



1. Project 1 - PageRank Algorithm and Social Networks

The initial project queries, along with additional information from the presentation slides, comprise this chapter. Notably, this section has less extensive documentation compared to others, primarily due to the limited details provided in the slides. For a more comprehensive grasp of this topic, it's advisable to consult both the slides and the relevant textbook. Contributions for enriching the content of this chapter are welcome, and interested individuals are encouraged to submit a pull request with their additions.

1. Understanding PageRank

The PageRank algorithm, a seminal creation of Google's co-founders Larry Page and Sergey Brin, fundamentally revolutionized web search technologies. It operates on a principle that evaluates the importance of web pages based solely on the network of links within the World Wide Web. When a specific query is made, Google employs the PageRank algorithm to sift through the vast web, identifying pages relevant to the query and ordering them based on their PageRank.

The core idea of PageRank can be visualized as a random surfer model. In this model, navigation across the web is simulated by randomly following outbound links from one web page to another. This process, however, encounters challenges such as dead ends (web pages lacking outbound links) and loops within tightly knit clusters of web pages. To mitigate these issues, PageRank introduces a randomization factor: at certain intervals, the surfer is directed to a randomly chosen web page, breaking free from loops and dead ends.

This model is mathematically conceptualized as a **Markov Chain** or **Markov Process**, where the probability of transitioning from one page to another is determined by the structure of the web's link network. The essence of PageRank lies in its recursive nature: a page is deemed significant and thus assigned a high rank if it is linked by other pages that are themselves of high rank. This recursive evaluation of page importance forms the backbone of the PageRank algorithm.

A matrix is called *col-stochastic* if it satisfies two conditions:

1. **Non-negativity:** every element in the matrix is non-negative. This means all entries a_{ij} of the matrix are greater than or equal to zero ($a_{ij} \geq 0$)
2. **Column Sum Equals One:** The sum of the elements in each column is equal to one.

2. Numerical Methods for Eigenvalue Problems

2.1. Power Method

The Power Method is an iterative algorithm used to compute the dominant eigenvalue and corresponding eigenvector of a matrix. It is especially effective for large matrices where determining all eigenvalues is not feasible.

Power Method is particularly used for estimating numerically largest or smallest Eigen Value and its corresponding Eigenvector.

The Power Iteration: An approximation method for computing the eigenvectors of a matrix. It uses the following recursive formula to find the principal eigenvector r of a matrix M :

$$r^{n+1} = \frac{M \cdot r^n}{||M \cdot r^n||}$$

In case of PR, M is col-stochastic, and r is a probability vector, therefore $||M \cdot r|| = 1$, and we reduce the power iteration to:

$$r^{n+1} = M \cdot r^n$$

Example:

$$A = \begin{bmatrix} 1 & -3 & 2 \\ 4 & 4 & -1 \\ 6 & 3 & 5 \end{bmatrix} x_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$
$$\text{Iteration 1: } Ax_0 = \begin{bmatrix} 1 & -3 & 2 \\ 4 & 4 & -1 \\ 6 & 3 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 7 \\ 14 \end{bmatrix}$$

Here the 14 is the largest eigenvalue and its corresponding eigenvector. So on, continue with the iteration using the eigenvector where we take out 14:

$$\text{Iteration 2: } Ax_1 = A \cdot \begin{bmatrix} 0 \\ 0.5 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1 \\ 6.5 \end{bmatrix} = 6.5 \begin{bmatrix} 0.07 \\ 0.15 \\ 1 \end{bmatrix}$$

In this iteration the 6.5 is the largest eigenvalue. And continue iteratively for n iterations. It will happen that same eigenvalues will repeat recursively and then happen when it converge.

Power Iteration method:

1. **Input:** Graph M with n nodes and m links.
2. **Initialize:** $r_0 = [1/n, \dots, 1/n]^T$
3. **While** $|r_k - r_{k-1}|_1 \geq \epsilon$ do:
4. $\tilde{r} = M \cdot r_{k-1}$
5. $r_k = \tilde{r} / ||\tilde{r}||$
6. **end while**

Flow equations in matrix form $r = M \cdot r$:

- x is an eigenvector with corresponding eigenvalue λ if: $Ax = \lambda x$
- r (the rank vector) is an eigenvector of the col. stochastic web adjacency matrix M
- Starting from any stochastic u , the rank r is the limit $M \cdot M \cdots Mu$
- The associated Eigenvalue $\lambda = 1$
- We can efficiently solve for r .

Convergence: The method converges to the eigenvector associated with the largest eigenvalue, assuming that the largest eigenvalue's magnitude is strictly greater than the second largest.

2.2. Inverse Iteration Method (Shift and Invert)

Inverse Iteration, often used with a shift (known as Shift and Invert), is a technique to find an eigenvalue and its corresponding eigenvector close to a given target μ .

- The power Iteration has cost $O(n^2)$ (matrix-vector product) per iteration.
- Convergence can be very slow if for eigenvalues $\lambda_1 > \lambda_2 > \dots \lambda_n$ of the matrix $A \in \mathbb{R}^{n \times n}$ we have $|\frac{\lambda_2}{\lambda_1}| \approx 1$
- Shift and invert technique improves the convergence rate, but solves a linear systems at each iteration with $O(n^3)$ cost.

Practical example:

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \quad b_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \mu = -0.5$$

$$b_{k+1} = \frac{(A - \mu I)^{-1} b_k}{\|(A - \mu I)^{-1} b_k\|}$$

$$(A - \mu I)^{-1} = \begin{bmatrix} 0.5 & 1 \\ 1 & 1.5 \end{bmatrix}^{-1} = \begin{bmatrix} -6 & 4 \\ 4 & -2 \end{bmatrix}$$

$$b_k = \begin{bmatrix} -6 & 4 \\ 4 & -2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 2 \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\vdots \text{ Several Iterations Later}$$

$$b_8 = \begin{bmatrix} -1.651180 \\ 1 \end{bmatrix} \quad \lambda = \frac{b_8^T A b_8}{b_8^T b_8} = 0.6180$$

Rayleigh Quotient Iteration

First do the *Inverse Power Method* with *Shift* + *Rayleigh Quotient* at each step:

$$\text{Inverse Power Method} \quad b_{k+1} = \frac{(A - \mu I)^{-1} b_k}{\|(A - \mu I)^{-1} b_k\|}$$

$$\text{Rayleigh Quotient} \quad \mu_k = \frac{b_k^T A b_k}{b_k^T b_k} \quad \text{Each Step}$$

Compute new $(A - \mu_k I)^{-1}$ each iteration.

Helps to find the convergence quickly but more expensive.

Algorithm Inverse Iteration:

1. **Input:** Graph A , Initial guess v_0 , shift α
2. **Initialize** $r_0 = [1/n, \dots, 1/n]^T$
3. **for** $k=1, 2, \dots$, until termination criteria *do*:
4. Solve $(A - \alpha I)\tilde{v} = v_{k-1}$
5. $v_k = \tilde{v}/\|v\|$
6. $\lambda^k = v_k^T A v_k$
7. **end for**

Shift and inverse technique, main idea:

- λ_i : eigs of A , then $\lambda_i - \alpha$: of $A\alpha I$, and $\mu_i = \frac{1}{\lambda_i - \alpha}$: of $B = (A - \alpha I)^{-1}$
- If $\alpha = \lambda_1 \Rightarrow$ Rapid Convergence
- Popular choices for $\alpha = |A|_1$ or $\alpha = \lambda^{k-1}$
- If α is fixed, inversion of B takes place outside the for loop.
- If v_0 normalized, $\lambda^0 = v_0^T A v_0$ and $\alpha = \lambda^{k-1}$, the Rayleigh quotient is approximated.
- $\lambda(v) = \frac{v^T A v}{v^T v}$
- Cubic convergence justifies increased cost/iteration.

Usage: This method is particularly useful when the largest eigenvalue is not of primary interest, or when the target is to find eigenvalues close to a specific value μ .

3. Comparison of Methods

- **Power Method:** Best suited for finding the largest eigenvalue. Simple to implement but may converge slowly if the gap between the largest and second-largest eigenvalues is narrow.
- **Inverse Iteration:** Ideal for locating eigenvalues near a specific target μ . Requires more computation due to matrix inversion or solving a system of linear equations at each step.

4. Theoretical Questions

4.1. Eigenvectors, Eigenvalues, and Eigenbasis

- **Eigenvector:** An eigenvector of a matrix A is a non-zero vector \mathbf{v} which, when multiplied by A , results in a scaled version of \mathbf{v} , i.e., $A\mathbf{x} = \lambda\mathbf{x}$ for some scalar λ .
- **Eigenvalue:** The scalar λ in the equation $A\mathbf{x} = \lambda\mathbf{x}$ is the eigenvalue corresponding to eigenvector \mathbf{x} .
- **Eigenbasis:** An eigenbasis of a square matrix A is a set of linearly independent eigenvectors of A that span the entire vector space.

4.2. Assumptions for Convergence of the Power Method

- The matrix should preferably be **symmetric**; asymmetric matrices may lead to inaccurate results.
- A **dominant eigenvalue** is necessary; if multiple eigenvalues have the same magnitude, the method may not converge properly.
- **Linearly independent eigenvectors** are required.
- The **initial guess** significantly impacts the convergence rate.

4.3. Shift and Invert Approach

The Shift and Invert approach is used when the Power Method's asymptotic error constant $\left(\frac{\lambda_1}{\lambda_2}\right)$ is small, and λ_2 is very close to λ_1 . This method involves shifting the matrix A by a scalar α and then inverting it. The choice of α is crucial, and the approach is computationally intensive, making it suitable for smaller problems.

4.4. Cost Difference: Power Method vs. Inverse Iteration

- **Power Method:** Involves matrix-vector multiplication with complexity $O(n^2)$.
- **Inverse Iteration:** Involves solving a linear system in each iteration with complexity $O(n^3)$.

4.5. Rayleigh Quotient

The Rayleigh Quotient for a vector \mathbf{V} is defined as $\mu(\mathbf{V}) = \frac{\mathbf{V}^T \mathbf{A} \mathbf{V}}{\mathbf{V}^T \mathbf{V}}$. It provides the associated eigenvalue for an eigenvector \mathbf{V} , or an approximation of an eigenvalue in a least-squares sense. The quotient $\mu(\mathbf{V}_k)$ gives an estimate for the eigenvalue λ_1 in the k -th iteration.

4.6. Definition of a Clique

In the context of graph theory, a **clique** is a subset of vertices within a graph such that every two distinct vertices are connected by an edge. In other words, a clique is a complete subgraph of a graph.

Formally, given a graph $G = (V, E)$, where V is the set of vertices and E is the set of edges, a clique in G is a subset $C \subseteq V$ such that for every pair of vertices $v, w \in C$, there is an edge $(v, w) \in E$.

Characteristics:

- In a clique, every vertex is adjacent to all other vertices in the clique.
- The number of edges in a clique of n vertices is $\frac{n(n-1)}{2}$, representing a fully connected subgraph.
- The size of the largest clique in a graph is known as the *clique number* of the graph.

Example: A simple example of a clique is a triangle in a network, where three nodes are all connected to each other. In a social network, a clique might represent a group of friends where every person is directly connected to every other person in the group.

5. Reverse Cuthill-McKee Ordering

The *Reverse Cuthill-McKee Ordering* is an algorithm used in numerical linear algebra to reorder the rows and columns of a sparse matrix with the goal of reducing its bandwidth.

The *Reverse Cuthill-McKee (RCM)* algorithm is a graph reordering technique primarily used to reduce the bandwidth of a sparse matrix.

The **Bandwidth** of a matrix is a measure of how far from the main diagonal the non-zero elements are located. Reducing the bandwidth is important in numerical analysis because it can make storage and computations more efficient, especially for sparse matrices. For example, if the furthest non-zero element from the diagonal is three rows and columns away, the bandwidth is three.

How RCM works:

RCM is a two-step process: First applying the *Cuthill McKee (CM)* and then reversing the resulting order.

- *Step 1: Cuthill-McKee Ordering*
- **Degree Of Nodes:** The algorithm starts by computing the degree of each node (number of non-zero connections) in the graph.
- **Starting Node:** Choose a node with the smallest degree as the starting node. If there are multiple, any can be chosen. This choice can significantly affect the final result.

- **Bredth-First Search(BFS):** Perform a BFS starting from the chosen node. in BFS, you visit a node, then visit all its neighbours before moving to the next level of neighbours.
- **Neighbor Ordering:** When visiting the neighbors of a node, order them by increasing degree. This step is crucial for reducing the bandwidth.
- **Continue Until all Nodes are visited**
- *Step 2: Reverse the Order:*
Once the CM ordering is complete, simply reverse this order. This reversed order is the final RCM order.

A practical example:

Consider a sparse matrix. Converting it to a graph, applying RCM, and then converting it back to a matrix form will often result in a Matrix where non-zero elements are clustered closer to the diagonal, reducing the bandwidth.

5.1. Bandwidth of a Matrix

The bandwidth of a matrix refers to the maximum distance of non-zero elements from the main diagonal. In matrices with low bandwidth, non-zero elements are concentrated close to the diagonal. Reducing the bandwidth of a matrix is often desirable for computational efficiency.

5.2. Cuthill-McKee Algorithm

Developed by Edward Cuthill and James McKee, this algorithm reorders a sparse matrix to minimize its bandwidth. The process involves:

- Starting from a vertex (row or column) with the lowest degree (fewest non-zero elements).
- Performing a breadth-first search to order the vertices.
- The resulting ordering of vertices is used to permute the rows and columns of the matrix.

5.3. Reverse Ordering

The *Reverse Cuthill-McKee Ordering* is an extension where the order determined by the Cuthill-McKee algorithm is then reversed. This reversal often further reduces the bandwidth, improving the efficiency of matrix computations.

6. Applications

- **Matrix Computations:** Particularly beneficial for the efficient computation of large sparse matrices in linear systems and eigenvalue problems.
- **Graph Theory:** Used in optimizing the representation of graphs as adjacency matrices.
- **Scientific Computing:** Crucial in areas like finite element analysis, where sparse matrices are common.

6.1. MATLAB Implementation

MATLAB implements this ordering through the `symrcm` function, which takes a sparse matrix and returns a permutation vector to rearrange the matrix, minimizing bandwidth.

7. Explanation of MATLAB Pagerank1 Function - POWER METHOD

7.1. Function Definition

```
1 function x = pagerank1(U, G, p)
```

- **Purpose:** Defines the function 'pagerank1' to compute the PageRank of web pages. - **Inputs:** - 'U': A cell array of URLs. - 'G': An adjacency matrix representing web page links. - 'p': The damping factor, typically set to 0.85.

7.2. Set Default Damping Factor

```
1 if nargin < 3, p = .85; end
```

- **Purpose:** Sets 'p' to 0.85 if not provided. - **Theory:** The damping factor models the probability of following links.

7.3. Calculate Out-Degrees and In-Degrees

```
1 [~, n] = size(G);  
2 c = sum(G, 1);  
3 r = sum(G, 2);
```

- **Purpose:** Determines the out-degree ('c') and in-degree ('r') for each page. - **Theory:** - 'c': Out-degree is the number of outbound links. - 'r': In-degree is the number of inbound links.

7.4. Create Diagonal Matrix D

```
1 k = find(c ~= 0);  
2 D = sparse(k, k, 1./c(k), n, n);
```

- **Purpose:** Creates a diagonal matrix 'D' with reciprocal out-degrees. - **Theory:** Normalizes 'G' for the transition probability matrix.

7.5. Initialize Power Method

```
1 e = ones(n, 1);  
2 G = p * G * D;  
3 z = ((1 - p) * (c ~= 0) + (c == 0)) / n;  
4 x = e / n;
```

- **Purpose:** Prepares the Power Method iteration. - **Theory:** Modifies 'G' for transition probabilities and handles dangling nodes.

7.6. Power Method Iteration

```
1 tolerance = 0.0001;  
2 x_old = 0;  
3 while norm(x - x_old, 1) > tolerance  
4     x_old = x;  
5     x = G * x + e * (z * x);  
6 end
```

- **Purpose:** Executes the Power Method for PageRank. - **Theory:** Iteratively refines the PageRank vector.

7.7. Normalize PageRank Vector

```
1 x = x / sum(x);
```

- **Purpose:** Ensures the sum of PageRank values is 1. - **Theory:** PageRank is a probability distribution.

7.8. Display Results

```
1 shg; bar(x); title('Page Rank');
```

- **Purpose:** Shows a bar graph of the PageRank values.

7.9. Print Dominant URLs

```
1 if nargin < 1
2     % ... code to print URLs and ranks ...
3 end
```

- **Purpose:** Prints PageRanks and URLs if no output argument is specified.

8. Explanation of MATLAB Pagerank2 Function - INVERSE

8.1. Function Definition

```
1 function x = pagerank1(U, G, p)
```

- **Purpose:** Defines the function 'pagerank1' to compute the PageRank of web pages. - **Inputs:** - 'U': A cell array of URLs. - 'G': An adjacency matrix representing web page links. - 'p': The damping factor, typically set to 0.85.

8.2. Set Default Damping Factor

```
1 if nargin < 3, p = .85; end
```

- **Purpose:** Sets 'p' to 0.85 if not provided. - **Theory:** The damping factor models the probability of following links.

8.3. Calculate Out-Degrees and In-Degrees

```
1 [~, n] = size(G);
2 c = sum(G, 1);
3 r = sum(G, 2);
```

- **Purpose:** Determines the out-degree ('c') and in-degree ('r') for each page. - **Theory:** - 'c': Out-degree is the number of outbound links. - 'r': In-degree is the number of inbound links.

8.4. Create Diagonal Matrix D

```
1 k = find(c ~= 0);
2 D = sparse(k, k, 1./c(k), n, n);
```

- **Purpose:** Creates a diagonal matrix 'D' with reciprocal out-degrees. - **Theory:** Normalizes 'G' for the transition probability matrix.

8.5. Initialize Inverse Method

```
1 z = ((1-p) * (c~=0) + (c==0)) /n;  
2 x = e/n;  
3 A = p * G * D + e * z;  
4 alpha = 0.99;  
5 tol = 0.0001;  
6  
7 prev_x = 0;
```

- Purpose: - Theory:

8.6. Inverse Method Iteration

```
1 while norm(x - prev_x, 1) >= tol  
2     prev_x = x;  
3     x = (A - alpha * I) \ x;  
4     x = x/norm(x);  
5 end
```

- Purpose: - Theory:

8.7. Normalize PageRank Vector

```
1 x = x / sum(x);
```

- Purpose: Ensures the sum of PageRank values is 1. - Theory: PageRank is a probability distribution.

8.8. Display Results

```
1 shg; bar(x); title('Page Rank');
```

- Purpose: Shows a bar graph of the PageRank values.

8.9. Print Dominant URLs

```
1 if nargin < 1  
2     % ... code to print URLs and ranks ...  
3 end
```

- Purpose: Prints PageRanks and URLs if no output argument is specified.

2. Project 2 - Graph Partitioning

The second project questions plus some content from the slides.

1. How to solve the partitioning problem?

- Partitioning **with** nodal coordinates each node has (x,y,z) coordinates (partition space):
 - Coordinate Bisection
 - Recursive Coordinate Bisection
 - Interlary Partitioning

- Partitioning **without** Nodal coordinates
 - Spectral Methods

2. Coordinate Bisection

Coordinate Bisection is a graph partitioning technique that involves dividing a graph into subgraphs based on hyperplanes orthogonal to the coordinate axes. This method is particularly useful when the vertices of the graph are embedded in a coordinate space, such as Euclidean space.

Hyperplanes and Orthogonality: In Coordinate Bisection, a hyperplane refers to a flat, $n-1$ dimensional subset of an n -dimensional space. For a graph embedded in 2D space, this hyperplane is a line, and in 3D space, it is a plane. The key feature of this hyperplane is its orthogonality to the coordinate axes, meaning it is perpendicular to one of the axes.

