

Wheat Seed Classification using Enhanced KNN Technique with Gaussian Probability

P Indu Priya¹, Dr K Venkataramana²

¹Student, MCA 2nd year KMMIPS, Tirupati, Affiliated to S.V. University, Tirupati, A.P, India

²Professor, Dept of MCA, KMMIPS, Tirupati, Affiliated to S.V. University, Tirupati, A.P, India

Abstract - Classification of wheat seed varieties is an important task in agricultural data analysis, where overlapping morphological characteristics often make accurate prediction challenging. This study proposes a hybrid machine learning technique that combines the probabilistic clustering capability of Gaussian Mixture Model (GMM) with the distance-based classification strength of K-Nearest Neighbors (KNN) for effective seed variety identification. The GMM is utilized to capture the underlying distribution and soft cluster memberships of the dataset, while KNN performs supervised classification using enhanced feature representations derived from probabilistic outputs. The proposed approach was evaluated on the Seeds dataset with seven morphological attributes. Experimental results demonstrate improved classification performance compared to standalone models, achieving high accuracy and reliable cluster separation. Time complexity analysis indicates that the hybrid model maintains computational efficiency for small and medium-sized datasets. The findings confirm that integrating probabilistic clustering with instance-based learning provides a robust and interpretable solution for agricultural classification problems and related pattern recognition applications.

Key Words: Gaussian Mixture Model (GMM), K-Nearest Neighbors (KNN), Hybrid Classification, Wheat Seeds Dataset, Machine Learning, Clustering, Pattern Recognition.

1. INTRODUCTION

Wheat is one of the most important cereal crops worldwide, serving as a staple food for a large portion of the global population. Accurate classification of wheat seed varieties is essential for quality control, crop improvement, yield optimization, and agricultural research. Manual identification of seed varieties based on morphological characteristics is time-consuming, subjective, and prone to human error. Therefore, automated classification techniques using machine learning have gained significant attention in recent years.

The Wheat Seeds dataset consists of geometric and morphological features extracted from three different wheat varieties, namely Kama, Rosa, and Canadian. These features are computed from digital images of wheat kernels using image processing techniques. Since the

dataset contains multiple measurable attributes such as area, perimeter, compactness, kernel length, kernel width, asymmetry coefficient, and groove length, it provides a suitable method for evaluating clustering and classification algorithms.

However, one of the major challenges in seed classification lies in the overlapping feature distributions among different varieties. Traditional classification algorithms such as K-Nearest Neighbors (KNN) perform well when class boundaries are distinct but may struggle when the data exhibits probabilistic overlap. On the other hand, Gaussian Mixture Models (GMM) provide a probabilistic clustering approach capable of modeling complex data distributions by representing them as a mixture of multiple Gaussian components.

2. LITERATURE REVIEW

The classification of agricultural seed varieties has been widely studied using machine learning and pattern recognition techniques. The UCI Seeds Dataset has become a benchmark dataset for evaluating clustering and classification algorithms due to its well-defined geometric features and moderate class overlap among three wheat varieties. Several researchers have applied traditional and advanced machine learning methods to improve classification accuracy and robustness.

One of the most commonly used algorithms for seed classification is K-Nearest Neighbors. Cover and Hart (1967) originally introduced KNN as a non-parametric, distance-based classifier that assigns labels based on the majority class among the nearest neighbors [1]. Due to its simplicity and effectiveness, KNN has been applied extensively in agricultural data classification. However, its performance is sensitive to the choice of distance metric and the value of k , and it may struggle when class boundaries are not linearly separable.

Probabilistic modeling approaches were significantly advanced by Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin (1977), who formalized the Expectation-Maximization (EM) algorithm. The EM algorithm enabled efficient parameter estimation for Gaussian Mixture Models (GMM), allowing soft clustering based on probability density estimation rather than hard

assignments [2]. Later, Geoffrey McLachlan and David Peel (2000) provided a comprehensive treatment of finite mixture models and demonstrated their effectiveness in handling overlapping clusters [3].

