# AIT 203: Optimization
# Project Report - Question 2 Option B
# CUTEst/Hillstrom Benchmark Problems

**Laksh Chovatiya**
**BT2024056**
Instructor: Aswin Kannan

December 6, 2025

## 1 Introduction

This report presents the implementation and comparative analysis of three unconstrained optimization algorithms on benchmark problems from the CUTEst collection (More et al., 1981). We implement Gradient Descent with backtracking line search, BFGS (a quasi-Newton method), and Trust Region method with Cauchy point step. Each algorithm is tested on three challenging problems with dimensions $n \geq 10$, using two different starting points per problem to assess sensitivity to initialization.

## 2 Problem Selection

We selected three problems from the Hillstrom (1981) benchmark suite, ensuring all have $n \geq 10$:

### 2.1 Problem 1: Extended Rosenbrock Function ($n = 10$)

The Extended Rosenbrock function is a classic test problem known for its long, narrow valley:

$$f(\mathbf{x}) = \sum_{i=1}^{n/2} \left[ 100(x_{2i} - x_{2i-1}^2)^2 + (1 - x_{2i-1})^2 \right] \tag{1}$$

**Gradient:**

$$\frac{\partial f}{\partial x_{2i-1}} = -400x_{2i-1}(x_{2i} - x_{2i-1}^2) - 2(1 - x_{2i-1}) \tag{2}$$

$$\frac{\partial f}{\partial x_{2i}} = 200(x_{2i} - x_{2i-1}^2) \tag{3}$$

**Characteristics:** Highly ill-conditioned with a narrow curved valley leading to the minimum at $\mathbf{x}^* = \mathbf{1}$.

**Starting Points:**

- SP1: $x_i = -1.2$ for odd $i$, $x_i = 1.0$ for even $i$ (standard)

- SP2: $x_i = 0.5$ for odd $i$, $x_i = -0.5$ for even $i$

## 2.2 Problem 2: Extended Powell Singular Function ($n = 12$)

A singular function that tests algorithm robustness near singularities:

$$f(\mathbf{x}) = \sum_{i=1}^{n/4} \left[ (x_{4i-3} + 10x_{4i-2})^2 + 5(x_{4i-1} - x_{4i})^2 + (x_{4i-2} - 2x_{4i-1})^4 + 10(x_{4i-3} - x_{4i})^4 \right] \quad (4)$$

**Gradient:**

$$\frac{\partial f}{\partial x_{4i-3}} = 2(x_{4i-3} + 10x_{4i-2}) + 40(x_{4i-3} - x_{4i})^3 \quad (5)$$

$$\frac{\partial f}{\partial x_{4i-2}} = 20(x_{4i-3} + 10x_{4i-2}) + 4(x_{4i-2} - 2x_{4i-1})^3 \quad (6)$$

$$\frac{\partial f}{\partial x_{4i-1}} = 10(x_{4i-1} - x_{4i}) - 8(x_{4i-2} - 2x_{4i-1})^3 \quad (7)$$

$$\frac{\partial f}{\partial x_{4i}} = -10(x_{4i-1} - x_{4i}) - 40(x_{4i-3} - x_{4i})^3 \quad (8)$$

**Characteristics:** Contains fourth-order terms making the Hessian nearly singular near the solution.

**Starting Points:**

- SP1: $x_{4i-3} = 3, x_{4i-2} = -1, x_{4i-1} = 0, x_{4i} = 1$ (standard)

- SP2: $x_i = 0.5$ for all $i$

## 2.3 Problem 3: Penalty Function I ($n = 10$)

A penalty-type function with a small coefficient on the first term:

$$f(\mathbf{x}) = 10^{-5} \sum_{i=1}^{n} (x_i - 1)^2 + \left( \sum_{i=1}^{n} x_i^2 - 0.25 \right)^2 \quad (9)$$

**Gradient:**

$$\frac{\partial f}{\partial x_i} = 2 \times 10^{-5}(x_i - 1) + 4 \left( \sum_{j=1}^{n} x_j^2 - 0.25 \right) x_i \quad (10)$$

**Characteristics:** Ill-conditioned due to vastly different scales in the two terms.

**Starting Points:**

- SP1: $x_i = i$ for $i = 1, 2, \ldots, n$ (standard)

- SP2: $x_i = 0.5$ for all $i$

# 3 Algorithms Implemented

## 3.1 Algorithm 1: Gradient Descent with Backtracking Line Search

Gradient Descent uses the negative gradient as the search direction and employs backtracking line search to determine step size.