Partitioning Process: The process of partitioning using Coordinate Bisection involves the following steps:

1. Determine the axis along which the graph will be bisected. This choice can be based on various criteria, such as the distribution of nodes or specific application requirements.
2. Identify the hyperplane orthogonal to the chosen axis. In 2D, this could be a vertical or horizontal line, depending on whether the bisection is along the x-axis or y-axis.
3. Divide the graph into two subgraphs based on the position of the vertices relative to the hyperplane. Nodes on one side of the hyperplane form one subgraph, and nodes on the other side form the second subgraph.

Advantages and Use Cases: Coordinate Bisection is advantageous in scenarios where the graph naturally resides in a geometric space, and the vertex coordinates are meaningful or related to the structure of the graph. Examples include spatial networks, sensor networks, and certain types of social networks.

Limitations: One limitation of this method is that it might not always produce optimally balanced partitions, especially in graphs where node positions are not uniformly distributed across the coordinate space. Additionally, the method's effectiveness can be limited if the graph's structure does not correlate well with the geometric layout of the nodes.

Algorithmic Complexity: The computational complexity of Coordinate Bisection is generally low, particularly because it relies on simple geometric calculations. This makes it an efficient method for partitioning large graphs, provided the graph structure is amenable to this approach. In summary, Coordinate Bisection is a geometrically intuitive method for graph partitioning, leveraging the orthogonal hyperplanes to the coordinate axes. Its effectiveness is highly dependent on the spatial distribution of the graph's vertices and the correlation between vertex coordinates and graph structure.

```

1  function [part1,part2] = bisection_coordinate(A,xy,picture)
2  % bisection_coordinate: Coordinate bisection partition of a mesh.
3  %
4  % [part1,part2] = bisection_coordinate(A,xy,picture) returns a list of the vertices on
   one side of a partition
5  %   obtained by bisection perpendicular to a coordinate axis. We try every
6  %   coordinate axis and return the best cut.
7  %   Input A is the adjacency matrix of the mesh;
8  %   each row of xy is the coordinates of a point in d-space.
9  %
10 % bisection_coordinate(A,xy,1) also draws a picture.
11
```

```

12 d = size(xy,2);
13 best_cut = inf;
14 for dim = 1:d
15     v = zeros(d,1);
16     v(dim) = 1;
17     [p1,p2] = partition(xy,v);
18     this_cut = cutsize(A,p1);
19     if this_cut < best_cut
20         best_cut = this_cut;
21         part1 = p1;
22         part2 = p2;
23     end
24 end
25
26 if picture == 1
27     clf reset
28     gplotpart(A,xy,part1);
29     title('Coordinate bisection')
30 end
31
32
33 end

```

3. Recursive Coordinate Bisection (RBC)

Recursive Coordinate Bisection (RBC) is an advanced technique for graph partitioning that extends the basic concept of Coordinate Bisection. It involves iteratively dividing a graph into smaller subgraphs using hyperplanes orthogonal to the coordinate axes, with each subsequent partition treated as a separate graph for further bisection.

Initial Bisection: RBC begins with a standard Coordinate Bisection:

- The graph is divided into two equal parts using a cutting plane orthogonal to one of the coordinate axes. This axis can be chosen based on various factors, such as node distribution or specific characteristics of the graph.
- The choice of the orthogonal cutting plane is crucial, as it influences the balance and characteristics of the resulting subgraphs.

Recursive Partitioning: After the initial partition, RBC applies the Coordinate Bisection recursively:

1. Each of the two subgraphs obtained from the initial partition is considered as an independent graph.
2. Coordinate Bisection is then reapplied to each of these subgraphs, again dividing them into two smaller subgraphs along a chosen coordinate axis.
3. This recursive process continues until a predefined condition is met, such as achieving subgraphs of a certain size or depth of recursion.

Advantages of RBC:

- **Balanced Partitions:** RBC is effective in creating balanced partitions, which is crucial for many applications, such as parallel computing.
- **Flexibility:** It allows for flexibility in choosing the cutting plane at each recursive step, potentially leading to better graph partitions.

Challenges and Considerations:

- **Choice of Axis:** The selection of the axis for the cutting plane at each recursive step can significantly affect the quality of partitioning.
- **Termination Criteria:** It is important to define appropriate termination criteria for the recursion to avoid excessively small or unbalanced subgraphs.

Applications: RBC is particularly useful in scenarios where a balanced partitioning of a large graph is necessary, such as in load balancing for parallel processing, data distribution in distributed systems, and network analysis.

In summary, Recursive Coordinate Bisection is a method that iteratively applies Coordinate Bisection to subgraphs, dividing the original graph into progressively smaller and more balanced partitions. Its recursive nature allows for a more refined partitioning of the graph compared to a single application of Coordinate Bisection.

```

1  function [map,sepj,sepA] = rec_bisection(method,levels,A,varargin)
2  % rec_bisection Separate a graph recursively.
3  %
4  % [map,sepj,sepA] = rec_bisection(method,levels,A,arg...) partitions the
5  % mesh or graph A recursively by a specified method. 'method' is the name
6  % of the 2-way edge separator function to call. levels is the number of
7  % levels of partitioning. A is the adjacency matrix of the graph. arg2,
8  % arg3, arg4 are optional additional arguments to the 2-way function. arg2
9  % is special: If it has the same number of rows as A, then the recursive
10 % calls use the appropriate subset of arg2 as well as of A. This is useful
11 % if the partitioner uses xy coords. map is a vector of integers from 0 to
12 % 2^levels-1, indexed by vertex, giving the partition number for each
13 % vertex. sepj is a list of separating edges; each row is an edge [i j].
14 % sepA is the adjacency matrix of the graph less the separating edges.
15
16 minpoints = 8; % Don't separate pieces smaller than this.
17
18 nargcall = nargin - 2;
19
20 n = size(A,1);
21
22 if n < minpoints || levels < 1
23     map = zeros(1,n);
24 else
25     % Call the partitioner to split the graph, giving one part.
26     [p1,p2] = feval(method, A, varargin{:});
27
28     % Call recursively.
29     vararginp1 = varargin;
30     vararginp2 = varargin;
31     if nargcall >= 2
32         if size(varargin{1},1) == n
33             vararginp1{1} = varargin{1}(p1,:);
34             vararginp2{1} = varargin{1}(p2,:);
35         end
36     end
37
38     mapa = rec_bisection(method,levels-1,A(p1,p1),vararginp1{:});
39     mapb = rec_bisection(method,levels-1,A(p2,p2),vararginp2{:});
40
41     % Set up the whole map.
42     map = zeros(1,n);
43     mapb = mapb + max(mapa) + 1;
44     map(p1) = mapa;
45     map(p2) = mapb;
46 end

```

```

47
48 % Set up the separating edge list and separated graph.
49 if nargout >= 2
50     [i,j] = find(A);
51     f = find(map(i) > map(j));
52     sepij = [i(f) j(f)];
53     f = find(map(i) == map(j));
54     sepA = sparse(i(f), j(f), 1, n, n);
55 end
56
57 end

```

4. Inertial Partitioning

Inertial Partitioning, a distinctive graph partitioning approach, relies on the geometric layout of a graph, particularly the coordinates assigned to its vertices. The main goal is to identify a hyperplane that efficiently bisects the graph.

In 2D Space: The method involves finding a line l in a 2-dimensional space that minimizes the sum of squared distances from the nodes to the line. This line is defined by a point $P = (X, Y)$ and a direction vector $u = (u_1, u_2)$, where $\|u\|_2 = \sqrt{u_1^2 + u_2^2}$. The line is represented as $l = \{P + \alpha u | \alpha \in \mathbb{R}\}$, and each node v_i has coordinates (x_i, y_i) .

Center of Mass: The center of mass of the nodes is calculated as:

$$x = \frac{1}{n} \sum_{i=1}^n x_i, \quad y = \frac{1}{n} \sum_{i=1}^n y_i, \quad (1)$$

where n is the number of nodes.

Minimizing Sum of Distances: The objective is to minimize:

$$\begin{aligned}
 \sum_{i=1}^n d_i^2 &= \sum_{i=1}^n (x_i - x)^2 + (y_i - y)^2 - (u_1(x_i - x) + u_2(y_i - y))^2 \\
 &= \left(\sum_{i=1}^n (x_i - x)^2 \right) u_1^2 + \left(\sum_{i=1}^n (y_i - y)^2 \right) u_2^2 + 2 \left(\sum_{i=1}^n (x_i - x)(y_i - y) \right) u_1 u_2 \\
 &= S_{xx} + S_{xy} + S_{xy}
 \end{aligned} \quad (2)$$

Matrix M and Eigenvector: The matrix M is formed from the sums S_{xx}, S_{xy}, S_{yy} and is given by:

$$u^T \begin{bmatrix} S_{xx} & S_{xy} \\ S_{xy} & S_{yy} \end{bmatrix} u = u^t M u. \quad (3)$$

The optimal direction vector u is the normalized eigenvector associated with the smallest eigenvalue of M . **Procedure:** The inertial partitioning algorithm can be summarized in the following pseudocode:

```

1 require: Graph G(V,E),  $P_i = (x_i, y_i)$ ,  $i = 1, \dots, n$  the coordinates of the nodes.
2
3 Ensure: A bisection of G
4 function INERTIALBISECTION(graph G(V,E), P_i)
5     Calculate the center of mass of the points.
6     Compute the eigenvector associated with the smallest eigenvalue of M.
7     Partition the nodes along the line H which goes through the center of mass is
        orthogonal to L.
8     return  $V_1, V_2$ 
9 end function

```

MATLAB Implementation:

```

1 centerOfMass = mean(xy);
2
3 Sxx = sum((xy(:,1) - centerOfMass(1)).^2);
4 Syy = sum((xy(:,2) - centerOfMass(2)).^2);
5 Sxy = sum((xy(:,1) - centerOfMass(1)) .* (xy(:,2) - centerOfMass(2)));
6
7 M = [Sxx Sxy; Sxy Syy];
8
9 [V, ~] = eigs(M, 1, 'smallestabs');
10
11 smallestVec = V(:, 1);
12
13 u = smallestVec;
14
15 Un = [-u(2), u(1)];
16
17 [part1, part2] = partition(xy, Un);

```

Generalization of Coordinate Partitioning: The approach can be extended to select a line L that is not necessarily orthogonal to the axes. The line $L : ax + by + c = 0$, where $a^2 + b^2 = 1$, is chosen such that it minimizes the squared distances of nodes from L . The partitioning is then performed using the median of the projections of nodes on L .

Inertial Algorithm: The algorithm involves assembling matrix M , finding the eigenvector u_1 , and computing line L . Nodes are then partitioned along line H , orthogonal to L and passing through the center of mass.

5. Spectral Partitioning

Spectral Partitioning is a method in graph theory, based on the work of Fiedler in the 1970s and popularized by Horst Simon in 1995. It leverages the eigenvectors and eigenvalues (the *spectrum*) of a graph's Laplacian matrix for partitioning.

Spectrum of a Graph: The spectrum of a graph refers to its eigenvectors v_i and their corresponding eigenvalues λ_i , typically ordered by the magnitude of the eigenvalues.

$$\Lambda = (\lambda_1, \lambda_2, \dots, \lambda_n) \quad \text{where} \quad \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \quad (4)$$

Theory Behind Spectral Partitioning: The spectral method employs linear algebra, particularly the spectral graph theorem, which enables the decomposition of a real symmetric matrix. Consider an undirected graph $G(V, E)$ with vertices $V = \{v_1, \dots, v_n\}$ divided into two subsets V_1, V_2 such that $V_1 \cup V_2 = V$ and $V_1 \cap V_2 = \emptyset$. The indicator vector

$$x \in R^n$$

is defined as:

$$x_i = \begin{cases} 1, & \text{if } i \in V_1 \\ 0, & \text{if } i \in V_2 \end{cases} \quad (5)$$

Graph Laplacian Matrix: The graph Laplacian matrix L is used, encoding the degree $d_i = \sum_{j=1}^n w_{ij}$ (number of connected edges) of each vertex on its diagonal and the negative weights $-w_{ij}$. The adjacency matrix W represents the connections, where $w_{ij} > 0$ if v_i is connected to v_j , and $w_{ij} = 0$ otherwise. The graph Laplacian is symmetric and positive semidefinite, with its smallest eigenvalue being $\lambda_1 = 0$.

Fiedler's Eigenvector for Partitioning: The eigenvector associated with the second smallest eigenvalue λ_2 , known as Fiedler's eigenvector, is crucial for spectral graph partitioning. Nodes v_i of the graph are associated with entries in this eigenvector w_2 . Thresholding around the median value of w_2 results in two balanced partitions.

Procedure: The spectral partitioning algorithm can be summarized in the following pseudocode:

```

1 require: Graph G(V,E)
2 Ensure: A bisection of G
3 function SpectralBisection(graph G(V,E))
4     Form the graph Laplacian matrix L.
5     Calculate the second smallest eigenvalue  $\lambda_2$  and its associated eigenvector
         $w_2$ .
6     Set 0 or the median of all components of  $w_2$  as the threshold.
7     Choose  $V_1 := \{v_i \in V \mid w_i < \text{thres}\}$ ,  $V_2 := \{v_i \in V \mid w_i \geq \text{thres}\}$ .
8     return V1, V2
9 end function

```

In this method, the graph is bisected by partitioning the vertices based on the sign or relative magnitude of their corresponding entries in the Fiedler vector w_2 .

In MatLab:

```

1 degreeMatrix = diag(sum(A,2));
2 laplacian = degreeMatrix - A;
3 [V, ~] = eigs(laplacian, 2, 'smallestabs');
4 fiedlerVector = V(:, 2);
5 threshold = 0;
6 part1 = find(fiedlerVector < threshold);
7 part2 = find(fiedlerVector >= threshold);

```

6. Questions from from the project 2

The comparison with recursive bisection to k-way partitioning result that is not one better from the other one. The results are influenced by factors such as the graph structure and number of partitions required.

3. Spectral Graph Clustering

Clustering algorithms:

- Flat Algorithms, that has no explicit structure that would relate clusters to each other and usually start with a random partitioning (k-means or model based clustering) and it is a Spectral Clustering.
- Hierarchical algorithms, bottom-up and agglomerative or top-down and divisive.

1. Hierarchical Clustering

Hierarchical clustering is more informative than the unstructured set of flat clusters and easier to decide on the number of clusters by looking at the dendrogram.

But it is impossible to undo the previous step, once the instances have been assigned to a cluster they can no longer be moved around. Its complexity is $O(n^2)$ not suitable for very large datasets. The order of the data has impact on the final results.

1.1. Agglomerative Clustering

Agglomerate to collect or gather into a cluster or mass. Treat each data entry as a singleton cluster and then successively merge pairs of clusters till all clusters have been merged into a single one. Usually visualized with a dendrogram. Does not require a prespecified number of clusters.

1.2. Divisive Clustering

Start at the top with all data in one cluster. The cluster is split using a flat algorithm. Apply this procedure recursively. It is more complex cause implement Hierarchical+ flat but it is more efficient and accurate.

2. Graph Theory

Information regarding the inner structure of a data set in terms of:

- Cliques, subsets of connected nodes, each pair of elements is connected.
- Clusters, highly connected groups of nodes
- centrality, important nodes, hubs
- Outlier, unimportant Nodes

2.1. Graph Construction & Graph Notation

For an undirected graph $G(V, E)$, with vertex set $V = v_1, \dots, v_n$ and a subset $A \subset V$:

Indicator Vector: $x \in \mathbb{R}^n$

$$x_i = \begin{cases} 1, & i \in A, \\ 0, & i \in \bar{A} \end{cases} \quad (6)$$

Shorthand notation: $i \in A \iff v_i \in A$

Weights:

$$w_{ij} = \begin{cases} > 0, & \text{if } v_i \text{ connected to } v_j \\ = 0 & \text{otherwise} \end{cases} \quad (7)$$

Degree:

$$d_i = \sum_{j=1}^n w_{ij}$$

This sum only runs over all vertices adjacent to v_i .

Degree Matrix: $D \in \mathbb{R}^{n \times n}$

$$D = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix}$$

For two sets $A, B \subset V$, the weighted adjacency matrix is defined as

$$W := \left[\sum_{\substack{i \in A \\ j \in B}} w_{ij} \right] \in \mathbb{R}^{n \times n},$$

where w_{ij} represents the weight of the edge from vertex i to vertex j . The matrix W can be expressed as

$$W = \begin{bmatrix} 0 & w_{12} & \cdots & w_{1n} \\ w_{21} & 0 & \cdots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n1} & w_{n2} & \cdots & 0 \end{bmatrix},$$

where each w_{ij} is the weight of the edge between vertices i and j , and the diagonal elements are typically zero (assuming no self-loops in the graph).

Measuring the size of a subset:

- The cardinality of set A , denoted as $|A|$, measures the size of A by its number of vertices.
- The volume of set A , denoted as $\text{vol}(A)$, is defined as $\text{vol}(A) := \sum_{i \in A} d_i$. It measures the size of A by summing the degrees (or weights of all edges attached) of vertices in A .

The unnormalized graph Laplacian matrix L is defined as $L = D - W$, where D is the degree matrix and W is the adjacency matrix. It can be expressed as

$$L = \begin{bmatrix} d_1 & -w_{12} & \cdots & -w_{1n} \\ -w_{21} & d_2 & \cdots & -w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ -w_{n1} & -w_{n2} & \cdots & d_n \end{bmatrix},$$

where d_i represents the degree of vertex i , and w_{ij} represents the weight of the edge between vertices i and j . In this matrix, the diagonal entries d_i are the degrees of the vertices, and the off-diagonal entries $-w_{ij}$ are the negated weights of the edges.

Properties of L : The graph Laplacian matrix L exhibits several important properties which are fundamental in graph theory and spectral graph analysis. These properties include:

- **Quadratic Form:** For every vector $x \in \mathbb{R}^n$, the quadratic form of L is given by:

$$x^T L x = \frac{1}{2} \sum_{i,j=1}^n w_{ij} (x_i - x_j)^2,$$

where w_{ij} is the weight of the edge between vertices i and j , and x_i, x_j are the components of the vector x .

- **Symmetry and Positive Semi-Definiteness:** The Laplacian matrix L is symmetric and positive semi-definite. This implies that all eigenvalues of L are non-negative and real.
- **Smallest Eigenvalue:** The smallest eigenvalue of L is 0, and the corresponding eigenvector is the constant vector $\mathbf{1}$ (a vector of all ones). This property reflects the fact that the sum of each row in L is zero.
- **Eigenvalue Spectrum:** The eigenvalues of L are non-negative and real-valued, ordered as $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$. These eigenvalues provide significant insights into the structure and connectivity of the graph.

These properties make the Laplacian matrix L a powerful tool in various applications, including spectral clustering, network analysis, and the study of graph dynamics.

Graph Construction - Connectivity Matrix G

The goal of constructing a connectivity matrix G is to model the local neighborhood relationships between data points. There are several approaches to building such a graph:

- **Epsilon-Neighborhood Graph:** In this approach, a threshold value ε is identified. An edge is included between two nodes if the distance between them is smaller than or equal to ε . This method emphasizes local connectivity based on proximity.
- **K-Nearest Neighbors (K-NN) Graph:** This method involves inserting edges between a node and its k -nearest neighbors. Here, each node is connected to at least k other nodes, ensuring that every node has a minimum degree of k . The K-NN graph is particularly useful for capturing local structures in the data.

- **Fully Connected Graph:** In a fully connected graph, there is an edge between every pair of nodes, resulting in a complete graph. This approach models the most dense possible connectivity, assuming every data point is connected to every other point.

Each of these methods has its own applications and implications on the structure and properties of the graph. The choice of method depends on the specific requirements of the problem and the nature of the data.

