

# Stellar Classification using Dataset SDSS17

Kristina Levina

08 september, 2022

## Introduction

Stellar classification is the classification of stars based on their spectral characteristics. Stellar Classification Dataset - SDSS17 [1] provides data for classification of stars, galaxies, and quasars based on their spectral characteristics. This dataset contains 17 attributes, one of which is *class*. Other 18 attributes include various spectral characteristics, object identifiers, and measurement-related parameters. Overall, the dataset contains 100000 observations of space taken by the Sloan Digital Sky Survey (SDSS).

The distribution of stars, quasars, and galaxies is shown in Fig. 1.

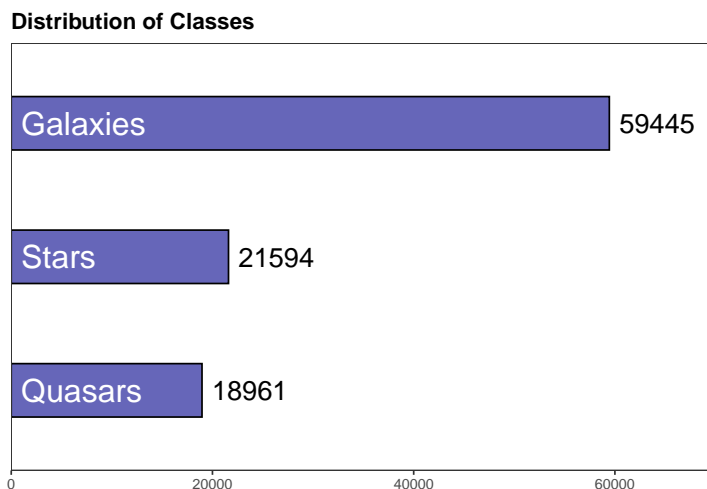


Figure 1: Distribution of stars, galaxies, and quasars in Stellar Classification Dataset - SDSS17

As one can observe from Fig. 1, the problem is quite imbalanced.

## Feature Selection

Full description of the attributes can be found in [1]. We need to select meaningful attributes to be used in our model using domain knowledge and exploratory analysis.

First, `obj_ID` is the object identifier, the unique value that identifies the object in the image catalog used by the CAS. Therefore, it should not be used in further analysis because it is not a spectral characteristic.

Second, `alpha` and `delta` specify the right ascension and declination angles, respectively. These angles specify the position of the object in the sky, making the use of the feature irrelevant for the model.

Third, `u`, `r`, `g`, `i`, `redshift`, and `z` are spectral characteristics. They should be included into the model.

Fourth, `run_ID`, `rerun_ID`, `cam_col`, `fiber_ID`, `plate`, `MJD`, and `field_ID` are features related to the measurement setup, which makes them irrelevant for stellar classification based on spectroscopic data.

Fifth, `spec_obj_ID` is related to spectroscopic observations. We will use it in our model.

Let us now calculate the correlation between chosen features for understanding their relationship between each other.

```
##           u      r      g      i      z spec_obj_ID redshift
## u          1.000 0.999 0.999 0.046 0.998         0.030   0.014
## r          0.999 1.000 1.000 0.056 0.999         0.039   0.023
## g          0.999 1.000 1.000 0.056 0.999         0.039   0.023
## i          0.046 0.056 0.056 1.000 0.056         0.662   0.492
## z          0.998 0.999 0.999 0.056 1.000         0.038   0.030
## spec_obj_ID 0.030 0.039 0.039 0.662 0.038         1.000   0.389
## redshift    0.014 0.023 0.023 0.492 0.030         0.389   1.000
```

We observe that `r`, `g`, `u`, and `z` are strongly correlated. However, further analysis of the model performance has shown that inclusion of all these features into the model improves the results.

## Data Preprocessing

The selected features were scaled so that they all have a mean of 0 and a standard deviation of 1. The dataset was split to 80% train and 20% test randomly.

## Problem Visualisation

We selected two meaningful features `spec_obj_ID` versus `red_shift` and visualised the distribution of classes. The results are shown in Fig. 2.