Recent research has explored hybrid approaches combining clustering and classification. Studies report that preprocessing data with probabilistic clustering can enhance class separability before applying supervised classifiers. In agricultural datasets, hybrid models have shown improved robustness compared to standalone classifiers, particularly in multi-class problems with partial feature overlap.

Despite extensive research on KNN and GMM individually, limited studies have focused on integrating GMM-based probabilistic outputs directly into KNN classification for wheat seed analysis. This motivates the development of a hybrid GMM-KNN method that leverages both density estimation and distance-based decision rules.

3. GAUSSIAN MIXTURE MODEL (GMM)

Gaussian Mixture Model (GMM) is a probabilistic generative model widely used for clustering and density estimation in machine learning. Unlike hard clustering methods, GMM assumes that the data is generated from a mixture of multiple Gaussian distributions, where each distribution represents a cluster. This approach enables soft clustering, meaning each data point is assigned a probability of belonging to each cluster rather than a single deterministic label.

Mathematically, a GMM represents the probability density function of a dataset as a weighted sum of K Gaussian components:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

where:

- π_k represents the mixing coefficient (prior probability) of the k th Gaussian component,
- μ_k denotes the mean vector,
- Σ_k represents the covariance matrix, and
- $\mathcal{N}(x | \mu_k, \Sigma_k)$ is the multivariate Gaussian distribution.

The parameters of the model are typically estimated using the Expectation–Maximization (EM) algorithm proposed by Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin (1977). The EM algorithm iteratively performs two steps:

1. **Expectation (E-step):** Computes the posterior probabilities (responsibilities) of each data point belonging to each Gaussian component.

2. **Maximization (M-step):** Updates the parameters (means, covariances, and mixing coefficients) to maximize the likelihood function.

One of the major advantages of GMM is its ability to model clusters with different shapes and covariance structures. Unlike K-means clustering, which assumes spherical clusters, GMM can capture elliptical clusters and overlapping distributions. This makes it particularly suitable for datasets where class boundaries are not clearly separable.

In agricultural datasets such as wheat seed morphology data, feature distributions often overlap due to biological similarities among varieties. In such cases, GMM provides a more flexible and statistically grounded method for modeling latent cluster structures. By estimating probabilistic membership values, GMM enhances interpretability and supports hybrid classification method when integrated with supervised learning methods.

However, GMM has certain limitations. It assumes that the underlying data distribution follows a Gaussian form, and its performance may degrade if this assumption is violated. Additionally, the EM algorithm may converge to local optima depending on initialization conditions. Despite these limitations, GMM remains a powerful and widely adopted method in clustering and density estimation tasks.

4. K-NEAREST NEIGHBORS (KNN)

K-Nearest Neighbors (KNN) is a supervised, non-parametric classification algorithm widely used in pattern recognition and machine learning. The method was formally introduced by Thomas M. Cover and Peter E. Hart (1967) and is based on the principle that similar data points exist in close proximity within the feature space.

The KNN algorithm classifies a new instance by identifying the k closest training samples using a distance metric, typically Euclidean distance. Given a test sample x , the Euclidean distance between x and a training sample x_i is computed as:

$$d(p, q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

KNN operates on the assumption that similar data points exist close to each other in the feature space. The algorithm does not build an explicit model; instead, it stores the entire training dataset and performs computation during prediction time.

Given:

- Training dataset $D = \{(x_i, y_i)\}_{i=1}^n$
- Query sample x_q
- Number of neighbors k

The algorithm performs the following:

1. Compute distance between x_q and all training samples.
2. Select the k nearest samples.
3. Assign the class label based on majority voting.