Graph Construction - Similarity Matrix S

The construction of a graph involves defining a similarity function on the dataset to ensure that the local neighborhoods, as delineated by this function, are meaningful. For instance, in the context of text documents, it's important to verify that documents with a high similarity score indeed belong to the same category or have similar content.

Similarity in Euclidean Space: A common scenario involves data points situated in Euclidean space. Here, the similarity between points can be quantified using various functions, such as the Gaussian similarity function.

Combining Similarity and Connectivity: The full similarity matrix S is constructed using the chosen similarity function. Concurrently, an epsilon-connectivity graph represented by matrix G is established, which captures the local neighborhood structure based on a threshold ε .

Element-wise Multiplication: The operation $S \cdot G = W$ signifies an element-wise multiplication of matrices S and G , resulting in a sparse matrix W . This matrix W retains elements only in positions where both S and G have significant entries, effectively filtering the similarity matrix with the graph's connectivity structure.

Resulting Sparse Matrix: The resulting sparse matrix W represents a weighted adjacency matrix of the graph, where the weights are influenced both by the similarity of data points and the epsilon-connectivity criterion. This approach ensures that the graph structure is both locally relevant and reflective of the underlying data's similarity.

2.2. Spanning Trees and Minimum Spanning Trees

The concepts of Spanning Trees and Minimum Spanning Trees are fundamental in graph theory, particularly in the context of undirected graphs. These structures have important applications in various fields, including network design and optimization.

- **Tree:** A tree is an undirected graph in which any two vertices are connected by exactly one path. It is a connected graph with no cycles, meaning there is a unique path between any pair of nodes.
- **Spanning Tree:** A spanning tree of a graph G is a subgraph that includes all the vertices of G , forming a tree structure. The spanning tree contains the minimum number of edges required to connect all the vertices, ensuring there are no cycles and that the graph remains connected.
- **Minimum Spanning Tree (MST):** The minimum spanning tree takes the concept of a spanning tree further by minimizing the total edge weight. In an MST, the sum of the weights of the edges is as small as possible, making it a spanning tree with the least total cost. This is particularly important in scenarios where costs or distances are associated with the edges of the graph, such as in network design or optimization problems.

3. Drawing graphs using eigenvectors

A graph $G(V, E)$ is a structure that models a relation E over a set of V entities. Relational information are visualized by graph drawing.

```

1 load meshes.mat
2 % Matrix calculations
3 L = Airfoil;
```

```

4  xy = Airfoilxy;
5  D = diag(diag(L));
6  W = D - L;
7
8  % Eigen-decomposition
9  [V, lambda] = eigs (L, 3, 'SA');
10
11 % Plot and compare
12 figure;
13 subplot(2,1,1);
14 gplot(w,xy);
15 subplot(2,1,2);
16 % locate vertex i at pos: xi = (v2(i),v3(i))
17 gplot(W, V(:, [2,3]));

```

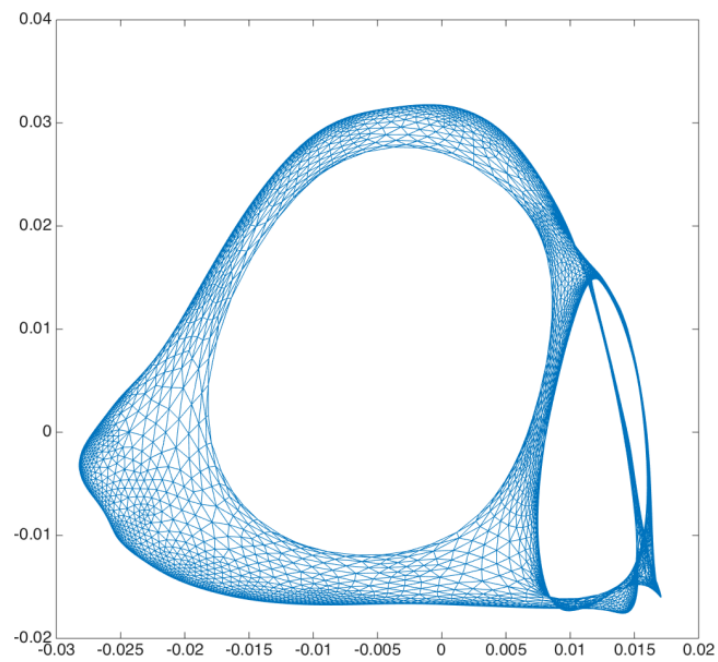


Figure 1: First plot with (w,xy).

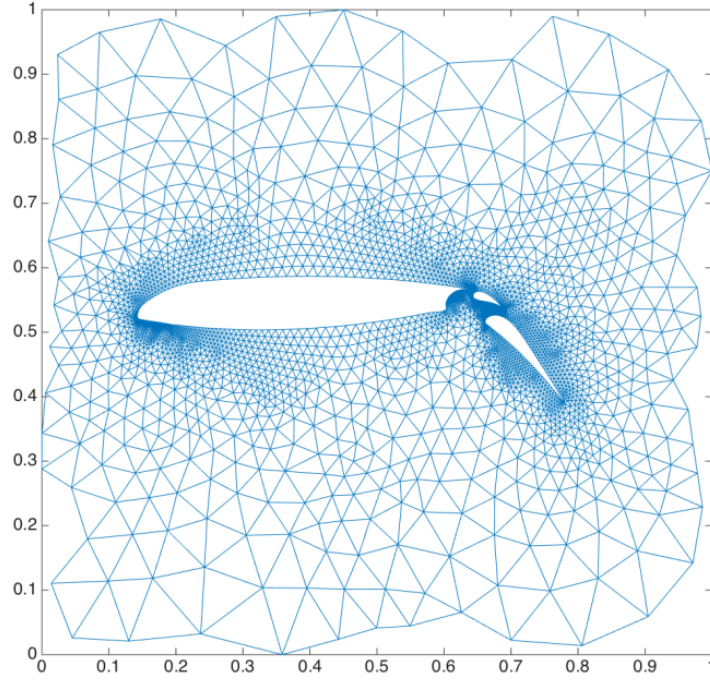


Figure 2: Second Plot with $(W, V(:, [2,3]))$

4. k-way Spectral Clustering

Recursive Bi-partition:

- *Method:* The method involves recursively applying a bi-partitioning algorithm in a hierarchical, divisive manner to the data set. This process continues until the desired number of clusters (k) is achieved.
- *Advantages:* The main advantage of this approach is its simplicity in terms of the algorithmic concept, making it easy to implement and understand.
- *Disadvantages:* However, recursive bi-partitioning can be inefficient and unstable, particularly in terms of maintaining the quality of the clusters as the number of partitions increases.

Clustering Multiple Eigenvectors:

- *Method:* This approach involves constructing a reduced space from the first k eigenvectors of the graph Laplacian. The data points are then clustered in this reduced eigenvector space.
- *Advantages:* Clustering multiple eigenvectors makes full use of the information provided by these eigenvectors. It generally has less computational complexity compared to other methods, and the clustering results are often quite satisfactory.
- *Disadvantages:* A key challenge is determining the appropriate number of eigenvectors to use, as not all eigenvectors contribute equally to the quality of the clustering. Often, only the first few eigenvectors provide satisfactory results, making the selection process critical and not straightforward.

In summary, k-way spectral clustering offers different approaches, each with its own set of advantages and challenges. The choice between recursive bi-partitioning and clustering multiple eigenvectors depends on the specific requirements and constraints of the clustering task at hand.

5. Unnormalized Spectral Clustering

Input: Given a similarity matrix $S \in \mathbb{R}^{n \times n}$ and a specified number of clusters k to construct.

Procedure:

1. *Construct a Similarity Graph:* Form a similarity graph based on the matrix S . Let W be its weighted adjacency matrix.
2. *Compute the Unnormalized Laplacian:* Calculate the unnormalized Laplacian matrix L using $L = D - W$, where D is the degree matrix of the graph.
3. *Eigenvector Computation:* Compute the first k eigenvectors u_1, \dots, u_k of the Laplacian matrix L .
4. *Form Matrix U :* Construct matrix $U \in \mathbb{R}^{n \times k}$ which contains the vectors u_1, \dots, u_k as columns. This step represents a dimensionality reduction from $n \times n$ to $n \times k$.
5. *Row Vector Formation:* For each $i = 1, \dots, n$, define $y_i \in \mathbb{R}^k$ as the vector corresponding to the i -th row of U .
6. *Clustering:* Apply the k-means algorithm to cluster the points $(y_i)_{i=1, \dots, n}$ in \mathbb{R}^k into k clusters C_1, \dots, C_k .

Outcome: This process results in the partitioning of the data into k clusters based on the spectral properties of the Laplacian matrix. The unnormalized spectral clustering approach is particularly effective in identifying underlying structures in data that are not immediately apparent in the original space.

6. K-means Clustering

K-means is one of the most well-known and widely used clustering algorithms. It clusters n objects based on attributes into k partitions, where $k < n$, effectively grouping the objects based on attributes/features into k distinct groups.

Algorithm Process:

1. *Initialization:* Start with some initial cluster centers.
2. *Iteration:*
 - Assign each example to the closest center.
 - Recalculate the centers as the mean of the points in a cluster.

Drawbacks of K-means:

1. *Euclidean Distance:* The use of Euclidean distance treats the data space as isotropic, meaning distances remain unchanged by translation or rotation. Data points in each cluster are modeled as lying within a sphere around the cluster centroid. However, this spherical assumption may not always hold, leading to non-intuitive behaviors in clearly identifiable clusters.

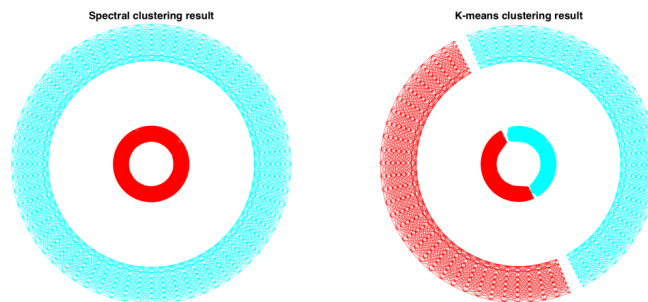


Figure 3: K-means clustering with non-intuitive behavior.

2. *Impact of Outliers:* In a linear Euclidean space, small changes in data result in proportionally small changes to the position of cluster centroids. This can be problematic when outliers are present, as they can disproportionately influence the centroid location.

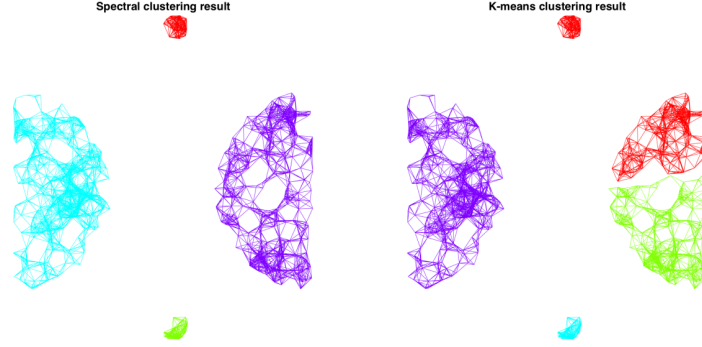


Figure 4: Behavior of K-means in the presence of outliers.

3. *Equal Volume Assumption:* K-means clusters data points based on their geometric closeness to the cluster centroid and does not account for the varying densities of each cluster. This implicit assumption that each cluster occupies the same volume in data space often leads to clusters containing an equal number of data points, which may not always be ideal.

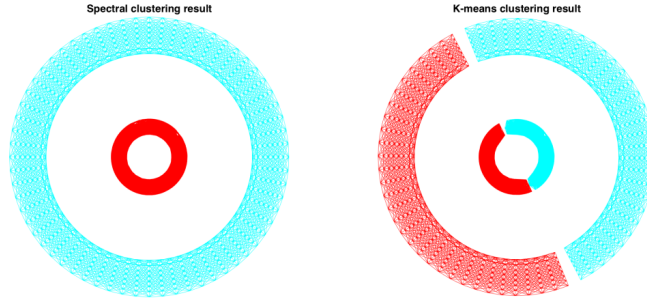


Figure 5: K-means clustering with the same quantity of points per cluster.

Despite these drawbacks, K-means remains a popular choice for clustering due to its simplicity and efficiency, particularly in scenarios where its underlying assumptions align well with the data distribution.

7. Balanced Graph Partitioning Analogy

The key insight in balanced graph partitioning is to separate points into different groups based on their similarities. The analogy involves partitioning a graph such that the edges between different groups have low weight, and edges within a group have high weight.

For a given number k of subsets, our goal is to find a partition A_1, \dots, A_k that minimizes:

$$\text{cut}(A_1, \dots, A_k) = \frac{1}{2} \sum_{i=1}^k W(A_i, \bar{A}_i) \quad (\text{Not balanced})$$

$$\text{RatioCut}(A_1, \dots, A_k) = \sum_{i=1}^k \frac{\text{cut}(A_1, \dots, A_k)}{|A_i|} \quad (\text{Balanced})$$

The minimum of $\sum_{i=1}^k (\frac{1}{|A_i|})$ coincides.

Example of RatioCut: Let $k = 2$. Our goal is to solve the minimization problem:

$$\min \text{RatioCut}(A, \bar{A})$$

Given a subset $A \subset V$, we define the vector $x = (x_1, \dots, x_n)^T \in \mathbb{R}^n$ with

$$x_i = \begin{cases} \sqrt{\frac{|\bar{A}|}{|A|}}, & \text{if } v_i \in A, \\ -\sqrt{\frac{|A|}{|\bar{A}|}}, & \text{if } v_i \in \bar{A}. \end{cases}$$

The RatioCut objective function can be rewritten using the graph Laplacian:

$$x^T L x = x^T D x - x^T W x = \sum_{i=1}^n d_i x_i^2 - \sum_{i,j=1}^n x_i x_j w_{ij} = \dots = \text{Rcut}(A, \bar{A}) \cdot |V|$$

Additionally, the vector x is orthogonal to the constant vector $\mathbf{1}$:

$$x^T \mathbf{1} = \sum_{i=1}^n x_i = |A| \sqrt{\frac{|\bar{A}|}{|A|}} - |\bar{A}| \sqrt{\frac{|A|}{|\bar{A}|}} = 0$$

And x measures the cardinality of the entire set:

$$\|x\|^2 = x^T x = \sum_{i=1}^n x_i^2 = |A| + |\bar{A}| = n = |V|$$

The minimization problem

$$\min \text{RatioCut}(A, \bar{A})$$

can equivalently be stated as

$$\min x^T L x \quad \text{subject to } x \perp \mathbf{1}, \|x\|^2 = n$$

The Rayleigh-Ritz theorem indicates that the solution of this problem is given by the Fiedler eigenvector (the eigenvector corresponding to the second smallest eigenvalue of L). However, to obtain a partition of the graph, we need to convert the real-valued solution vector x of the relaxed problem into a discrete indicator vector. In spectral partitioning, we use the sign of x as an indicator function:

$$v_i \in \begin{cases} A, & \text{if } x_i \geq 0, \\ \bar{A}, & \text{if } x_i < 0. \end{cases}$$

Spectral clustering algorithms instead regard the coordinates x_i as points in \mathbb{R} and cluster them into two groups C and \bar{C} using a flat clustering algorithm (k-means, in this case). The clustering results are then associated with the underlying data points:

$$v_i \in \begin{cases} A, & \text{if } x_i \in C, \\ \bar{A}, & \text{if } x_i \in \bar{C}. \end{cases}$$

This procedure is referred to as unnormalized spectral clustering for $k = 2$ clusters.

7.1. Computational Complexity

The overall computational complexity of the spectral clustering algorithm is determined by the complexity of its individual components, which include the computation of eigenvalues and eigenvectors, construction of the similarity matrix, and the K-means clustering step.

- **Eigenvalue and Eigenvector Computation:** The computational complexity of calculating eigenvalues and eigenvectors of a matrix is $O(n^3)$, where n is the number of nodes in the graph. This step is often the most computationally intensive part of the process.
- **Similarity Matrix Construction:** Constructing the similarity matrix has a computational complexity of $O(n^2)$. This involves calculating the pairwise similarities between all nodes, which scales quadratically with the number of nodes.

- **K-means Clustering:** The K-means step in spectral clustering has a complexity of $O(nldk)$, where l is the number of iterations in the K-means algorithm, d is the dimensionality of the input data, and k is the number of clusters to be formed. The complexity depends on both the size of the dataset and the dimensions of the feature space.

Understanding these complexities is crucial in evaluating the feasibility and efficiency of spectral clustering, particularly for large-scale datasets. Optimizations and approximations may be necessary in practice to handle large n , high d , or when a large k is desired.

8. More on Clusters

Clustering, particularly in the context of graph theory, can be effectively formulated as a graph cut problem. This approach is especially beneficial in scenarios involving non-convex clustering, where traditional clustering methods might struggle.

- **Graph Cut Formulation:** In this formulation, the objective is to partition a graph in such a way that the cut between different clusters (subgraphs) is minimized. This involves finding a balance between minimizing the edges between different clusters while maximizing the connectivity within each cluster.
- **Eigenvalue Decomposition of the Laplacian Matrix:** A key step in solving the graph cut problem is the eigendecomposition of the Laplacian matrix associated with the graph. This decomposition provides a representation of the data in a lower-dimensional space, which is more amenable to clustering.
- **Lower-Dimensional Representation:** By representing the data in a reduced dimensionality, it becomes easier to identify clusters, as the intrinsic structure of the data is more apparent in this space. This approach leverages the spectral properties of the Laplacian matrix to reveal the underlying clustering structure.
- **Empirical Success:** Empirically, this method of clustering has been highly successful, particularly in complex datasets where traditional clustering methods may fail to capture the true data structure. It is especially adept at uncovering the intrinsic grouping in data that is not linearly separable.

Overall, clustering as a graph cut problem, solved via eigendecomposition of the Laplacian matrix, offers a powerful and flexible approach for identifying clusters in a wide range of data types and structures.

