

A Performance Comparison between Isolation Forest and K-Means Isolation Forest on Anomaly Detection Tasks

Răzvan-Nicolae Căpățină Mihnea-Vicențiu Bucă



Faculty of Mathematics and Informatics
University of Bucharest

February 2, 2026

- 1 Isolation Forest (IF)
- 2 k-Means-Based Isolation Forest (kMIF)
- 3 Elbow Rule
- 4 Parallelization
- 5 Experiments: Synthetic Data
- 6 Experiments: Real-World Data
- 7 Future Work

Isolation Forest (IF) - Implementation

Core Idea

Anomalies are **few and different** → easier to isolate

Tree Construction (Binary Splitting)

- ① **Feature Selection:** Choose random dimension q (uniform)
- ② **Split Value:** Choose random value $v \in [\min_q, \max_q]$
- ③ **Partition:**
 - Left child: $D_L = \{x \mid x_q < v\}$
 - Right child: $D_R = \{x \mid x_q \geq v\}$
- ④ **Recurse** until termination conditions

Termination Conditions

- Maximum depth / reached
- Node contains ≤ 1 sample
- All samples identical

Isolation Forest - Anomaly Score

Expected Path Length

$$h(x) = \begin{cases} d(x) & \text{if } x \text{ is a natural leaf (not max-depth blocked)} \\ d(x) + c(m(x)) & \text{if } x \text{ is a leaf due to max-depth pruning} \end{cases}$$

$m(x)$ = # of trained samples for leaf x

Anomaly Score Formula

$$s(x, n) = 2^{-\frac{\mathbb{E}[h(x)]}{c(n)}} \quad \text{where} \quad c(n) = 2[\ln(n-1) + 0.5772] - \frac{2(n-1)}{n}$$

- $\mathbb{E}[h(x)]$ = mean path length across all trees
- $c(n)$ = normalization factor (avg. BST search length)
- 0.5772 = Euler-Mascheroni constant

Interpretation

Score $\approx 1 \rightarrow$ Anomaly — Score $\approx 0 \rightarrow$ Normal



Isolation Forest - Ensemble

Ensemble Construction

Build t independent Isolation Trees, each on a random subsample

Recommended Parameters (Original Paper)

- **Sub-sampling size:** $n = 256$
- **Maximum depth:** $l = \lceil \log_2 n \rceil = 8$
- **Ensemble size:** $t = 100$ trees

Anomaly Score Aggregation

- ① For each tree i : compute path length $h_i(x)$
- ② Average across ensemble: $\bar{h}(x) = \frac{1}{t} \sum_{i=1}^t h_i(x)$
- ③ Final score: $s(x) = 2^{-\bar{h}(x)/c(n)}$

Threshold

Use contamination rate or default threshold = 0.5

Outline

- 1 Isolation Forest (IF)
- 2 k-Means-Based Isolation Forest (kMIF)
- 3 Elbow Rule
- 4 Parallelization
- 5 Experiments: Synthetic Data
- 6 Experiments: Real-World Data
- 7 Future Work

k -Means-Based Isolation Forest (kMIF)

Key Modification

Replace **binary splits** with **multi-branch k -Means clustering**

Tree Construction

- ① **Feature Selection:** Choose random dimension q (same as IF)
- ② **k -Means Clustering:** Apply to values along dimension q
- ③ **Optimal k :** Determined dynamically via **Elbow Rule**
- ④ **Partition:** Split into k subsets (children nodes)
- ⑤ **Store:** Cluster centers and radii for scoring

Advantages

- **Data-aware partitioning** adapts to local density
- More children in dense regions, fewer in sparse regions
- Theoretically better fit to complex data structures

Normality Score at Each Node

$$s_j(x) = \max \left(0, 1 - \frac{d(x, c_q)}{d(c_l, c_q)} \right) \in [0, 1]$$

- $d(x, c_q)$ = distance from x to its cluster center
- $d(c_l, c_q)$ = cluster radius (distance to boundary)
- Clamped to $[0, 1]$ for out-of-sample data

Final Anomaly Score

$$\alpha(x) = 1 - \frac{1}{t} \sum_{i=1}^t \sum_{j=1}^{M_i} s_j(x)$$

M_i = path depth in tree i , score depends on **cluster distance**, not path length

Ensemble Construction

Build t multi-branch trees, each using k -Means splits

Parameters

- Max depth: $l = 9$
- Sub-sample: $n = 256$
- Trees: $t = 100$
- k candidates: $\{2, 3, 5, 7\}$