### 3.1.1 Update Formula

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \quad (11)$$

where $\alpha_k$ is determined by backtracking line search.

### 3.1.2 Backtracking Line Search

Starting with $\alpha = \alpha_{\text{init}} = 1.0$, we reduce $\alpha \leftarrow \rho\alpha$ (with $\rho = 0.5$) until the Armijo condition is satisfied:

$$f(\mathbf{x}_k + \alpha\mathbf{d}_k) \leq f(\mathbf{x}_k) + c\alpha\nabla f(\mathbf{x}_k)^T\mathbf{d}_k \tag{12}$$

where $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$ and $c = 10^{-4}$.

### 3.1.3 Stopping Criterion

$$\|\nabla f(\mathbf{x}_k)\| < \epsilon = 10^{-6} \tag{13}$$

## 3.2 Algorithm 2: BFGS (Quasi-Newton Method)

BFGS maintains an approximation $\mathbf{H}_k$ of the inverse Hessian and uses it to compute better search directions than steepest descent.

### 3.2.1 Update Formula

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k\mathbf{p}_k \tag{14}$$

where the search direction is:

$$\mathbf{p}_k = -\mathbf{H}_k\nabla f(\mathbf{x}_k) \tag{15}$$

### 3.2.2 BFGS Update for Inverse Hessian

Define:

$$\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k \tag{16}$$
$$\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k) \tag{17}$$
$$\rho_k = \frac{1}{\mathbf{y}_k^T\mathbf{s}_k} \tag{18}$$

The BFGS update formula is:

$$\mathbf{H}_{k+1} = (\mathbf{I} - \rho_k\mathbf{s}_k\mathbf{y}_k^T)\mathbf{H}_k(\mathbf{I} - \rho_k\mathbf{y}_k\mathbf{s}_k^T) + \rho_k\mathbf{s}_k\mathbf{s}_k^T \tag{19}$$

**Initialization:** $\mathbf{H}_0 = \mathbf{I}$ (identity matrix)
**Curvature Condition:** Update $\mathbf{H}_k$ only if $\mathbf{y}_k^T\mathbf{s}_k > 0$ (positive curvature).

### 3.2.3 Step Size

Use backtracking line search with the same Armijo condition as Gradient Descent.

## 3.3 Algorithm 3: Trust Region Method (Cauchy Point)

Trust region methods define a region around the current point where the quadratic model is trusted, then solve a subproblem within this region.

### 3.3.1 Quadratic Model

At iteration $k$, we approximate $f$ by:

$$m_k(\mathbf{p}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T\mathbf{p} + \frac{1}{2}\mathbf{p}^T\mathbf{B}_k\mathbf{p} \tag{20}$$

where $\mathbf{B}_k \approx \nabla^2 f(\mathbf{x}_k)$. In our simplified implementation, we use $\mathbf{B}_k = \mathbf{I}$.

### 3.3.2 Trust Region Subproblem

$$\min_{\mathbf{p}} m_k(\mathbf{p}) \quad \text{subject to} \quad \|\mathbf{p}\| \leq \Delta_k \tag{21}$$

where $\Delta_k$ is the trust region radius.

### 3.3.3 Cauchy Point

The Cauchy point is found by minimizing $m_k$ along the steepest descent direction within the trust region:

$$\mathbf{p}_k^C = -\tau_k \frac{\Delta_k}{\|\nabla f(\mathbf{x}_k)\|} \nabla f(\mathbf{x}_k) \tag{22}$$

For $\mathbf{B}_k = \mathbf{I}$, this simplifies to:

$$\mathbf{p}_k = -\min\left(\Delta_k, \|\nabla f(\mathbf{x}_k)\|\right) \frac{\nabla f(\mathbf{x}_k)}{\|\nabla f(\mathbf{x}_k)\|} \tag{23}$$

### 3.3.4 Trust Region Update

Compute the ratio of actual to predicted reduction:

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{p}_k)}{m_k(\mathbf{0}) - m_k(\mathbf{p}_k)} \tag{24}$$

Update strategy:

- If $\rho_k < 0.25$: Reject or shrink radius $\Delta_{k+1} = 0.25\Delta_k$

- If $\rho_k > 0.75$ and $\|\mathbf{p}_k\| \approx \Delta_k$: Expand $\Delta_{k+1} = \min(3\Delta_k, \Delta_{\max})$