9. Project 3

The `epsilonSimGraph.mat` for generate the similarity matrix of the epsilon-similarity graph.

```

1 function [G] = epsilonSimGraph(epsilon,Pts)
2 % Construct an epsilon similarity graph
3 % Input
4 % epsilon: size of neighborhood (calculate from Prim's Algorithm)
5 % Pts    : coordinate list of the sample
6 %
7 % Output
8 % A      : the epsilon similarity matrix
9 S = similarityfunc(Pts, epsilon);
10 threshold = 0.5;
11 G = S >= threshold;
12 end

```

CreateLapl.m for create Laplacian Matrix:

```

1 function [L,Diag] = CreateLapl(W)
2 % Create the Laplacian matrix of a graph
3 % Input
4 % W : Adjacency matrix
5 % Output
6 % L : Laplacian
7 % Degree matrix
8 Diag = zeros(size(W,1));
9 for i = 1:size(W,1)
10     Diag(i,i) = sum(W(:,i));
11 end
12 % Construct the Laplacian
13 L = Diag - W;
14 end

```

Create the gaussian similarity function

```

1 function [S] = GausSimilarityfunc(Pts, sigma)
2 % Create the similarity matrix S from the coordinate list of the input points
3 if nargin < 2
4     n = length(Pts(:,1));
5     sigma = log(n);
6 end
7
8 S = squareform(pdist(Pts));
9 S = exp(-S.^2 ./ (2*sigma^2));
10 end

```

Function to perform k-means clustering:

```

1 function [D,x] = kmeans_mod(Y,K,n)
2 % Function to perform k-means clustering
3 % -----
4 % Input
5 % Y: data matrix
6 % K: # of clusters
7 % n: number of points
8 % Output
9 % D: centroids
10 % x: assignment(indicator) vector
11
12
13 cand = unique(round(10000*Y)/10000,'rows');
14 c = datasample(cand,K,'Replace',false); % pick K random input pts as initial centroids
15 D = c;
16 x = zeros(size(Y,1),1);
17 D_old = inf*ones(size(D));
18 count = 1;
19 while norm(D - D_old) > 0.00001 & count < 500
20     D_old = D;
21     % Assign points to clusters
22     for i = 1:n
23         min = inf;
24         for k = 1:K
25             dist = norm(Y(i,:) - D(k,:));
26             if dist < min
27                 min = dist;
28                 x(i) = k;
29             end
30         end

```

```

31     end
32
33     % Update centroid locations
34     for k = 1:K
35         ind = find(x == k);
36         if length(ind) == 1
37             D(k,:) = Y(ind,:);
38         else
39             D(k,:) = mean(Y(ind,:));
40         end
41     end
42
43     count = count + 1;
44 end
45
46 for i = 1:n
47     min = inf;
48     for k = 1:K
49         dist = norm(Y(i,:) - D(k,:));
50         if dist < min
51             min = dist;
52             x(i) = k;
53         end
54     end
55 end
56 end

```

KNN simalrity graph function:

```

1  function [G] = kNNSimGraph(Pts)
2  % Construct a k-nearest neighbors similarity graph
3  % Input
4  % k      : # of neighbors
5  % Pts    : coordinate list of the sample
6  %
7  % Output
8  % G      : the kNN similarity matrix
9
10
11 n = length(Pts(:,1));
12 kn = ceil(2*log(n));
13
14
15 fprintf('kNN similarity graph\n');
16 for i = 1:n
17     s = repmat(Pts(i,:),n,1);
18     d = Pts - s;
19     e = sum(d.^2,2);
20     [~,ind] = sort(e);
21     ind(1) = [];
22     nbrs = ind(1:kn);
23     G(i,nbrs) = 1;
24     G(nbrs,i) = 1;
25 end
26 end

```

4. Image Deblurring and Conjugate Gradient Solver

We explore solving a system of linear equations $Ax = b$, where $A \in \mathbb{R}^{n \times n}$ is a nonsingular coefficient matrix, $b \in \mathbb{R}^n$ is a constant vector, and $x \in \mathbb{R}^n$ is the unknown solution vector.

1. Direct Methods

Direct methods are designed to compute the exact solution in a finite number of steps. Their characteristics include:

- **Computational Cost:** The cost for classical Gaussian elimination is $O(n^3)$. For sparse matrices, this cost reduces to $O(n^{1.5})$ or $O(n^2)$.
- **Memory Consumption:** These methods have a high memory requirement due to additional fill-in elements created during the elimination process.
- **Limitation:** The complexity of Gaussian elimination becomes prohibitively high for large systems of linear equations.

This complexity leads to the consideration of alternative approaches, particularly iterative methods.

2. Iterative Methods

Iterative methods start with an arbitrary initial guess and iteratively compute an approximate solution:

- **Computational Cost:** The cost for the conjugate gradient algorithm ranges between $O(n)$ and $O(n^{3/2})$, making it more scalable for large systems.
- **Memory Efficiency:** These methods do not require additional memory for fill-ins, which is beneficial for large systems.
- **Convergence:** The convergence of iterative methods depends on the properties of matrix A and typically occurs after a few iterations. This makes them suitable for problems where an exact solution is not necessary, or where the exact solution is computationally infeasible.

In summary, while direct methods offer exact solutions, their computational and memory demands make them less viable for large systems. Iterative methods, such as the conjugate gradient algorithm, provide a practical alternative by offering approximate solutions with lower computational and memory requirements.

2.1. Sketch of an Iterative Method

The algorithm for an iterative method can be structured as follows:

1. **Initialization:** Start with $m = 0$ and set x^m as the initial guess for the solution.
2. **Iteration:** Repeat the following steps until the error estimate is less than a predetermined threshold ϵ :
 - a) Compute a new solution x^{m+1} .
 - b) Update the error estimate based on x^{m+1} .

2.2. Definition of Error and Residual

In iterative methods, two important concepts are *Error* and *Residual*, which help in measuring the accuracy of the solution at each iteration.

- **Error:** The error at step m quantifies the deviation of the approximate solution $x^{(m)}$ from the exact solution x . It is defined as:

$$e^{(m)} = x - x^{(m)} = A^{-1}b - x^{(m)}$$

Here, $e^{(m)}$ represents the difference between the true solution and the current approximation.

- **Residual:** The residual is a measure of the real error and is relatively easier to compute compared to the error. It is defined as:

$$r^{(m)} = b - Ax^{(m)} = Ae^{(m)}$$

This expression shows that the residual $r^{(m)}$ is equivalent to the error transformed by the matrix A (i.e., $Ae^{(m)}$). The significance of this relationship will be evident in subsequent discussions.

- **Residual Norm:** The norm of the residual, denoted as $\|r^{(m)}\|$, quantifies the magnitude of the residual. Additionally, the relative residual norm is defined as $\frac{\|r^{(m)}\|}{\|b\|}$, which normalizes the residual norm by the norm of the constant vector b in the system $Ax = b$.

Understanding error and residual, and their respective norms, is crucial in assessing the convergence and accuracy of iterative methods in solving systems of linear equations. $x^{(m)}$ is the approximation in the m -th iteration.

3. Minimization Problem in 1D Case

The basic idea in this context is to transform the linear equation $a \cdot x = b$ into an equivalent minimization problem. This approach involves defining a function $f(x)$ whose minimum we seek to find.

- **Minimization Function:** The function to be minimized is given by:

$$f(x) := \frac{1}{2}a \cdot x^2 - b \cdot x \quad (1)$$

where a is a scalar coefficient and b is a constant term.

- **Precondition:** We assume that a is symmetric and positive definite. In the context of a scalar coefficient, this simply means $a > 0$

$$f'(x) = 2 \cdot \frac{1}{2}x - b = ax - b$$

- **Solving the Minimization Problem:** Setting the derivative equal to zero, $f'(x) = 0$, leads us back to the original problem:

$$a \cdot x = b$$

This shows that solving the minimization problem $f(x)$ is equivalent to solving the linear equation $a \cdot x = b$, under the precondition that a is positive.

3.1. Minimization Problem

- **Basic Idea**

Instead of $Ax = b$ solve an equivalent minimization problem

$$f(x) := \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$$

- **Precondition**, let A be a symmetric positive definite matrix:

$$\langle Ax, x \rangle > 0, \text{ if } x \neq \text{zerovector}$$

if we take the derivative of (3) it will lead to:

$$f'(x) = \frac{1}{2} A^T x + \frac{1}{2} Ax - b = Ax - b \quad \text{condition : } A^t = A$$

- Linear equations with solution at the intersection of both lines.
- Quadratic Form $f(x)$ and gradients $f'(x)$

3.2. Method of Steepest Descent

The method of steepest descent is an iterative approach used for finding the minimum of a function. It involves taking steps proportional to the negative of the gradient of the function at the current point.

- **Basic Idea:** At each iteration, take a step in the direction of the negative gradient at the current point $x^{(m)}$. This direction is given by $-f'(x^{(m)})$.
- **Step Length:** Choose a step length such that the function f is minimized along this search direction.

Algorithm:

1. Initialize the residual: $r^{(0)} = b - Ax^{(0)}$.
2. Repeat until the norm of the residual $\|r^{(0)}\|$ is less than a threshold ϵ :
 - a) Find a scalar α that minimizes $f(x^{(m)} + \alpha r^{(m)})$.
 - b) Update the solution: $x^{(m+1)} = x^{(m)} + \alpha r^{(m)}$.
 - c) Update the residual: $r^{(m+1)} = b - Ax^{(m+1)}$.

In this method, the residual not only measures the error but also serves as the search direction for the next iteration.

The goal is to minimize the function f along a chosen search vector $p^{(m)}$.

Approach: To find the minimum of f along $p^{(m)}$, we consider the function at the point $x^{(m+1)} = x^{(m)} + \alpha p^{(m)}$. The condition for f to be minimized at this point is given by setting its derivative with respect to α to zero:

$$\frac{d}{d\alpha} f(x^{(m+1)}) = 0$$

Using the Chain Rule: By applying the chain rule, we can express this derivative as:

$$\frac{d}{d\alpha} f(x^{(m+1)}) = \left\langle f'(x^{(m+1)}), \frac{d}{d\alpha} (x^{(m)} + \alpha p^{(m)}) \right\rangle$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product. Simplifying further, we get:

$$\frac{d}{d\alpha} f(x^{(m+1)}) = \langle r^{(m+1)}, p^{(m)} \rangle$$

Steepest Descent Specifics: For the method of steepest descent, we specifically choose $p^{(m)} = r^{(m)}$, where $r^{(m)}$ is the residual at the m -th iteration. This choice leads to a descent direction that is steepest in terms of the current position in the function's domain. We search for an α , such that $r^{(m+1)} \perp p^{(m)}$. Rearranging for α :

$$\langle r^{(m+1)}, p^{(m)} \rangle = 0 \dots \alpha = \frac{\langle r^{(m)}, p^{(m)} \rangle}{\langle Ap^{(m)}, p^{(m)} \rangle}$$

The algorithm for the method of steepest descent is summarized as follows:

1. **Initialization:** Start with the initial residual computed as $r^{(0)} = b - Ax^{(0)}$.
2. **Iterative Process:** Repeat the following steps until the norm of the residual $\|r^{(m)}\|$ is less than a predetermined threshold ϵ :

$$\begin{aligned}\alpha^{(m)} &= \frac{\langle r^{(m)}, r^{(m)} \rangle}{\langle Ar^{(m)}, r^{(m)} \rangle} \\ x^{(m+1)} &= x^{(m)} + \alpha^{(m)} r^{(m)} \\ r^{(m+1)} &= r^{(m)} - \alpha^{(m)} Ar^{(m)}\end{aligned}$$

3. **Efficiency Note:** The computation of $r^{(m+1)}$ does not require the matrix-vector product $Ax^{(m)}$. Instead, it utilizes the expression $r^{(m+1)} = r^{(m)} - \alpha^{(m)} Ar^{(m)}$, enhancing computational efficiency.

This algorithm effectively implements the steepest descent method by iteratively updating the solution vector $x^{(m)}$ and the residual $r^{(m)}$ using the calculated step size $\alpha^{(m)}$.

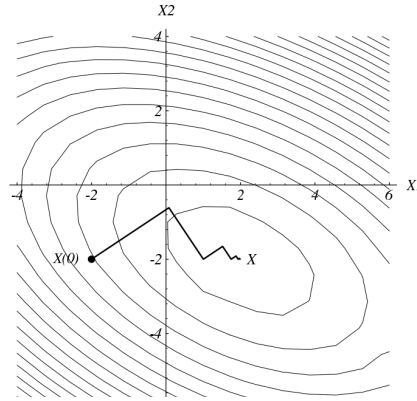


Figure 6: Method of steepest descent with $x^{(0)} = (-2 - 2)^T$

Consider the function $f(x) = \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$, which we aim to minimize to compute the solution of $Ax = b$.

Example for $n = 2$: Suppose matrix A has two orthogonal eigenvectors v_1 and v_2 with associated eigenvalues λ_1 and λ_2 . Any vector x can be expressed as a linear combination of these eigenvectors:

$$x = a_1 \cdot v_1 + a_2 \cdot v_2, \quad a_1, a_2 \in \mathbb{R}$$

For the case where $b = 0$, the function $f(x)$ can be expanded as follows:

$$\begin{aligned} f(x) &= \frac{1}{2}x^T A x \\ &= \frac{1}{2}(a_1 v_1 + a_2 v_2)^T (a_1 \lambda_1 v_1 + a_2 \lambda_2 v_2) \\ &= \frac{1}{2}(a_1^2 \lambda_1 + a_2^2 \lambda_2) \end{aligned}$$

This expression represents the function $f(x)$ in terms of the eigenvalues and eigenvectors of A .

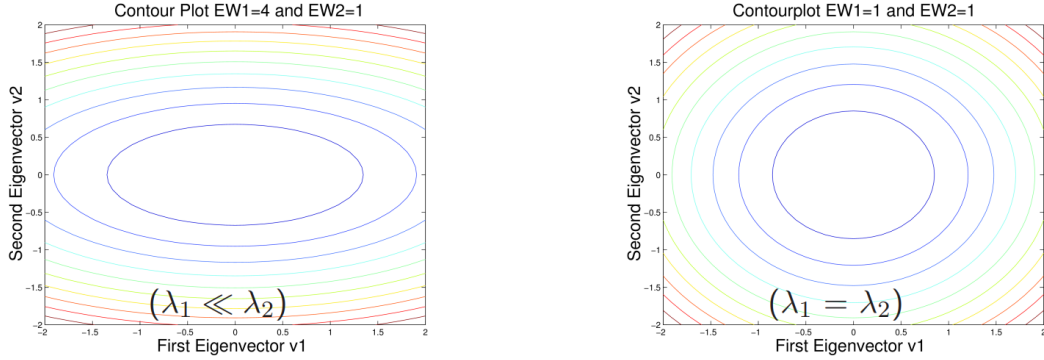


Figure 7: Contour lines of the functional.

3.3. Impact on Convergence

The condition number of a matrix significantly impacts the convergence behavior of iterative methods like steepest descent. The influence can be summarized as follows:

- **Large Condition Number with Good Starting Point:** If the condition number of the matrix is large and the initial guess is by chance close to the true solution, the steepest descent method may converge quickly. However, this fast convergence is more accidental and not guaranteed.
- **Large Condition Number with Bad Starting Point:** Conversely, if the condition number is large and the starting point is not close to the solution, the convergence of the steepest descent method can be very slow. This scenario underscores the sensitivity of the method to initial guesses.
- **Small Condition Number:** A matrix with a small condition number generally leads to good convergence properties of the steepest descent method, regardless of the starting vector. This consistent performance is one of the reasons why preconditioning (which aims to reduce the condition number) is a valuable approach in iterative methods.

4. Conjugate Gradient Method

Basic Idea: The Conjugate Gradient Method improves upon the steepest descent method by considering optimal search steps in each iteration.

Approach:

- Define n orthogonal search directions, $p^{(0)}, p^{(1)}, \dots, p^{(n-1)}$, allowing for a maximum of n steps.
- The minimization criterion is based on the orthogonality of the error $e^{(m+1)}$ with the search direction $p^{(m)}$, such that $p^{(m)} \perp e^{(m+1)}$.

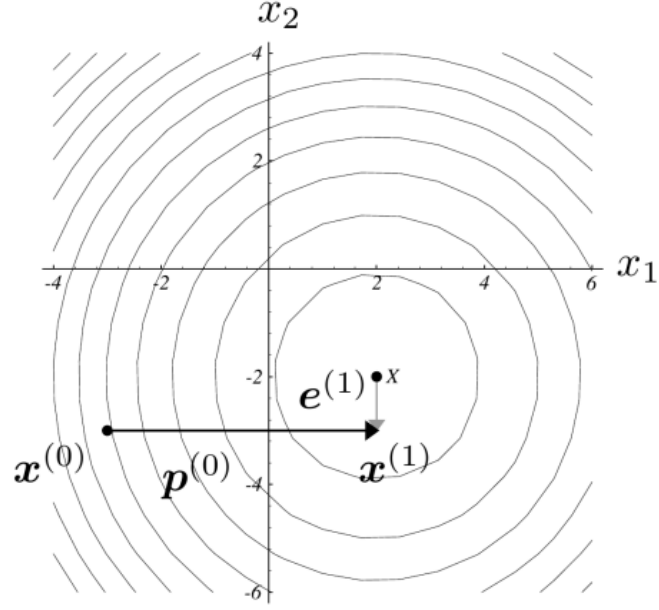


Figure 8: Illustration of the orthogonality in the Conjugate Gradient Method.

Error Elimination: Each step eliminates the error component in the current search direction. If $p^{(m)} \perp e^{(m+1)}$, then the scalar product must be zero:

$$\langle p^{(m)}, e^{(m+1)} \rangle = 0$$

$$\langle p^{(m)}, e^{(m)} + \alpha^{(m)} p^{(m)} \rangle = 0$$

$$\alpha^{(m)} = -\frac{\langle p^{(m)}, e^{(m)} \rangle}{\langle p^{(m)}, p^{(m)} \rangle}$$

However, $e^{(m)}$ is typically unknown.

A-Orthogonality: To overcome this, the method uses A -orthogonal (or conjugate) vectors. Two vectors a and b are A -orthogonal if:

$$a^T A b = \langle a, A b \rangle = 0 \quad \Leftrightarrow \quad a \perp_A b$$

For $p^{(m)} \perp_A e^{(m+1)}$ and $r^{(m)} = A e^{(m)}$:

$$\langle p^{(m)}, A e^{(m+1)} \rangle = 0$$

$$\langle p^{(m)}, A e^{(m)} + \alpha^{(m)} A p^{(m)} \rangle = 0$$

$$\alpha^{(m)} = -\frac{\langle p^{(m)}, A e^{(m)} \rangle}{\langle p^{(m)}, A p^{(m)} \rangle} = \frac{\langle p^{(m)}, r^{(m)} \rangle}{\langle p^{(m)}, A p^{(m)} \rangle}$$

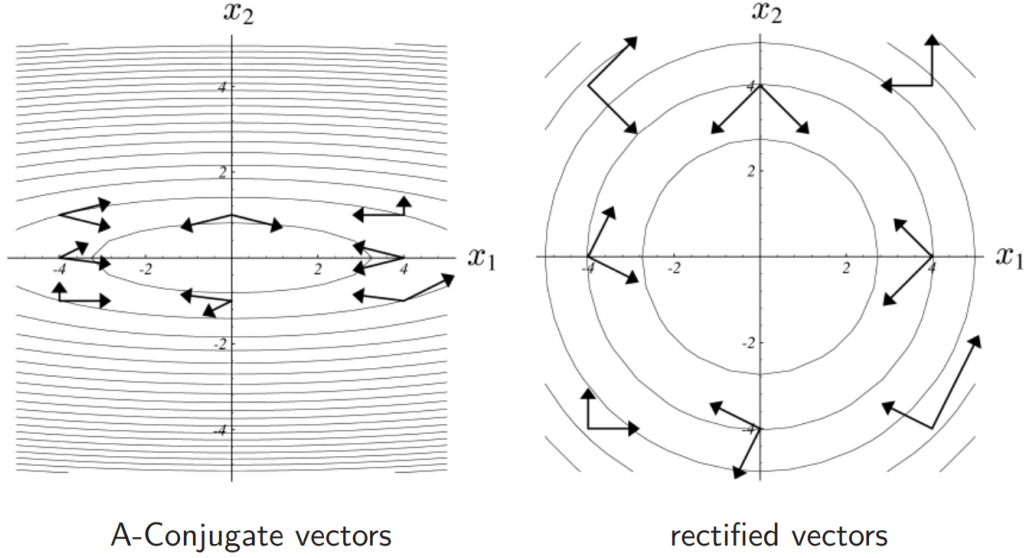


Figure 9: Further demonstration of A -orthogonality in the method.

5. Conjugate Directions

Conditions: For the minimization of $f'(x^{(m)} + \alpha p^{(m)})$, the algorithm iteratively follows these steps:

1. **Initialization:** Start with $r^{(0)} = b - Ax^{(0)}$.
2. **Iterative Process:** For all $m = 1, \dots, n$:
 - Select a search vector $p^{(m)}$ that is conjugate to all previously computed $p^{(l)}$ for $l < m$.
 - Compute $\alpha^{(m)}$ and update $x^{(m+1)}$ and $r^{(m+1)}$ as:

$$\alpha^{(m)} = \frac{\langle r^{(m)}, p^{(m)} \rangle}{\langle Ap^{(m)}, p^{(m)} \rangle}$$

$$x^{(m+1)} = x^{(m)} + \alpha^{(m)} p^{(m)}$$

$$r^{(m+1)} = r^{(m)} - \alpha^{(m)} Ap^{(m)}$$

5.1. Conjugate Gradient Method

Basic Idea: The Conjugate Gradient Method iteratively constructs conjugate search directions using residuals.

