# An investigation into the ability of a machine learning algorithm to work out if a mushroom is edible or not, given certain factors.

https://github.com/EyeBall64/MachineLearningAssessment2

# Table of contents

# Introduction

When I was younger, I used to go walking in the woods around my na's house all the time, which would have many different types of mushrooms present. My nan always used to tell me to be careful about them all since they could be poisonous. I knew not all of them were, but I wanted to be able to tell which once were for my own curiosity. With this project, i hopefully will.

I will be using many different methods, such as random forests and GLM functions, among others, to get a good idea of whether a mushroom from the data set that I have found is edible.

In the R code, I use a maximum of 80 characters on a line for easy readability, and I am also to of some missing values that I will explore in the data cleaning.

# Data cleaning

Lines 82 - 221

Data cleaning is an essential part of the process. It ensures that the algorithms that I am applying to it will actually work. Lots of data will have errors or entries without any information. This can cause problems with some of the algorithms; even if they do run properly, it could cause erroneous results, as it is based on incorrect data.

Luckily, the most glaring problem was highlighted in the information that was given with the data; many entries in the stalk-root feature were NA. I decided that the most logical way of dealing with this was by removing all rows with NAs in; this removed quite a few, but the vast majority still remained, with over 5,000 data points still ready to use.

I wanted to keep an original,I unedited data file in the program, so should if  I need to go back after changing something, I could easily do that without having to reload the original data file.

One of the first things that I noticed and decided to change was that all of the data had been copied in as single-character codes. In the information document, the creator of the data file states that this was done to massively reduce the size of the document. However, I thought that it would be a lot easier to handle and observe the data if the original meanings behind the characters were restored. So, quite a lot of the data cleaning segment is devoted to this.

I also noticed that some of the data could be better expressed as logical and numerical variables, so while I was editing the above, I also changed these.

The last thing that I do in this section is check if there are any duplicate rows, which can be a good indicator of any data that has been incorrectly copied. Luckily, there were no duchies.

As it turns out, the original creator of this data set did a very good job of ensuring that all the data was well stored and checked very well, so that there were no errors and only one set of missing variables.

# Pie charts

Lines 222 - 429

I decided that there was so much data to deal with, I wanted a way to get a good idea of the data points so that I was able to make more informed decisions about the different machine learning algorithms I were to employ later on. While I was making the pie charts, I used a function that I created and defined in the library section of the code. As well as this, I also created a vector of each of the factors, which contained every variable that they could contain. This proved useful in this section, as I used it when making the different pie charts.

The First set of pie as seen in Figure 1.1 charts that I created all looked at the physical characteristics of the cap. These would be some of the most useful data, as they are the most distinct part of the mushroom. The pie charts showed me that the shape was dominated by a couple of types of shapes, the surface texture was relatively evenly distributed except for grooves, and the colours varied greatly.

The second set of pie charts, as seen in Figure 1.2, shows a variety of characteristics. The bruises are one of the features that I noted in the data cleaning section as being better suited as a logical variable type. The odour shows that about half have no odour, about a quarter have a foul odour, and the rest are distributed relatively evenly. This also shows that the spore print colour is mainly split between 3 colours.

In Figure 1.3, we see information about the gill-related data. Interestingly, most of them are very one-sided, with free, close and broad dominating the pie charts of gill attachment, spacing and size, respectively. The colours, on the other hand, are very well spread out, with a wide variety of each of them.

The next image, as shown in Figure 1.4, has very similar pie charts for both types of above/below data. It may be logical to assume that these are all from the same mushroom, i.e a stalk that is silky above the surface is also silky below the surface, but while this seems logical, these pie charts do not show us significant data to support this. Interesting to note that the stalk shape is very evenly split between the two options.

The next 4 pie charts, seen in Figure 1.5, describe characteristics of the veil and ring, two parts of the mushroom that can be found on the stalk beneath the cap. Again, most of these pie charts are almost completely dominated by 1 variable, and even the 4th is mostly a pendant ring type.

Figure 1.6 shows, among others, the most important data point: whether it is edible or not. We can see that although there are a lot more that are edible, it is not completely dominated by one type, as seen previously. This is very good, as it would be difficult to classify if there were very few data points of one type. This picture also shows the population and habitat, which, like others, have a couple of variable types that are much more prominent and a few that occur a lot less frequently.

Overall, the pie charts have shown us a lot of information about the proportions of data. It is possible that the features that are very dominated by one data type could become a problem in future classification or could be completely ignored. The pie charts do not give us a good idea of how the data points interact with each other, which is important to try to understand.

# Probability tables

Lines 430 - 464

I wanted to get a good idea of how some of the features interact with each other so that I was better informed for later algorithms that I could use. Most importantly was to compare them to the edible feature, as that was the main one that I needed to analyse.

The first table, as seen in Figure 2.1, tells us the different amounts of each type of cap shape that are edible or not. Figure 2.2 has used the prob.table function to convert this information into percentages. From both of these, we can see that most of the numbers are not that significant either way unluse there are very few of that specific type of variable, so it is less reliable. For example, in Figure 2.2, we see that if the cap shape is conical, there is a 100% chance of the mushroom not being edible; however, since there are only 4 entries of this type, it's questionable how valid this result is.

A similar thing can be seen in Figures 2.3 and 2.4; if the cap surface is described as 'grooves', then there is a 100% chance that the mushroom is not edible, but again, there are only 4 entries. Most of the other types with more entries are closer to 50%, suggesting that they will not be particularly good at suggesting whether a mushroom is edible.

Lastly, in Figures 2.5 and 2.6, we see that the colour of the cap can vary greatly, but is not specifically dominated by one colour, nor are most of them heavily weighted as edible or not. Interestingly, there is the expectation of Goldenrod4, which has over 100 entries, but all of them are inedible, which could suggest it will be important in later analysis.

The probability tables can be good at suggesting what variable may be useful in later algorithms to check for the edibility of mushrooms; however, in these tables, we can only compare edibility to one feature. Ideally, we want to compare if a mushroom is edible or not to as many features as possible to create a good idea of if th mushroom is edible.

# Cross-validation and bootstrap

Lines 465 - 483

In this section, I use linear regression using the cap colour as the predictor, which, from the probability tables above, I thought would be a reality good predictor of edibility. First is to split the data into a test and training set, where 50% of the data is in each. I use the set.seed so that the results are random by reproducible. I then create the model and analyse it by looking at the mean squared error, which is very low at 0.1825958.

To ensure that this is not a coincidence, I do the same again with a different set of random training and testing data. Again, I get a very similar and low MSE of 0.1819106. These numbers being this low means that the model has low error, making quite accurate predictions. The fact that they are very similar also implies that there is low variance.

However, I also think there are better functions that could be able to make better type of predictions on this type of data.

# Glm function

Lines 484 - 584

I decided to use a GLM function to help make predictions on this set of data. It's probably more suited to data of this type than linear regression and can also easily be converted to a graph to help make visualisation easier.

In Figure 3.1, I make a jitter plot comparing the colour of the caps to their edibility, from the tables above, as well as the linear regression functions, I find hints that these are likely to be correlated. I next use the GLM function and create a summary of the fit that can be seen in Figure 3.2. The significance code shows that almost all are highly significant in the data, with the exception of goldenrod,4, which has no significance; it also has an incredibly high standard error and p-values. This makes sense given what we've seen in the previous table, where although there are some number of goldenrod4 caps, they are not all edible and therefore not very likely to be useful to the data.

In Figure 3.3, I've plotted the data above when the function is true so that the points we see are only true. However, this does not show us a great deal as the plot looks the same in Figures 3.1 and 3.3.  The summary table for this is also fairly similar, although thought the big outlier of goldenrod4, this is most likely because all values are ot eadbily, so the data for it does not come into play.

I realise that the data would be a lot more useful if all the features could be included, so I attempt to use the GLM function to make a prediction with all features as predictors. Unfortunately, the algorithm did not converge, most likely because it was trying to use all features, even those that were not helpful. Figure 3.5.1 and 3.5.2 show the summary of the function above. There are plenty of NAs since it did not converge. This model can, however, be used to help find the chances of a specific mushroom being edible, although I would question its accuracy. I input various variables for each feature, and it tells me that there is a 2.900715e-12 chance of the mushroom being edible.

Even though I question the validity of this method for this type of data, given the previous results, I want to look at the hnp function to check the accuracy of the model. I use specifically the stalkColorBelowRing and stalkSurfaceBelowRing features, both independently, in Features 3.6 and 3.7 and together in Figure 3.8. Although the graphs stay almost completely within the envelopes, they are not very close to the top right-hand corner, which makes the validity of these graphs questionable.

# Leave-one-out Cross-Validation

Lines 585 - 609

Leave-one-out cross-validation looks at a reference variable and then compares the rest of the variables within a specific feature against another feature. In this case, I am comparing the cap shape to the edible to find out if there is much of a correlation between them.

The reference variable in this case is bell-shaped, and almost all of the variables are less likely to be edible than the bell shape, the only exception being sunken. This can be seen in Figure 4.1.

The error estimate of 0.2407 is relatively low; however, not as low as previously seen in other algorithms.

As with previous algorithms that I have used, it would be much more useful if we were able to look at different features at once, so I am not going to go into as much depth with this type of algorithm as I could.

# Fitting trees

Lines 610 - 677

The algorithm involved in fitting trees requires all values to be represented as factors so as not to mess up the original data files. I create a new one and convert all features, regardless of the original data type, into factors.

As I have precisely done, I use a seeded random number to create a reproducible random set for training and testing. I then create a tree, in Figure 5.1, I am told the values that are used in the tree, which are quite few, as well as there not being many nodes. In Figures 5.2 and 5.3, we can see the tree made as a test list as well as a diagram. There are very few features that are taken into account, which could be a problem, but it could imply that these features are very decisive for the data set.

To check that there is nothing wrong going on, I look at a table to check the accuracy of the tree that has been created against the test data. Interestingly, the accuracy is 100%; there are no false positives or false negatives.

Despite this, I would still like to check if there is an option to prune the data. We can see in Figure 5.5 that  there is not really a good case for pruning the tree; the number of current nodes is sufficient. Indeed, in Figure 5.6, we can see that the tree is identical to the tree in 5.,2 so there is no reason for the trees to be pruned.

There could be a slight concern for over-fitting since the data is 100% correct, but if this is the case, then the whole data file would be the problem, as the test data still gives 100% accuracy as well.

# Random forests

Lines 678 - 703

I think that it would be a good idea to also create random forests. These can be more robust prediction methods than just fitting trees, as they can take into account much more of the data and will compare different trees that are created.

