

COMPARATIVE STUDY

This table presents a comprehensive comparison of various machine learning models and their performance with and without Principal Component Analysis (PCA) on the Wisconsin Breast Cancer classification dataset.

| S no. | Predictive Models | | |
|-------|--------------------------------|-------------------|----------------------|
| | Technique | Accuracy With PCA | Accuracy Without PCA |
| 1. | Artificial Neural Network | 99.12 | 97.37 |
| 2. | AdaBoost | 97.37 | 97.37 |
| 3. | AutoML(TPOT) | 98.25 | 99.12 |
| 4. | Bagging Classifier | 95.61 | 97.37 |
| 5. | Convolutional Neural Network | 94.74 | 97.37 |
| 6. | Categorical Boosting | 97.37 | 96.49 |
| 7. | Decision Trees | 94.74 | 94.74 |
| 8. | Ensemble Methods | 96.49 | 95.61 |
| 9. | Gaussian Mixture Model | 07.02 | 05.26 |
| 10. | Extreme Gradient Boosting | 95.61 | 98.25 |
| 11. | K-Means Clustering | 37.71 | 97.71 |
| 12. | K-Nearest Neighbors | 94.74 | 95.61 |
| 13. | Light Gradient Boosting | 96.49 | 94.74 |
| 14. | Logistic Regression | 97.37 | 98.25 |
| 15. | Naive Bayes | 96.49 | 92.11 |
| 16. | Neural Networks | 96.49 | 98.25 |
| 17. | Random Forest Model | 96.49 | 94.74 |
| 18. | Support Vector Machines | 95.61 | 96.49 |
| 19. | MobileNetV2 | 95.61 | - |
| 20. | EfficientNetB0 | 96.49 | - |
| 21. | DenseNet121 | 53.51 | - |
| 22. | Generative Adversarial Network | 62.28 | - |

1. Artificial Neural Network (ANN)

- a. Advantages: With an accuracy of 99.12% with PCA, ANNs are exceptional at capturing complex patterns and interactions in the data, making them highly effective for this classification task.

- b. Disadvantages: They can be computationally intensive and may overfit if not properly regularized.

2. AdaBoost

- a. Advantages: AdaBoost is a robust classifier that combines multiple weak learners to improve performance. Its consistent accuracy (97.37%) indicates good generalization.
- b. Disadvantages: It can be sensitive to noisy data and outliers, which might affect performance if the dataset is not preprocessed effectively.

3. AutoML (TPOT)

- a. Advantages: TPOT automates the machine learning pipeline, efficiently finding a well-performing model. Its performance is quite impressive, especially without PCA (99.12%).
- b. Disadvantages: It can be time-consuming to run and requires computational resources.

4. Bagging Classifier

- a. Advantages: It reduces variance and helps to avoid overfitting, which is evident from its fairly consistent performance.
- b. Disadvantages: It doesn't handle the bias of the base estimators, so if they're biased, the ensemble will be too.

5. Convolutional Neural Network (CNN)

- a. Advantages: CNNs can capture spatial hierarchies in the data. While more commonly used in image data, they can be adapted for structured datasets.
- b. Disadvantages: They can be complex and computationally intensive with a lot of parameters to tune.

6. Categorical Boosting (CatBoost)

- a. Advantages: It provides powerful binary classification and can handle categorical data well. It shows reasonable performance across the board.
- b. Disadvantages: It can be slower to train compared to other models and may overfit without proper parameter tuning.

7. Decision Trees

- a. Advantages: They are easy to interpret and don't require data scaling. They work well for non-linear relationships.
- b. Disadvantages: Prone to overfitting and don't generalize well without ensemble methods.

8. Ensemble Methods

- a. Advantages: Combining predictions from multiple models usually improves accuracy and reduces the chance of an individual model's biases affecting the outcome.
- b. Disadvantages: They can be complex and computationally intensive.

9. Gaussian Mixture Model (GMM)

- a. Advantages: GMM is good for density estimation and is a flexible model for classification.
- b. Disadvantages: Its very low accuracy indicates it might not be the right choice for this particular dataset, possibly due to the assumption about the data distribution not holding.

10. Extreme Gradient Boosting (XGBoost)

- a. Advantages: XGBoost provides a robust way to handle a variety of data types and relationships, with impressive performance.
- b. Disadvantages: It can be prone to overfitting and is sometimes sensitive to the specific hyperparameters used.

11. K-Means Clustering

- a. Advantages: It's an unsupervised method, useful for identifying structures in the data.
- b. Disadvantages: Its low accuracy in this context suggests it's not suitable for direct classification tasks without additional processing.

12. K-Nearest Neighbors (KNN)

- a. Advantages: KNN is simple and effective, performing reasonably well in this case. It's also intuitive and doesn't make assumptions about the data's distribution.
- b. Disadvantages: Its performance can degrade with high-dimensional data, and it can be computationally intensive for large datasets.

13. Light Gradient Boosting (LGBM)

- a. Advantages: LGBM is efficient and fast, even on large datasets, and handles imbalanced data well.
- b. Disadvantages: Like other gradient boosting methods, it can overfit if not carefully tuned.

14. Logistic Regression

- a. Advantages: It's a simple and explainable model, performing quite well in this scenario.
- b. Disadvantages: It can struggle with complex, non-linear relationships unless feature engineering is applied.

15. Naive Bayes

- a. Advantages: It's fast and performs decently well. It works well with small datasets and can be used for both binary and multiclass classification.
- b. Disadvantages: Its assumption of feature independence rarely holds true in real-world data, which can limit its performance.

16. Neural Networks

- a. Advantages: Similar to ANNs, they are powerful and flexible, capturing complex relationships in the data.
- b. Disadvantages: They require significant data and computational power and are prone to overfitting without regularization.

17. Random Forest Model

- a. Advantages: Random forests are robust, handle overfitting well, and provide variable importance measures.
- b. Disadvantages**: They can be less interpretable and slower to predict compared to some other models.

18. Support Vector Machines (SVM)

- a. Advantages: SVMs are effective in high-dimensional spaces and versatile with different kernel functions.
- b. Disadvantages: They can be memory-intensive and tricky to tune due to the importance of picking the right kernel.

19. MobileNetV2, EfficientNetB0, DenseNet121

- a. Advantages: These are advanced deep learning models, usually used for image classification tasks. They can capture complex patterns and hierarchies.
- b. Disadvantages: They are typically overkill for structured datasets like the one discussed here and require significant computational resources.

20. Generative Adversarial Network (GAN)

- a. Advantages: GANs are powerful for generating new data instances and capturing data distributions.
- b. Disadvantages: Their low accuracy here indicates they're not suitable for direct classification tasks and are complex to train.

In conclusion, while the best model depends on the specific needs and constraints of the task, models like ANNs, AutoML (TPOT), and XGBoost generally show strong performance on the Wisconsin Breast Cancer classification dataset. However, consideration of computational resources, the need for interpretability, and the specific nuances of the dataset should guide the final choice.