Search Direction:

- Construct the next search direction $p^{(m+1)}$ using the residual $r^{(m+1)}$ and a scaling factor $\beta^{(m+1)}$:

$$p^{(m+1)} = r^{(m+1)} + \beta^{(m+1)} p^{(m)}$$

- The factor $\beta^{(m+1)}$ is chosen to ensure A -orthogonality:

$$\beta^{(m+1)} = -\frac{\langle r^{(m+1)}, Ap^{(m)} \rangle}{\langle p^{(m)}, Ap^{(m)} \rangle} = \frac{\langle r^{(m+1)}, r^{(m+1)} \rangle}{\langle r^{(m)}, r^{(m)} \rangle}$$

- The initial search vector is $r^{(0)}$.

Krylov-Subspace: The method constructs a Krylov-subspace U_m defined as:

$$U_m := \text{span}\{p^{(0)}, Ap^{(0)}, A^2p^{(0)}, \dots, A^m p^{(0)}\}$$

Algorithm Steps:

1. **Initialization:** Start with $r^{(0)} := b - Ax^{(0)}$, where $x^{(0)}$ is arbitrary, and set $p^{(0)} := r^{(0)}$.

2. **Iteration:** For $m = 0, 1, \dots, n - 1$:

- Calculate step length $\alpha^{(m)}$ and update solution $x^{(m+1)}$ and residual $r^{(m+1)}$:

$$\alpha^{(m)} = \frac{\langle r^{(m)}, p^{(m)} \rangle}{\langle Ap^{(m)}, p^{(m)} \rangle}$$

$$x^{(m+1)} = x^{(m)} + \alpha^{(m)} p^{(m)}$$

$$r^{(m+1)} = r^{(m)} - \alpha^{(m)} Ap^{(m)}$$

- Compute $\beta^{(m+1)}$ and update search direction $p^{(m+1)}$:

$$\beta^{(m+1)} = \frac{\langle r^{(m+1)}, r^{(m+1)} \rangle}{\langle r^{(m)}, r^{(m)} \rangle}$$

$$p^{(m+1)} = r^{(m+1)} + \beta^{(m+1)} p^{(m)}$$

Elements of the Algorithm:

- $x^{(m)}$: Approximate solution at step m .
- $r^{(m)}$: Residual at step m .
- $p^{(m)}$: Search direction at step m .
- $\alpha^{(m)}$: Step length in search direction $p^{(m)}$ for next iterate $x^{(m+1)}$.
- $\beta^{(m+1)}$: Factor to compute new search direction $p^{(m+1)}$ from $r^{(m+1)}$ and $p^{(m)}$, ensuring A -orthogonality.

The exact solution is theoretically obtained after n iterations, making the Conjugate Gradient method a direct method. However, in practice, fewer iterations are often sufficient for an approximate solution.

6. Preconditioning in Conjugate Gradient Method

Problem: The rate of convergence of the Conjugate Gradient (CG) method is heavily influenced by the condition number of matrix A , denoted as $k(A)$. The convergence relation is given by:

$$\|x^{(m)} - x^*\|_A \leq 2 \left(\frac{\sqrt{k(A)} - 1}{\sqrt{k(A)} + 1} \right)^m \|x_0 - x^*\|_A$$

where the condition number $k(A)$ is defined as:

$$k(A) = \|A\|_2 \|A^{-1}\|_2 = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$$

If $k(A)$ is significantly greater than 1, the convergence rate is slow.

Basic Idea of Preconditioning: To improve the condition number of A by applying a preconditioner M , such that $k(M^{-1}A)$ is much smaller than $k(A)$. The equivalent problem becomes $M^{-1}Ax = M^{-1}b$, where:

- M should be easy to invert.
- The resulting matrix $M^{-1}A$ must remain symmetric positive definite (SPD) for the CG method to be applicable.
- Decompose M as $EE^T = M$ to transform $Ax = b$ into $E^{-1}AE^{-T}\hat{x} = E^{-1}b$, where $\hat{x} = E^Tx$ and $E^{-1}AE^{-T}$ is SPD.

Methods of Preconditioning:

1. **Diagonal Preconditioning:** Choose $M = \text{diag}(A)$.
2. **Incomplete LU-decomposition:** Choose $M = LU \approx A - R$, where it is crucial that the fill-in and the time for factorization are minimized.

Effective preconditioning enhances the performance of the CG method by reducing the condition number, thereby accelerating the convergence rate.

6.1. Jacobi Iteration Method

The Jacobi Iteration Method is an iterative process for solving a system of linear equations by updating each variable independently in a systematic manner.

Iterative Formula: In the Jacobi Iteration Method, each component x_j of the solution vector is updated using the current approximation of the other components:

$$x_j^{(m+1)} = \frac{1}{a_{jj}} \left(b_j - \sum_{k \neq j} a_{jk} x_k^{(m)} \right)$$

This method treats each equation independently, allowing for parallel computation.

Matrix Formulation: The matrix A is decomposed into its diagonal components D and off-diagonal components E :

$$A = D + E$$

The iterative process can then be expressed in matrix form as:

$$Ax = b \Rightarrow Dx = -Ex + b \Rightarrow x = -D^{-1}Ex + D^{-1}b$$

$$x^{(m+1)} = R_j x^{(m)} + c_j$$

where $R_j = -D^{-1}E$ and $c_j = D^{-1}b$.

Spectral Radius and Convergence: The convergence of the Jacobi method is governed by the spectral radius of the matrix R_j :

$$\rho(R_j) = \max\{|\lambda| : \lambda \text{ is an eigenvalue of } R_j\}$$

The Jacobi iteration converges if and only if the spectral radius $\rho(R_j) < 1$.

Geometric Interpretation: Consider a symmetric, real $n \times n$ matrix A with eigenvectors x_k corresponding to eigenvalues λ_k . The multiplication of A with a vector x is a linear transformation involving scaling and rotation:

$$A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}, \quad x = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad y = Ax = \begin{bmatrix} 5 \\ 4 \end{bmatrix}$$

The eigenvectors of A form an orthonormal basis, and any vector can be represented as a linear combination of these eigenvectors.

Eigenvector Representation: A symmetric real matrix has eigenvalues λ_k with associated orthonormal eigenvectors x_k . These eigenvectors form a complete system of orthonormal vectors:

$$\langle x_i, x_j \rangle = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

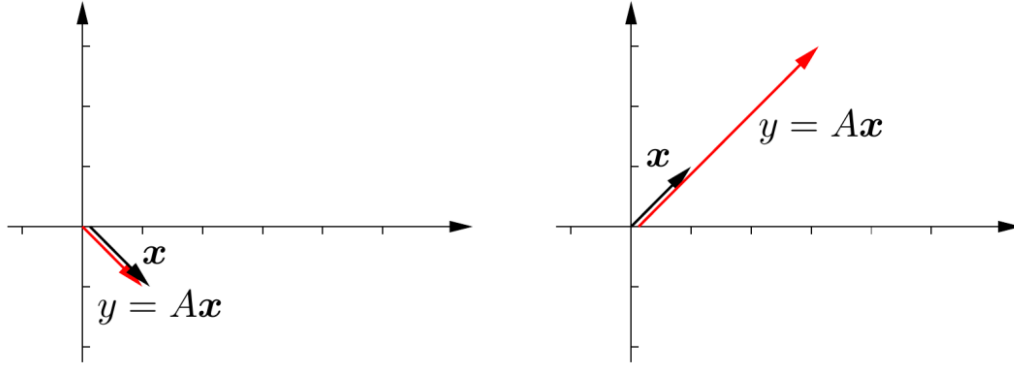


Figure 10: Graphical representation of the Jacobi Iteration Method.

Theorem and Convergence Criteria: The Jacobi iteration converges for all starting vectors if the spectral radius $\rho(R_j) < 1$. Good convergence is typically observed when the matrix A is close to the identity matrix or when the spectral condition number is close to 1.

6.2. Convergence of the Jacobi Iteration Method

The convergence of the Jacobi iteration method is crucial for its effectiveness in solving linear systems. This convergence is determined by the spectral properties of the matrix involved in the iterations.

Convergence Theorem: The Jacobi iteration defined as $x^{(m+1)} = R_j x^{(m)} + c_j$ converges to the solution x^* for any initial vector $x^{(0)}$ if and only if the spectral radius $\rho(R_j) < 1$.

Explanation: Any vector in \mathbb{R}^n can be expressed as a linear combination of the eigenvectors of R_j , provided R_j is not singular. If any eigenvalue of R_j is greater than 1, the Jacobi iteration will not converge.

General Criteria for Good Convergence: - The condition number k should be close to 1.
- The matrix A should be close to the identity matrix I . - These methods will converge in one iteration if $A = \beta I$.

Effect on the Error Term: The iteration affects the error term $e^{(m)}$ as follows:

$$\begin{aligned} x^{(m+1)} &= R_j x^{(m)} + c_j = R_j(x + e^{(m)}) + c_j = R_j x + c_j + R_j e^{(m)} = x + R_j e^{(m)} \\ e^{(m+1)} &= R_j e^{(m)} \quad \text{where} \quad e^{(m+1)} = x - x^{(m+1)} \end{aligned}$$

If $\rho(R_j) < 1$, then the error $e^{(m)}$ converges to zero as $m \rightarrow \infty$.

Comparison with Gauss-Seidel Iteration: While the Jacobi method uses only the components from the previous iteration, the Gauss-Seidel iteration method immediately incorporates newly computed components $x_j^{(m+1)}$:

$$x_j^{(m+1)} = \frac{1}{a_{jj}} \left(b_j - \sum_{k=1}^{j-1} a_{jk} x_k^{(m+1)} - \sum_{k=j+1}^n a_{jk} x_k^{(m)} \right)$$

This approach can lead to faster convergence than the Jacobi method.

7. Project 4

7.1. Simple Image Deblurring - Problem Definition

Suppose we have a blurred image B which we want to refine (or deblur) and we know the transformation that lead to this blurred image. Then how we can efficiently reconstruct the original image?

Let $X \in \mathbb{R}^{n \times n}$ be the original square, grayscale image matrix, where each matrix entry corresponds to one pixel value. Hence X corresponds to an image with n pixels per side. We can convert the original image matrix X into a column vector $x \in \mathbb{R}^{n^2}$, where the entries of X get stacked into a vector, either row-or-column-wise. This process is referred to as *Vectorization*. If we now define as b the vectorized form of the blurred image B , we can write the following system of equations:

$$Ax = b$$

Here, $A \in \mathbb{R}^{n^2 \times n^2}$ indicates the transformation matrix coming from the repeated application of what is referred to as the *Image Kernel* which in our case produces a blurring effect.

The goal is to solve the linear system $Ax = b$ for the original vectorized image x , where the transformation matrix A and the blurred image B (and hence b) are known. By solving the system, we can recover the original image, getting rid of the blurring effect applied on top of it.

Image Kernels Image kernels, also known as convolution kernels, play a significant role in image processing and machine learning. They are small matrices that are used to apply various effects to images. The application of these kernels to an image is known as *convolution*.

Applications of Image Kernels:

- **Blurring:** Softening the image by reducing the intensity variations between adjacent pixels.
- **Sharpening:** Enhancing the edges in the image, making it appear more crisp.
- **Outlining:** Detecting the edges within the image by highlighting the boundaries of objects.
- **Embossing:** Creating a 3D effect by emphasizing edges and texture in the image.

Role in Feature Extraction: In the realm of machine learning, image kernels are utilized for *Feature Extraction*. This process involves identifying and isolating the most important parts of an image, which is essential for tasks such as image recognition and classification. The application of kernels in this context is a critical step in preparing the image data for machine learning algorithms.

Convolution Process: Convolution refers to the process of applying the kernel to the image. This involves sliding the kernel over the image, usually from top-left to bottom-right, and at each position, multiplying the overlapping pixel values with the corresponding kernel values, then summing these products. The result of this sum replaces the pixel value at the current position. This process is repeated for every pixel in the image, thus applying the desired effect comprehensively.

7.2. Simple Image Deblurring - Visualization and Vectorization of Black and White Images

In image processing, particularly with black and white images, two key steps are visualization and vectorization. These processes allow us to manipulate and analyze the image data effectively.

Visualization with *imshow*: The black and white image under consideration is stored as a matrix, where each element represents the intensity of a pixel. To visualize this image matrix, MATLAB's *imshow* function is used. This function interprets the matrix values as pixel intensities and displays the corresponding image. The command *imshow(ImageMatrix)* renders the image on the screen, allowing us to observe the effects of any processing like blurring or filtering.

Vectorization Process: Vectorization is the process of converting the 2D image matrix into a 1D vector. This is crucial for certain image processing techniques, particularly when interfacing with algorithms that expect input data in a linear format. In MATLAB, this can be achieved using the $(:)$ operator. For a given matrix B , representing the blurred image, vectorization is done by $B(:)$. This operation stacks the columns of B into a single column vector.

Customizing Vectorization Order: The default vectorization in MATLAB is column-wise.

However, if a different ordering is required (for example, row-wise), it can be achieved by first transposing the matrix. For instance, if we need to vectorize the image matrix B row by row, we first transpose B using $B = B'$, and then apply the vectorization: $b = B(:, :)$. This sequence of operations converts the 2D image matrix into a 1D vector, where the elements are ordered according to the rows of the original matrix.

Importance of Consistent Vectorization: It is essential to ensure that the vectorization ordering is consistent with the format expected by any algorithms or operations applied to the image data. Inconsistent vectorization can lead to incorrect data interpretation and processing results.

7.3. Simple Image Deblurring - Understanding the Transformation Matrix

In image processing, deblurring involves reversing the effect of a blur, which is typically represented by a kernel matrix. This process can be better understood by examining the roles of the kernel and the transformation matrix.

Kernel Matrix (K) in Image Blurring: - The kernel matrix $K \in \mathbb{R}^{d \times d}$ represents the blurring effect. Each element in this matrix specifies the weight given to neighboring pixels. - When an image is blurred, each pixel's new value becomes a weighted average of its own and its neighbors' values, based on K .

Construction of Transformation Matrix (A): - The transformation matrix A models the blurring process over the entire image. While K deals with local pixel neighborhoods, A applies the blur effect across the entire image. - A is constructed so that each row corresponds to a pixel in the vectorized image, and the non-zero elements in each row of A are derived from the kernel K . - The size of A is much larger than K , specifically $n^2 \times n^2$ for an image with $n \times n$ pixels.

Characteristics of Transformation Matrix (A):

- *Borders Handling:* A is designed to handle image borders appropriately. It ignores contributions from outside the image borders, indicated by conditions like $\max(i, j) > n$.
- *Banded Matrix:* Since A is derived from a smaller kernel K , it is a d^2 -banded matrix, implying sparsity. Most elements in A are zero, reflecting the localized influence of each pixel.
- *Memory Considerations:* The sparsity of A is crucial for memory efficiency. Defining a dense matrix of size $n^2 \times n^2$ is often impractical due to high memory requirements.

Example Kernel for Assignment: In this assignment, the provided transformation matrix A is generated from a normalized blur image kernel K , as shown below:

$$K = \frac{1}{605} \cdot \begin{bmatrix} 100 & 9 & 9 & 9 & 9 & 9 & 100 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 100 & 9 & 9 & 9 & 9 & 9 & 100 \end{bmatrix}$$

This kernel is applied to a row-vectorized image matrix. The transformation matrix A that you will use is derived from this kernel and reflects the combined effect of these weights on the entire image.

7.4. Understanding the Conjugate Gradient Method

The Conjugate Gradient (CG) method is an algorithm for efficiently solving systems of linear equations of the form $Ax = b$, where the matrix A is symmetric and positive-definite.

Characteristics of the Conjugate Gradient Method:

- *Direct vs. Iterative Method:* Theoretically, CG can be seen as a direct method since it computes the exact solution within a finite number of iterations, not exceeding the size of A . However, in practice, due to round-off errors and other perturbations, the directions in the algorithm may not remain precisely conjugate, and the exact solution is seldom reached.
- *Iterative Nature:* Despite this, CG is highly effective as an iterative method. It produces increasingly accurate approximations x_k to the exact solution, often reaching an acceptable level of accuracy in a number of iterations much smaller than the problem size.
- *Rate of Convergence:* The convergence rate is generally linear and heavily influenced by the condition number $k(A)$ of the matrix A . A higher $k(A)$ typically results in slower convergence.

Algorithmic Steps: The CG algorithm begins with an initial guess x_0 and iteratively refines this guess through a series of operations until the desired tolerance is met. The essential steps of the CG algorithm are as follows:

```
1 % Initialize variables
2 r = b - A * x0; % Residual
3 d = r; % Direction
4 p_old = dot(r, r);
5
6 % Iterative process
7 for i = 0, 1, ..., n-1
8     s = A * d;
9     alpha = p_old / dot(d, s);
10    x = x + alpha * d;
11    r = r - alpha * s;
12    p_new = dot(r, r);
13    beta = p_new / p_old;
14    d = r + beta * d;
15    p_old = p_new;
16 end
```

Each iteration of the algorithm performs matrix-vector and vector-vector operations to progressively refine the solution. The variables r , d , and x represent the residual, the direction, and the current solution estimate, respectively. The coefficients α and β are computed to ensure the efficiency and accuracy of the algorithm.

7.5. Conjugate Gradient - Symmetrix And Positivity

In the specific case of image deblurring, the transformation matrix A we encounter is symmetric and of full rank. However, it may not be positive-definite, which is a prerequisite for applying the Conjugate Gradient (CG) method directly. To address this, we can transform the original system into a form that meets the requirements for CG.

Augmented System Approach: We construct an augmented system that is positive-definite by pre-multiplying the original equation $Ax = b$ by the transpose of A . This results in a new system of equations:

$$A^T Ax = A^T b$$

In this formulation, the new transformation matrix $\tilde{A} = A^T A$ is symmetric and positive-definite. The right-hand side of the equation becomes $\tilde{b} = A^T b$. This augmented system is suitable for applying the CG method.

Characteristics of the Augmented System:

- *Symmetry*: The matrix \tilde{A} is symmetric since it is a product of A^T and A .
- *Positive-Definiteness*: The matrix \tilde{A} is guaranteed to be positive-definite as it is the product of a matrix and its transpose.
- *Suitability for CG*: With \tilde{A} being symmetric and positive-definite, the augmented system $\tilde{A}x = \tilde{b}$ is now suitable for the Conjugate Gradient method.

By solving the augmented system instead of the original one, we effectively utilize the CG method's advantages even when the initial matrix A does not satisfy all the necessary conditions for the method's direct application.

7.6. Conjugate Gradient - Condition Number

The effectiveness of the Conjugate Gradient (CG) method in solving a linear system $Ax = b$ is significantly influenced by the system's condition number, denoted as $k(A)$. The condition number essentially measures the sensitivity of the solution x to changes in b . In scenarios where minor variations in b lead to substantial fluctuations in x , the system is termed *ill-conditioned*, characterized by a large condition number.

Condition Number and CG Method: - *Effect on CG Convergence*: A high condition number adversely affects the convergence rate of the CG method. Although the total number of iterations required for an exact solution correlates to the count of unique eigenvalues of A , the condition number impacts the rate of convergence within these iterations.

- *Definition and Calculation*: For a square matrix A , the condition number is defined as:

$$k(A) = \frac{\sigma_{\max}}{\sigma_{\min}}$$

Here, σ represents the singular values of A . In the context of a real symmetric matrix like A , these singular values are equivalent to the absolute values of the eigenvalues, i.e., $\sigma = |\lambda|$.

Specifics in Our Case: - *Real Symmetric Matrix*: Given that A is symmetric and real, the condition number is determined by the ratio of the largest to the smallest eigenvalue magnitudes.