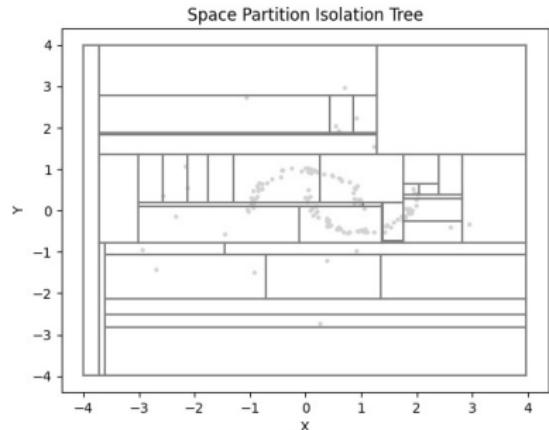
Score Aggregation

- ➊ Sum normality scores per tree
- ➋ Average across ensemble
- ➌ Anomaly = $1 - \text{avg. normality}$

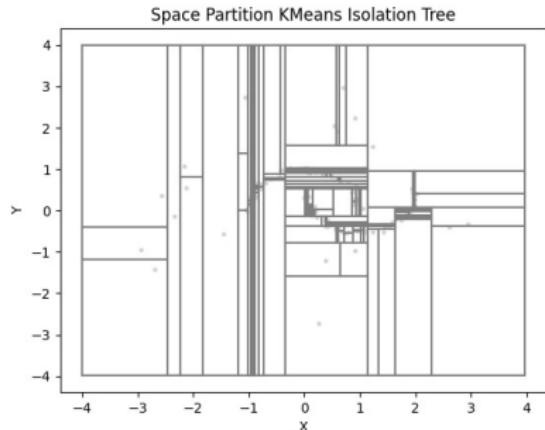
Key Difference from IF

Score depends more on **cluster membership distance** than path length

ITree vs k-Means ITree



ITree Partition on Two Moons Dataset



k-Means ITree Partition on Two Moons Dataset

Outline

1 Isolation Forest (IF)

2 k-Means-Based Isolation Forest (kMIF)