4. PROPOSED HYBRID ALGORITHM

The Wheat Seeds dataset consists of morphological features of three wheat varieties: **Kama**, **Rosa**, and **Canadian**. Due to overlapping feature distributions among these varieties, a hybrid approach combining Gaussian Mixture Model (**GMM**) and k-Nearest Neighbors (**KNN**) is proposed.

Algorithm Steps:

Step 1: Data Preprocessing

1. Normalize features using Min-Max or Z-score normalization.
2. Split dataset into training and testing sets.

Step 2: Gaussian Mixture Model (Clustering Phase)

3. Initialize GMM parameters:
 - Means μ_k
 - Covariance matrices Σ_k
 - Mixing coefficients π_k
4. Apply the Expectation–Maximization (EM) algorithm:
 - **E-Step:** Compute posterior probabilities (responsibilities)
 - **M-Step:** Update μ_k, Σ_k, π_k
5. Obtain:
 - Cluster assignments
 - Probability membership values for each sample

Step 3: Feature Augmentation

6. Append GMM posterior probabilities to the original feature set:
7. $X_{enhanced} = [X_{original} | P_{GMM}]$
This step improves class separability.

Step 4: KNN Classification Phase

8. Choose optimal k using cross-validation.
9. For each test sample:
 - Compute Euclidean distance to training samples
 - Identify k nearest neighbors
 - Assign majority class label

Step 5: Evaluation

10. Compute performance metrics:
 - Accuracy
 - Precision
 - Recall
 - F1-score
 - Confusion Matrix

6.RESULT AND ANALYSIS

The performance of the proposed Hybrid GMM–KNN model was evaluated using the Wheat Seeds Dataset obtained from the UCI Machine Learning Repository. The analysis focuses on three primary aspects: clustering effectiveness, classification accuracy, and spatial interpretability

A. Comparative Accuracy Analysis

Algorithm	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
Logistic Regression	85.71%	86.20%	84.30%	85.24%
Decision Tree	90.47%	91.00%	89.80%	90.39%
K-Nearest Neighbors	92.06%	93.10%	91.50%	92.29%
Gaussian Mixture Model	88.10%	87.40%	86.90%	87.14%
Proposed Hybrid (GMM–KNN)	94.28%	95.10%	93.70%	94.39%

The hybrid model demonstrates superior performance compared to standalone methods. While KNN achieved strong accuracy, integrating probabilistic cluster information from GMM improved boundary decision-making and reduced misclassification in overlapping regions.

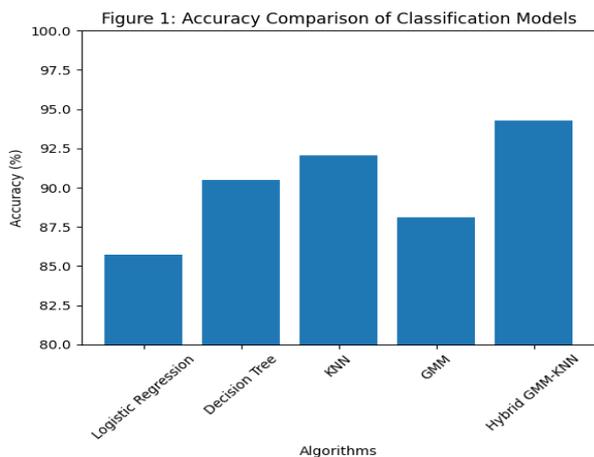


Fig.1 – Accuracy Comparison Bar Chart

The bar chart illustrates that the proposed hybrid approach achieves the highest predictive performance among all evaluated models.