- *Estimating Condition Number*: To estimate A 's condition number in MATLAB, the command `condest` (or `cond` for dense matrices) is utilized. However, this computation can be time-consuming for large matrices.

- *Augmented System*: In our image deblurring scenario, we are dealing with the augmented system $A^T Ax = A^T b$, where the condition number is squared, $k(A^T A) = k(A)^2$, potentially exacerbating the issue of slow convergence.

- *Preconditioning as a Solution*: Despite these challenges, preconditioning techniques can be employed to partially mitigate the effects of a high condition number and enhance the convergence rate of the CG method in our deblurring application.

7.7. Conjugate Gradient - Preconditioning

Preconditioning is a crucial technique in the Conjugate Gradient (CG) method, particularly when dealing with a system where the matrix \tilde{A} is not ideally conditioned. The goal is to transform the original system into one that has a more favorable condition number, thereby enhancing the

efficiency of the CG method.

Concept of Preconditioning:

- *Objective:* We introduce a symmetric positive-definite preconditioner P with the aim of approximating the inverse of \tilde{A} , such that $P^{-1}\tilde{A} \approx I$. In essence, P^{-1} should closely resemble \tilde{A}^{-1} .
- *Cholesky Factorization:* The preconditioner P can be decomposed into $P = LL^T$, where L is the Cholesky factor. This decomposition facilitates solving the preconditioned augmented system.

Preconditioned Augmented System:

- *Transformation:* We aim to solve the system $P^{-1}\tilde{A}x = P^{-1}\tilde{b}$. By applying the Cholesky decomposition, this system can be re-expressed as:

$$(L^{-1}\tilde{A}L^{-T})(Lx) = L^{-1}\tilde{b}$$

Simplifying this, we arrive at a transformed system $Ax = b$, where A and b have been preconditioned to improve the condition number. - *Impact:* This transformation reduces the condition number $k(\tilde{A})$ and narrows the range of eigenvalues, thereby enhancing the convergence rate of the CG method.

Incomplete Cholesky (IC) Factorization: - *Approach:* A common strategy for preconditioning involves using Incomplete Cholesky factorization. In this method, we only compute the Cholesky factors for the non-zero elements of \tilde{A} . This approach yields a cost-effective approximation of \tilde{A} , represented by the preconditioner $P = F^T F$, where F is the sparse IC factor.

- *Challenges and Solutions:* However, the existence of F is not always guaranteed due to its enforced sparsity. To overcome this, a diagonal shift is often applied to P , ensuring its positive-definiteness and making F computable.

- *Implementation:* The MATLAB function *ichol* can be used for performing IC factorization, providing an effective and computationally efficient preconditioner.

In summary, preconditioning in the CG method is a powerful tool to address ill-conditioned systems, particularly in applications like image deblurring, where the transformation matrix may have an unfavorable condition number. By wisely choosing and implementing a preconditioner, we can significantly enhance the performance of the CG algorithm.

7.8. Project - General Questions

7.9. Project - Determining the Size of the Matrix A

In the context of image processing, particularly for operations like blurring, it's essential to understand the dimensions of the involved matrices.

Background: - *Matrix B :* Consider B as an $n \times n$ matrix representing a blurred image. - *Matrix A :* Matrix A is the transformation matrix responsible for the blurring effect applied on the original image.

Vectorization and Dimensionality:

- *Vectorized Form:* In the equation $Ax = b$, both x and b are vectorized forms of the matrices representing the original image (X) and the blurred image (B), respectively. - *Size of x and b :* Since X and B are $n \times n$ matrices, their vectorized forms, x and b , will each have a size of n^2 . This is due to the flattening of the $n \times n$ matrix into a single column vector.

Size of Matrix A : - *Calculation:* Given that $Ax = b$ and both x and b are vectors of length n^2 , the matrix A must therefore be an $n^2 \times n^2$ matrix to conform with the rules of matrix multiplication. - *Practical Computation:* In MATLAB, you can compute the dimension of A using the size of B . The code snippet below demonstrates this calculation:

```

1   dim = size(B, 1);
2   dim = dim * dim;

```

This code first retrieves the size of one dimension of B and then squares it to find the total number of elements in the vectorized form, which corresponds to the dimensions of A .

Implications: - The resulting dimensionality of A highlights the computational complexity involved in image processing tasks, such as deblurring. The significantly large size of A in practical applications necessitates efficient computational techniques and algorithms for image restoration.

7.10. Project - Determining the Number of Diagonal Bands in Matrix A

Understanding the structure of the transformation matrix A , particularly in terms of its diagonal bands, is crucial for image processing applications like blurring.

Kernel Matrix and Transformation Matrix: - *Kernel Matrix K :* The kernel matrix K used for blurring the image is normalized and has dimensions 7×7 . It is represented as:

$$K = \frac{1}{605} \cdot \begin{bmatrix} 100 & 9 & 9 & 9 & 9 & 9 & 100 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 9 & 1 & 1 & 1 & 1 & 1 & 9 \\ 100 & 9 & 9 & 9 & 9 & 9 & 100 \end{bmatrix}$$

- *Transformation Matrix A :* Matrix A is derived from K , where each non-zero element of K contributes to a diagonal in A .

Number of Diagonal Bands: - *Calculation:* The number of diagonal bands in A is determined by the dimension of K . Given that K is a 7×7 matrix, the transformation matrix A will have d^2 diagonal bands, where d is the dimension of K . - *Result:* Therefore, with $d = 7$, matrix A will have $7^2 = 49$ diagonal bands.

Implications: - The structure of A as a banded matrix is significant because it implies that A is sparse. This sparsity is beneficial for computational efficiency, especially when dealing with large images and matrices in image processing tasks. - Each diagonal band in A corresponds to a particular offset in the image blur effect, as defined by the kernel K .

7.11. Project - What is the length of the vectorized blurred image b ?

Understanding the size of the vectorized blurred image b is fundamental in image processing, especially when working with linear systems in the form $Ax = b$.

Matrix-to-Vector Conversion: - *Original Blurred Image Matrix B :* The blurred image is initially represented as a matrix B , which is typically a square matrix of size $n \times n$. - *Vectorization Process:* The process of vectorization transforms the matrix B into a vector b . This involves rearranging all elements of B into a single column vector.

Calculating the Length of b : - *Formula:* The length of the vector b is determined by the total number of elements in the matrix B . For a matrix of size $n \times n$, this is calculated as n^2 . - *Example:* If B is a 250×250 matrix, the length of b will be $250^2 = 62500$.

Implications in Linear Algebra: - In the context of solving the linear system $Ax = b$ for image deblurring, understanding the size of b is crucial. It directly influences the dimensions of matrix A and the solution vector x . - The size of b also has implications for computational requirements, especially in terms of memory allocation and processing power needed for matrix operations.

7.12. Project - If A is not symmetric, how would this affect \tilde{A} ?

When dealing with image deblurring and solving the corresponding linear system, it is important to consider the properties of the transformation matrix A , especially its symmetry.

Symmetry in Transformation Matrices: - *Non-Symmetric Matrix A :* If A is not symmetric, this attribute in itself does not hinder the formulation of the augmented transformation matrix \tilde{A} . - *Augmented Matrix \tilde{A} :* The augmented matrix \tilde{A} , defined as $A^T A$, inherently becomes symmetric regardless of the symmetry of A . This is a direct result of the matrix multiplication properties.

Characteristics of \tilde{A} : - *Symmetry:* The matrix multiplication $A^T A$ results in a symmetric matrix since $(A^T A)^T = A^T (A^T)^T = A^T A$, demonstrating the symmetry of \tilde{A} . - *Positive-Definiteness:* Although symmetry is ensured, the positive-definiteness of \tilde{A} is not guaranteed by the non-symmetry of A alone. Additional properties of A may influence this aspect.

Implications for Image Deblurring: - In the context of image deblurring, if A is originally not symmetric, the constructed \tilde{A} will still be symmetric, which is favorable for certain numerical methods, such as the Conjugate Gradient method. - However, the solver's effectiveness and the nature of the solution may be influenced by the characteristics of A and subsequently \tilde{A} , particularly regarding its positive-definiteness and condition number.

7.13. Project - Effect of Non-Symmetric A on Augmented Matrix \tilde{A}

When dealing with image deblurring and solving the corresponding linear system, it is important to consider the properties of the transformation matrix A , especially its symmetry.

Symmetry in Transformation Matrices: - *Non-Symmetric Matrix A :* If A is not symmetric, this attribute in itself does not hinder the formulation of the augmented transformation matrix \tilde{A} . - *Augmented Matrix \tilde{A} :* The augmented matrix \tilde{A} , defined as $A^T A$, inherently becomes symmetric regardless of the symmetry of A . This is a direct result of the matrix multiplication properties.

Characteristics of \tilde{A} : - *Symmetry:* The matrix multiplication $A^T A$ results in a symmetric matrix since $(A^T A)^T = A^T (A^T)^T = A^T A$, demonstrating the symmetry of \tilde{A} . - *Positive-Definiteness:* Although symmetry is ensured, the positive-definiteness of \tilde{A} is not guaranteed by the non-symmetry of A alone. Additional properties of A may influence this aspect.

Implications for Image Deblurring: - In the context of image deblurring, if A is originally not symmetric, the constructed \tilde{A} will still be symmetric, which is favorable for certain numerical methods, such as the Conjugate Gradient method. - However, the solver's effectiveness and the nature of the solution may be influenced by the characteristics of A and subsequently \tilde{A} , particularly regarding its positive-definiteness and condition number.

7.14. Project - Explain why $Ax = b$ for x is equivalent to minimizing $\frac{1}{2}x^T Ax - b^T x$ over x , assuming that A is symmetric positive-definite.

Objective Function: - Consider the quadratic form $f(x) = \frac{1}{2}x^T Ax - b^T x$ as our objective function. - The goal is to find the value of x that minimizes this function.

Differentiation and Critical Points: - To find the minimum, we differentiate $f(x)$ with respect to x and set the derivative to zero. - The derivative is $f'(x) = \frac{1}{2}(A^T + A)x - b$.

Symmetry of A : - Since A is symmetric, $A^T = A$. - Therefore, $f'(x)$ simplifies to $Ax - b$.

Finding the Minimum: - Setting $f'(x) = 0$, we get the equation $Ax - b = 0$ or equivalently $Ax = b$. - The critical point where $f'(x) = 0$ corresponds to the solution of the linear system $Ax = b$.

Conclusion: - Since A is positive-definite, the quadratic form $f(x)$ is convex, and any critical point is a global minimum. - Thus, solving $Ax = b$ is equivalent to finding the minimum of $\frac{1}{2}x^T Ax - b^T x$ over x when A is symmetric and positive-definite.

7.15. Conjugate Gradient - Conjugate Gradient Solver

Function CG solver $[x, rvec] = myCG(A, b, x0, maxitr, tol]$ where x and $rvec$ are, respectively the solution value and a vector containing the residual at every iteration.

```

1  function [x, rvec] = myCG(coeff_matrix, right_hand_side, initial_guess, max_iterations,
    tolerance)
2      x = initial_guess;
3      residual = right_hand_side - coeff_matrix * x;
4      search_direction = residual;
5      rvec = zeros(max_iterations, 1);
6
7      for k = 1:max_iterations
8          alpha = (residual' * residual) / (search_direction' * coeff_matrix *
                search_direction);
9          x = x + alpha * search_direction;
10         residual_old = residual;
11         residual = residual - alpha * coeff_matrix * search_direction;
12         beta = (residual' * residual) / (residual_old' * residual_old);
13         search_direction = residual + beta * search_direction;
14         rvec(k) = norm(residual);
15
16         if norm(residual) < tolerance
17             break;
18         end
19     end
20     rvec = rvec(1:k);
21 end

```

5. Project 5 - Linear Programming and the Simplex Method

1. Inequalities and System of Inequalities

We are given the following system of inequalities:

$$\begin{cases} x + 2y \leq 8 \\ 3x - 5y \leq 13 \\ x \geq 0, y \geq 0 \end{cases} \quad (8)$$

To determine the solution region on the plane, we need to find the set of all points (x, y) that satisfy these inequalities.

Rearranging Inequalities: - The first inequality, $x + 2y \leq 8$, can be rearranged to $y \leq -\frac{1}{2}x + 4$.
- The second inequality, $3x - 5y \leq 13$, can be rearranged to $y \geq \frac{3}{5}x - \frac{13}{5}$.

Graphical Interpretation: - These inequalities represent regions in the xy -plane. - The first inequality represents the area below the line $y = -\frac{1}{2}x + 4$. - The second inequality represents the area above the line $y = \frac{3}{5}x - \frac{13}{5}$.

Non-Negative Constraints: - The constraints $x \geq 0$ and $y \geq 0$ limit the solution to the first

quadrant of the plane.

Identifying the Solution Region: - The solution region is the intersection of these areas, considering all inequalities. - It is a bounded region, defined by the overlapping areas of the inequalities.

Possible Cases: - In some cases, systems of inequalities might not have a solution (incompatible constraints). - In other cases, the solution region could be unbounded, allowing any values within the constraints.

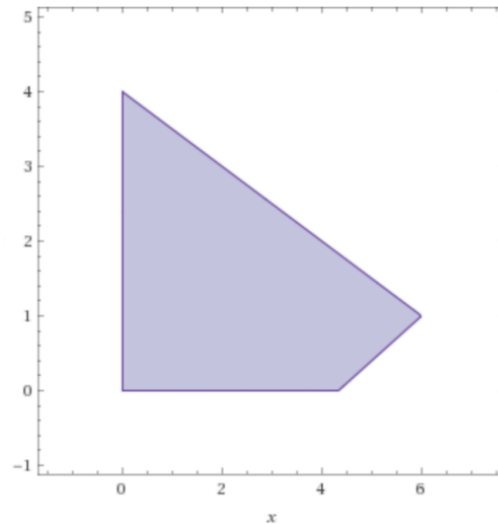


Figure 11: Graphical representation of the solution region for the given system of inequalities.

2. Linear Programming

Linear programming is a method used in mathematical optimization, where the goal is to maximize or minimize a linear objective function. This method finds applications in diverse fields for optimizing various processes, such as minimizing costs or maximizing profits.

Objective Function: The objective function in linear programming is a linear equation of the form:

$$z = c_1x_1 + c_2x_2 + \dots + c_nx_n$$

where c_i are coefficients and x_i are variables.

Constraints: The optimization is subject to a set of linear constraints, typically expressed as inequalities. These can be represented as:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \dots + a_{1,n}x_n &\geq h_1 \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \dots + a_{m,n}x_n &\leq h_m \end{aligned}$$

Matrix Notation: To express the problem in a compact form, we define the following quantities:

$$\begin{aligned} c &\in \mathbb{R}^n : \text{Coefficient vector} \\ A &\in \mathbb{R}^{m \times n} : \text{Constraint coefficients matrix} \\ h &\in \mathbb{R}^m : \text{Constraint values vector} \\ x &\in \mathbb{R}^n : \text{Solution vector} \end{aligned}$$

2.1. Standard Form of Linear Programming:

For convenience, we convert all inequalities to either upper bounds (for maximization problems) or lower bounds (for minimization problems), and include non-negativity constraints. The standard forms are:

Maximization Problem:

$$\begin{aligned}
&\text{maximize} && z = c^T x \\
&\text{s.t.} && Ax \leq h \\
&&& x \geq 0
\end{aligned}$$

Minimization Problem:

$$\begin{aligned}
&\text{minimize} && z = c^T x \\
&\text{s.t.} && Ax \geq h \\
&&& x \geq 0
\end{aligned}$$

In these forms, the non-negativity constraint $x \geq 0$ is explicitly stated, ensuring that the solution vector x contains only non-negative values.

2.2. Fundamental Theorem of Linear Programming

The Fundamental Theorem of Linear Programming provides critical insights into the solution of linear programming problems. It states:

- If a linear programming problem has a solution, this solution is guaranteed to occur at a vertex (corner point) of the feasible solution set. This feasible solution set is defined by the constraints of the problem.
- In the scenario where multiple solutions exist, at least one of these solutions will be found at a vertex of the feasible solution set. This implies that while multiple solutions might span across the edges of the feasible region, they all intersect at the vertices.
- Regardless of whether there is a single solution or multiple solutions, the value of the objective function at these optimal solutions remains unique.

A graphical solution method can be effectively used for linear programming problems involving two variables. The steps are as follows:

1. Solve the system of inequalities defined by the constraints of the problem.
2. Identify the feasible region on a plane. This region is a polygonal area where all points satisfy the given constraints. According to the Fundamental Theorem of Linear Programming, at least one optimal solution will be located at a vertex of this feasible region.
3. Evaluate the objective function at each vertex of the feasible region. Choose the vertex that optimizes (maximizes or minimizes) the value of the objective function. If the feasible region is bounded, both the maximum and minimum values of the objective function can be determined.

Consider the following example of a maximization problem:

$$\begin{aligned}
&\text{maximize} && z = 4x + 6y \\
&\text{s.t.} && \\
&&& -x + y \leq 11, \\
&&& x + y \leq 27, \\
&&& 2x + 5y \leq 90, \\
&&& x, y \geq 0.
\end{aligned}$$

This system of inequalities can be rewritten as:

$$\begin{cases} y \leq x + 11, \\ y \leq -x + 27, \\ y \leq -\frac{2}{5}x + 18, \\ x, y \geq 0. \end{cases}$$

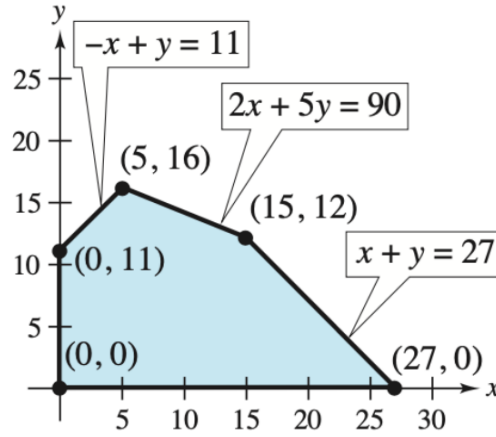


Figure 12: Graphical representation of the feasible region.

To find the optimal solution, evaluate the objective function at each vertex of the feasible region:

- At Vertex $(0,0)$: $z = 0$,
- At Vertex $(0,11)$: $z = 66$,
- At Vertex $(5,16)$: $z = 116$,
- At Vertex $(15,12)$: $z = 132$ (Optimal),
- At Vertex $(27,0)$: $z = 108$.

The optimal value of the objective function, z^* , is achieved at the vertex $(15,12)$.

2.3. Formulating a Problem as Linear Program

The liquid portion of a diet is to provide at least 300 calories, 36 units of vitamin A, and 90 units of vitamin C daily. A cup of dietary drink X provides 60 calories, 12 units of vitamin A, and 10 units of vitamin C. A cup of dietary drink Y provides 60 calories, 6 units of vitamin A, and 30 units of vitamin C. Now, suppose that dietary drink X cost \$0.12 per cup and drink Y cost \$0.15 per cup. How many cups of each drink should be consumed each day to minimise the cost and still meet the state daily requirements?

We obtain the minimization problem:

$$\begin{aligned} \min \quad & z = 0.12x + 0.15y \\ \text{s.t.} \quad & 60x + 60y \geq 300, \\ & 12x + 6y \geq 36, \\ & 10x + 30y \geq 90, \\ & x, y \geq 0. \end{aligned}$$

$$\begin{cases} 60x + 60y \geq 300 & \Rightarrow & y \geq -x + 5 \\ 12x + 6y \geq 36 & \Rightarrow & y \geq -2x + 6 \\ 10x + 30y \geq 90 & \Rightarrow & y \geq -\frac{1}{3}x + 3, y \geq 0. \end{cases}$$

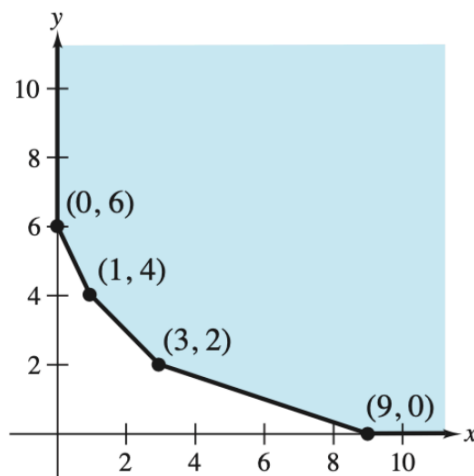


Figure 13: Graphical representation of the feasible region.