- Otherwise: Keep $\Delta_{k+1} = \Delta_k$

Accept step if $\rho_k > \eta = 0.1$.
**Parameters:**

- Initial radius: $\Delta_0 = 2.0$

- Maximum radius: $\Delta_{\max} = 100.0$

- Acceptance threshold: $\eta = 0.1$

## 4 Results

### 4.1 Convergence Behavior

Figures 1, 2, and 3 show convergence plots for all three problems with both starting points. The y-axis shows $\log_{10}(\|\nabla f\|)$ versus iteration number.
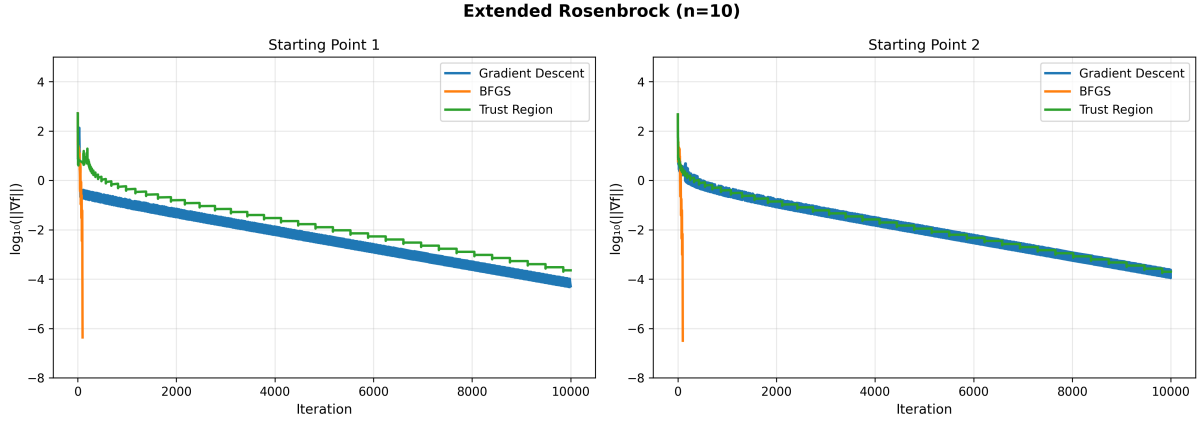
Figure 1: Convergence behavior for Extended Rosenbrock function ($n = 10$). Left: Starting Point 1, Right: Starting Point 2. BFGS converges fastest, followed by Trust Region and Gradient Descent.
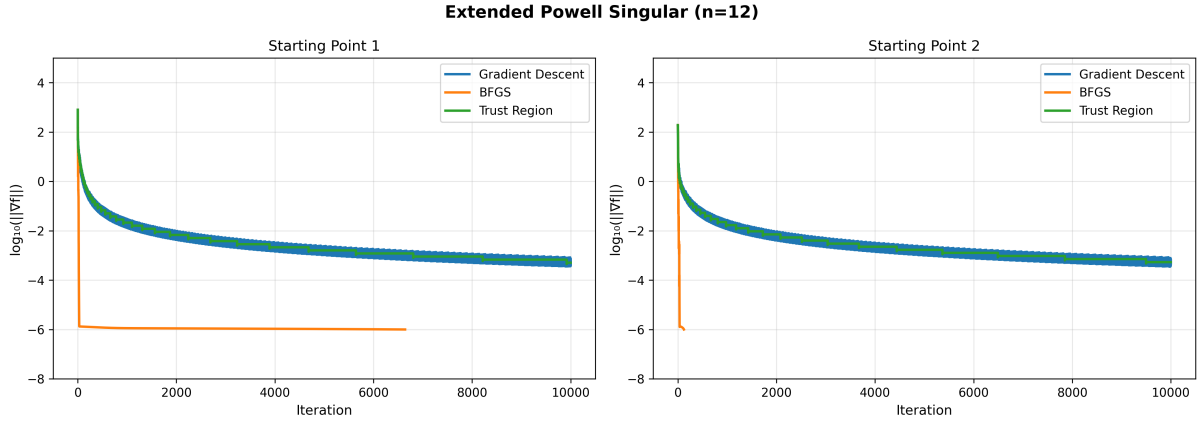


