32     % Plot 1: Total FLOPs vs n
33     subplot(2, 2, 1);
34     loglog(n_vals, total_flops_vec, 'o-', 'LineWidth', 2);
35     xlabel('Matrix size n');
36     ylabel('Total FLOPs');
37     title('Log-Log: Total FLOPs vs. n');
38     grid on;

```

```
37     hold on;
38
39     log_n = log(n_vals);
40     log_flops = log(total_flops_vec);
41     coeffs_flops = polyfit(log_n, log_flops, 1);
42     b = coeffs_flops(1);
43     a = exp(coeffs_flops(2));
44     fit_flops = a * n_vals.^b;
45     loglog(n_vals, fit_flops, '--', 'LineWidth', 1.5);
46     text(n_vals(round(end * 0.3)), fit_flops(round(end * 0.3)),
47
48 % Plot 2: Relative Error vs Iteration
49 subplot(2, 2, 2);
50 hold on;
51 for i = 1:length(n_vals)
52     rel_errors = rel_error_mat{i};
53     k = 0:length(rel_errors)-1;
54     semilogy(k, rel_errors, 'DisplayName', sprintf('n = %d',
55     n_vals(i)));
56
57 end
58 xlabel('Iteration');
59 ylabel('Relative Error (log scale)');
60 title('Semi-Log: Relative Error vs. Iteration');
61 legend('show', 'Location', 'northeastoutside');
62 grid on;
63
64 % Plot 3: |rho| vs. n
65 subplot(2, 2, 3);
66 loglog(n_vals, abs(rho_vec), 'o-', 'LineWidth', 2);
67 xlabel('Matrix size n');
68 ylabel('|Convergence Rate|');
69 title('Log-Log: | rho | vs. n');
70 grid on;
71
72 % Plot 4: Fit |rho| vs. n
73 subplot(2, 2, 4);
74 loglog(n_vals, abs(rho_vec), 'o-', 'LineWidth', 2);
75 xlabel('Matrix size n');
76 ylabel('|Convergence Rate|');
77 title('Log-Log: Fitted | rho | vs. n');
78 grid on;
79 hold on;
80
81 log_rho = log(abs(rho_vec));
82 coeffs_rho = polyfit(log_n, log_rho, 1);
83 slope = coeffs_rho(1);
84 a_rho = exp(coeffs_rho(2));
```

```
83     fit_rho = a_rho * n_vals.^slope;
84     loglog(n_vals, fit_rho, '--', 'LineWidth', 1.5);
85     text(n_vals(round(end * 0.3)), fit_rho(round(end * 0.3)),
86          'FontSize', 10, 'BackgroundColor', 'white');
87
88     % User defined function
89     function [num_iter, total_flops, rel_errors] = grad_descent(
90         n, tolerance)
91         A = generate_tridiagonal_matrix(n);
92         iterations = 0;
93         x_star = (1:n)';
94         x_current = zeros(n,1);
95         eigenvalues = eig(A);