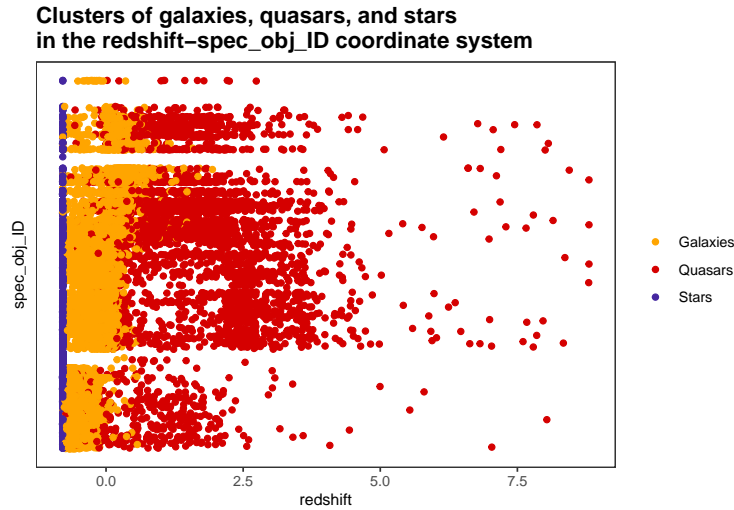


Figure 2: Visualised distribution of classes in the coordinate system of `spec-obj-ID` and `red-shift`

Figure 2 shows three distinct clusters that should be easily identified using logistic regression. We see that the red shift value that corresponds to quasars is the highest because these are the most distant objects we can observe. Stars have the lowest red shift value. The value of the red shift of galaxies is in between.

## Model and Results

We use logistic regression model with `multinom` function from `nnet` package. To evaluate the model performance, we use one-versus rest approach, and investigate the following six cases: 1) Train data, galaxies-versus-rest classification; 2) Train data, quasars-versus-rest classification; 3) Train data, stars-versus-rest classification; 4) Test data, galaxies-versus-rest classification; 5) Test data, quasars-versus-rest classification; 6) Test data, stars-versus-rest classification.

For the classification of stars versus other objects, the obtained train and test accuracies are 99.4% and 99.5%, respectively. Figure 3 shows the corresponding confusion matrices. The accuracy metric is a good choice in this case because the number of misclassification cases is low. Hence, stars are easy to distinguish from galaxies and quasars.

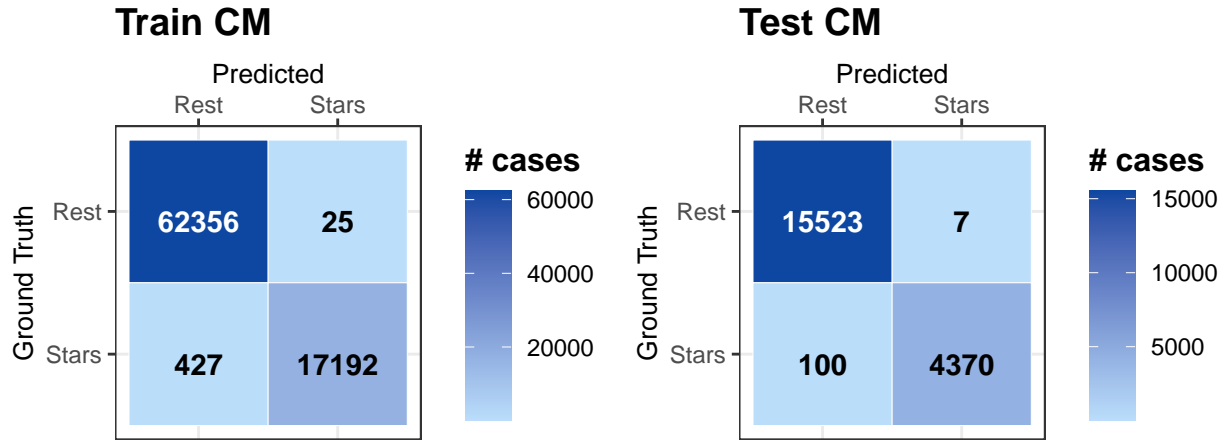


Figure 3: Confusion matrices for train and test data for classification of stars versus other classes.

