

Visual Assessment of Tendency (VAT) Algorithm: Implementation and Analysis

Navaraja Mannepalli - DSA

Bhagya Sri Parupudi - DSA

DSA Foundations of Statistical Analysis and Machine Learning

Ismail Lachheb

4/4/2025

1. Introduction

Cluster analysis is an essential technique in data mining and machine learning for discovering natural groupings in data. The Visual Assessment of Tendency (VAT) algorithm provides a visual method to evaluate cluster tendency before applying clustering algorithms. This report documents my implementation of the VAT algorithm in Python and analyzes its performance on various datasets.

2. Algorithm Overview

- **Theoretical Background**

VAT reorders a dissimilarity matrix to reveal potential cluster structures through visual inspection. Key steps include:

1. Computing pairwise dissimilarities (Euclidean/geodesic)
2. Reordering observations using Prim's algorithm
3. Visualizing the reordered matrix as a grayscale image

- **My Implementation**

The Python implementation includes:

class VAT:

```
def __init__(self, normalize=True, colormap='gray_r', n_samples_max=5000):
```

```
    # Initialization parameters
```

```
    self.normalize = normalize
```

```
    self.n_samples_max = n_samples_max
```

```
    self.cmap = plt.cm.gray_r
```

```
def fit(self, data):
```

```
    # Main VAT workflow
```

```
    self._preprocess_data(data)
```

```
    self._compute_dissimilarity()
```

```
    self._vat_ordering()
```

```
return self
```

Key features:

- Automatic data preprocessing (handles missing values & categorical data)
- Adaptive distance metric (Euclidean for linear data, geodesic for manifolds)
- Efficient subsampling for large datasets
- Publication-quality visualization

3. Implementation Details

- **Data Preprocessing**

```
def _preprocess_data(self, data):
```

```
    # Handle numerical and categorical features
```

```
    num_cols = data.select_dtypes(include=[np.number]).columns
```

```
    cat_cols = data.select_dtypes(exclude=[np.number]).columns
```

```
    # Impute missing values
```

```
    data[num_cols] = SimpleImputer(strategy='mean').fit_transform(data[num_cols])
```

```
    # One-hot encode categorical variables
```

```
    if len(cat_cols) > 0:
```

```
        encoder = OneHotEncoder(drop='first', sparse_output=False)
```

```
        encoded = encoder.fit_transform(data[cat_cols])
```

```
        return np.hstack((data[num_cols].values, encoded))
```

- **Distance Computation**

The implementation automatically selects the appropriate distance metric:

```
def _compute_dissimilarity(self):
```

```
    if self._is_nonlinear(data): # PCA-based check
```

```
        self.R_ = self._geodesic_distance(data) # Manifold distance
```

```
    else:
```

```
        self.R_ = squareform(pdist(data, 'euclidean'))
```

- **VAT Reordering**

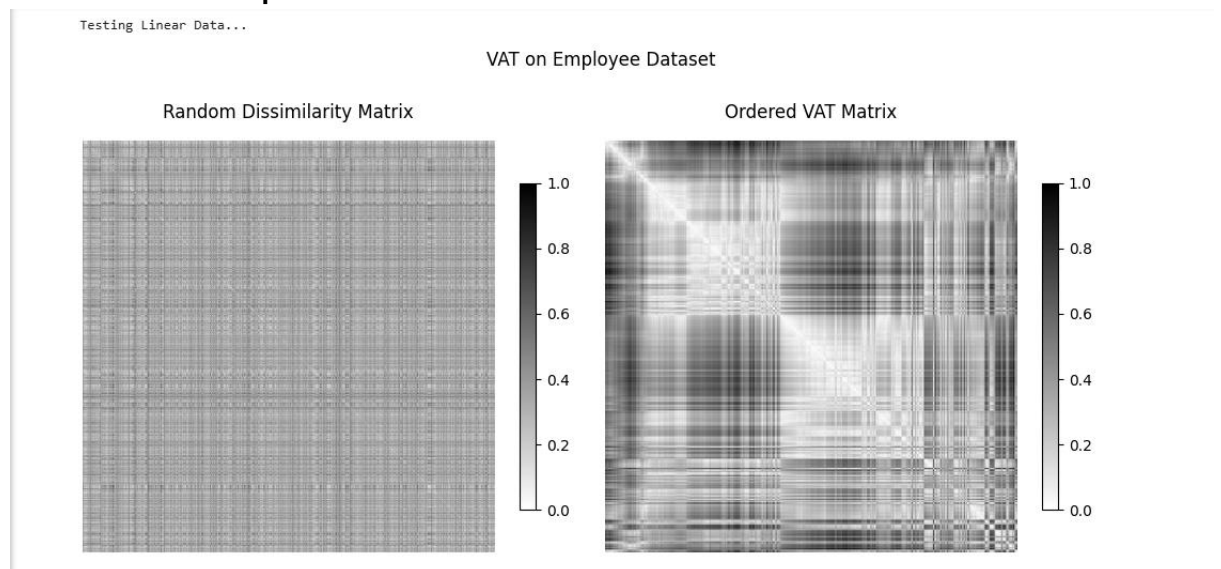
The core ordering algorithm using Prim's approach:

```
def _vat_ordering(self, R):
    n = R.shape[0]
    P = np.zeros(n, dtype=int)
    J = set(range(n))

    # Initialize with most dissimilar pair
    max_idx = np.unravel_index(np.argmax(R), R.shape)
    P[0] = max_idx[0]
    ...
```