96         s = 2 / (min(eigenvalues) + max(eigenvalues));
97         rel_errors = [];
98         norm_x_star = norm(x_star);
99         rel_errors(end+1) = norm(x_current - x_star) /
100             norm_x_star;
101         tic;
102         while true
103             iterations = iterations + 1;
104             grad = A * (x_current - x_star);
105             x_next = x_current - s * grad;
106
107             % Track relative error
108             rel_error = norm(x_next - x_star) / norm_x_star;
109             rel_errors(end+1) = rel_error;
110
111             if rel_error < tolerance
112                 break;
113             end
114             x_current = x_next;
115         end
116         flops_per_iteration = n^2;
117         num_iter = iterations;
118         total_flops = flops_per_iteration * num_iter;
119     end
```

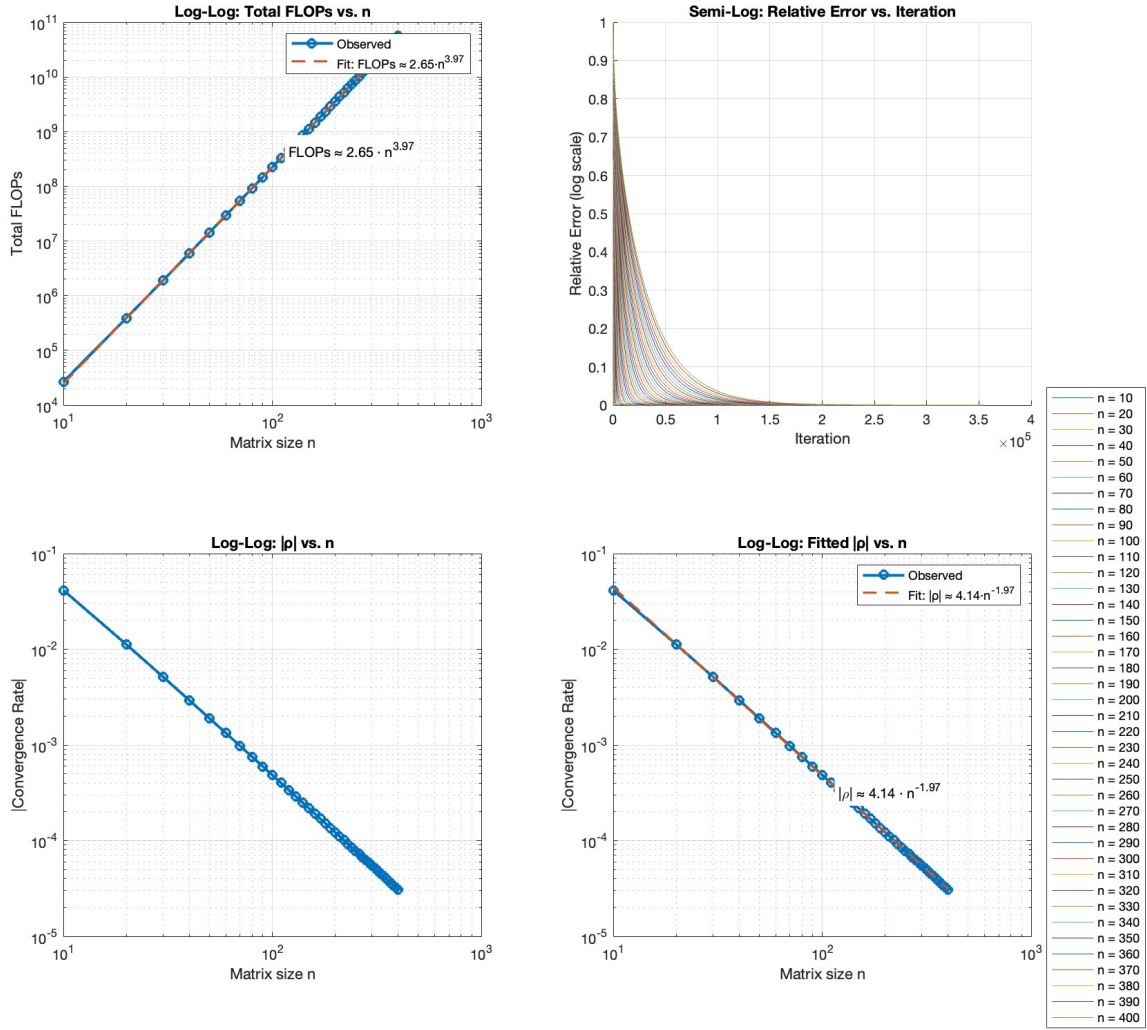


Figure 1: Gradient descent graphs

We can see that the computational complexity of this algorithm is $\mathcal{O}(n^4)$ from the top right graph. This implies that the rate of iteration increasing is $\mathcal{O}(n^2)$ as within the actual iteration, we do a matrix vector operation, $\mathcal{O}(n^2)$. Thus, the complexity for each step is $\mathcal{O}(n^2)$ the bottom right graph, we know that convergence rate is $\mathcal{O}(n^2)$ as we can extract the rate of convergence from the negative of the exponent. This is proved in the following, and let e_k denote the relative error:

$$\rho = \max_i |1 - s \lambda_i(A)| = \frac{\lambda_{\max}(A) - \lambda_{\min}(A)}{\lambda_{\max}(A) + \lambda_{\min}(A)} = \frac{\kappa - 1}{\kappa + 1} = 1 - \Theta\left(\frac{1}{n^2}\right).$$

Hence the error after k steps satisfies

$$\|x^{(k)} - x^*\| \leq \rho^k \|x^{(0)} - x^*\|.$$

To drive $\|x^{(k)} - x^*\| \leq \epsilon$, solve

$$\rho^k \leq \frac{\epsilon}{\|x^{(0)} - x^*\|} \implies k \geq \frac{\ln(\epsilon/\|x^{(0)} - x^*\|)}{\ln(\rho)} = O(n^2 \ln \frac{1}{\epsilon}).$$

since we are using a fixed ϵ , it suffices to conclude that our algorithm converges in $\mathcal{O}(n^2)$ time. \square

Problem 3. Repeat Problem 2 by implementing a conjugate gradient method.