For the classification of quasars versus other objects, the obtained train and test  $\phi$  coefficients are 89.2% and 89.2%, respectively. Figure 4 shows the corresponding confusion matrices. The  $\phi$  coefficient metric is a good choice in this case because of the considerable imbalance of quasars with respect to other classes and difficulty in distinguishing quasars from galaxies based on the studied spectral characteristics.

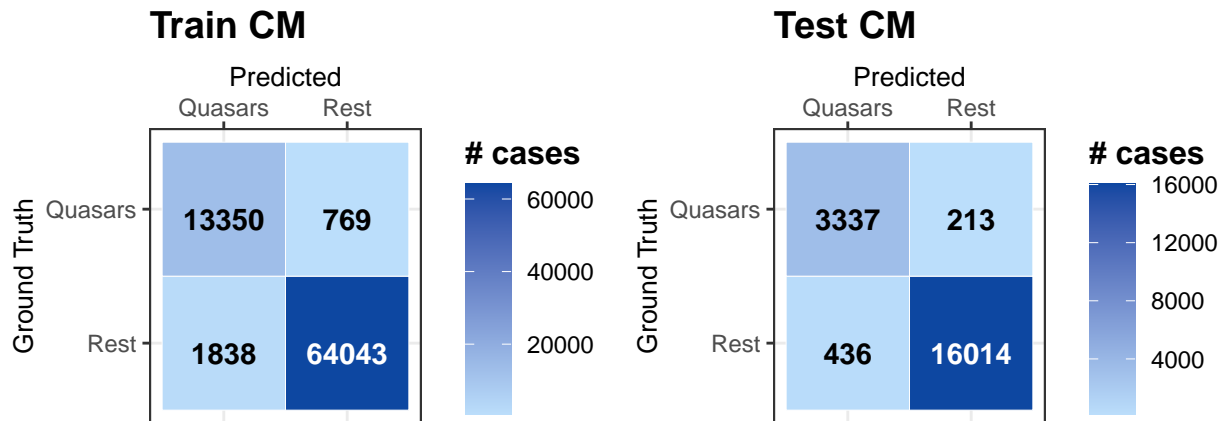


Figure 4: Confusion matrices for train and test data for classification of quasars versus other classes.

Considering the considerable imbalance, the model performance is decent in distinguishing quasars from other classes.

Finally, we assess the model performance in distinguishing galaxies from other classes.

For the classification of galaxies versus other objects, the obtained train and test accuracies are 96.2% and 96.2%, respectively. Figure 5 shows the corresponding confusion matrices. The accuracy metric is a good choice in this case because the number of false positives and false negatives is similar and the number of galaxies is quite balanced with respect to the number of other classes.