To find the optimal solution, evaluate the objective function at each vertex of the feasible region:

- At Vertex $(0, 6)$: $z = 0.90$,
- At Vertex $(1, 4)$: $z = 0.72$,
- At Vertex $(3, 2)$: $z = 0.66$,
- At Vertex $(9, 0)$: $z = 1.08$,

The optimal is at Vertex $(3, 2)$: $z^* = 0.66$.

2.4. The Simplex Method

The Simplex Method, introduced by George Dantzig in 1947, is an algorithm designed for efficiently solving linear programming problems, especially useful for scenarios with a large number of variables and constraints. However, in the worst-case scenario, the Simplex Method can have exponential complexity due to the potentially enormous number of candidate solutions that need to be examined to find the optimal value.

Formulating the Problem: To utilize the Simplex Method, the linear programming problem must first be expressed in standard form. This involves transforming all inequality constraints into equalities by introducing slack variables in a maximization problem or surplus variables in a minimization problem. These additional variables, which have zero coefficients in the objective function, represent the deviation from the constraint's right-hand side value in the optimal solution.

Dimensionality and Variables: The coefficient vector is now expanded to $c \in \mathbb{R}^{n+m}$, where m is the number of constraints. Similarly, the unknowns vector expands to $x \in \mathbb{R}^{n+m}$, encompassing both the original variables x_1, \dots, x_n and the slack/surplus variables s_1, \dots, s_m . The coefficient matrix A also changes accordingly to $A \in \mathbb{R}^{m \times (n+m)}$.

Algorithm Overview: The Simplex Method begins by identifying a basis of the extended coefficient matrix A , which is a square matrix with full rank comprising a mix of original and newly added variables. Variables outside this basis are termed *nonbasic variables*. By setting nonbasic variables to zero, a *basic solution* is derived. It is crucial to verify that all basic solutions satisfy the non-negativity condition to ensure feasibility.

According to the Fundamental Theorem of Linear Programming, the existence of a feasible solution implies the existence of a feasible basic solution, and likewise, an optimal solution implies the existence of an optimal basic solution. The Simplex algorithm iterates through basic solutions, starting

from a feasible one, and progressing towards the optimum. The worst-case number of iterations equals the total number of basic solutions, given by:

$$N = \frac{(m+n)!}{m! \cdot n!}.$$

This number grows exponentially with the increase in variables and constraints, highlighting the potential complexity issue of the Simplex Method. Despite this, under typical conditions, the Simplex Method often reaches the optimal solution in a practical number of iterations.

In summary, the Simplex Method is a powerful tool for solving linear programming problems but requires careful formulation of the problem and an understanding of its potential complexity in scenarios with a large number of variables and constraints.

2.5. The Simplex Method - Maximization Problem

For the maximization problem in standard linear programming form, we can introduce slack variables to convert inequality constraints into equalities. From the original set of constraints:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,n}x_n &\leq h_1 \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \cdots + a_{m,n}x_n &\leq h_m \end{aligned}$$

We transform them to:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,n}x_n + s_1 &= h_1 \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \cdots + a_{m,n}x_n + s_m &= h_m \end{aligned}$$

Here, s_1, s_2, \dots, s_m are the slack variables. The problem is then reformulated as:

$$c = \begin{bmatrix} c_1 \\ \vdots \\ c_n \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \\ s_1 \\ \vdots \\ s_m \end{bmatrix}, \quad A = \begin{bmatrix} a_{1,1} & \cdots & a_{1,n} & 1 & 0 & \cdots & 0 \\ a_{2,1} & \cdots & a_{2,n} & 0 & 1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & \cdots & a_{m,n} & 0 & 0 & \cdots & 1 \end{bmatrix}, \quad h = \begin{bmatrix} h_1 \\ \vdots \\ h_m \end{bmatrix}$$

Additionally, non-negativity constraints are applied to both the decision variables (x_1, \dots, x_n) and the slack variables (s_1, \dots, s_m).

2.6. The Simplex Method - Minimization Problem

In the minimization version of a linear programming problem, we address standard form inequalities by introducing *surplus variables*. This transformation converts inequality constraints into equalities, shifting from the original constraints:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,n}x_n &\geq h_1, \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \cdots + a_{m,n}x_n &\geq h_m \end{aligned}$$

to the following system:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,n}x_n - s_1 &= h_1, \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \cdots + a_{m,n}x_n - s_m &= h_m \end{aligned}$$

Here, s_1, s_2, \dots, s_m are the surplus variables. The problem is then reformulated with the following parameters:

$$c = \begin{bmatrix} c_1 \\ \vdots \\ c_n \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \\ s_1 \\ \vdots \\ s_m \end{bmatrix}, \quad A = \begin{bmatrix} a_{1,1} & \cdots & a_{1,n} & -1 & 0 & \cdots & 0 \\ a_{2,1} & \cdots & a_{2,n} & 0 & -1 & \cdots & 0 \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{m,1} & \cdots & a_{m,n} & 0 & 0 & \cdots & -1 \end{bmatrix}, \quad h = \begin{bmatrix} h_1 \\ \vdots \\ h_m \end{bmatrix}$$

Additionally, non-negativity constraints are imposed on both the decision variables (x_1, \dots, x_n) and the surplus variables (s_1, \dots, s_m) .

2.7. Basic and Nonbasic Variables

In linear programming, the matrix A can be rearranged by freely swapping its rows and corresponding elements of vector h , and by swapping its columns and corresponding elements of vector x . We partition A into two sub-matrices:

$$A = [B \quad D]$$

Here, $B \in \mathbb{R}^{m \times m}$ consists of m linearly independent columns from A , ensuring full rank and invertibility. $D \in \mathbb{R}^{m \times (n-m)}$ contains the remaining columns. Similarly, vectors x and c are partitioned as:

$$x = \begin{bmatrix} x_B \\ x_D \end{bmatrix}, \quad c = \begin{bmatrix} c_B \\ c_D \end{bmatrix}$$

with $x_B, c_B \in \mathbb{R}^m$ and $x_D, c_D \in \mathbb{R}^{n-m}$.

The constraint system $Ax = h$ can be rewritten using this decomposition:

$$Bx_B + Dx_D = h \quad \Rightarrow \quad Bx_B = h - Dx_D$$

Variables in x_B are known as *basic variables*, while those in x_D are *nonbasic variables*.

Inverting the matrix B (which is always nonsingular due to its full rank), we obtain:

$$x_B = B^{-1}h - B^{-1}Dx_D$$

This equation represents a general solution of the problem, with the values of basic variables x_B depending on the values assigned to the nonbasic variables x_D . Note that x_B can include both the original unknowns x_1, \dots, x_n and the slack variables s_1, \dots, s_m .

Basic Solution:

A *basic solution* is obtained by setting $x_D = 0$:

$$x_B = B^{-1}h$$

If the non-negativity condition $x_B \geq 0$ is satisfied, the basic solution is feasible. If one or more components of x_B are zero, the solution is degenerate. Basic solutions correspond to vertices of the feasible region polytope, defined by the intersections of constraints.

2.8. The Simplex Method: Algorithm

The Simplex method is a systematic procedure for solving linear programming problems. The algorithm can be summarized as follows:

1. Convert the given problem into its standard form by adding slack or surplus variables as needed.

2. Identify a feasible starting basic solution. This can often be achieved by solving an auxiliary problem if a clear starting solution is not immediately apparent.
3. **While** the optimality criterion is not met:
 - a) Apply the iteration rule, which involves exchanging one of the basic variables with a nonbasic variable. This step moves the solution from one vertex of the feasible region to another, seeking to improve the objective function value.
4. **End While**
5. Return the basic solution that satisfies the optimality criterion.

This method iteratively moves towards the optimal solution by navigating the vertices of the feasible region defined by the constraints.

2.9. Optimal Basic Solution

The existence of a feasible solution, as stated by the Fundamental Theorem of Linear Programming, guarantees the existence of a feasible basic solution. Similarly, if an optimal solution exists, there must also be an optimal basic solution. The Simplex method navigates through the vertices (corners) of the polytope defined by the constraints of the program until the optimality criterion is met.

Downside: The potential number of basic solutions increases exponentially with the number of variables and constraints. The maximum possible number of iterations can be calculated as:

$$N = \frac{(m+n)!}{m!n!}$$

where m is the number of constraints and n is the number of variables.

The optimality condition at each iteration is evaluated using the **Reduced Cost Coefficients**:

$$r_D = c_D^T - c_B^T B^{-1} D$$

Here, c_D represents the coefficients of the nonbasic variables in the objective function, c_B represents the coefficients of the basic variables, B^{-1} is the inverse of the basic matrix, and D is the matrix of nonbasic columns.

The optimality conditions for the Simplex method are:

- **Maximization:** $r_D \leq 0$
- **Minimization:** $r_D \geq 0$

Iterative Rule: If the optimality condition is not satisfied, select the variable with the highest (for maximization) or lowest (for minimization) reduced cost coefficient r_D to enter the basis. Determine the departing variable by calculating the ratio:

$$\frac{B^{-1}h}{B^{-1}D}$$

for the column associated with the entering variable and select the variable with the smallest **positive** value. This step ensures the maintenance of a feasible solution while seeking improvement in the objective function value.

2.10. Feasible Starting Basic Solution: Auxiliary Problem

To initiate the Simplex method, a feasible starting basic solution is required, which is typically obtained by solving an **auxiliary problem**. The initial step involves ensuring that all the elements of the right-hand side vector h are positive. This condition is achieved by either adding slack variables or subtracting surplus variables to/from the original constraints, ensuring $h_i \geq 0$ for all i .

The constraints are modified as follows:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,n}x_n \pm s_1 &= h_1, \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \cdots + a_{m,n}x_n \pm s_m &= h_m, \end{aligned}$$

where the sign of s_i is chosen to ensure non-negativity of h_i .

To construct the auxiliary problem, we add artificial variables u_1, \dots, u_m to each constraint:

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,n}x_n \pm s_1 + u_1 &= h_1, \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \cdots + a_{m,n}x_n \pm s_m + u_m &= h_m. \end{aligned}$$

The auxiliary problem's objective function is then defined as the sum of these artificial variables:

$$z_{aux} = u_1 + u_2 + \cdots + u_m.$$

This auxiliary problem is inherently a **minimization** problem. The initial feasible basic solution corresponds to the case where $z_{aux} = 0$. If z_{aux} cannot be minimized to zero, it implies that neither the auxiliary problem nor the original linear programming problem has a feasible solution.

2.11. Simplex Method: Minimization Example

Consider the following linear programming problem for minimization:

$$\begin{aligned} \min \quad & z = 9x_1 + 11x_2 + 8x_3, \\ \text{s.t.} \quad & \\ & 2x_1 + 2x_2 + 3x_3 \geq 210, \\ & 3x_1 + 4x_2 + 2x_3 \geq 360, \\ & x_1, x_2, x_3 \geq 0. \end{aligned}$$

To transform inequalities into equalities, surplus variables are subtracted:

$$\begin{aligned} 2x_1 + 2x_2 + 3x_3 - s_1 &= 210, \\ 3x_1 + 4x_2 + 2x_3 - s_2 &= 360, \\ x_1, x_2, x_3, s_1, s_2 &\geq 0. \end{aligned}$$

Vectors and matrices are defined as:

$$c = \begin{bmatrix} 9 \\ 11 \\ 8 \\ 0 \\ 0 \end{bmatrix}, \quad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ s_1 \\ s_2 \end{bmatrix}, \quad A = \begin{bmatrix} 2 & 2 & 3 & -1 & 0 \\ 3 & 4 & 2 & 0 & -1 \end{bmatrix}, \quad h = \begin{bmatrix} 210 \\ 360 \end{bmatrix}.$$

Separating into basic (B) and non-basic (D) components:

$$c_B = \begin{bmatrix} 8 \\ 0 \end{bmatrix}, \quad x_B = \begin{bmatrix} x_3 \\ s_1 \end{bmatrix}, \quad B = \begin{bmatrix} 3 & -1 \\ 2 & 0 \end{bmatrix}, \quad h = \begin{bmatrix} 210 \\ 360 \end{bmatrix},$$

$$c_D = \begin{bmatrix} 9 \\ 11 \\ 0 \end{bmatrix}, \quad x_D = \begin{bmatrix} x_1 \\ x_2 \\ s_2 \end{bmatrix}, \quad D = \begin{bmatrix} 2 & 1 & 0 \\ 3 & 4 & -1 \end{bmatrix}.$$

Computing the basic solution with $x_D = 0$:

$$x_B = B^{-1}h = \begin{bmatrix} 0.5 & -0.25 \\ -0.5 & 0.75 \end{bmatrix} \begin{bmatrix} 210 \\ 360 \end{bmatrix} = \begin{bmatrix} 180 \\ 330 \end{bmatrix}.$$

This feasible solution has an objective function value of $z = c_B^T x_B = 1440$. To determine if it is optimal, compute the reduced cost coefficients:

$$r_D = c_D^T - c_B^T B^{-1} D = \begin{bmatrix} -3 & -5 & 4 \end{bmatrix}.$$

Since there are negative coefficients in r_D , the solution is not optimal for a minimization problem. The entering variable is x_2 , corresponding to the smallest value in r_D . For the departing variable, calculate the ratios:

$$\frac{B^{-1}h}{B^{-1}D} = \begin{bmatrix} 90 \\ 66 \end{bmatrix}.$$

Selecting the smallest positive value indicates s_1 as the departing variable. Reformulating the basic and non-basic variables:

$$c_B = \begin{bmatrix} 8 \\ 11 \end{bmatrix}, \quad x_B = \begin{bmatrix} x_3 \\ x_2 \end{bmatrix}, \quad B = \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix},$$

$$c_D = \begin{bmatrix} 9 \\ 0 \\ 0 \end{bmatrix}, \quad x_D = \begin{bmatrix} x_1 \\ s_1 \\ s_2 \end{bmatrix}, \quad D = \begin{bmatrix} 2 & -1 & 0 \\ 3 & 0 & -1 \end{bmatrix}.$$

The new basic solution yields $x_B = \begin{bmatrix} 48 \\ 66 \end{bmatrix}$ with $z = 1100$. Further iterations involve computing new r_D , identifying entering and departing variables, and recalculating x_B until all reduced cost coefficients are non-negative, indicating an optimal solution has been found.

3. Project 5

3.1. Project 5 - Graphical Solution of Linear Programming Problems

(1)

$$\begin{aligned} \min \quad & z = 4x + y \\ \text{s.t.} \quad & x + 2y \leq 40 \\ & x + y \geq 30 \\ & 2x + 3y \geq 72 \\ & x, y \geq 0 \end{aligned}$$

- (2) A tailor plans to sell two types of trousers, with production costs of 25 CHF and 40 CHF, respectively. The former type can be sold for 85 CHF, while the latter for 110 CHF. The tailor estimates a total monthly demand of 265 trousers. Find the number of units of each type of trousers that should be produced in order to maximise the net profit of the tailor, if we assume that the he cannot spend more than 7000 CHF in raw materials.

Start by writing problem (2) as a linear programming problem. Then complete the following tasks:

- Solve the system of inequalities.
- Plot the feasible region identified by the constraints.
- Find the optimal solution and the value of the objective function in that point.

4. Solve the system of inequalities.

1.

$$\begin{aligned} \min \quad & z = 4x + y \\ \text{s.t.} \quad & x + 2y \leq 40 \\ & x + y \geq 30 \\ & 2x + 3y \geq 72 \\ & x, y \geq 0 \end{aligned}$$

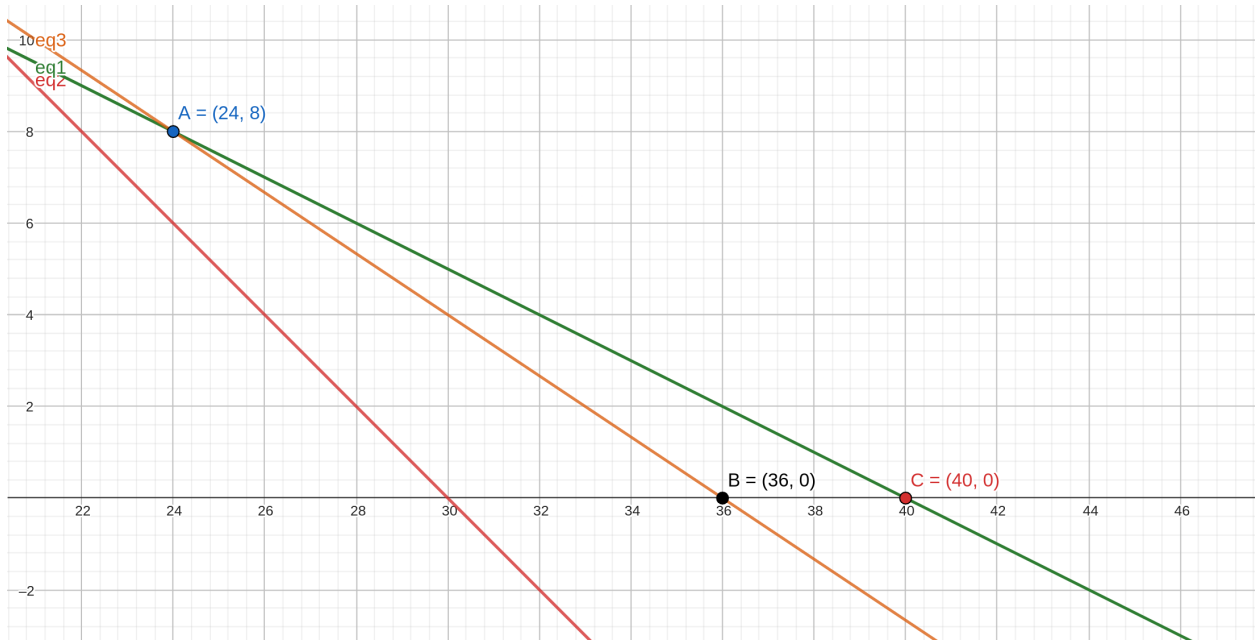


Figure 14: Graphical Illustration of the equations with geogebra.

In Figure 14, the graphical representation of each inequality constraint has been transformed into its equivalent equality for visualization purposes. The line for the inequality $x + 2y \leq 40$ is accompanied by shading in the area below it, whereas for the inequalities $x + y \geq 30$ and $2x + 3y \geq 72$, shading is applied above their respective lines.

The intersections of these lines form the vertices of the feasible region, identifiable as points:

- $A = (24, 8)$
- $B = (36, 0)$
- $C = (40, 0)$

To determine the minimum value of z , the equation $z = 4x + y$ is solved at each of these points.

For $A = (24, 8)$

$$\begin{aligned} z_1 &= 4(24) + 1(8) \\ z_1 &= 96 + 8 \\ z_1 &= 104 \end{aligned}$$

For $A = (36, 0)$

$$\begin{aligned} z_1 &= 4(36) + 1(0) \\ z_1 &= 144 + 0 \\ z_1 &= 144 \end{aligned}$$

For $A = (40, 0)$

$$\begin{aligned}z_1 &= 4(40) + 1(0) \\z_1 &= 160 + 0 \\z_1 &= 160\end{aligned}$$

The minimum value achieved for z is $z_1 = 104$.

2.

Consider x as the quantity of the first type of trousers and y as the quantity of the second type. The tailor's objective is to maximize the net profit, which is calculated as the difference between the selling price and the production cost.