Solution:

```
1      n_vals = 10:10:2000;
2      tolerance = 1.5e-5;
3
4      iters_vec = zeros(length(n_vals), 1);
5      total_flops_vec = zeros(length(n_vals), 1);
6      rho_vec = zeros(length(n_vals), 1);
7      rel_error_mat = cell(length(n_vals), 1);
8
9      for i = 1:length(n_vals)
10         disp("conj");
11         disp(i)
12         n = n_vals(i);
13         [num_iter, total_flops, rel_errors] = conj_grad_descent(
14             n, tolerance);
15
16         iters_vec(i) = num_iter;
17         total_flops_vec(i) = total_flops;
18         rel_error_mat{i} = rel_errors;
19
20         % Estimate convergence rate from slope of log(error) vs.
21         % iteration
22         k = 0:length(rel_errors)-1;
23         log_err = log(rel_errors(:));
24         coeffs = polyfit(k', log_err, 1);
25         rho_vec(i) = coeffs(1); % Negative convergence slope
26     end
27
28     % ---- Start Combined Plot ----
29     figure;
30
31     % ---- Plot 1: Log-Log Total FLOPs vs. n ----
32     subplot(2,2,1);
33     loglog(n_vals, total_flops_vec, 'o-', 'LineWidth', 2);
34     xlabel('Matrix size n');
35     ylabel('Total FLOPs');
36     title('Log-Log: Total FLOPs vs. n');
37     grid on;
38     hold on;
39
40     log_n = log(n_vals);
41     log_flops = log(total_flops_vec);
42     coeffs_flops = polyfit(log_n, log_flops, 1);
43     b = coeffs_flops(1);
44     a = exp(coeffs_flops(2));
```



```
43     fit_flops = a * n_vals.^b;
44     loglog(n_vals, fit_flops, '--', 'LineWidth', 1.5);
45     text(n_vals(round(end*0.3)), fit_flops(round(end*0.3)), ...
46
47
48     % ---- Plot 2: Semi-log Relative Error vs. Iteration ----
49     subplot(2,2,2);
50     hold on;
51     for i = 1:3:length(n_vals) % Plot fewer curves for
52         readability
53         rel_errors = rel_error_mat{i};
54         k = 0:length(rel_errors)-1;
55         semilogy(k, rel_errors, 'DisplayName', sprintf('n = %d',
56             n_vals(i)));
57     end
58     xlabel('Iteration');
59     ylabel('Relative Error (log scale)');
60     title('Semi-Log: Relative Error vs. Iteration');
61     legend('show', 'Location', 'northeastoutside');
62     grid on;
63
64     % ---- Plot 3: Log-Log |rho| vs. n ----
65     subplot(2,2,3);
66     loglog(n_vals, abs(rho_vec), 'o-', 'LineWidth', 2);
67     xlabel('Matrix size n');
68     ylabel('|Convergence Rate|');
69     title('Log-Log: |rho| vs. n');
70     grid on;
71
72     % ---- Plot 4: Fit |rho| approx a n^b ----
73     subplot(2,2,4);
74     loglog(n_vals, abs(rho_vec), 'o-', 'LineWidth', 2);
75     xlabel('Matrix size n');
76     ylabel('|Convergence Rate|');
77     title('Log-Log: Fitted |rho| vs. n');
78     grid on;
79     hold on;
80
81     log_rho = log(abs(rho_vec));
82     coeffs_rho = polyfit(log_n, log_rho, 1);
83     slope = coeffs_rho(1);
84     a_rho = exp(coeffs_rho(2));
85     fit_rho = a_rho * n_vals.^slope;
86     loglog(n_vals, fit_rho, '--', 'LineWidth', 1.5);
87     text(n_vals(round(end*0.3)), fit_rho(round(end*0.3)), ...