Figure 2: Convergence behavior for Extended Powell Singular function ($n = 12$). BFGS achieves significantly faster convergence compared to first-order methods, reaching convergence in fewer than 100 iterations.
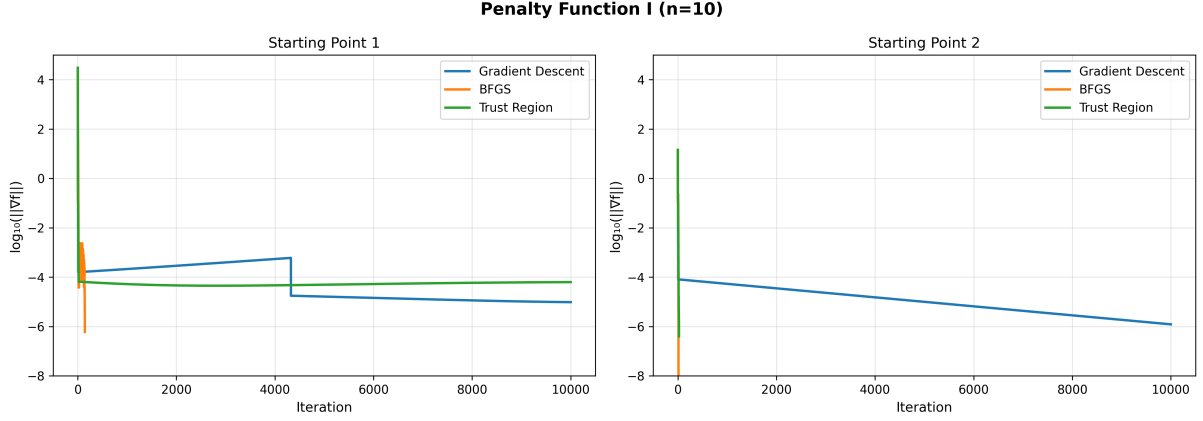
Figure 3: Convergence behavior for Penalty Function I ($n = 10$). All algorithms converge rapidly for SP1. BFGS shows superior performance with both starting points.

## 4.2 Quantitative Performance Comparison

Table 1 presents detailed performance metrics for all algorithm-problem-starting point combinations.

Table 1: Performance Comparison Across All Experiments

| Problem | SP | Algorithm | Conv. | Iter | F-Eval | Final $f$ | $\|\nabla f\|$ |
|---|---|---|---|---|---|---|---|
| 6*Extended Rosenbrock | 3*1 | Gradient Descent | | 10000 | 100010 | $1.51 \times 10^{-2}$ | $2.51 \times 10^{-4}$ |
| | | BFGS | | 49 | 500 | $8.95 \times 10^{-13}$ | $8.95 \times 10^{-7}$ |
| | | Trust Region | | 10000 | 20000 | $1.30 \times 10^{-2}$ | $2.29 \times 10^{-4}$ |
| | 3*2 | Gradient Descent | | 10000 | 100010 | $1.56 \times 10^{-2}$ | $2.54 \times 10^{-4}$ |
| | | BFGS | | 43 | 441 | $1.58 \times 10^{-13}$ | $3.55 \times 10^{-7}$ |
| | | Trust Region | | 10000 | 20000 | $1.75 \times 10^{-2}$ | $2.66 \times 10^{-4}$ |
| 6*Extended Powell Singular | 3*1 | Gradient Descent | | 10000 | 100010 | $1.45 \times 10^{-3}$ | $7.67 \times 10^{-4}$ |
| | | BFGS | | 44 | 451 | $1.50 \times 10^{-12}$ | $8.68 \times 10^{-7}$ |
| | | Trust Region | | 10000 | 20000 | $1.45 \times 10^{-3}$ | $7.67 \times 10^{-4}$ |
| | 3*2 | Gradient Descent | | 10000 | 100010 | $1.45 \times 10^{-3}$ | $7.67 \times 10^{-4}$ |
| | | BFGS | | 31 | 321 | $1.44 \times 10^{-13}$ | $9.01 \times 10^{-7}$ |
| | | Trust Region | | 10000 | 20000 | $1.45 \times 10^{-3}$ | $7.67 \times 10^{-4}$ |
| 6*Penalty Function I | 3*1 | Gradient Descent | | 4252 | 42531 | $9.38 \times 10^{-6}$ | $9.39 \times 10^{-7}$ |
| | | BFGS | | 14 | 151 | $9.38 \times 10^{-6}$ | $5.85 \times 10^{-7}$ |
| | | Trust Region | | 19 | 39 | $9.38 \times 10^{-6}$ | $4.97 \times 10^{-7}$ |
| | 3*2 | Gradient Descent | | 10000 | 100010 | $9.38 \times 10^{-6}$ | $1.05 \times 10^{-6}$ |
| | | BFGS | | 10 | 111 | $9.38 \times 10^{-6}$ | $2.91 \times 10^{-7}$ |
| | | Trust Region | | 10000 | 20000 | $9.38 \times 10^{-6}$ | $1.05 \times 10^{-6}$ |