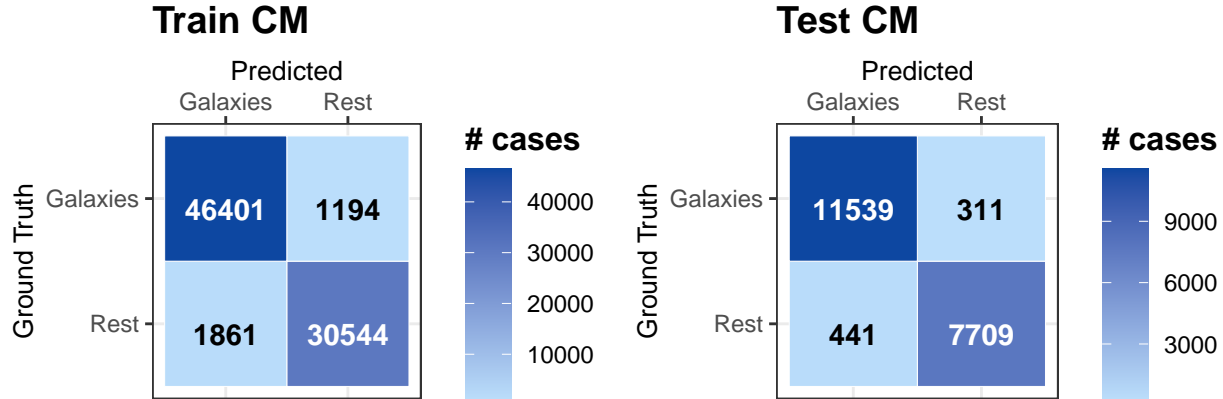


Figure 5: Confusion matrices for train and test data for classification of galaxies versus other classes.

## Conclusion

**Summary:** In this study, Stellar Classification Dataset - SDSS17 was explored. This is a multi-class classification problem with three classes: galaxies, stars, and quasars. The logistic regression model with `multinom` function from `nnet` package was utilised to solve this problem. The obtained accuracy for distinguishing stars from other classes is 99.4%. The obtained  $\phi$  coefficient for distinguishing quasars from other classes is 89.2%. The obtained accuracy for distinguishing galaxies from other classes is 96.2%.

**Strengths and limitations:** We can observe that the model performance is excellent for distinguishing stars from other classes. The model performance is slightly worse but still good for distinguishing galaxies from other classes. In contrast, the model performance is bad for distinguishing quasars from other classes.

**Future scope:** We will attempt to increase the performance of quasars' classification without compromising the performance of stars' and galaxies' classification by improving the model. Furthermore, the performance of the model selected herein will be compared with those of other models, and the best model will be used. In addition, the features used in the model should be investigated further; for example, basis expansion can be used, and correlated features should be explored further. Finally, the model evaluation criteria will be adjusted; for example, ROC and AUC curves will be explored.

## References

[1] <https://www.kaggle.com/datasets/fedesoriano/stellar-classification-dataset-sdss17> (The data released by the SDSS is under public domain.)

## Appendix

### Full code

```
#####
# LIBRARIES
#####
```

```

library(nnet)
library(ggplot2)
#####
# READING THE DATA
#####
data <- read.csv("data/star_classification.csv", stringsAsFactors = TRUE)
#####
# FEATURE SELECTION
#####
# Now, let us see the correlation between features
features <- data.frame(u = data$u, r = data$g, g = data$g, i = data$i,
                      z = data$z, spec_obj_ID = data$spec_obj_ID,
                      redshift = data$redshift)

target <- data$class

# Compute the correlation matrix
cor_mat <- cor(features)
print(cor_mat)

# We observe that r, g, u, and z are strongly correlated.

#####
# DATA PREPROCESSING
#####
features <- scale(features)
data_prep <- as.data.frame(cbind(target, features))
n <- dim(data_prep)[1]
# Data division
set.seed(12345)
tr_ind <- sample(1:n, floor(0.8*n))
tr <- data_prep[tr_ind, ]
te <- data_prep[-tr_ind, ]
# logistic regression with multinom
#####
# DATA VISUALISATION
#####

# Let us visualise the distribution of data points as spec_obj_ID versus red_shift
# colored by class

plot <- ggplot(data = te) + theme_bw() +
  geom_point(mapping = aes(x = redshift, y = spec_obj_ID,
                          color = as.factor(target))) +
  scale_color_manual(values = c("orange", "#d50000", "#4527a0"),
                    name = NULL,
                    labels = c("Galaxies", "Quasars", "Stars")) +
  labs(title = paste0("Clusters of galaxies, quasars, and stars\nin the",
                      "redshift-spec_obj_ID coordinate system")) +
  theme(panel.grid.major = element_blank(),
        panel.grid.minor = element_blank(),
        axis.text.y = element_blank(),
        axis.ticks.y = element_blank(),
        title = element_text(size = 12, face = 'bold'),