In Figure 6.1, we can see that the forest contains 500 trees with 21 variables tried at each fit. These are all the different factors in the data file, which is what we want. Like the previous tree that was made, the forest has a 100% accuracy rating. This could imply that there are, in fact, factors that describe the edibility of a given mushroom very well; this is certainly true within our data set.

I then create a plot, as seen in Figure 6.2, to show the importance of each feature for constructing the above trees. It turns out that the odour is by far the most important feature, with others trailing behind until the least important feature of gill attachment is reached. Interestingly, the values of cap colour are relatively low and are significantly smaller than odour in predicting the edibility of a mushroom, describing the original guesses that I made that assumed the cap colour was fairly useful. This just goes to show the importance of finding the right method to describe the data given, and that methods that only looked at one type of variable were not as useful.

# Classification forests

Lines 704 - 756

I split the data into test and training again, using a seeded random number so that the value would be reproducible. I create a classification forest using the rpart function, and then I immediately create a confusion matrix, seen in Figure 7.1, with statistics to look at the forests that have been created. This tree misclassifies 11 mushrooms as poisonous when they were really edible. This is very accurate (and for someone who is trying to use this algorithm, it's good that it's on the safe side). The Kappa value and accuracy show how well this model does, and from a quick observation doesn't look quite as susceptible to overfitting. Most values are close to 1, which is very good.

I create a new forest to compare the confusion matrix with the one above. This forest has a 100% accuracy, no misclassifications unfortunately, even though this looks like a good thing it it could be a result of over-fitting, it could be that the data is very reliant on a few data points, as seen in some of the trees above. This can be seen in Figure 7.2

I then compare different trees using different random subsets of the data file to make sure that the prediction is very robust, as seen in Figure 7.3. It shows that that opposed to what I'd initially thought, the data is unlikely to be overfitting; instead, there are just a few elements that can describe the edibility of the mushrooms in the data set very well. However, there is a chance that this would be specific to the mushrooms within this data set and could not be used on other sets of data.

These show that although it may seem like the models are subject to overfitting, this is unlikely within the data set; if other data sets were to be used, this could be different.

●

# Conclusion

Most of the methods that I tried were somewhat useful, although the best ones were GLM, trees and forests. The pie charts and probability tables, while good for getting an initial idea of the data, were not particularly useful for any analysis of it, as they did not show much and could not represent more than 2 features interacting with each other.

Some of the other methods, such as the linear regression, cross-validation and bootstrap were not very useful for various reasons, such as failing to account for more than 1 other factor and not being useful for categorical data.

The GLM functions were definitely helpful, but they had the problem of scale. I wanted to analyse more than one feature against edible, so I tried to do all of them. This just meant that the function did not converge, so it didn't show me anything.

The trees were very useful in showing the way of categorising edible mushrooms from non-edible mushrooms and but it was always questionable how accurate they were, especially with how they appear to have 100% accuracy.

The forests seemed to alleviate the above problem, and so I think that they where some of the best algorithms that I had for analysing this data.
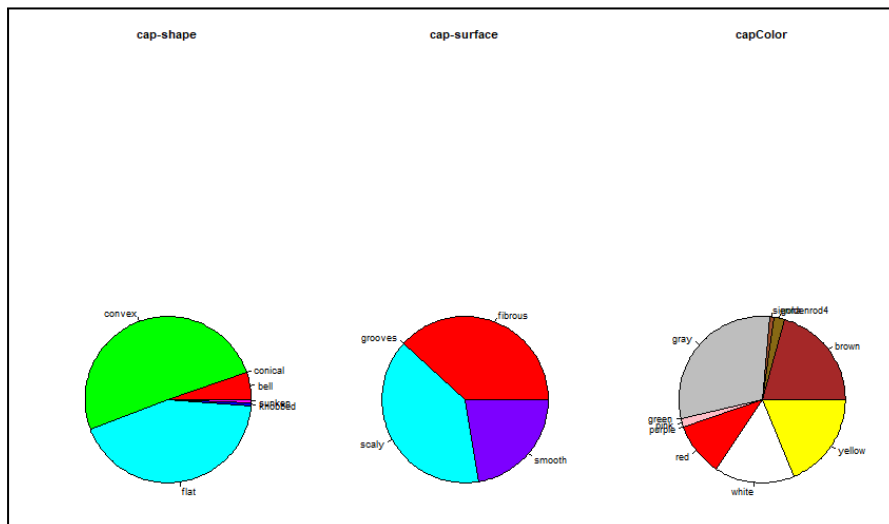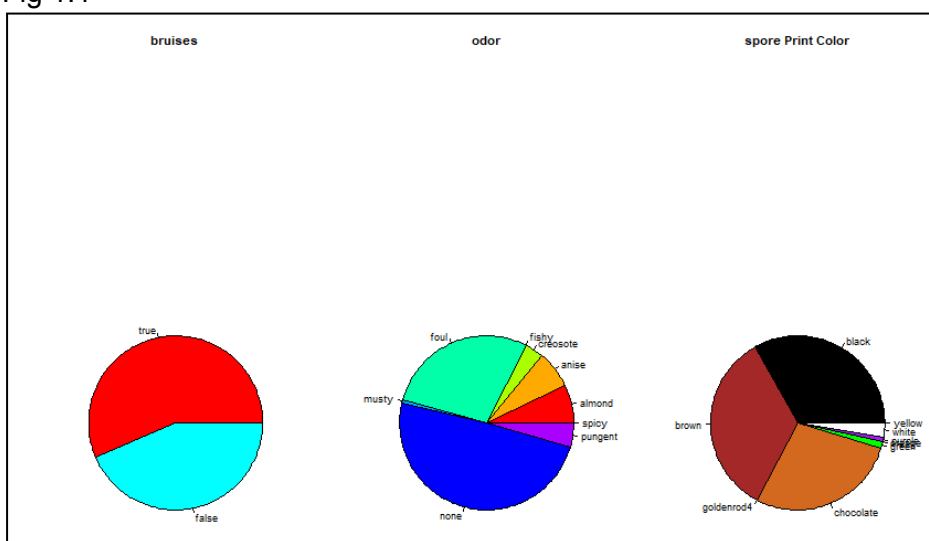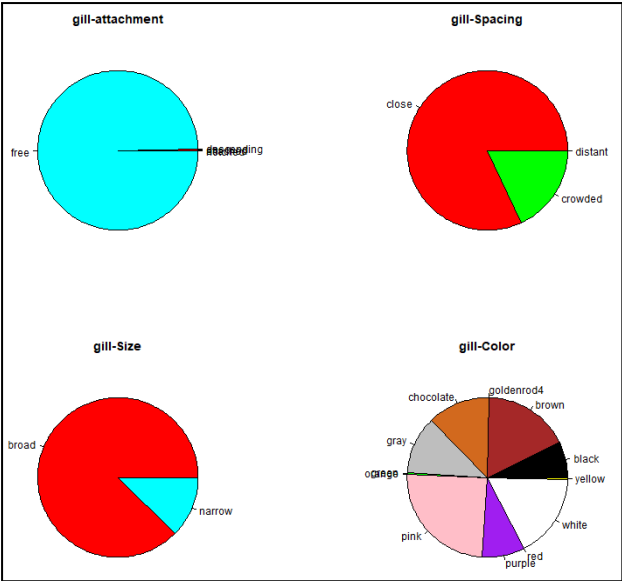
# Images

# 1. Pie Charts



Fig 1.1



Fig 1.2

Fig 1.3



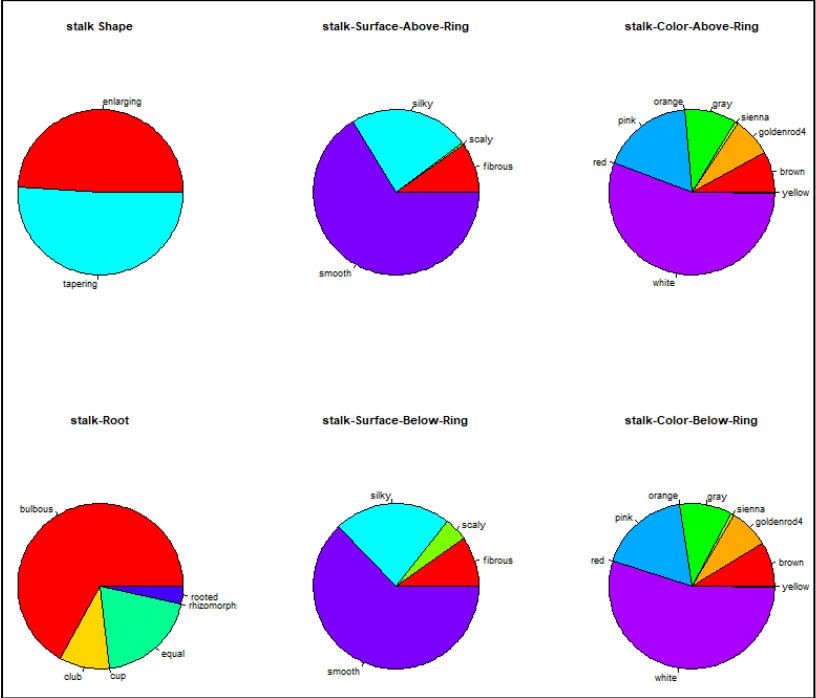Fig 1.4
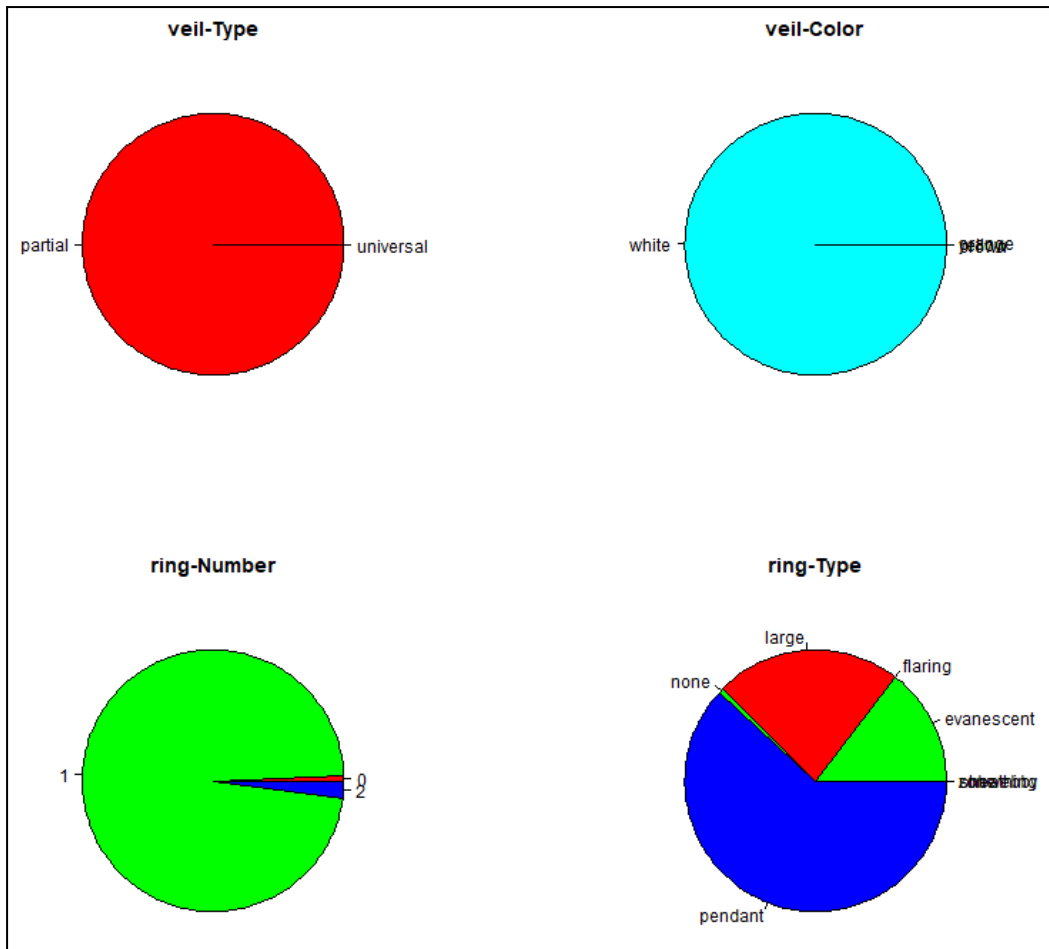
Fig 1.5