88         'FontSize', 10, 'BackgroundColor', 'white');
```

```
88     sgtitle('Conjugate Gradient Convergence Analysis (
89         Tridiagonal A)');
90
91     % User defined function
92
93     function A = generate_tridiagonal_matrix(n)
94         % Generates an n x n symmetric tridiagonal matrix with:
95         % 2 on the diagonal and -1 on the sub- and super-
96         % diagonals
97
98         e = ones(n,1);
99         A = 2*diag(e) - diag(e(1:end-1),1) - diag(e(1:end-1),-1)
100         ;
101     end
102
103     function [num_iter, total_flops, rel_errors] =
104         conj_grad_descent(m, tolerance)
105         A = generate_tridiagonal_matrix(m);
106
107         initial_guess = zeros(m, 1);
108         x_star = (1:m)';
109         b = A * x_star;
110
111         residual_old = b - A * initial_guess;
112         direction = residual_old;
113         rs_old = residual_old' * residual_old;
114
115         iter = 0;
116         rel_errors = [];
117         norm_x_star = norm(x_star);
118
119         while sqrt(rs_old) > tolerance
120             Ap = A * direction;
121             alpha = rs_old / (direction' * Ap);
122
123             initial_guess = initial_guess + alpha * direction;
124             residual_old = residual_old - alpha * Ap;
125             rel_error_k = norm(initial_guess - x_star) /
126                 norm_x_star;
127             rel_errors = [rel_errors, rel_error_k];
128
129             rs_new = residual_old' * residual_old;
130             beta = rs_new / rs_old;
131
132             direction = residual_old + beta * direction;
133             rs_old = rs_new;
```

```

130         iter = iter + 1;
131     end
132
133     num_iter = iter;
134     total_flops = num_iter * m^2; % dense matrix assumption
135     ; for tridiagonal use ~5m
136 end

```

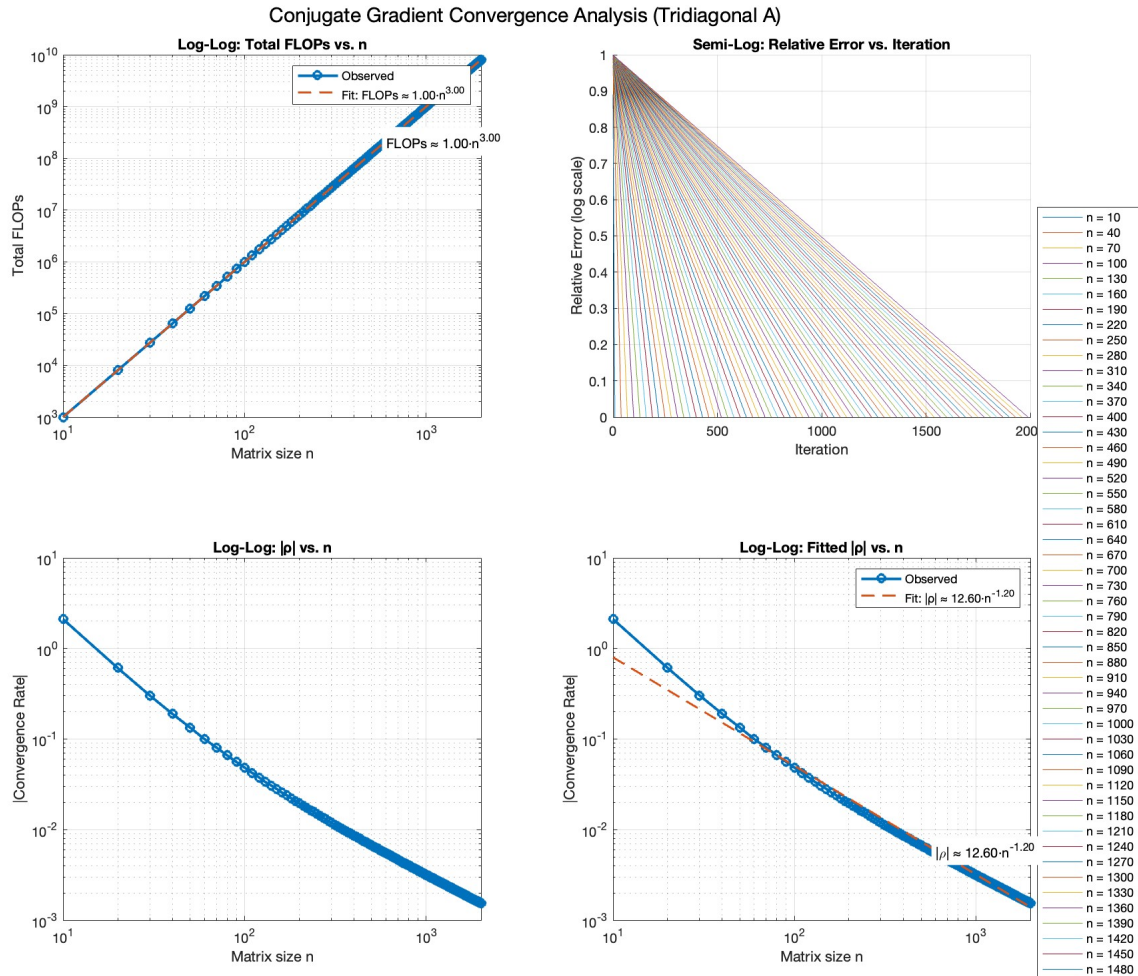


Figure 2: Conjugate Gradient descent method