```

```

axis.title = element_text(size = 10, face = "plain"))

print(plot)

#####
# LOGISTIC REGRESSION
#####
logreg <- multinom(formula = target ~ ., data = tr)

pred_tr <- predict(logreg, tr, type = "class")
pred_te <- predict(logreg, te, type = "class")
#####
# PERFORMANCE EVALUATION
#####
CM_tr <- table(pred_tr, tr$target, deparse.level = 0)
CM_te <- table(pred_te, te$target, deparse.level = 0)
rownames(CM_te) <- c("GAL_pr", "QUA_pr", "STA_pr")
colnames(CM_te) <- c("GAL_tr", "QUA_tr", "STA_tr")

cat("The confusion matrix for train data is\n"); print(CM_tr)
cat("The confusion matrix for test data is\n"); print(CM_te)
# 1 is GALAXY, 2 is QUASAR, 3 is STAR

# Let us introduce this renewed stars-versus-rest CM
CM_tr_stars_vs_rest <- matrix(c(CM_tr[1, 1] + CM_tr[1, 2] + CM_tr[2, 1] +
                                CM_tr[2, 2],
                                CM_tr[2, 3] + CM_tr[1, 3],
                                CM_tr[3, 1] + CM_tr[3, 2],
                                CM_tr[3, 3]), byrow = TRUE, ncol = 2)

colnames(CM_tr_stars_vs_rest) <- c("rest_tr", "STA_tr")
rownames(CM_tr_stars_vs_rest) <- c("rest_pr", "STA_pr")

print("For train data")
print(CM_tr_stars_vs_rest)
GT <- factor(c("Rest", "Stars", "Rest", "Stars"))
predicted <- factor(c("Rest", "Rest", "Stars", "Stars"))
values <- c(CM_tr_stars_vs_rest)
df <- data.frame(GT, predicted, values)

# Visualise this CM
plot <- ggplot(data = df, mapping = aes(x = predicted, y = GT)) +
  geom_tile(aes(fill = values), color = "white") +
  geom_text(data = subset(df, values > 30000),
            aes(label = values,
                vjust = 1, color = "white", fontface = "bold", family = "arial") +
  geom_text(data = subset(df, values < 30000),
            aes(label = values,
                vjust = 1, color = "black", fontface = "bold", family = "arial") +
  scale_fill_gradient(low = "#bbdefb", high = "#0d47a1") +
  theme_bw() +
  scale_x_discrete(position = "top") +
  scale_y_discrete(limits=rev) +

```

```

labs(x = "Predicted", y = "Ground Truth", fill = "# cases",
     title = "Train CM for stars-versus-rest case") +
theme(title = element_text(size = 12, face = 'bold'),
      axis.title = element_text(size = 10, face = "plain"))

print(plot)

CM_te_stars_vs_rest <- matrix(c(CM_te[1, 1] + CM_te[1, 2] + CM_te[2, 1] +
                              CM_te[2, 2],
                              CM_te[2, 3] + CM_te[1, 3],
                              CM_te[3, 1] + CM_te[3, 2],
                              CM_te[3, 3]), byrow = TRUE, ncol = 2)

colnames(CM_te_stars_vs_rest) <- c("rest_tr", "STA_tr")
rownames(CM_te_stars_vs_rest) <- c("rest_pr", "STA_pr")

print("For test data")
print(CM_te_stars_vs_rest)
GT <- factor(c("Rest", "Stars", "Rest", "Stars"))
predicted <- factor(c("Rest", "Rest", "Stars", "Stars"))
values <- c(CM_te_stars_vs_rest)
df <- data.frame(GT, predicted, values)

# Visualise this CM
plot <- ggplot(data = df, mapping = aes(x = predicted, y = GT)) +
  geom_tile(aes(fill = values), color = "white") +
  geom_text(data = subset(df, values > 10000),
            aes(label = values),
            vjust = 1, color = "white", fontface = "bold", family = "arial") +
  geom_text(data = subset(df, values < 10000),
            aes(label = values),
            vjust = 1, color = "black", fontface = "bold", family = "arial") +
  scale_fill_gradient(low = "#bbdefb", high = "#0d47a1") +
  theme_bw() +
  scale_x_discrete(position = "top") +
  scale_y_discrete(limits=rev) +
  labs(x = "Predicted", y = "Ground Truth", fill = "# cases",
       title = "Test CM for stars-versus-rest case") +
  theme(title = element_text(size = 12, face = 'bold'),
        axis.title = element_text(size = 10, face = "plain"))

print(plot)

cat("Train accuracy for distinguishing between stars and other classes\n")
MCR_tr <- (CM_tr_stars_vs_rest[1, 1] + CM_tr_stars_vs_rest[2, 2]) /
  sum(CM_tr_stars_vs_rest)
cat(round(MCR_tr*100, 1), "%", sep="")

cat("Test accuracy for distinguishing between stars and other classes\n")
MCR_te <- (CM_te_stars_vs_rest[1, 1] + CM_te_stars_vs_rest[2, 2]) /
  sum(CM_te_stars_vs_rest)
cat(round(MCR_te*100, 1), "%", sep="")