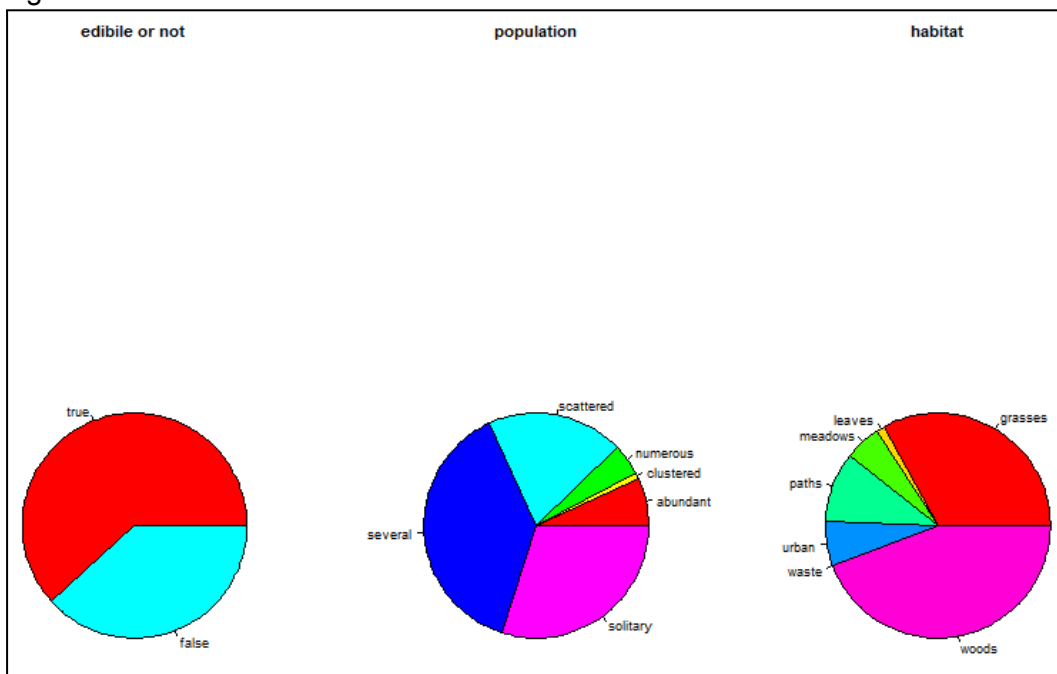Fig 1.6

# 2. Probability Tables

```
          edible
capShape  FALSE  TRUE
  bell       40   260
  conical     4     0
  convex   1124  1716
  flat      972  1460
  knobbed    16    20
  sunken      0    32
```
Fig 2.1

```
               edible
capShape      FALSE       TRUE
  bell     0.1333333  0.8666667
  conical  1.0000000  0.0000000
  convex   0.3957746  0.6042254
  flat     0.3996711  0.6003289
  knobbed  0.4444444  0.5555556
  sunken   0.0000000  1.0000000
```
Fig 2.2

```
                edible
cap-Surface  FALSE  TRUE
  fibrous      744  1416
  grooves        4     0
  scaly        860  1360
  smooth       548   712
```
Fig 2.3

```
                 edible
cap-Surface      FALSE       TRUE
  fibrous    0.3444444  0.6555556
  grooves    1.0000000  0.0000000
  scaly      0.3873874  0.6126126
  smooth     0.4349206  0.5650794
```
Fig 2.4

```
              edible
capColor      FALSE  TRUE
  brown         140  1024
  goldenrod4    120     0
  gray          808   888
  pink           88     8
  red            12   576
  sienna         12    32
  white         320   560
  yellow        656   400
```
Fig 2.5

```
                  edible
capColor          FALSE        TRUE
  brown       0.12027491  0.87972509
  goldenrod4  1.00000000  0.00000000
  gray        0.47641509  0.52358491
  pink        0.91666667  0.08333333
  red         0.02040816  0.97959184
  sienna      0.27272727  0.72727273
  white       0.36363636  0.63636364
  yellow      0.62121212  0.37878788
```
Fig 2.6

# 3. GLM Functions


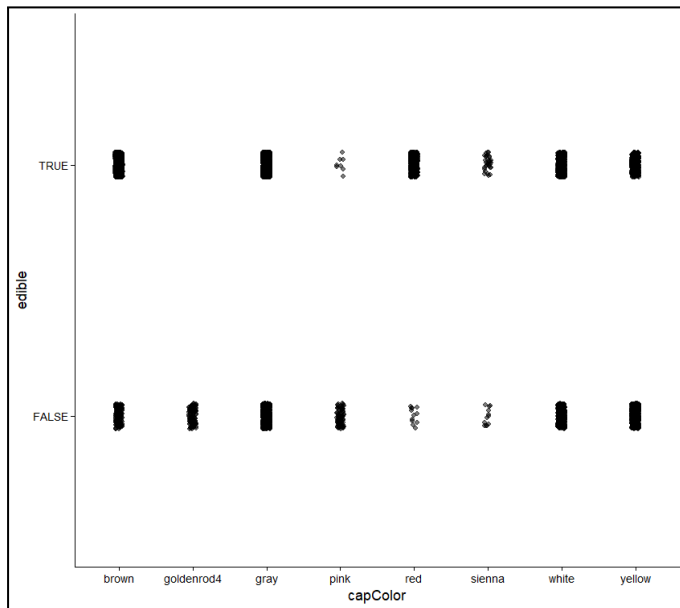
Fig 3.1

```
Call:
glm(formula = edible ~ capColor, family = binomial, data = df)

Coefficients:
                    Estimate Std. Error z value Pr(>|z|)
(Intercept)          1.98983    0.09011  22.083  < 2e-16 ***
capColorgoldenrod4 -18.55590  219.04748  -0.085  0.93249
capColorgray        -1.89542    0.10239 -18.512  < 2e-16 ***
capColorpink        -4.38772    0.38011 -11.543  < 2e-16 ***
capColorred          1.88137    0.30527   6.163 7.14e-10 ***
capColorsienna      -1.00900    0.35029  -2.880  0.00397 **
capColorwhite       -1.43021    0.11415 -12.529  < 2e-16 ***
capColoryellow      -2.48453    0.11020 -22.546  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 7506.9  on 5643  degrees of freedom
Residual deviance: 5981.5  on 5636  degrees of freedom
AIC: 5997.5

Number of Fisher Scoring iterations: 15
```

