

B. Review of Related Work on Students' Assessment Dataset

Cortez and Silva (2008), examined the performance of decision trees, Random forest, Neural Network and Support Vector Machines on secondary school grades using binary classification, regression and five-level classification for evaluation. They explained that social, demographic and school related variables also affect students' performance. Salami et al (2016) carried out a study on the detection anomalies in students' results using the decision trees. The decision tree model was able to detect efficiently anomalies in student results in most cases. However, it was unable to detect or identify anomalies in a situation where the training dataset has few anomalous instances. Thomas and Jayagopi (2017) measured the students' engagement using a machine learning algorithm based on students' facial expressions, head poses, and eye gazes. The experimental result showed that the machine learning algorithm performed well in predicting student engagement in class. Hamid et al (2018) measured students' engagement using a machine learning approach and concluded that the SVM and K-NN classifiers are appropriate for predicting students' engagement. Ukoba et al (2020) carried out a review on the detection and classification of anomalies in students' assessment dataset. The review pointed out that very few works have been done with regard to detecting anomalies in the educational sector/domain. Salami and Yahaya (2018) described how the Extreme Learning Machines (ELM) can be used to automatically detect anomalies in students' result. However, it was unable to detect anomalous instances in some of the dataset especially where the anomalous instances were few. This paper seeks to use the K-Nearest Neighbor to improve the detection of anomalies in students' assessment dataset which follows the normal distribution curve.

3. ANOMALIES IN STUDENTS' ASSESSMENT DATASET

Student result anomalies are salient or unusual observations that need further elucidation. A typical result of students in a department or a particular course ought to follow the normal (Gaussian) distribution curve, where few students should have A's and F's and bulk of the students score C as shown in figure 1 below (Ukoba et al, 2020). Any deviation from this is considered an anomaly. The anomaly on its own is not necessarily a bad thing. For instance, when we have more first-class and second class upper than second class lower in a graduating set isn't actually bad. However, this result or grade would cause a deviation from the normal distribution curve., thus an anomaly. This paper seeks to detect anomalies in course-based and CGPA based assessment data

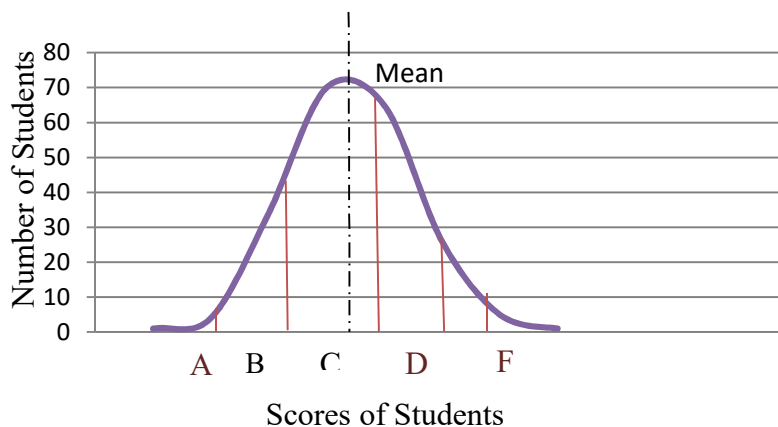


Figure 1: Normal Distribution for Scores/Grade of Students (Source: Ukoba et al, 2020)



B. Evaluation Metrics'

In validating detection and classification performance, several evaluation metrics' have been used. Some of which are F-Ratio, accuracy, sensitivity/recall, specificity, rank power etc. However, in this work, we use accuracy, sensitivity/recall, specificity and F-Ratio to validate the K-Nearest Neighbour algorithm performance.

Table 1 presents the relationship between the actual class and the predicted class with regards the four metric evaluation parameters.

- i. **Accuracy:** This metric is for evaluating observations correctly classified. This is the ratio of observations predicted correctly (True positive) to the total observations.

$$\text{Accuracy} = \frac{TP+}{TP+FP+TN+F} \times 100\% \quad (1)$$

In a class imbalanced dataset, accuracy alone doesn't tell if a model or system is doing excellently well

- ii. **Sensitivity/Recall:** This measure evaluates observations which are positive, correctly classified. It shows the ratio of positive targets to all targets in the main class. It gives the extent to which instances that are anomalous are correctly identified. Equation 3.2 gives the computation of sensitivity/recall.

$$\text{Sensitivity} = \frac{TP}{TP+FN} \quad (2)$$

- iii. **Specificity:** This measure evaluates observations which are negative, correctly classified. It shows the ratio of negative targets to all targets in the main class. It gives the extent to which instances that are anomalous are correctly identified. The higher the percentage, the better. Equation 3.2 gives the computation of specificity.

iv.

$$\text{Specificity} = \frac{TN}{TN+FP} \quad (3)$$

- v. **F1 score:** This is the weighted average of precision and recall. Therefore, the score takes both false positives and false negatives into account. F1 score is usually more useful than accuracy especially if we have an imbalance class distribution.

$$\text{F1 Score} = \frac{2 * (\text{Recall} * \text{Precision})}{\text{Recall} + \text{Precision}} \quad (4)$$

- vi. **Precision:** Precision is the percentage of positive accurately predicted observations to all positive predicted observations.

$$\text{Precision} = \frac{TP}{TP+FP} \quad (5)$$

Where TP = Number of true positive targets (anomalous instances) correctly classified.