3 Elbow Rule

4 Parallelization

5 Experiments: Synthetic Data

6 Experiments: Real-World Data

7 Future Work

Elbow Rule for Optimal k

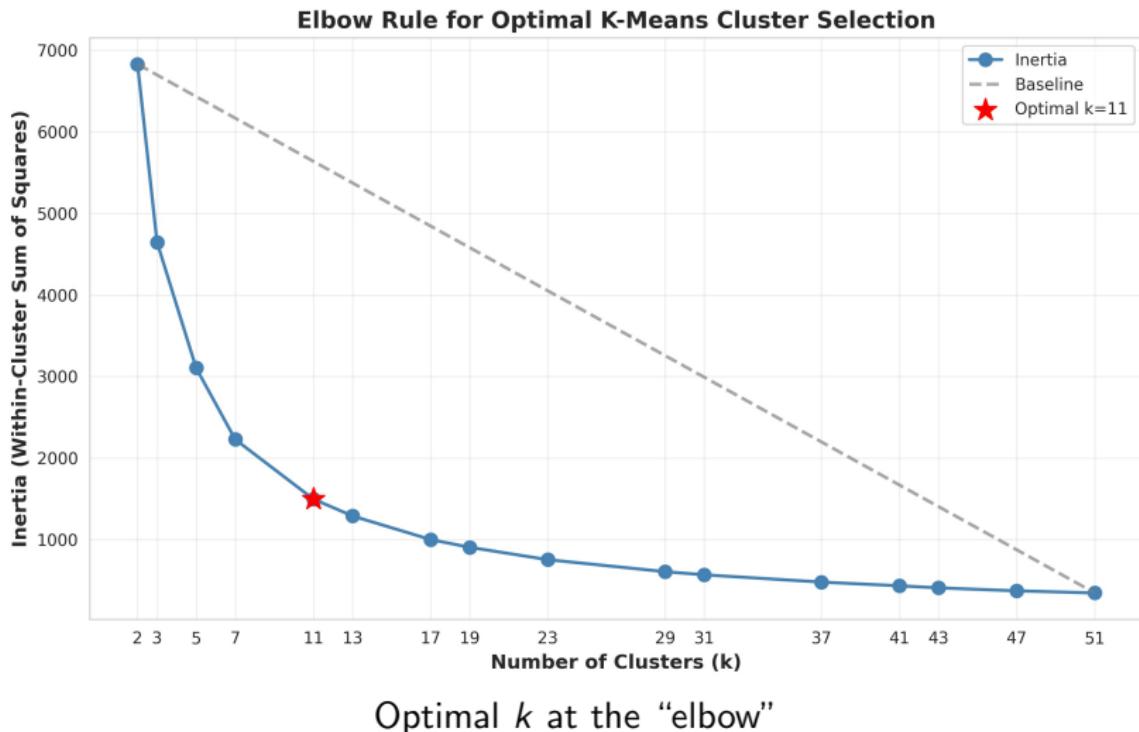
How It Works

- ① Run k -Means for $k \in \{2, 3, 5, 7\}$
- ② Track **WCSS** (Within-Cluster Sum of Squares)
- ③ Find the “elbow” point

Finding the Elbow

- Draw baseline from first to last point
- Select k with max perpendicular distance

Elbow Rule for Optimal k



Cost

k -Means at **every node** = major overhead

Outline

- 1 Isolation Forest (IF)
- 2 k-Means-Based Isolation Forest (kMIF)
- 3 Elbow Rule
- 4 Parallelization
- 5 Experiments: Synthetic Data
- 6 Experiments: Real-World Data
- 7 Future Work

Parallelization Strategy

Key Insight

Each tree in the ensemble is **independent** → parallelizable

Training Phase

- Distribute t trees across cores
- Independent tree construction
- Unique seeds for reproducibility

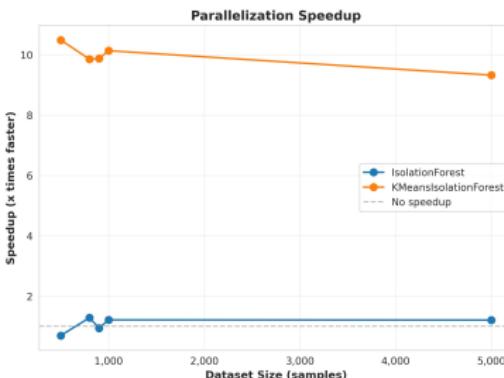
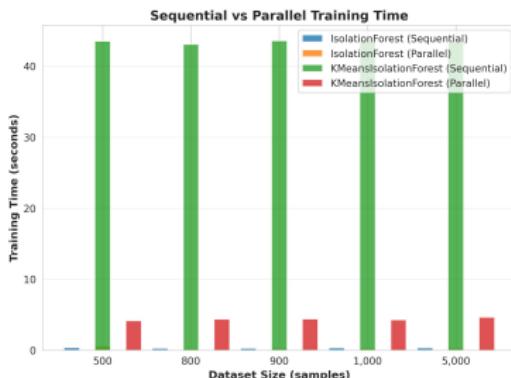
Scoring Phase

- Distribute scoring across cores
- Process tree subsets in parallel
- Aggregate scores at the end

Why Essential for kMIF?

IF: $O(1)$ node split → already fast — kMIF: $O(P \cdot k \cdot n')$ → needs parallelization