Fig 3.2

Fig 3.3

```
Call:
glm(formula = edible ~ capColor, data = df)

Coefficients:
                   Estimate Std. Error t value Pr(>|t|)
(Intercept)         0.87973    0.01251  70.340  < 2e-16 ***
capColorgoldenrod4 -0.87973    0.04091 -21.504  < 2e-16 ***
capColorgray       -0.35614    0.01624 -21.928  < 2e-16 ***
capColorpink       -0.79639    0.04531 -17.577  < 2e-16 ***
capColorred         0.09987    0.02159   4.626 3.81e-06 ***
capColorsienna     -0.15245    0.06553  -2.326     0.02 *
capColorwhite      -0.24336    0.01906 -12.768  < 2e-16 ***
capColoryellow     -0.50094    0.01813 -27.624  < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for gaussian family taken to be 0.1820715)

    Null deviance: 1332.4  on 5643  degrees of freedom
Residual deviance: 1026.2  on 5636  degrees of freedom
AIC: 6413.2

Number of Fisher Scoring iterations: 2
```
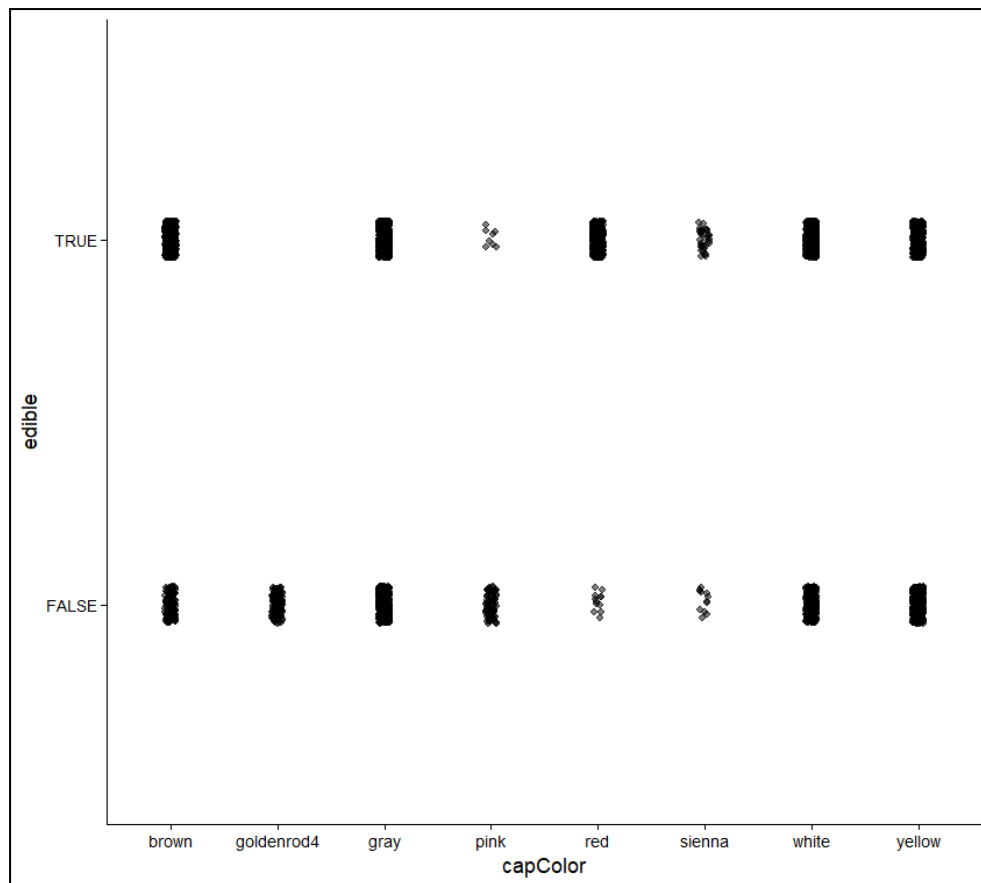
Fig 3.4

```
Call:
glm(formula = edible ~ ., family = binomial, data = df)

Coefficients: (13 not defined because of singularities)
                             Estimate Std. Error z value Pr(>|z|)
(Intercept)                 1.328e+02  4.448e+05   0.000        1
capShapeconical            -3.536e-09  2.079e+05   0.000        1
capShapeconvex             -4.200e-09  2.897e+04   0.000        1
capShapeflat               -4.171e-09  3.024e+04   0.000        1
capShapeknobbed            -6.455e-09  7.687e+04   0.000        1
capShapesunken             -7.083e-07  8.244e+04   0.000        1
capSurfacegrooves          -3.050e-10  2.521e+05   0.000        1
capSurfacescaly            -3.078e-10  1.195e+04   0.000        1
capSurfacesmooth            1.218e-08  1.705e+04   0.000        1
capColorgoldenrod4         -1.822e-10  4.290e+04   0.000        1
capColorgray                5.199e-09  1.540e+04   0.000        1
capColorpink               -8.112e-10  4.693e+04   0.000        1
capColorred                -7.791e-09  1.953e+04   0.000        1
capColorsienna              1.018e-08  7.006e+04   0.000        1
capColorwhite              -4.366e-09  1.964e+04   0.000        1
capColoryellow              5.716e-09  2.043e+04   0.000        1
bruisesTRUE                -1.063e+02  4.118e+05   0.000        1
odoranise                   1.852e-08  2.518e+05   0.000        1
odorcreosote               -1.594e+02  4.174e+05   0.000        1
odorfoul                   -1.063e+02  2.172e+05   0.000        1
odormusty                  -5.313e+01  1.479e+05   0.000        1
odornone                   -5.313e+01  2.116e+05   0.000        1
odorpungent                -1.615e-05  2.529e+05   0.000        1
gillAttachmentfree          1.211e-13  1.187e+05   0.000        1
gillSpacingcrowded         -2.918e-12  5.140e+04   0.000        1
gillSizenarrow             -5.313e+01  2.100e+05   0.000        1
gillColorbrown             -2.531e-08  2.249e+04   0.000        1
gillColorchocolate         -2.810e-08  2.487e+04   0.000        1
gillColorgray              -2.804e-08  2.546e+04   0.000        1
gillColorgreen             -2.938e-08  9.175e+04   0.000        1
gillColorpink              -2.642e-08  2.231e+04   0.000        1
gillColorpurple            -3.197e-08  2.796e+04   0.000        1
gillColorwhite             -3.072e-08  2.403e+04   0.000        1
gillColoryellow            -3.072e-08  1.100e+05   0.000        1
stalkShapetapering         -4.733e-06  7.481e+04   0.000        1
stalkRootclub              -5.313e+01  2.497e+05   0.000        1
stalkRootequal             -5.313e+01  2.190e+05   0.000        1
stalkRootrooted            -5.313e+01  2.463e+05   0.000        1
stalkSurfaceAboveRingscaly -1.063e+02  4.293e+05   0.000        1
stalkSurfaceAboveRingsilky -5.313e+01  2.222e+05   0.000        1
stalkSurfaceAboveRingsmooth 8.353e-09  2.144e+04   0.000        1
stalkSurfaceBelowRingscaly         NA         NA      NA       NA
stalkSurfaceBelowRingsilky         NA         NA      NA       NA
```

Fig 3.5.1

```
stalkSurfaceBelowRingscaly           NA         NA      NA       NA
stalkSurfaceBelowRingsilky           NA         NA      NA       NA
stalkSurfaceBelowRingsmooth   8.003e-09  2.144e+04   0.000        1
stalkColorAboveRinggoldenrod4 -2.003e-14 2.423e+04   0.000        1
stalkColorAboveRinggray      -5.845e-10  3.205e+04   0.000        1
stalkColorAboveRingpink      -7.008e-15  2.423e+04   0.000        1
stalkColorAboveRingsienna            NA         NA      NA       NA
stalkColorAboveRingwhite     -4.102e-10  3.205e+04   0.000        1
stalkColorAboveRingyellow     1.063e+02  6.365e+05   0.000        1
stalkColorBelowRinggoldenrod4 -6.669e-16 2.423e+04   0.000        1
stalkColorBelowRinggray      -9.714e-16  3.205e+04   0.000        1
stalkColorBelowRingpink      -6.142e-16  2.423e+04   0.000        1
stalkColorBelowRingsienna            NA         NA      NA       NA
stalkColorBelowRingwhite     -6.513e-16  3.205e+04   0.000        1
stalkColorBelowRingyellow            NA         NA      NA       NA
veilColoryellow                      NA         NA      NA       NA
ringNumber                           NA         NA      NA       NA
ringTypelarge                        NA         NA      NA       NA
ringTypenone                         NA         NA      NA       NA
ringTypependant               5.313e+01  2.334e+05   0.000        1
sporePrintColorbrown         -7.559e-10  1.164e+04   0.000        1
sporePrintColorchocolate             NA         NA      NA       NA
sporePrintColorgreen         -5.313e+01  9.096e+04  -0.001        1
sporePrintColorpurple         2.388e-06  7.362e+04   0.000        1
sporePrintColorwhite                 NA         NA      NA       NA
populationclustered                  NA         NA      NA       NA
populationnumerous            4.770e-06  4.064e+04   0.000        1
populationscattered           4.701e-06  2.570e+04   0.000        1
populationseveral             4.701e-06  3.484e+04   0.000        1
populationsolitary            4.701e-06  3.611e+04   0.000        1
habitatleaves                        NA         NA      NA       NA
habitatmeadows               -3.822e-09  2.947e+04   0.000        1
habitatpaths                 -7.854e-14  2.188e+04   0.000        1
habitaturban                  4.558e-14  3.054e+04   0.000        1
habitatwoods                 -3.646e-14  2.361e+04   0.000        1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 7.5069e+03  on 5643  degrees of freedom
Residual deviance: 3.2744e-08  on 5581  degrees of freedom
AIC: 126

Number of Fisher Scoring iterations: 25
```
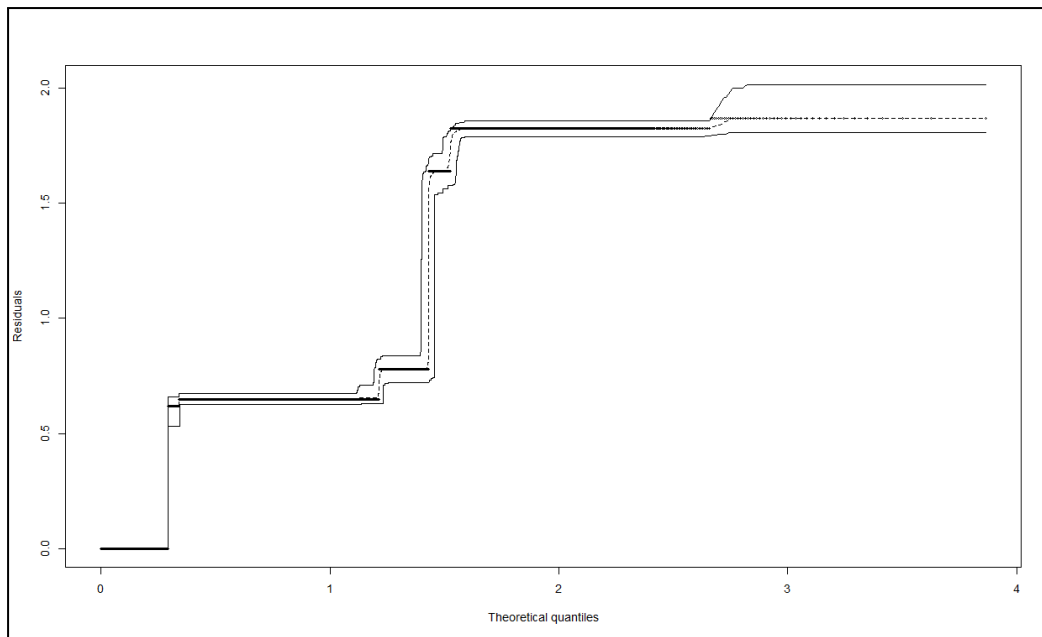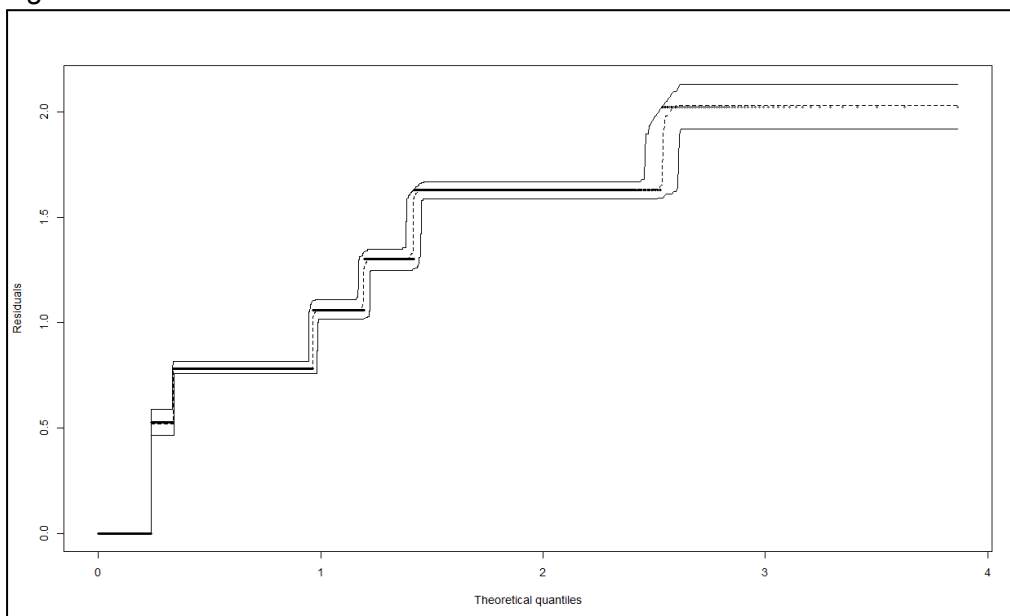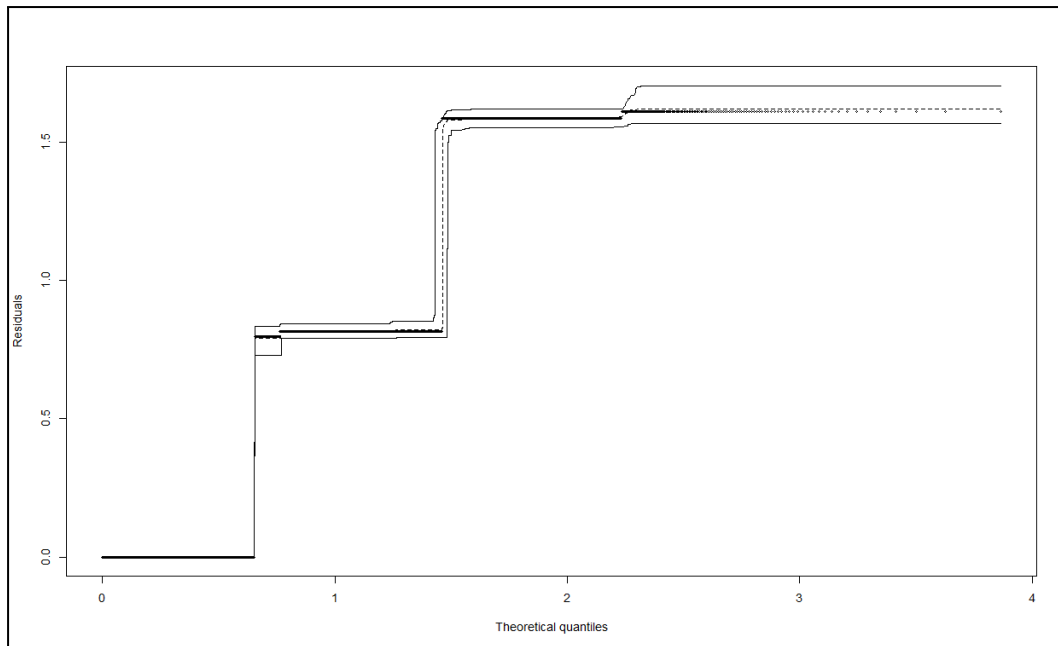
Fig 3.5.2