```

```

# Now, let us consider distinguishing quasars from other data. The confusion
# matrix is as follows:
CM_tr_quasars_vs_rest <- matrix(c(CM_tr[1, 1] + CM_tr[1, 3] + CM_tr[3, 1] +
                                CM_tr[3, 3],
                                CM_tr[1, 2] + CM_tr[3, 2],
                                CM_tr[2, 1] + CM_tr[2, 3],
                                CM_tr[2, 2]), byrow = TRUE, ncol = 2)

colnames(CM_tr_quasars_vs_rest) <- c("rest_tr", "QUA_tr")
rownames(CM_tr_quasars_vs_rest) <- c("rest_pr", "QUA_pr")

print("For train data")
print(CM_tr_quasars_vs_rest)

GT <- factor(c("Rest", "Quasars", "Rest", "Quasars"))
predicted <- factor(c("Rest", "Rest", "Quasars", "Quasars"))
values <- c(CM_tr_quasars_vs_rest)
df <- data.frame(GT, predicted, values)

# Visualise this CM
plot <- ggplot(data = df, mapping = aes(x = predicted, y = GT)) +
  geom_tile(aes(fill = values), color = "white") +
  geom_text(data = subset(df, values > 40000),
            aes(label = values,
                vjust = 1, color = "white", fontface = "bold", family = "arial") +
  geom_text(data = subset(df, values < 40000),
            aes(label = values,
                vjust = 1, color = "black", fontface = "bold", family = "arial") +
  scale_fill_gradient(low = "#bbdefb", high = "#0d47a1") +
  theme_bw() +
  scale_x_discrete(position = "top") +
  scale_y_discrete(limits=rev) +
  labs(x = "Predicted", y = "Ground Truth", fill = "# cases",
       title = "Train CM for quasars-versus-rest case") +
  theme(title = element_text(size = 12, face = 'bold'),
        axis.title = element_text(size = 10, face = "plain"))

print(plot)

CM_te_quasars_vs_rest <- matrix(c(CM_te[1, 1] + CM_te[1, 3] + CM_te[3, 1] +
                                CM_te[3, 3],
                                CM_te[1, 2] + CM_te[3, 2],
                                CM_te[2, 1] + CM_te[2, 3],
                                CM_te[2, 2]), byrow = TRUE, ncol = 2)

colnames(CM_te_quasars_vs_rest) <- c("rest_tr", "QUA_tr")
rownames(CM_te_quasars_vs_rest) <- c("rest_pr", "QUA_pr")

print("For test data")
print(CM_te_quasars_vs_rest)

GT <- factor(c("Rest", "Quasars", "Rest", "Quasars"))
predicted <- factor(c("Rest", "Rest", "Quasars", "Quasars"))