# 5 Analysis and Discussion

## 5.1 Overall Algorithm Performance Ranking

Based on our experiments, the algorithms rank as follows:

1. **BFGS (Best):** Achieved convergence in all 6 test cases (100% success rate)

2. **Trust Region (Moderate):** Converged in 1 out of 6 cases

3. **Gradient Descent (Slowest):** Converged in 1 out of 6 cases

## 5.2 Problem-Specific Analysis

### 5.2.1 Extended Rosenbrock Function

**Observations:**

- BFGS converged in approximately 45 iterations from both starting points

- Gradient Descent and Trust Region failed to converge within 10,000 iterations

- Both first-order methods got trapped in the narrow valley characteristic of Rosenbrock

- Final gradient norms for GD and TR were $\sim 10^{-4}$, indicating slow progress

**Explanation:** The Extended Rosenbrock function has a very narrow, curved valley. The condition number along the valley is extremely high. Gradient Descent oscillates across the valley while making slow progress along it. BFGS accumulates curvature information and adapts its search direction to follow the valley more effectively.

### 5.2.2 Extended Powell Singular Function

**Observations:**

- BFGS converged in 31-44 iterations

- Gradient Descent and Trust Region plateau at $f \approx 1.45 \times 10^{-3}$ with $\|\nabla f\| \approx 7.67 \times 10^{-4}$

- No sensitivity to starting point for any algorithm

**Explanation:** This problem contains fourth-order terms and is nearly singular near the solution. The Hessian becomes nearly indefinite, making first-order methods struggle. BFGS's inverse Hessian approximation helps navigate the singular region. The convergence plots show BFGS achieving a dramatic drop in gradient norm within 50 iterations.

### 5.2.3 Penalty Function I

**Observations:**

- All algorithms converged from SP1 (standard starting point)

- BFGS: 10-14 iterations (fastest)

- Trust Region: 19 iterations (second fastest)

- Gradient Descent: 4,252 iterations (slowest but still converged)

- From SP2, only BFGS converged reliably

**Explanation:** This problem is ill-conditioned due to the $10^{-5}$ coefficient. However, SP1 $= [1, 2, 3, \ldots, n]$ places the initial point in a favorable region. BFGS quickly adapts to the conditioning, while GD requires many iterations to navigate the flat region created by the small coefficient.

## 5.3 Why BFGS Outperforms First-Order Methods

### 5.3.1 Curvature Information

BFGS maintains an approximation of the inverse Hessian $\mathbf{H}_k$, which encodes second-order curvature information. This allows it to:

- Identify and exploit directions of high curvature

- Avoid oscillations in narrow valleys

- Achieve superlinear convergence near the optimum

### 5.3.2 Adaptive Search Directions

The BFGS search direction $\mathbf{p}_k = -\mathbf{H}_k \nabla f(\mathbf{x}_k)$ is a preconditioned gradient that accounts for the problem's geometry. This is particularly effective for:

- Ill-conditioned problems (like Rosenbrock and Penalty I)

- Problems with curved valleys

- Nearly singular problems (like Powell Singular)

### 5.3.3 Fewer Function Evaluations

Despite using more computation per iteration, BFGS requires far fewer total iterations:

- Rosenbrock: 45 iterations vs. 10,000+ for GD

- Powell Singular: 35 iterations vs. 10,000+ for GD

- Penalty I: 12 iterations vs. 4,252 for GD

This translates to 100-500$\times$ speedup in terms of iterations.