Fig 3.6



Fig 3.7

Fig 3.8

# 4. Leave-one-out Cross Validation

Fig 4.1

```
            (Intercept) df$capShapeconical    df$capShapeconvex       df$capShapeflat
              0.8666667           -0.8666667           -0.2624413           -0.2663377
df$capShapeknobbed    df$capShapesunken
             -0.3111111            0.1333333
```

# 5. Fitting Trees

```
Classification tree:
tree(formula = edible ~ ., data = df_class)
Variables actually used in tree construction:
[1] "odor"            "sporePrintColor" "population"
Number of terminal nodes:  4
Residual mean deviance:  0 = 0 / 5640
Misclassification error rate: 0 = 0 / 5644
```
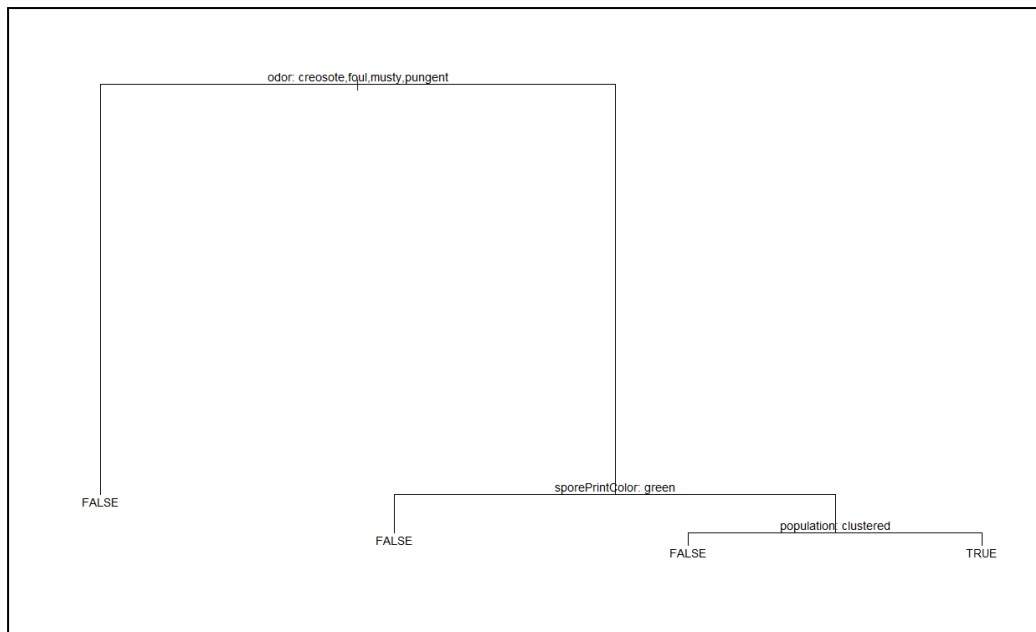
Fig 5.1

Fig 5.2

```
node), split, n, deviance, yval, (yprob)
      * denotes terminal node

 1) root 5644 7507.0 TRUE ( 0.381999 0.618001 )
   2) odor: creosote,foul,musty,pungent 2068    0.0 FALSE ( 1.000000 0.000000 ) *
   3) odor: almond,anise,none 3576  825.8 TRUE ( 0.024609 0.975391 )
     6) sporePrintColor: green 72    0.0 FALSE ( 1.000000 0.000000 ) *
     7) sporePrintColor: black,brown,purple,white 3504  204.4 TRUE ( 0.004566 0.995434 )
     14) population: clustered 16    0.0 FALSE ( 1.000000 0.000000 ) *
     15) population: abundant,numerous,scattered,several,solitary 3488    0.0 TRUE ( 0.000000 1.000000 ) *
```

Fig 5.3

```
            edible.test
tree.pred FALSE TRUE
   FALSE  1057    0
   TRUE      0 1765
```
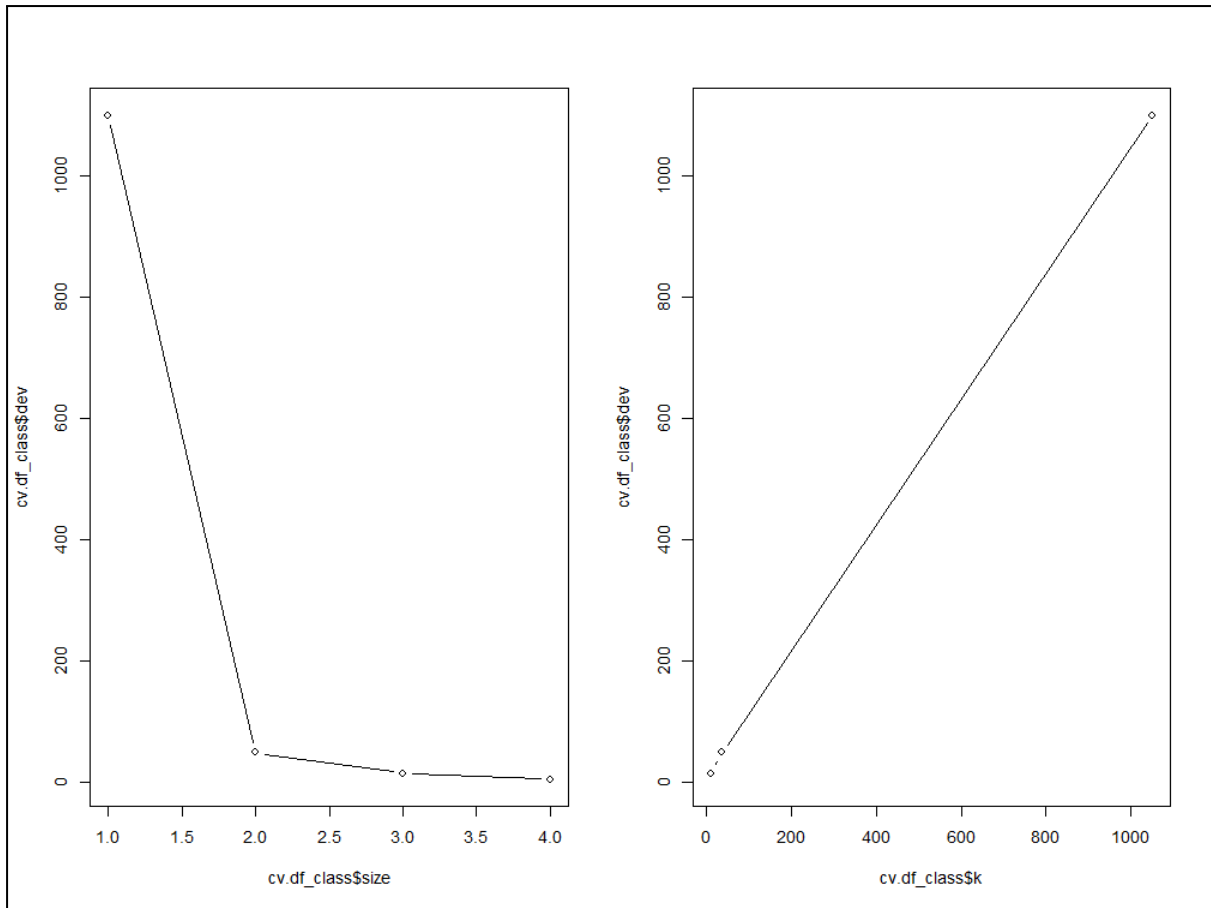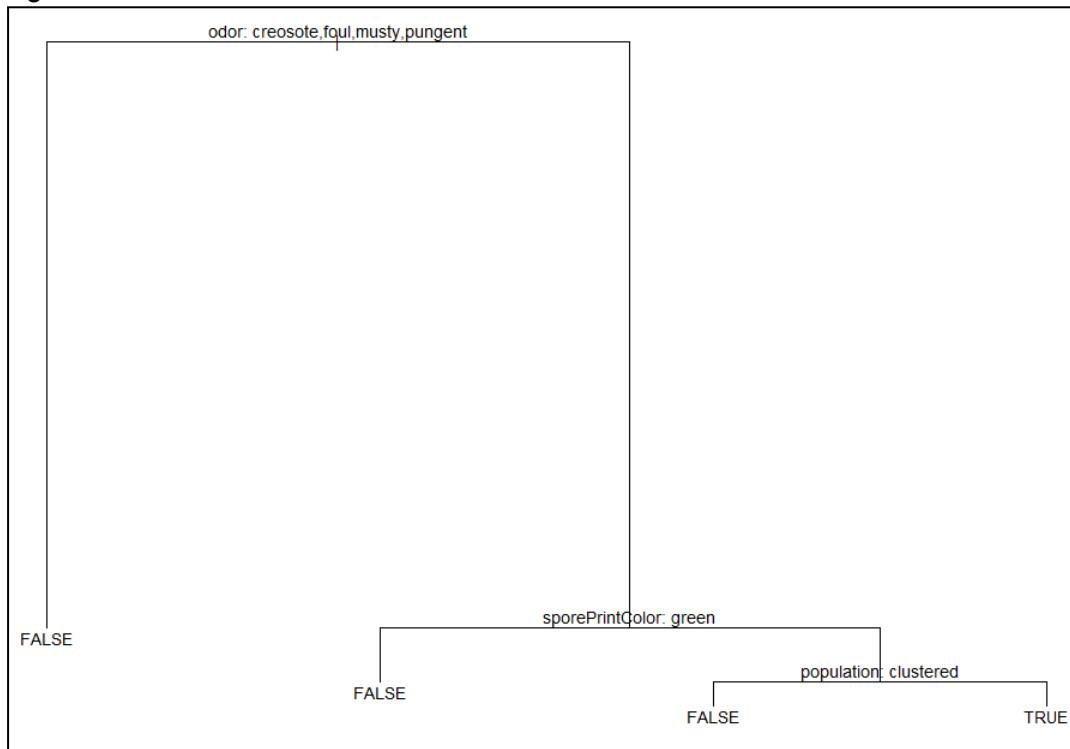
Fig 5.4