```



```

values <- c(CM_te_quasars_vs_rest)
df <- data.frame(GT, predicted, values)

# Visualise this CM
plot <- ggplot(data = df, mapping = aes(x = predicted, y = GT)) +
  geom_tile(aes(fill = values), color = "white") +
  geom_text(data = subset(df, values > 6000),
    aes(label = values),
    vjust = 1, color = "white", fontface = "bold", family = "arial") +
  geom_text(data = subset(df, values < 6000),
    aes(label = values),
    vjust = 1, color = "black", fontface = "bold", family = "arial") +
  scale_fill_gradient(low = "#bbdefb", high = "#0d47a1") +
  theme_bw() +
  scale_x_discrete(position = "top") +
  scale_y_discrete(limits=rev) +
  labs(x = "Predicted", y = "Ground Truth", fill = "# cases",
    title = "Test CM for quasars-versus-rest case") +
  theme(title = element_text(size = 12, face = 'bold'),
    axis.title = element_text(size = 10, face = "plain"))

print(plot)

# Here, we observe the large number of misclassification cases. To assess the
# model performance, considering the data imbalance (low number of quasars wrt
# other classes, as shown in Figure X), we use the Phi coefficient, which is a
# recommended choice for binary classification in the case of imbalanced data

phi_coef <- function(M){
  numerator <- (M[1, 1]*M[2, 2] - M[1, 2]*M[2, 1])

  MCC <- numerator/sqrt(M[1, 1] + M[1, 2])
  MCC <- MCC/sqrt(M[2, 1] + M[2, 2])
  MCC <- MCC/sqrt(M[1, 1] + M[2, 1])
  MCC <- MCC/sqrt(M[1, 2] + M[2, 2])
  return(MCC)
}

cat("Train Phi coeff for distinguishing between quasars and other classes\n",
  round(phi_coef(CM_tr_quasars_vs_rest)*100, 1), "%", sep = "")
cat("Test Phi coeff for distinguishing between quasars and other classes\n",
  round(phi_coef(CM_te_quasars_vs_rest)*100, 1), "%", sep = "")

# Considering the considerable imbalance, the model performance is decent in
# distinguishing quasars from other classes

# Finally, we want to assess the model performance in distinguishing galaxies
# from other classes.

CM_tr_galaxies_vs_rest <- matrix(c(CM_tr[2, 2] + CM_tr[2, 3] + CM_tr[3, 2] +
  CM_tr[3, 3],
  CM_tr[1, 2] + CM_tr[1, 3],
  CM_tr[2, 1] + CM_tr[3, 1],

```

```

CM_tr[1, 1]), byrow = TRUE, ncol = 2)

colnames(CM_tr_galaxies_vs_rest) <- c("rest_tr", "GAL_tr")
rownames(CM_tr_galaxies_vs_rest) <- c("rest_pr", "GAL_pr")

print("For train data")
print(CM_tr_galaxies_vs_rest)
GT <- factor(c("Rest", "Galaxies", "Rest", "Galaxies"))
predicted <- factor(c("Rest", "Rest", "Galaxies", "Galaxies"))
values <- c(CM_tr_galaxies_vs_rest)
df <- data.frame(GT, predicted, values)

# Visualise this CM
plot <- ggplot(data = df, mapping = aes(x = predicted, y = GT)) +
  geom_tile(aes(fill = values), color = "white") +
  geom_text(data = subset(df, values > 30000),
    aes(label = values),
    vjust = 1, color = "white", fontface = "bold", family = "arial") +
  geom_text(data = subset(df, values < 30000),
    aes(label = values),
    vjust = 1, color = "black", fontface = "bold", family = "arial") +
  scale_fill_gradient(low = "#bbdefb", high = "#0d47a1") +
  theme_bw() +
  scale_x_discrete(position = "top") +
  scale_y_discrete(limits=rev) +
  labs(x = "Predicted", y = "Ground Truth", fill = "# cases",
    title = "Train CM for galaxies-versus-rest case") +
  theme(title = element_text(size = 12, face = 'bold'),
    axis.title = element_text(size = 10, face = "plain"))

print(plot)

CM_te_galaxies_vs_rest <- matrix(c(CM_te[2, 2] + CM_te[2, 3] + CM_te[3, 2] +
  CM_te[3, 3],
  CM_te[1, 2] + CM_te[1, 3],
  CM_te[2, 1] + CM_te[3, 1],
  CM_te[1, 1]), byrow = TRUE, ncol = 2)

colnames(CM_te_galaxies_vs_rest) <- c("rest_tr", "GAL_tr")
rownames(CM_te_galaxies_vs_rest) <- c("rest_pr", "GAL_pr")

print("For test data")
print(CM_te_galaxies_vs_rest)
GT <- factor(c("Rest", "Galaxies", "Rest", "Galaxies"))
predicted <- factor(c("Rest", "Rest", "Galaxies", "Galaxies"))
values <- c(CM_te_galaxies_vs_rest)
df <- data.frame(GT, predicted, values)

# Visualise this CM
plot <- ggplot(data = df, mapping = aes(x = predicted, y = GT)) +
  geom_tile(aes(fill = values), color = "white") +
  geom_text(data = subset(df, values > 6000),
    aes(label = values),