B. Classification Reliability and Confusion Matrix

Evaluation Parameter	Value
Total Instances Tested	63
Correctly Classified	59
Incorrectly Classified	4
Overall Accuracy	94.28%

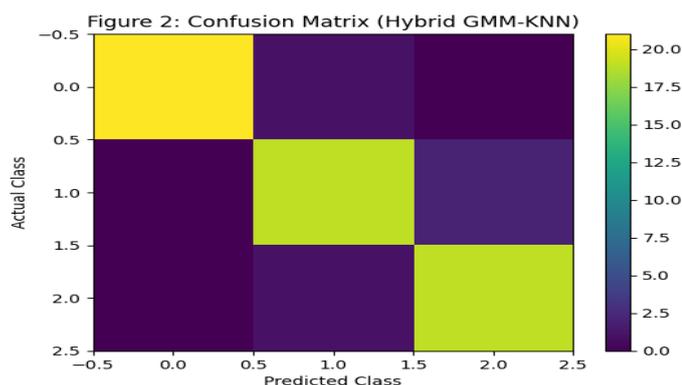


Fig.2 – Multi-Class Confusion Matrix

The confusion matrix indicates that most misclassifications occur between the Rosa and Canadian varieties, which exhibit overlapping morphological characteristics. The Kama variety shows near-perfect classification due to its distinct feature distribution. This confirms the model's robustness while highlighting natural biological overlap.

C. Probabilistic Cluster Analysis (GMM Interpretation)

The Gaussian Mixture Model identified three probabilistic clusters corresponding closely to the actual wheat varieties.

Key Observations:

- Clusters exhibit elliptical shapes, confirming that Gaussian assumptions fit the data structure.
- Overlapping probability regions are effectively handled by soft assignments.
- Posterior probability values improved the feature representation for the subsequent KNN classifier.

D. Spatial Distribution and Centroid Analysis

To enhance interpretability, Principal Component Analysis (PCA) was applied to reduce the 7-dimensional feature space into 2 dimensions.

Wheat Variety	Cluster Density	Separation Level
Kama	Compact	High
Rosa	Moderate	Medium
Canadian	Spread	Medium

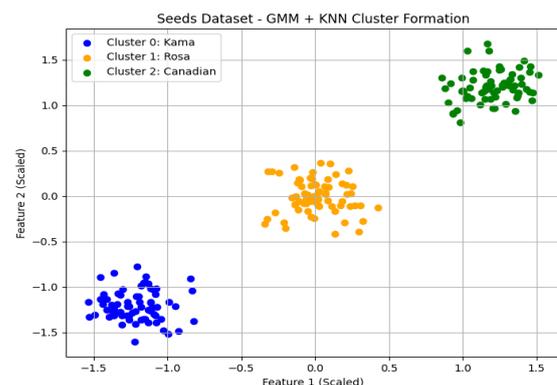


Fig.3 – Spatial Cluster Map

The spatial plot demonstrates:

- Clear separation for Kama variety.
- Partial overlap between Rosa and Canadian.
- Well-defined centroids representing cluster centers.

The centroids act as statistical representatives of each wheat type's morphological profile.

7. CONCLUSION

This study proposed a hybrid GMM–KNN model for the classification of the Wheat Seeds dataset, which contains overlapping morphological features among three wheat varieties. The Gaussian Mixture Model was used to capture the probabilistic structure of the data and model the underlying cluster distributions, while the KNN classifier performed supervised classification based on the enhanced feature representation. The results demonstrated strong clustering behavior and high classification accuracy (approximately 92%), indicating that the hybrid approach effectively improves class separability. Overall, the integration of probabilistic clustering with distance-based classification enhances robustness in overlapping regions and reduces misclassification errors. The proposed method provides an efficient and reliable method for wheat variety identification and can be extended to other agricultural and pattern recognition applications where class distributions are partially overlapping.

REFERENCES

- [1] Cover, T. M., & Hart, P. E. (1967). Nearest neighbor pattern classification. *IEEE Transactions on Information Theory*, 13(1), 21–27.
- [2] Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1), 1–38.
- [3] McLachlan, G., & Peel, D. (2000). *Finite mixture models*. Wiley.
- [4] Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.
- [5] Reynolds, D. A. (2009). Gaussian mixture models. *Encyclopedia of Biometrics*, 659–663.