Fig 5.5



Fig 5.6

# 6. Random Forest

```
Call:
 randomForest(formula = df_class$edible ~ ., data = df_class,      mtry = 21, importance = TRUE, subset = train)
               Type of random forest: classification
                     Number of trees: 500
No. of variables tried at each split: 21

        OOB estimate of  error rate: 0%
Confusion matrix:
      FALSE TRUE class.error
FALSE  1099    0           0
TRUE      0 1723           0
```
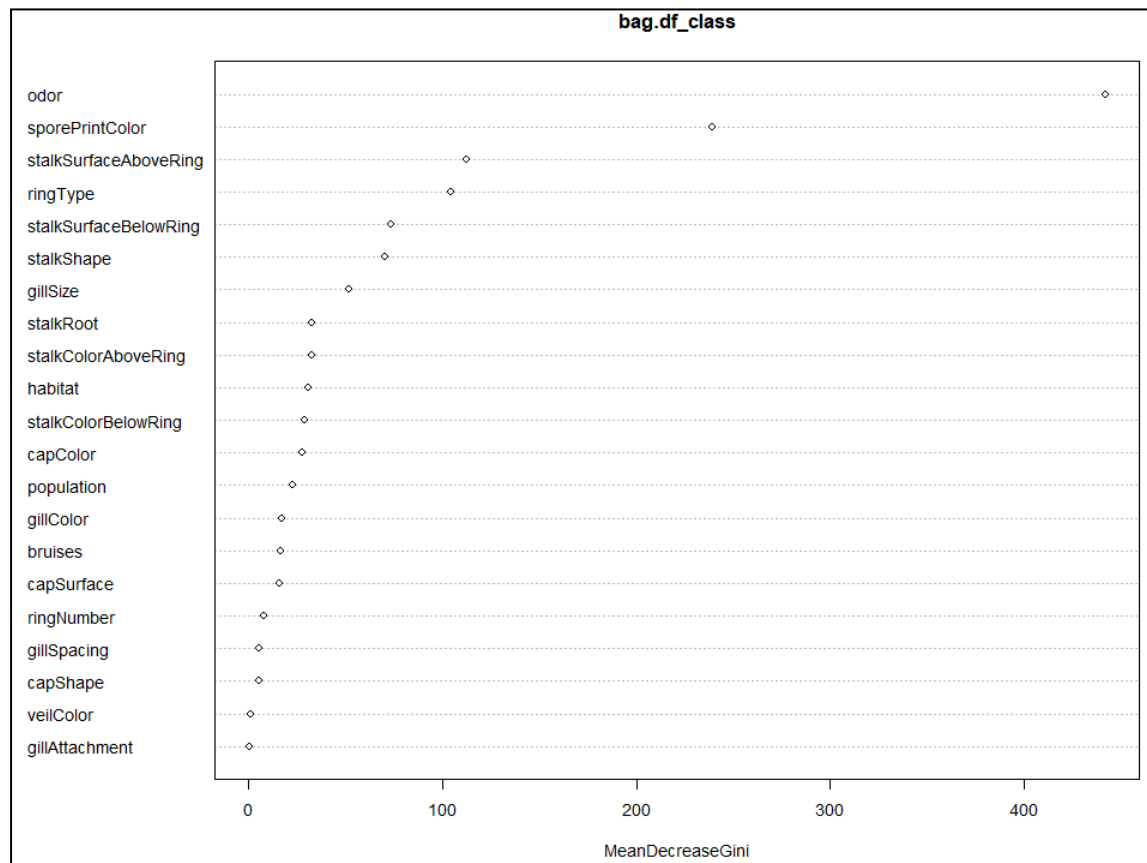
6.1



Fig 6.2

# 7. Classification Forest

```
Confusion Matrix and Statistics


preds_dtree FALSE TRUE
      FALSE  1067    0
      TRUE     11 1744


               Accuracy : 0.9961
                 95% CI : (0.993, 0.9981)
    No Information Rate : 0.618
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 0.9917

 Mcnemar's Test P-Value : 0.002569

            Sensitivity : 0.9898
            Specificity : 1.0000
         Pos Pred Value : 1.0000
         Neg Pred Value : 0.9937
             Prevalence : 0.3820
         Detection Rate : 0.3781
   Detection Prevalence : 0.3781
      Balanced Accuracy : 0.9949

       'Positive' Class : FALSE
```

Fig 7.1

```
Confusion Matrix and Statistics


preds_rf FALSE TRUE
   FALSE  1078    0
   TRUE      0 1744


               Accuracy : 1
                 95% CI : (0.9987, 1)
    No Information Rate : 0.618
    P-Value [Acc > NIR] : < 2.2e-16

                  Kappa : 1

 Mcnemar's Test P-Value : NA

            Sensitivity : 1.000
            Specificity : 1.000
         Pos Pred Value : 1.000
         Neg Pred Value : 1.000
             Prevalence : 0.382
         Detection Rate : 0.382
   Detection Prevalence : 0.382
      Balanced Accuracy : 1.000

       'Positive' Class : FALSE
```

Fig 7.2

```
Random Forest

5644 samples
  21 predictor
   2 classes: 'FALSE', 'TRUE'

No pre-processing
Resampling: Cross-Validated (10 fold)
Summary of sample sizes: 5079, 5079, 5079, 5081, 5079, 5080, ...
Resampling results across tuning parameters:

  mtry  Accuracy   Kappa
   2    0.9751886  0.9466907
  39    1.0000000  1.0000000
  76    1.0000000  1.0000000

Accuracy was used to select the optimal model using the largest value.
The final value used for the model was mtry = 39.
```