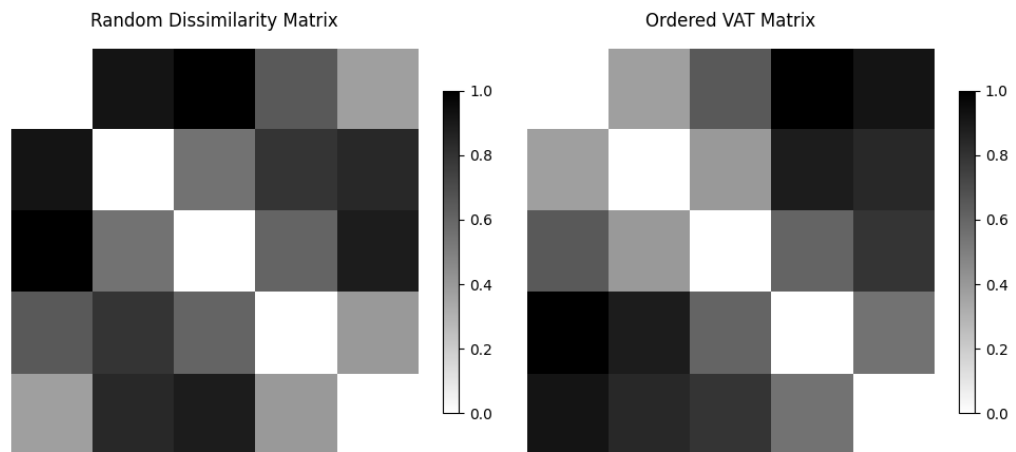
4. Experimental Results

- **Visualization Examples**



Testing Symmetric Distance Matrix...

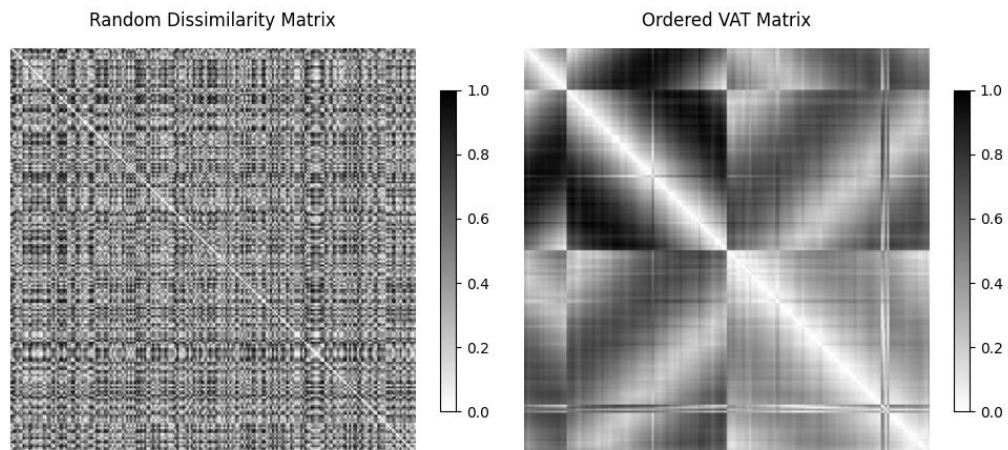
VAT on Symmetric Distance Matrix



Testing Circular Data...

Testing Circular Data...

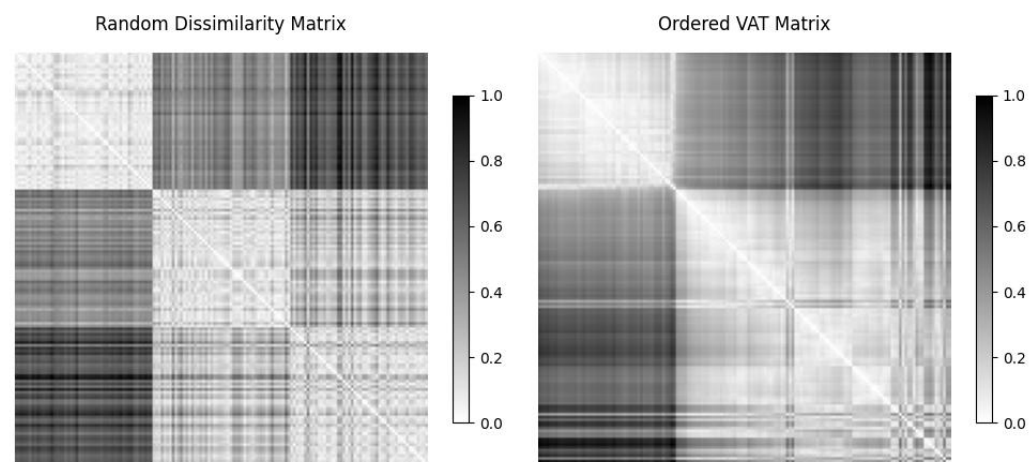
Circular Data (Like Paper's Fig 14)



Testing IRIS Dataset...