Parallelization Results



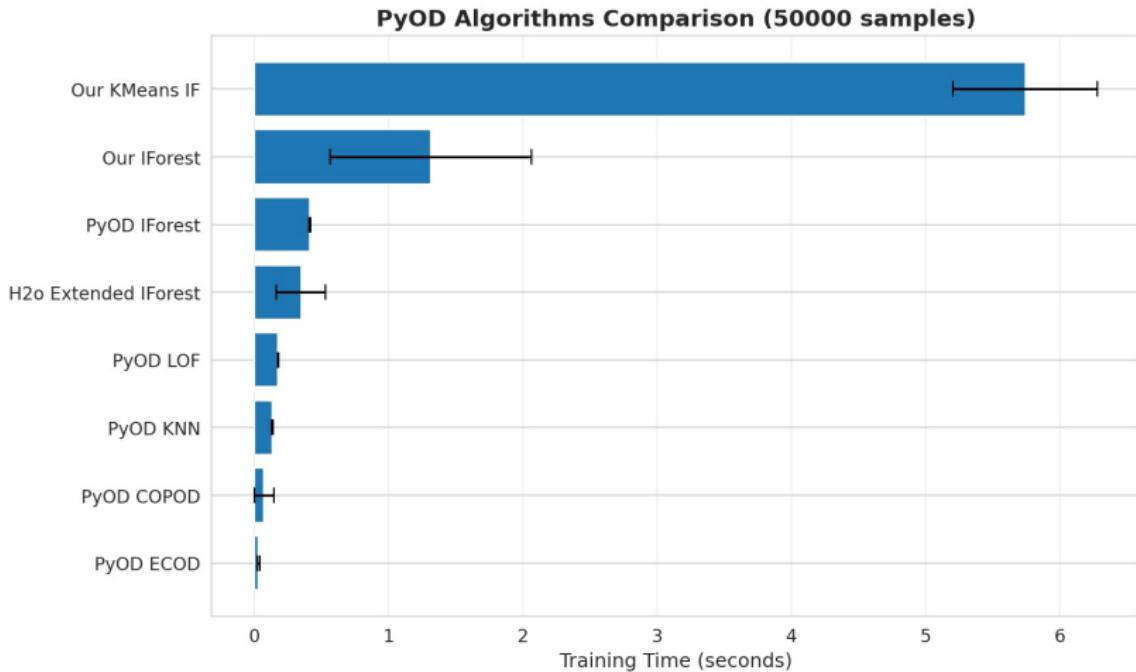
kMIF Speedup

- Sequential: > 43s
- Parallel: ~ 4s
- **9-10x speedup**

IF Speedup

- Already fast (< 0.4s)
- Minimal parallel benefit
- Overhead > gains

Parallelization Comparison



Training time comparison of our custom algorithms against H2O EIF and various PyOD models on 50,000 samples.

Outline

1 Isolation Forest (IF)

2 k-Means-Based Isolation Forest (kMIF)

3 Elbow Rule

4 Parallelization

5 Experiments: Synthetic Data

6 Experiments: Real-World Data

7 Future Work

Experimental Setup

- **Ensemble:** $t = 100$ trees, $n = 256$ samples, $l = 9$ depth
- **Metrics:** Accuracy, F1-Score, Precision, Recall, ROC-AUC, PR-AUC

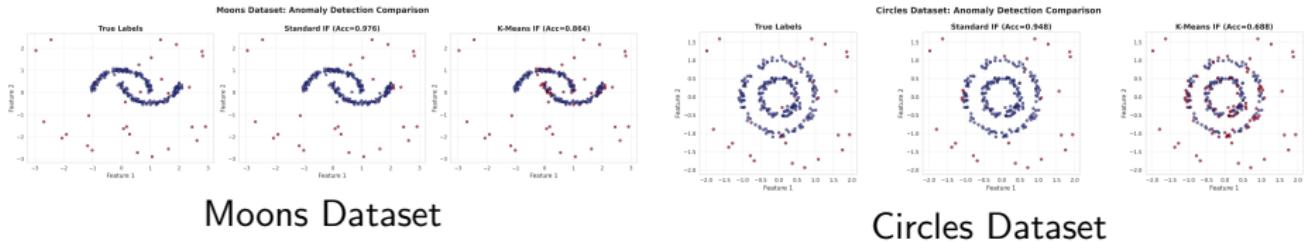
Synthetic Datasets Tested

- ① **9-Rectangle Dataset** - Rectilinear, axis-aligned clusters
- ② **Circles Dataset** - Concentric, non-convex clusters
- ③ **Moons Dataset** - Interlocking, non-convex manifolds
- ④ **S-Curve Dataset** - Complex, high-dimensional manifold
- ⑤ **High-Dim Synthetic** - 10D and 20D datasets