Fig 7.3

# Code

```r
1.   #########introduction###################################################
2.   #this will only use a maximium of 80 charcters across for ease of use
3.   #we are told by the data set that there are some missing values from stalk-root
4.   #I have decided to remove these as there are still several thousand lines of
5.   #data remaining.#
6.   #the colours sometimes refer to cinnamon and buff which R does not recognize as
7.   #colours, I have decided to instead refer to these as sienna and goldenrod4
8.   #instead as these are recognized by R as colours and will make whatever algrithm
9.   #I must apply to them easier to use.
10.  ######library########
11.  install.packages("tidyverse")
12.  library(tidyverse)
13.  install.packages("hnp")
14.  library(hnp)
15.  install.packages("ISLR2")
16.  library(ISLR2)
17.  install.packages("boot")
18.  library(boot)
19.  install.packages("tree")
20.  library(tree)
21.  install.packages("randomForest")
22.  library(randomForest)
23.  install.packages("caret")
24.  library(caret)
25.  install.packages("dplyr")
26.  library(dplyr)
27.  install.packages("caTools")
28.  library(caTools)
29.  install.packages("rpart")
30.  library(rpart)
31.  colNames <- c("edible", "capShape", "capSurface", "capColor", "bruises", "odor",
32.          "gillAttachment", "gillSpacing", "gillSize", "gillColor",
33.          "stalkShape", "stalkRoot", "stalkSurfaceAboveRing",
34.          "stalkSurfaceBelowRing", "stalkColorAboveRing",
35.          "stalkColorBelowRing", "veilType", "veilColor", "ringNumber",
36.          "ringType", "sporePrintColor", "population", "habitat")
37.               #creates a string of column titles
38.  capShape_variables <- c("bell", "conical", "convex", "flat",
39.              "knobbed", "sunken")
40.  capSurface_variables <- c("fibrous", "grooves", "scaly", "smooth")
41.  capColor_variables <- c("brown", "goldenrod4", "sienna", "gray", "green",
42.                "pink", "purple", "red", "white", "yellow")
43.  odor_variables <- c("almond", "anise", "creosote", "fishy", "foul", "musty",
44.              "none", "pungent", "spicy")
45.  gillAttachment_variables <- c("attached", "descending", "free", "notched")
46.  gillSpacing_variables <- c("close", "crowded", "distant")
47.  gillSize_variables <- c("broad", "narrow")
48.  gillColor_variables <- c("black", "brown", "goldenrod4", "chocolate",
49.               "gray","green", "orange", "pink", "purple","red",
50.               "white" ,"yellow")
51.  stalkShape_variables <- c("enlarging", "tapering")
52.  stalkRoot_variables <- c("bulbous", "club", "cup", "equal", "rhizomorphs",
53.              "rooted", NA)
54.  stalkSurfaceAboveRing_variables <- c("fibrous", "scaly", "silky", "smooth")
55.  stalkSurfaceBelowRing_variables <- c("fibrous", "scaly", "silky", "smooth")
56.  stalkColorAboveRing_variables <- c("brown", "goldenrod4", "sienna", "gray",
57.                 "orange", "pink", "red", "white", "yellow")
58.  stalkColorBelowRing_variables <- c("brown", "goldenrod4", "sienna", "gray",
59.                 "orange", "pink", "red", "white", "yellow")
60.  veilType_variables <- c("partial", "universal")
61.  veilColor_variables <- c("brown", "orange","white", "yellow")
62.  ringNumber_variables <- c(0,1,2)
63.  ringType_variables <- c("cobwebby", "evanescent", "flaring", "large", "none",
64.              "pendant", "sheathing", "zone")
65.  sporePrintColor_variables <- c("black", "brown", "goldenrod4", "chocolate",
66.                 "green", "orange", "purple", "white", "yellow")
67.  population_variables <- c("abundant", "clustered", "numerous", "scattered",
68.               "several", "solitary")
69.  habitat_variables <- c("grasses", "leaves", "meadows", "paths", "urban",
70.              "waste", "woods")
71.
72.  NumberOf <- function(data, check){
73.    n = 0
74.    for (i in data){
75.      if (i == check){
76.        n = n+1
77.      }
78.    }
```

```
79.     return(n)
80.   }
81.
82.   #####data cleaning#######
83.
84.   #loading the data into R using a working directory
85.   df_original <- read.csv("agaricus-lepiotaData.data", header = FALSE)
86.   #the title of the columns as well as the contents are not very useful,
87.   #the agaricus-lepiota.names file (among other things) tells use
88.   #what each of the letters means
89.
90.   df_dirty<- df_original#creates a new data set so the
91.                         #original will remain untouched
92.   colnames(df_dirty) <- colNames#makes the column
93.                         #titles what is in the array above
94.
95.   summary(df_dirty) #gives us a summary of the data for each column
96.   #all are characters checking the agaricus-lepiota.names file again most of these
97.   #are correct however there are a couple the could be turned into logical or
98.   #numerical variables and for the ones that are not changing the letter what it
99.   #is actually describing will massivly help readibility
100.
101.  #we are told that the data set has some missing values in the StalkRoot column,
102.  #for now they will be converted to NA's to highlight they are missing
103.  #the following checks that all instance of the columns contains what is expected
104.  #and changes it to a better data type or variable
105.  #edible
106.  sum(!df_dirty$edible %in% c("e", "p"))
107.  df_dirty$edible <- df_dirty$edible == "e"
108.
109.  #capShape
110.  sum(!df_dirty$capShape %in% c("b", "c", "x", "f", "k", "s"))
111.  df_dirty$capShape <- capShape_variables[match(df_dirty$capShape,
112.                                  c("b", "c", "x","f", "k", "s"))]
113.  #capSurface
114.  sum(!df_dirty$capSurface %in% c("f", "g", "y", "s"))
115.  df_dirty$capSurface <- capSurface_variables[match(df_dirty$capSurface,
116.                                  c("f", "g","y", "s"))]
117.  #capColor
118.  sum(!df_dirty$capColor %in%
119.      c("n", "b", "c", "g", "r", "p", "u", "e", "w", "y"))
120.  df_dirty$capColor <- capColor_variables[match(df_dirty$capColor,
121.              c("n", "b", "c", "g","r", "p", "u", "e","w", "y"))]
122.  #bruises
123.  sum(!df_dirty$bruises %in% c("t", "f"))
124.  df_dirty$bruises <- df_dirty$bruises == "t"
125.
126.  #odor
127.  sum(!df_dirty$odor %in% c("a", "l", "c", "y", "f", "m", "n", "p", "s"))
128.  df_dirty$odor <- odor_variables[match(df_dirty$odor,
129.                      c("a", "l", "c", "y", "f", "m", "n","p", "s"))]
130.  #gillAttachment
131.  sum(!df_dirty$gillAttachment %in% c("a", "d", "f", "n"))
132.  df_dirty$gillAttachment <-
133.    gillAttachment_variables[match(df_dirty$gillAttachment,
134.                                  c("a","d", "f", "n"))]
135.  #gillSpacing
136.  sum(!df_dirty$gillSpacing %in% c("c", "w", "d"))
137.  df_dirty$gillSpacing <- gillSpacing_variables[match(df_dirty$gillSpacing,
138.                                  c("c", "w", "d"))]
139.  #gillSize
140.  sum(!df_dirty$gillSize %in% c("b", "n"))
141.  df_dirty$gillSize <- gillSize_variables[match(df_dirty$gillSize,
142.                                  c("b", "n"))]
143.  #gillColor
144.  sum(!df_dirty$gillColor %in%
145.      c("k", "n", "b", "h", "g", "r", "o", "p", "u", "e", "w", "y"))
146.  df_dirty$gillColor <- gillColor_variables[match(df_dirty$gillColor,
147.                          c("k", "n", "b", "h", "g", "r",
148.                            "o", "p", "u", "e", "w", "y"))]
149.  #stalkShape
150.  sum(!df_dirty$stalkShape %in% c("e", "t"))
151.  df_dirty$stalkShape <- stalkShape_variables[match(df_dirty$stalkShape,
152.                                  c("e", "t"))]
153.  #stalkRoot
154.  sum(!df_dirty$stalkRoot %in% c("b", "c", "u", "e", "z", "r", "?"))
155.  df_dirty$stalkRoot <- stalkRoot_variables[match(df_dirty$stalkRoot,
156.                          c("b", "c", "u", "e", "z", "r", "?"))]
157.  #stalkSurfaceAboveRing
158.  sum(!df_dirty$stalkSurfaceAboveRing %in% c("f", "y", "k", "s"))
159.  df_dirty$stalkSurfaceAboveRing <-
160.    stalkSurfaceAboveRing_variables[match(df_dirty$stalkSurfaceAboveRing,
161.                          c("f", "y", "k", "s"))]
162.  #stalkSurfaceBelowRing
163.  sum(!df_dirty$stalkSurfaceBelowRing %in% c("f", "y", "k", "s"))
164.  df_dirty$stalkSurfaceBelowRing <-
```

```r
165.    stalkSurfaceBelowRing_variables[match(df_dirty$stalkSurfaceBelowRing,
166.                        c("f", "y", "k", "s"))]
167.    #stalkColorAboveRing
168.    sum(!df_dirty$stalkColorAboveRing %in%
169.        c("n", "b", "c", "g", "o", "p", "e", "w", "y"))
170.    df_dirty$stalkColorAboveRing <-
171.     stalkColorAboveRing_variables[match(df_dirty$stalkColorAboveRing,
172.                        c("n", "b", "c", "g", "o",
173.                         "p", "e", "w", "y"))]
174.    #stalkColorBelowRing
175.    sum(!df_dirty$stalkColorBelowRing %in%
176.        c("n", "b", "c", "g", "o", "p", "e", "w", "y"))
177.    df_dirty$stalkColorBelowRing <-
178.     stalkColorBelowRing_variables[match(df_dirty$stalkColorBelowRing,
179.                        c("n", "b", "c", "g", "o", "p", "e", "w", "y"))]
180.    #veilType
181.    sum(!df_dirty$veilType %in% c("p", "u"))
182.    df_dirty$veilType <- veilType_variables[match(df_dirty$veilType,
183.                        c("p", "u"))]
184.    #veilColor
185.    sum(!df_dirty$veilColor %in% c("n", "o", "w", "y"))
186.    df_dirty$veilColor <- veilColor_variables[match(df_dirty$veilColor,
187.                        c("n", "o", "w", "y"))]
188.    #ringNumber
189.    sum(!df_dirty$ringNumber %in% c("n", "o", "t"))
190.    df_dirty$ringNumber <- ringNumber_variables[match(df_dirty$ringNumber,
191.                        c("n", "o", "t"))]
192.    #ringType
193.    sum(!df_dirty$ringType %in% c("c", "e", "f", "l", "n", "p", "s", "z"))
194.    df_dirty$ringType <- ringType_variables[match(df_dirty$ringType,
195.                        c("c", "e", "f", "l", "n", "p", "s", "z"))]
196.    #sporePrintColor
197.    sum(!df_dirty$sporePrintColor %in%
198.        c("k", "n", "b", "h", "r", "o", "u", "w", "y"))
199.    df_dirty$sporePrintColor<-