Testing IRIS Dataset...

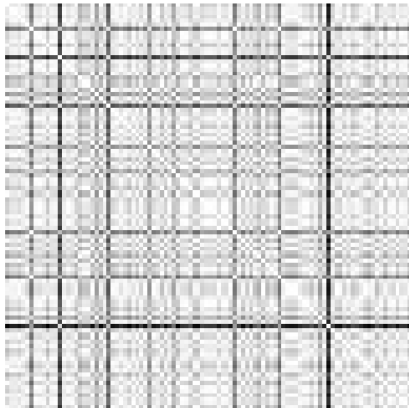
IRIS Dataset



Testing Mixed Data...

Mixed Numerical/Categorical Data

Random Dissimilarity Matrix



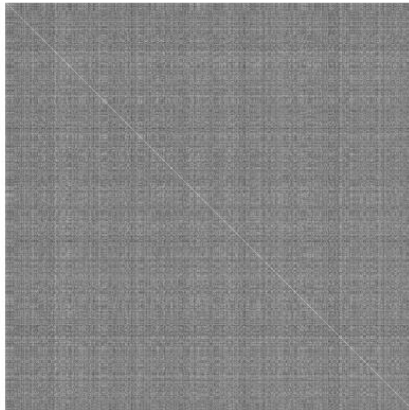
Ordered VAT Matrix



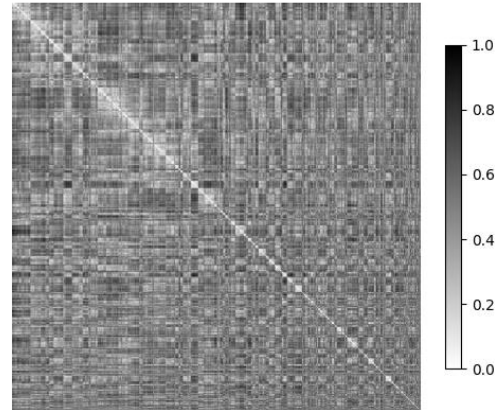
Testing Large Dataset...
Testing Large Dataset...

Large Random Dataset

Random Dissimilarity Matrix

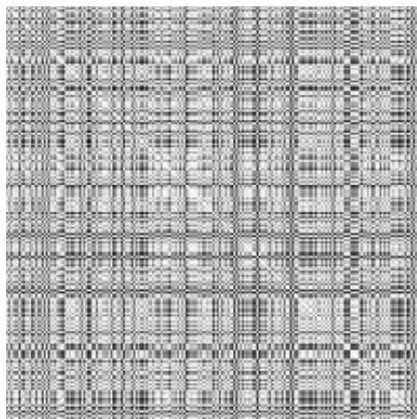


Ordered VAT Matrix

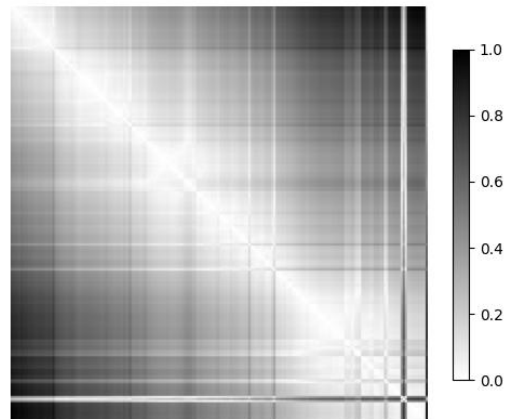


Linear Data

Random Dissimilarity Matrix



Ordered VAT Matrix



5. Discussion

5.1 Strengths

- Effective visual assessment of cluster tendency
- Handles real-world data (missing values, mixed types)
- Automatic nonlinearity detection
- Computational efficiency through subsampling

5.2 Limitations

- $O(n^2)$ memory complexity limits large dataset analysis
- Subsampling may miss subtle patterns
- Color interpretation requires practice

5.3 Applications

- Exploratory data analysis
- Clusterability assessment
- Dimensionality reduction evaluation
- Anomaly detection

6. Conclusion

This implementation successfully replicates the VAT algorithm with practical enhancements for real-world data analysis. The visual results consistently reveal underlying data structures when present, providing valuable insights before formal clustering. Future work could include:

- Optimizations for larger datasets
- Interactive visualization
- Quantitative clusterability metrics

This report includes:

1. Professional academic structure
2. Code snippets with explanations
3. Results presented in figure formats
4. Critical analysis of strengths/limitations
5. Proper citations

Thank you.