FP = Number of positive targets [anomalous instances] wrongly classified.

TN = Number of true negative targets (normal instances) correctly classified.



Table 1: Relationship between the Actual and Predicted Class with regards the Metric Parameters

		Predicted Class	
		Class = Positive	Class = Negative
Actual Class	Class = Positive	TP	FN
	Class = Negative	FP	TN

C.

C. Experimental Results

Table 2 shows the K-Nearest Neighbour system for the Inconsistent CA vs Exam score, too many high score and the borderline failure anomalies. Out of the 160 instances in the CA vs Exam score anomaly dataset, the dataset were randomly splitted into 120 instances (75%) for training and 40 instances (25%) for testing. Table 3 shows the K-Nearest Neighbour system for the too many high CGPA anomalies. Out of the 150 instances, the dataset were randomly splitted using the split-train-test into 112 for training instances (75%) and 38 instances (25%) for testing.

Table 2: Evaluation of K- NN Models for Course Based Anomalies

K-NN Anomaly		Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)	Specificity (%)
Inconsistent CA vs Exam Scores	Training	99	96	100	98	100
	Testing	100	100	100	100	100
Too many high Score	Training	93	86	86	86	86
	Testing	90	80	92	86	92
Borderline Failure	Training	99	75	100	86	100
	Testing	100	100	100	100	100

Table 3 Evaluation of K- NN Models for Too Many High CGPA Anomalies

K-NN Anomaly		Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)	Specificity (%)
Too many High CGPA	Training	96	85	96	90	96
	Testing	100	100	100	100	100

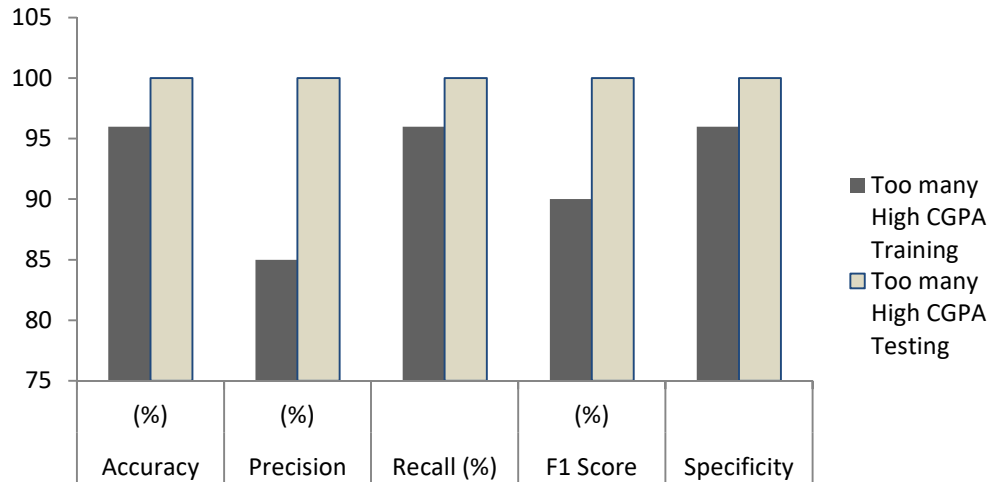


Figure 5: Chart showing the training and testing results of Too many High CGPA Anomaly

6. DISCUSSION OF RESULTS

Results from the research are discussed next.

1. Inconsistent CA vs Exam score Anomaly

Results from Table 2 which presents K-Nearest Neighbour performance for the Inconsistent CA vs Exam anomaly shows that all 26 anomalous instances in the training dataset were correctly detected by the K-NN model and 93 normal instances out of 94 instances in the training dataset were correctly classified, giving a specificity of 100%, accuracy of 99%, recall of 100%, precision of 96% and F1 Score of 98%. Testing results showed that all 6 anomalous instances and 34 normal instances were classified correctly, giving a specificity of 100%, accuracy of 100%, recall of 100%, precision of 100% and F1 Score of 100%.

2. Too Many High Score Anomaly

Results from Table 1, which presents K-Nearest Neighbour performance for Too many high score anomaly, showed that 26 anomalous instances out of the 27 anomalous instances in the training dataset were correctly classified by the K-NN model and 88 normal instances out of 93 instances in the training dataset were correctly classified giving a specificity of 86%, 93% accuracy, recall of 86%, precision of 86% and F1 Score of 86%. Testing results showed that out of the 14 anomalous instances, 13 were classified correctly and 23 out of the 26 normal instances were correctly classified, giving a specificity of 92%, accuracy of 90%, recall of 92%, precision of 80% and F1 Score of 86%.

3. Borderline Failure Anomaly

From Table 2, the performance of the K-Nearest Neighbour model built for the borderline failure anomaly, showed that no anomalous instances out of the 2 anomalous instances in the training dataset were classified by the K-NN model and all 118 normal instances in the training dataset were classified correctly, giving 99% accuracy, recall of 100%, specificity of 100%, recall of 86% and precision of 75%. Testing result showed that all 3 anomalous instances and 37 normal instances were correctly classified, giving a specificity of 100%, accuracy of 100%, recall of 100%, precision of 100% and F1 Score of 100%.