```

```

      vjust = 1, color = "white", fontface = "bold", family = "arial") +
geom_text(data = subset(df, values < 6000),
      aes(label = values),
      vjust = 1, color = "black", fontface = "bold", family = "arial") +
scale_fill_gradient(low = "#bbdefb", high = "#0d47a1") +
theme_bw() +
scale_x_discrete(position = "top") +
scale_y_discrete(limits=rev) +
labs(x = "Predicted", y = "Ground Truth", fill = "# cases",
      title = "Test CM for galaxies-versus-rest case") +
theme(title = element_text(size = 12, face = 'bold'),
      axis.title = element_text(size = 10, face = "plain"))

print(plot)
# For assessing the model performance for distinguishing galaxies from other classes,
# we can use the accuracy metric because the imbalance is not a problem here.

cat("Train accuracy for distinguishing between galaxies and other classes\n")
MCR_tr <- (CM_tr_galaxies_vs_rest[1, 1] + CM_tr_galaxies_vs_rest[2, 2]) /
  sum(CM_tr_galaxies_vs_rest)
cat(round(MCR_tr*100, 1), "%", sep="")

cat("Test accuracy for distinguishing between galaxies and other classes\n")
MCR_te <- (CM_te_galaxies_vs_rest[1, 1] + CM_te_galaxies_vs_rest[2, 2]) /
  sum(CM_te_galaxies_vs_rest)
cat(round(MCR_te*100, 1), "%", sep="")
#####
# CLASS DISTRIBUTION
#####
# Visualise the distribution of classes
classes = c("GALAXY", "STAR", "QSO")
counts = numeric(3)
for (i in 1:3){
  counts[i] = length(which(data$class == classes[i]))
}

df_classes <- data.frame(classes = classes, counts = counts)

plot <- ggplot(data = df_classes, mapping = aes(x = counts,
      y = reorder(classes, counts))) +
  theme_bw() +
  geom_bar(stat = "identity", fill = "darkblue", alpha = .6, width = .4,
    color = "black") +
  labs(title = "Distribution of Classes", x = "", y = "") +
  theme(panel.grid.major = element_blank(),
    panel.grid.minor = element_blank(),
    axis.text.y = element_blank(),
    axis.ticks.y = element_blank(),
    title = element_text(size = 12, face = 'bold')) +
  scale_x_continuous(expand = c(0, 10), limits = c(0, 70000)) +
  geom_text(
    aes(1000, y = classes, label = c("Galaxies", "Stars", "Quasars")),
    hjust = 0,

```

```
    nudge_x = 0,  
    color = "white",  
    family = "Calibri",  
    size = 7  
  ) +  
  geom_text(  
    aes(counts + 1000, y = classes, label = counts),  
    hjust = 0,  
    nudge_x = 0.3,  
    colour = "black",  
    family = "Calibri",  
    size = 6  
  )  
print(plot)
```