The net profit per unit for the first type of trousers is calculated as:

$$85 - 25 = 60 \text{ CHF per unit}$$

Similarly, the net profit per unit for the second type of trousers is:

$$110 - 40 = 70 \text{ CHF per unit}$$

Consequently, the objective function to be maximized is:

$$z = 60x + 70y$$

The constraints:

- Estimate total monthly demand is: $x + y \leq 265$
- Budget for raw materials: $25x + 40y \leq 7000 \text{ CHF}$
- Non-negativity: $x \geq 0 \quad y \geq 0$

$$\begin{aligned}\max \quad & z = 60x + 70y \\ \text{s.t.} \quad & x + y \leq 265 \\ & 25x + 40y \leq 7000 \\ & x, y \geq 0\end{aligned}$$

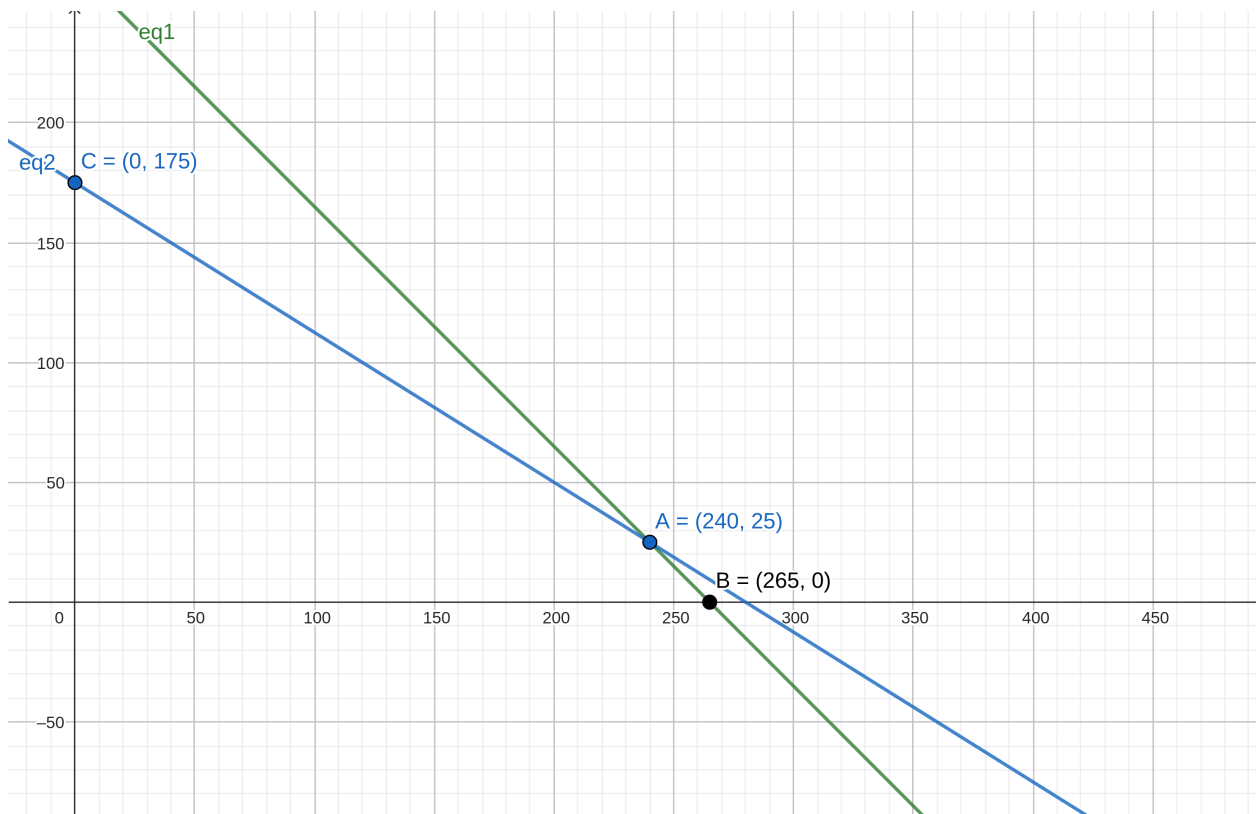


Figure 15: Graphical Illustration of the equations with geogebra.

It is possible to get three points on the vertices of the feasible region:

- $A = (240, 25)$
- $B = (265, 0)$
- $C = (0, 175)$

It is possible to find the minimum z from all the points trough solving the $z = 60x + 70y$.

For $A = (240, 25)$

$$\begin{aligned} z_1 &= 60(240) + 70(25) \\ z_1 &= 14400 + 1750 \\ z_1 &= 16150 \end{aligned}$$

For $A = (265, 0)$

$$\begin{aligned} z_1 &= 60(265) + 70(0) \\ z_1 &= 15900 \end{aligned}$$

For $A = (0, 175)$

$$\begin{aligned} z_1 &= 60(0) + 70(175) \\ z_1 &= 12250 \end{aligned}$$

The Minimum of Z is $z_1 = 12250$.

5. Project 5 - Implementation of the Simplex Method

SimplexSolve.m - which solves a maximisation or minimisation problem using the simplex method.

```
1  function [x_B,c_B,index_B] = simplexSolve
    (type,B,D,c_B,c_D,h,x_B,x_D,index_B,index_D,itMax)
2  % Solving a maximization problem with the simplex method
3
4  % Initialize the number of iterations
5  nIter = 0;
6
7  % Compute  $B^{-1} \cdot D$  and  $B^{-1} \cdot h$ 
8  BiD = B \ D;
9  Bih = B \ h;
10
11 % Compute the reduced cost coefficients
12 r_D = c_D - c_B * BiD;
13
14 tol = max(size(r_D)); % the optimality condition is satisfied if all reduced cost
    coefficients are positive or negative (depending on the problem)
15
16 % Check the optimality condition, in order to skip the loop if the solution is already
    optimal)
17 if(strcmp(type,'max'))
18     optCheck = sum(r_D <= 0);
19 elseif(strcmp(type,'min'))
20     optCheck = sum(r_D >= 0);
21 else
22     error('Incorrect type specified. Choose either a maximisation (max) or minimisation
        (min) problem.')
23 end
24
25 while(optCheck~=tol)
26
27     % TODO: Find the index of the entering variable
28     if(strcmp(type,'max'))
29         [~, idxIN]= max(r_D);
30     elseif(strcmp(type,'min'))
31         [~, idxIN] = min(r_D);
32     else
33         error('Incorrect type specified. Choose either a maximisation (max) or
            minimisation (min) problem.')
34     end
35
36     in = D(:,idxIN);
37     c_in = c_D(1,idxIN);
38     index_in = index_D(1,idxIN);
39
40     % Evaluate the coefficients ratio for the column corresponding to the entering
        variable
41     ratio = Bih ./ BiD(:, idxIN);
42
43     % Find the smallest positive ratio
44     idxOUT = find(ratio == min(ratio(ratio >= 0)), 1);
45     out = B(:,idxOUT);
46     c_out = c_B(1,idxOUT);
47     index_out = index_B(1,idxOUT);
48
49     % Update the matrices by exchanging the columns
50     B(:,idxOUT) = in;
51     D(:,idxIN) = out;
```

```

52     c_B(1,idxOUT) = c_in;
53     c_D(1,idxIN) = c_out;
54     index_B(1,idxOUT) = index_in;
55     index_D(1,idxIN) = index_out;
56
57     % Compute  $B^{-1} \cdot D$  and  $B^{-1} \cdot h$ 
58     BiD = B\D;
59     Bih = B\h;
60
61     % Compute the reduced cost coefficients
62     r_D = c_D - c_B * BiD;
63
64     % Check the optimality condition
65     if(strcmp(type,'max'))
66         optCheck = sum(r_D <= 0);
67     elseif(strcmp(type,'min'))
68         optCheck = sum(r_D >= 0);
69     else
70         error('Incorrect type specified. Choose either a maximisation (max) or
              minimisation (min) problem.')

```

simplex.m - is a wrapper which calls all the functions necessary to find a solution to the linear program.

```

1  function [z,x_B,index_B] = simplex (type,A,h,c,sign)
2  % Simplex method solver for a linear programming problem
3  % Input arguments:
4  %   type = 'max' for maximization, 'min' for minimization
5  %   A     = matrix holding the constraints coefficients
6  %   h     = coefficients of the constraints inequalities (rhs vector)
7  %   c     = coefficients of the objective function
8  %   sign = vector holding information about the constraints if the system
9  %           needs to be standardized (-1: less or equal, 0: equal, 1:vgreater or equal)
10
11  m = size(A,1);
12  n = size(A,2);
13
14  % Compute the maximum number of basic solutions of the original problem (i.e., the
      maximum number of iterations necessary to solve the problem)
15  itMax = nchoosek(m+n, n);
16  % Writing the problem in standard form
17  [A_aug,h,c_aug] = standardize(type,A,h,c,m,sign);

```

```

18
19 % Determining a starting feasible initial basic solution
20 [B,D,c_B,c_D,x_B,x_D,index_B,index_D] = auxiliary(A_aug,c_aug,h,m,n);
21
22 % Solving the linear programming problem with the Simplex method
23 [x_B,c_B,index_B] = simplexSolve(type,B,D,c_B,c_D,h,x_B,x_D,index_B,index_D,itMax);
24
25 % Compute the value of the objective function
26 z = round(c_B * x_B);
27
28 % Output of the solution
29 [x_B,index_B] = printSol(z,x_B,index_B,m,n);
30
31 end

```

standardize.m - which writes a maximisation or minimisation input problem in standard form.

```

1 function [A_aug,h,c_aug] = standardize(type,A,h,c,m,sign)
2 % return arguments are:
3 % (1) A_aug = augmented matrix A, containing also the surplus and slack variables
4 % (2) c_aug = augmented coefficients vector c (check compatibility of dimensions with A)
5 % (3) h, right hand side vector (remember to flip the signs when changing the
    inequalities)
6
7 aug_matrix = eye(m); % matrix corresponding to the slack/surplus variables
8
9 if(strcmp(type,'max'))
10     for i = 1:m
11         if(sign(i)==1)
12             % Using a surplus instead of a slack variable
13             aug_matrix(i,:) = -aug_matrix(i,:);
14         end
15     end
16 elseif(strcmp(type,'min'))
17     for i = 1:m
18         if(sign(i)==-1)
19             % Using a slack instead of a surplus variable
20             aug_matrix(i,:) = -aug_matrix(i,:);
21         end
22     end
23 else
24     error('Incorrect type specified. Choose either a maximisation (max) or minimisation
        (min) problem.')
25 end
26
27 % Correction on the sign of h for auxiliary problem (we want to ensure that h>=0, but we
    need to flip all the signs)
28 for i = 1:m
29     if(h(i)<0)
30         A(i,:) = -A(i,:);
31         h(i,:) = -h(i);
32         aug_matrix(i,:) = - aug_matrix(i,:);
33     end
34 end
35
36 c_aug = [c, zeros(1,m)];
37 if(strcmp(type,'max'))
38     % Extend matrix A by adding the slack variables
39     A_aug = [A, aug_matrix];
40 elseif(strcmp(type,'min'))

```

```

41 % Extend matrix A by adding the surplus variables
42 A_aug = [A, - aug_matrix];
43 else
44     error('Incorrect type specified. Choose either a maximisation (max) or minimisation
45         (min) problem.')
```

```

46 end
47 end
```

auxiliary.m - solves the auxiliary problem to find a feasible starting basic solution of the linear program.

```

1 function [B,D,c_B,c_D,x_B,x_D,index_B,index_D] = auxiliary (A_aug,c_aug,h,m,n)
2 % The auxiliary problem is always a minimization problem
3
4 % The output will be: B and D (basic and nonbasic matrices), c_B and c_D (subdivision of
5 % the coefficient vector in basic and nonbasic parts), x_B
6 % and x_D (basic and nonbasic variables) and index_B and index_D (to keep track of the
7 % variables indices)
8
9 % Redefine the problem by introducing the artificial variables required by
10 % the auxiliary problem (the objective function has to reach value 0)
11 A_aug = [A_aug, eye(m)];
12 c_aug = [c_aug, zeros(1,m)]; % original objective function coefficients
13 c_aux = [zeros(1,n+m), ones(1,m)]; % auxiliary c of minimization problem
14 index = 1:n+2*m; % index to keep track of the different variables
15
16 % Defining the basic elements and the nonbasic elements
17 B = A_aug(:,(n+m+1):(n+2*m)); % basic variables
18 D = A_aug(:,1:(n+m)); % nonbasic variables
19 c_Baux = c_aux(1,(n+m+1):(n+2*m));
20 c_Daux = c_aux(1,1:(n+m));
21 c_B = c_aug(1,(n+m+1):(n+2*m));
22 c_D = c_aug(1,1:(n+m));
23 x_B = h;
24 x_D = zeros((n+m),1);
25 index_B = index(1,(n+m+1):(n+2*m));
26 index_D = index(1,1:(n+m));
27
28 nIter = 0;
29 z = c_Baux*x_B;
30 itMax = factorial(2*m+n)/(factorial(n+m)*factorial(m));
31
32 % Compute  $B^{-1}D$  and  $B^{-1}h$ 
33 BiD = B\D;
34 Bih = B\h;
35
36 % Compute reduced cost coefficients
37 r_D = c_Daux - c_Baux*BiD;
38
39 while(z~=0)
40
41     % Find nonnegative index
42     idxIN = find(r_D==min(r_D));
43     % Using Bland's rule to avoid cycling
44     if(size(idxIN,2)>1)
45         idxIN = min(idxIN);
46     end
47     in = D(:,idxIN);
48     c_inaux = c_Daux(1,idxIN);
```

```

47     c_in = c_D(1,idxIN);
48     index_in = index_D(1,idxIN);
49
50     % Evaluating the coefficients ratio
51     inRatio = BiD(:,idxIN);
52     ratio = Bih./inRatio;
53
54     % Find the smallest ratio
55     for i = 1:size(ratio,1) % Eliminating negative ratios
56         if(ratio(i,1)<0)
57             ratio(i,1) = Inf;
58         end
59     end
60     idxOUT = find(ratio==min(ratio));
61     % Using Bland's rule to avoid cycling
62     if(size(idxOUT,1)>1)
63         idxOUT = min(idxOUT);
64     end
65     out = B(:,idxOUT);
66     c_outaux = c_Baux(1,idxOUT);
67     c_out = c_B(1,idxOUT);
68     index_out = index_B(1,idxOUT);
69
70     % Update the matrices by exchanging the columns
71     B(:,idxOUT) = in;
72     D(:,idxIN) = out;
73     c_Baux(1,idxOUT) = c_inaux;
74     c_Daux(1,idxIN) = c_outaux;
75     c_B(1,idxOUT) = c_in;
76     c_D(1,idxIN) = c_out;
77     index_B(1,idxOUT) = index_in;
78     index_D(1,idxIN) = index_out;
79
80     % Compute  $B^{-1}*D$  and  $B^{-1}*h$ 
81     BiD = B\D;
82     Bih = B\h;
83
84     % Compute reduced cost coefficients
85     r_D = c_Daux - c_Baux*BiD;
86
87     % Detect inefficient loop if nIter > total number of basic solutions
88     nIter = nIter + 1;
89     if(nIter>itMax)
90         error('The original LP problem does not admit a feasible solution.');
```

```

91     end
92
93     x_B = Bih - BiD*x_D;
94     z = c_Baux*x_B;
95
96 end
97
98 check = index_D<(n+m+1);
99 D = D(:,check);
100 index_D = index_D(1,check);
101 c_D = c_D(1,check);
102 x_D = x_D(check,1);
103
104 end

```

6. Project 5 - Applications to Real-life Example Cargo Aircraft

7. Linear program, objective function.

Profit Calculations:

For Cargo C_1 (135 CHF/t) :

$$p_{11} = 135 \cdot 1 = 135, \quad p_{12} = 135 \cdot 1.1 = 148.5,$$

$$p_{13} = 135 \cdot 1.2 = 162, \quad p_{14} = 135 \cdot 1.3 = 175.5$$

For Cargo C_2 (200 CHF/t) :

$$p_{21} = 200 \cdot 1 = 200, \quad p_{22} = 200 \cdot 1.1 = 220,$$

$$p_{23} = 200 \cdot 1.2 = 240, \quad p_{24} = 200 \cdot 1.3 = 260$$

For Cargo C_3 (410 CHF/t) :

$$p_{31} = 410 \cdot 1 = 410, \quad p_{32} = 410 \cdot 1.1 = 451,$$

$$p_{33} = 410 \cdot 1.2 = 492, \quad p_{34} = 410 \cdot 1.3 = 533$$

For Cargo C_4 (520 CHF/t) :

$$p_{41} = 520 \cdot 1 = 520, \quad p_{42} = 520 \cdot 1.1 = 572,$$

$$p_{43} = 520 \cdot 1.2 = 624, \quad p_{44} = 520 \cdot 1.3 = 676$$

Objective Function:

$$\begin{aligned} \text{Maximize } Z = & 135x_{11} + 148.5x_{12} + 162x_{13} + 175.5x_{14} + \\ & 200x_{21} + 220x_{22} + 240x_{23} + 260x_{24} + \\ & 410x_{31} + 451x_{32} + 492x_{33} + 533x_{34} + \\ & 520x_{41} + 572x_{42} + 624x_{43} + 676x_{44} \end{aligned}$$

Constraint - Compartment Weight Capacity:

$$x_{11} + x_{21} + x_{31} + x_{41} \leq 18 \quad (\text{Capacity of } S_1)$$

$$x_{12} + x_{22} + x_{32} + x_{42} \leq 32 \quad (\text{Capacity of } S_2)$$

$$x_{13} + x_{23} + x_{33} + x_{43} \leq 25 \quad (\text{Capacity of } S_3)$$

$$x_{14} + x_{24} + x_{34} + x_{44} \leq 17 \quad (\text{Capacity of } S_4)$$

Constraint - Compartment Volume Capacity:

$$320x_{11} + 510x_{21} + 630x_{31} + 125x_{41} \leq 11930 \quad (\text{Volume Capacity of } S_1)$$

$$320x_{12} + 510x_{22} + 630x_{32} + 125x_{42} \leq 22552 \quad (\text{Volume Capacity of } S_2)$$

$$320x_{13} + 510x_{23} + 630x_{33} + 125x_{43} \leq 11209 \quad (\text{Volume Capacity of } S_3)$$

$$320x_{14} + 510x_{24} + 630x_{34} + 125x_{44} \leq 5870 \quad (\text{Volume Capacity of } S_4)$$

Constraint - Cargo Weight Availability:

$$x_{11} + x_{12} + x_{13} + x_{14} \leq 16 \quad (\text{Availability of Cargo } C_1)$$

$$x_{21} + x_{22} + x_{23} + x_{24} \leq 32 \quad (\text{Availability of Cargo } C_2)$$

$$x_{31} + x_{32} + x_{33} + x_{34} \leq 40 \quad (\text{Availability of Cargo } C_3)$$

$$x_{41} + x_{42} + x_{43} + x_{44} \leq 28 \quad (\text{Availability of Cargo } C_4)$$

Constraint - Non-Negativity:

$$x_{ij} \geq 0 \quad \text{for all } i = 1, 2, 3, 4 \text{ and } j = 1, 2, 3, 4.$$

Calculated the adjusted profit for each type of cargo when placed in different compartments, the profit per ton for each cargo increases depending on the storage compartment.

The objective function, denoted as Z is the maximum total profit to maximize. It is the sum of

products of the profit per ton of each cargo in each compartment and the corresponding amount of cargo placed in that compartment:

$$\sum_{i=1}^4 \sum_{j=1}^4 p_{ij} x_{ij}$$

Where p_{ij} is the profit per ton for Cargo i in Compartment j and x_{ij} is the total of Cargo i in Compartment j .