We can see from the top right graph, the complexity for the whole algorithm is $\mathcal{O}(n^3)$. Since we know that in conjugate gradient, we eliminate the n components from the residual vector, the per step complexity must be $\mathcal{O}(n^2)$, which matches theory as we have a matrix vector product. Following a similar logic above, we know that convergence rate is approximately $\mathcal{O}(n)$. \square

Problem 4. Suppose in a gradient descent scheme, an error ϵ_k occurs:

$$x^{(k+1)} = x^{(k)} - s_k \nabla f(x^{(k)}) + \epsilon_k$$

where $\|\epsilon_k\| \leq \epsilon$, and $f = \frac{1}{2}(x - x_*)^T A(x - x_*)$ with A positive definite. Show the scheme can converge with a suitable choice of s_k . Make assumptions on ϵ if necessary.

Solution: Note that $\nabla f = A(x - x^*)$, we begin with the following.

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &= \|x^k - s_k(A(x^k - x^*)) - x^* + \epsilon_k\|^2 \\ &= \|(I - s_k A)(x^k - x^*) + \epsilon_k\|^2 \\ &\leq (\|(I - s_k A)(x^k - x^*)\| + \|\epsilon_k\|)^2 \\ &\leq (\|(I - s_k A)(x^k - x^*)\|)^2 + 2\|\epsilon_k\| \|(I - s_k A)(x^k - x^*)\| + \epsilon^2 \end{aligned}$$

As proven in class, we know that if $s_k \in (0, \frac{2}{\lambda_k})$, by definition of the spectral norm, we can prove that

$$\|I - s_k A\| = \max_i |\lambda_i|$$

thus, we see that $\|I - s_k A\| < 1$. For notational sake, let $\rho = \|I - s_k A\|$ and let $r_k = \|x^k - x^*\|$, thus, note that

$$(\|(I - s_k A)(x^k - x^*)\|)^2 + 2\|\epsilon_k\| \|(I - s_k A)(x^k - x^*)\| + \epsilon^2 \leq r_k^2 \rho^2 + 2\epsilon \rho r_k + \epsilon^2$$

we can see that we are left with bound:

$$r_{k+1}^2 \leq r_k^2 \rho^2 + 2\epsilon \rho r_k + \epsilon^2$$

We aim to find a bound $R \geq r_i, \forall i \in \{0, 1, 2, \dots, k\}$. Thus, it suffices to find a value such that:

$$\xi^2 = \xi^2 \rho^2 + 2\epsilon \rho \xi + \epsilon^2 \iff (1 - \xi^2) \rho^2 + 2\epsilon \rho \xi + \epsilon^2 = 0$$

Using the quadratic formula and simplifying:

$$\xi = \frac{\epsilon(\rho + 1)}{1 - \rho^2} = \frac{\epsilon}{1 - \rho}$$

Note that this above term is a constant dependent on the maximal error term. Thus, we can see that given a $s \in (0, \frac{2}{\lambda_{max}})$, the algorithm converges for a sufficiently small epsilon. \square