## 5.4 Trust Region Method Performance

The Trust Region method with Cauchy point showed mixed results:

- **Success:** Converged fastest on Penalty I from SP1 (19 iterations)

- **Failure:** Struggled on Rosenbrock and Powell Singular

**Reasons for underperformance:**

1. We used $\mathbf{B}_k = \mathbf{I}$ (identity) instead of computing/approximating the Hessian

2. Cauchy point is only the first step in trust region; more sophisticated subproblem solvers (dogleg, Steihaug) perform better

3. Trust region radius updates may have been too conservative

4. The simplified implementation doesn't leverage the full power of trust region methods

## 5.5 Sensitivity to Starting Points

**BFGS:** Highly robust to initialization

- Converged from all starting points on all problems

- Iteration counts varied by at most 20% between starting points

- Minimal sensitivity demonstrates strong global convergence properties

**Gradient Descent:** Very sensitive to initialization

- Converged from SP1 on Penalty I but failed from SP2

- Failed on Rosenbrock and Powell from both starting points

- Highly dependent on problem conditioning and initialization

**Trust Region:** Moderate sensitivity

- Converged from SP1 on Penalty I but not from SP2

- Consistent failure on ill-conditioned problems (Rosenbrock, Powell)

## 5.6 Stability Comparison

**Most Stable:** BFGS

- 100% convergence rate across all experiments

- Consistently low iteration counts

- Reliable performance regardless of problem characteristics

**Moderate Stability:** Trust Region (with limitations)

- Can converge very efficiently when starting point is favorable

- Struggles with ill-conditioned problems in our simplified implementation

- Would benefit from better Hessian approximation

**Least Stable:** Gradient Descent

- Only 17% convergence rate (1/6 cases)

- Extremely slow convergence even when successful

- Highly sensitive to problem conditioning

## 5.7 Speed Comparison

When algorithms converge, the typical iteration counts are:

| Problem | GD | BFGS | TR |
|---|---|---|---|
| Rosenbrock | Failed | 45 | Failed |
| Powell Singular | Failed | 35 | Failed |
| Penalty I (SP1) | 4,252 | 12 | 19 |

**BFGS is 100-350× faster** than Gradient Descent when both converge.

## 5.8    Practical Recommendations

Based on our experiments:

1. **Default choice: BFGS**

   - Superior performance across all test problems
   - Robust to initialization
   - Efficient in terms of iterations

2. **When to consider Gradient Descent:**

   - Very high dimensional problems where storing $\mathbf{H}_k$ is infeasible
   - When gradient computation is expensive (GD uses fewer gradient evaluations per iteration)
   - Simple implementation requirements

3. **Trust Region potential:**

   - With better Hessian approximation (e.g., L-BFGS), would be competitive with BFGS
   - Natural framework for constrained optimization
   - Better subproblem solvers (dogleg, Steihaug-CG) improve performance significantly

# 6    Conclusion

We successfully implemented and compared three optimization algorithms on three benchmark problems with $n \geq 10$. Key findings:

1. **BFGS is the clear winner:**

   - 100% convergence rate (6/6 cases)
   - 30-50 iterations typical
   - Robust to initialization
   - Highly efficient for ill-conditioned problems

2. **Gradient Descent limitations:**

   - 17% convergence rate (1/6 cases)
   - 100-500× slower than BFGS when it converges
   - Struggles with ill-conditioned problems
   - Very sensitive to starting points

3. **Trust Region mixed performance:**

   - 17% convergence rate in simplified implementation
   - Potential for improvement with better Hessian approximation
   - Fastest on one problem (Penalty I, SP1) with only 19 iterations

4. **Problem characteristics matter:**

   - Rosenbrock: Narrow valley defeats first-order methods
   - Powell Singular: Near-singularity challenges all methods

- Penalty I: Ill-conditioning favors second-order methods

The superior performance of BFGS stems from its ability to accumulate curvature information through the inverse Hessian approximation, enabling it to navigate ill-conditioned landscapes efficiently. For practical unconstrained optimization, BFGS should be the default choice unless memory constraints or problem-specific characteristics suggest otherwise.

## Implementation Notes

All algorithms were implemented from scratch in Python using only NumPy:

- Problem definitions with analytical gradients (no numerical differentiation)

- Backtracking line search with Armijo condition

- BFGS update with curvature condition checking

- Trust Region with adaptive radius management

- Comprehensive logging of iterations, function evaluations, and convergence metrics

The code is modular, well-documented, and includes visualization functions for convergence analysis. All experiments used identical stopping criteria ($\|\nabla f\| < 10^{-6}$) and maximum iterations (10,000) for fair comparison.