Goal

Test how each model handles different geometric challenges

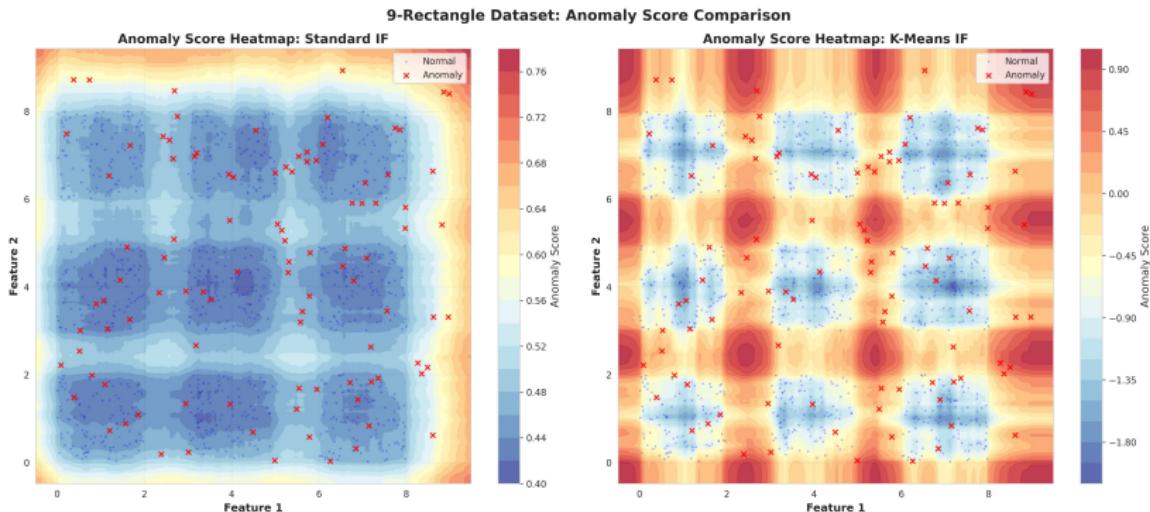
Synthetic Data Results



Key Findings

- **IF** outperforms kMIF on non-convex data (Moons F1: 0.88 vs 0.51)
- **kMIF** misclassifies dense core points as anomalies
- Both perform well on high-dimensional data

Synthetic Data Results



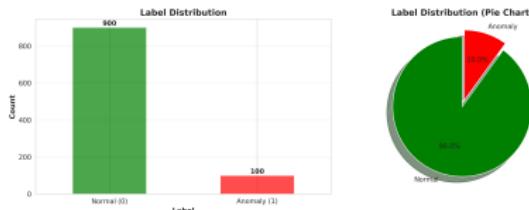
The KMeansIF method shows a clearer distinction and lower anomaly scores within the nine distinct, convex normal regions (blue), compared to the Standard IF, which exhibits more uniform scores across the feature space.

Outline

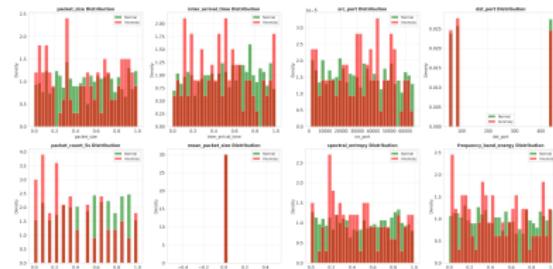
- 1 Isolation Forest (IF)
- 2 k-Means-Based Isolation Forest (kMIF)
- 3 Elbow Rule
- 4 Parallelization
- 5 Experiments: Synthetic Data
- 6 Experiments: Real-World Data
- 7 Future Work

Embedded System Network Security Dataset

- Real network traffic data
- **Extreme class imbalance:** 90% Normal, 10% Anomaly
- High feature overlap between classes

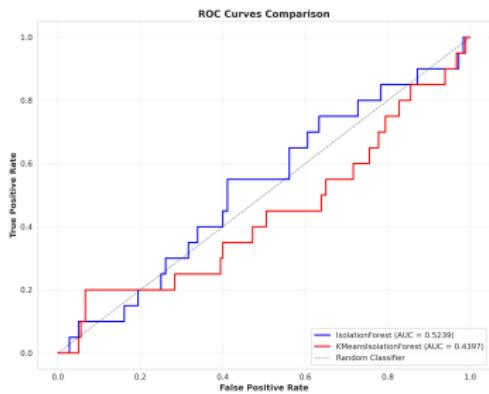
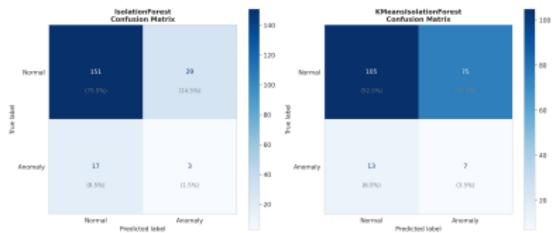


Class Distribution



Feature Overlap

ESNSDP Results



Performance

IF: Precision=0.094, ROC-AUC=0.524 — **kMIF:** Recall=0.35, 75 FPs

Cost

kMIF: 2x slower training, 4.5x slower prediction

Outline

1 Isolation Forest (IF)

2 k-Means-Based Isolation Forest (kMIF)

3 Elbow Rule

4 Parallelization

5 Experiments: Synthetic Data

6 Experiments: Real-World Data

7 Future Work

Conclusions & Future Work

Key Takeaways

- IF wins on non-convex data
- kMIF: complexity without gains
- Parallelization: 9-10x speedup
- IF: best for limited resources

Future Work

- Hybrid models
- Optimize kMIF overhead
- Alternative scoring
- More real-world datasets

Recommendation

Classic Isolation Forest remains optimal for general-purpose anomaly detection