# A Comparison of Classification Algorithms:

## Classification of celestial objects: stars, galaxies and quasars

| Classifiers | Accuracy Score | Strengths | Weaknesses |
|---|---|---|---|
| **1. Logistic Regression**<br><br>**cross validation logistic regression**<br><br>**with regularization**<br><br>important featrures: psfMag_u, psfMag_g , petromag_g , gr, ri, ug | 0.97<br><br>0.968<br><br><br>0.9693333 | 1.Most interpretable machine learning algorithms<br><br>2.Regularized to avoid overfitting | Underperform when there are multiple or non-linear decision boundaries |
| **2. SVM classifier**<br><br>**Using "OneVsRestClassifier"** | 0.954 | 1.Non-linear decision boundaries<br>2.Robust against overfitting, especially in high-dimensional space<br><br>3. Best classification performance (accuracy) on the training data. | 1.Don't scale well to larger datasets<br><br>2.Random forests are usually preferred over SVM's. |
| **3. KMeans** | - | Fast, simple, and surprisingly flexible | If the true underlying clusters in the data are not globular, then K-Means will produce poor clusters |
| **4. KNN classifier** | 0.904 | 1.Robust to noisy training data<br>2.Effective for large training data | 1.It is costly and lazy,<br>2.Requires full training data plus depends on the value of k<br>3. Has the issue of dimensionality because of the distance |
| **5. Ensemble Classifications**<br><br>• **Random Forest Classifier**<br><br>• **XGB classifier** | <br><br><br>0.97466667<br><br>0.9727 | 1.Perform well in practice<br>2.Robust to outliers,<br>3. scalable,<br>4. Naturally model non-linear decision boundaries<br>5.Overfitting is less | 1.Analysing theoretically is difficult<br>2.Large number of decision trees can slow down the algorithm in making real-time predictions.<br><br>3.If the data consists of categorical variables with different number of levels, then |

| Important Features:<br>Ug, iz, ri, psfFlux_u | | 6.Fast but not in all cases<br>7.Most effective and versatile<br>8.More robust to noise.<br>9.Can be grown in parallel.<br>10.Runs efficiently on large databases.<br><br>11.Has higher accuracy | the algorithm gets biased in favour of those attributes |
|---|---|---|---|
| **7. Decision Tree classifier** | 0.9384 | 1.can handle missing values nicely<br>2.best suited when the target function has discrete output values | 1.The more the number of decisions in a tree, less is the accuracy<br>2.do not fit well for continuous variables and result in instability and classification plateaus.<br>3.creating large decision trees that contain several branches is a complex and time-consuming task. |