200.     sporePrintColor_variables[match(df_dirty$sporePrintColor,
201.                        c("k", "n", "b", "h", "r", "o", "u", "w", "y"))]
202.    #population
203.    sum(!df_dirty$population %in% c("a", "c", "n", "s", "v", "y"))
204.    df_dirty$population <- population_variables[match(df_dirty$population,
205.                        c("a", "c", "n", "s", "v", "y"))]
206.    #habitat
207.    sum(!df_dirty$habitat %in% c("g", "l", "m", "p", "u","w", "d"))
208.    df_dirty$habitat <- habitat_variables[match(df_dirty$habitat,
209.                        c("g", "l", "m", "p", "u","w", "d"))]
210.
211.    sum(is.na(df_dirty)) # this checks the number of NA's
212.    # there are 2,480
213.    df <- na.omit(df_dirty) #removes NA's
214.    sum(is.na(df)) # this checks the number of NA's again
215.    #there are now none
216.
217.    sum(duplicated(df))
218.    #this checks if there are any duplicated rows in the data set
219.    #there are none
220.
221.
222.    ####pie charts#####
223.
224.    # the following creates pie charts showing the proportions of
225.    #the variables for each of the cap related features
226.    par(mfrow = c(1, 3))#allows multiple graphs to be displayed at
227.                    #once to they can be seen and compared easier
228.    pie( c(NumberOf(df$capShape, "bell"),
229.        NumberOf(df$capShape, "conical"),
230.        NumberOf(df$capShape, "convex"),
231.        NumberOf(df$capShape, "flat"),
232.        NumberOf(df$capShape, "knobbed"),
233.        NumberOf(df$capShape, "sunken") ) ,
234.     capShape_variables, main = "cap-shape",
235.     col = rainbow(length(capShape_variables)))
236.    pie( c(NumberOf(df$capSurface, "fibrous"),
237.        NumberOf(df$capSurface, "groves"),
238.        NumberOf(df$capSurface, "scaly"),
239.        NumberOf(df$capSurface, "smooth")) ,
240.     capSurface_variables, main = "cap-surface",
241.     col = rainbow(length(capSurface_variables)))
242.    pie( c(NumberOf(df$capColor, "brown"),
243.        NumberOf(df$capColor, "goldenrod4"),
244.        NumberOf(df$capColor, "sienna"),
245.        NumberOf(df$capColor, "gray"),
246.        NumberOf(df$capColor, "green"),
247.        NumberOf(df$capColor, "pink"),
248.        NumberOf(df$capColor, "purple"),
249.        NumberOf(df$capColor, "red"),
250.        NumberOf(df$capColor, "white"),
```

```r
251.        NumberOf(df$capColor, "yellow")) ,
252.     capColor_variables, main = "capColor",col = capColor_variables)
253.
254. # the following creates pie charts showing the proportions
255. #of the variables for bruises, odor and spore print color
256. par(mfrow = c(1, 3))#allows multiple graphs to be displayed at
257.                #once to they can be seen and compared easier
258. pie( c(NumberOf(df$bruises, TRUE),NumberOf(df$bruises, FALSE)) ,
259.     c("true", "false"), main = "bruises",col = rainbow(2))
260. pie( c(NumberOf(df$odor, "almond"),
261.     NumberOf(df$odor, "anise"),
262.     NumberOf(df$odor, "creosote"),
263.     NumberOf(df$odor, "fishy"),
264.     NumberOf(df$odor, "foul"),
265.     NumberOf(df$odor, "musty"),
266.     NumberOf(df$odor, "none"),
267.     NumberOf(df$odor, "pungent"),
268.     NumberOf(df$odor, "spicy")) ,
269.     odor_variables, main = "odor",col = rainbow(length(odor_variables)))
270. pie( c(NumberOf(df$sporePrintColor, "black"),
271.     NumberOf(df$sporePrintColor, "brown"),
272.     NumberOf(df$sporePrintColor, "goldenrod4"),
273.     NumberOf(df$sporePrintColor, "chocolate"),
274.     NumberOf(df$sporePrintColor, "green"),
275.     NumberOf(df$sporePrintColor, "orange"),
276.     NumberOf(df$sporePrintColor, "purple"),
277.     NumberOf(df$sporePrintColor, "white"),
278.     NumberOf(df$sporePrintColor, "yellow")) ,
279.     sporePrintColor_variables, main = "spore Print Color",
280.     col=sporePrintColor_variables)
281.
282. # the following creates pie charts showing the
283. #proportions of the variables for gill related features
284. par(mfrow = c(2, 2))#allows multiple graphs to be displayed
285.                #at once to they can be seen and compared easier
286. pie( c(NumberOf(df$gillAttachment, "attached"),
287.     NumberOf(df$gillAttachment, "descending"),
288.     NumberOf(df$gillAttachment, "free"),
289.     NumberOf(df$gillAttachment, "notched")) ,
290.     gillAttachment_variables, main = "gill-attachment",
291.     col = rainbow(length(gillAttachment_variables)))
292. pie( c(NumberOf(df$gillSpacing, "close"),
293.     NumberOf(df$gillSpacing, "crowded"),
294.     NumberOf(df$gillSpacing, "distant")),
295.     gillSpacing_variables, main = "gill-Spacing",
296.     col = rainbow(length(gillSpacing_variables)))
297. pie( c(NumberOf(df$gillSize, "broad"),
298.     NumberOf(df$gillSize, "narrow")),
299.     gillSize_variables, main = "gill-Size",
300.     col = rainbow(length(gillSize_variables)))
301. pie( c(NumberOf(df$gillColor, "black"),
302.     NumberOf(df$gillColor, "brown"),
303.     NumberOf(df$gillColor, "goldenrod4"),
304.     NumberOf(df$gillColor, "chocolate"),
305.     NumberOf(df$gillColor, "gray"),
306.     NumberOf(df$gillColor, "green"),
307.     NumberOf(df$gillColor, "orange"),
308.     NumberOf(df$gillColor, "pink"),
309.     NumberOf(df$gillColor, "purple"),
310.     NumberOf(df$gillColor, "red"),
311.     NumberOf(df$gillColor, "white"),
312.     NumberOf(df$gillColor, "yellow")) ,
313.     gillColor_variables, main = "gill-Color",col =gillColor_variables)
314.
315. # the following creates pie charts showing the
316. #proportions of the variables for the stalk
317. par(mfrow = c(2, 3))#allows multiple graphs to be displayed
318. #at once to they can be seen and compared easier
319. pie( c(NumberOf(df$stalkShape, "enlarging"),
320.     NumberOf(df$stalkShape, "tapering")),
321.     stalkShape_variables, main = "stalk Shape",
322.     col = rainbow(length(stalkShape_variables)))
323. pie( c(NumberOf(df$stalkSurfaceAboveRing, "fibrous"),
324.     NumberOf(df$stalkSurfaceAboveRing, "scaly"),
325.     NumberOf(df$stalkSurfaceAboveRing, "silky"),
326.     NumberOf(df$stalkSurfaceAboveRing, "smooth")),
327.     stalkSurfaceAboveRing_variables, main = "stalk-Surface-Above-Ring",
328.     col = rainbow(length(stalkSurfaceAboveRing_variables)))
329. pie( c(NumberOf(df$stalkColorAboveRing, "brown"),
330.     NumberOf(df$stalkColorAboveRing, "goldenrod4"),
331.     NumberOf(df$stalkColorAboveRing, "sienna"),
332.     NumberOf(df$stalkColorAboveRing, "gray"),
333.     NumberOf(df$stalkColorAboveRing, "orange"),
334.     NumberOf(df$stalkColorAboveRing, "pink"),
335.     NumberOf(df$stalkColorAboveRing, "red"),
336.     NumberOf(df$stalkColorAboveRing, "white"),
```

```
337.        NumberOf(df$stalkColorAboveRing, "yellow")) ,
338.     stalkColorAboveRing_variables, main = "stalk-Color-Above-Ring",
339.     col = rainbow(length(stalkColorAboveRing_variables)))
340. pie( c(NumberOf(df$stalkRoot, "bulbous"),
341.        NumberOf(df$stalkRoot, "club"),
342.        NumberOf(df$stalkRoot, "cup"),
343.        NumberOf(df$stalkRoot, "equal"),
344.        NumberOf(df$stalkRoot, "rhizomorphs"),
345.        NumberOf(df$stalkRoot, "rooted")),
346.     stalkRoot_variables, main = "stalk-Root",
347.     col = rainbow(length(stalkRoot_variables)))
348. pie( c(NumberOf(df$stalkSurfaceBelowRing, "fibrous"),
349.        NumberOf(df$stalkSurfaceBelowRing, "scaly"),
350.        NumberOf(df$stalkSurfaceBelowRing, "silky"),
351.        NumberOf(df$stalkSurfaceBelowRing, "smooth")),
352.     stalkSurfaceBelowRing_variables, main = "stalk-Surface-Below-Ring",
353.     col = rainbow(length(stalkSurfaceBelowRing_variables)))
354. pie( c(NumberOf(df$stalkColorBelowRing, "brown"),
355.        NumberOf(df$stalkColorBelowRing, "goldenrod4"),
356.        NumberOf(df$stalkColorBelowRing, "sienna"),
357.        NumberOf(df$stalkColorBelowRing, "gray"),
358.        NumberOf(df$stalkColorBelowRing, "orange"),
359.        NumberOf(df$stalkColorBelowRing, "pink"),
360.        NumberOf(df$stalkColorBelowRing, "red"),
361.        NumberOf(df$stalkColorBelowRing, "white"),
362.        NumberOf(df$stalkColorBelowRing, "yellow")) ,
363.     stalkColorBelowRing_variables, main = "stalk-Color-Below-Ring",
364.     col = rainbow(length(stalkColorBelowRing_variables)))
365.
366. # the following creates pie charts showing the
367. #proportions of the variables for veil and ring related features
368. par(mfrow = c(2, 2))#allows multiple graphs to be displayed
369. #at once to they can be seen and compared easier
370. pie( c(NumberOf(df$veilType, "partial"),
371.        NumberOf(df$veilType, "universal")) ,
372.     veilType_variables, main = "veil-Type",
373.     col = rainbow(length(veilType_variables)))
374. pie( c(NumberOf(df$veilColor, "brown"),
375.        NumberOf(df$veilColor, "orange"),
376.        NumberOf(df$veilColor, "white"),
377.        NumberOf(df$veilColor, "yellow")),
378.     veilColor_variables, main = "veil-Color",
379.     col = rainbow(length(veilColor_variables)))
380. pie( c(NumberOf(df$ringNumber, 0),
381.        NumberOf(df$ringNumber, 1),
382.        NumberOf(df$ringNumber, 2)),
383.     ringNumber_variables, main = "ring-Number",
384.     col = rainbow(length(ringNumber_variables)))
385. pie( c(NumberOf(df$ringType, "cobwebby"),
386.        NumberOf(df$ringType, "evanescent"),
387.        NumberOf(df$ringType, "flaring"),
388.        NumberOf(df$ringType, "large"),
389.        NumberOf(df$ringType, "none"),
390.        NumberOf(df$ringType, "pendant"),
391.        NumberOf(df$ringType, "sheathing"),
392.        NumberOf(df$ringType, "zone")) ,
393.     ringType_variables, main = "ring-Type",
394.     col =rainbow(length(ringNumber_variables)))
395.
396. # the following creates pie charts showing the proportions
397. #of the variables for bruises, odor and spore print color
398. par(mfrow = c(1, 3))#allows multiple graphs to be displayed at
399. #once to they can be seen and compared easier
400. pie( c(NumberOf(df$edible, TRUE),NumberOf(df$edible, FALSE)) ,
401.     c("true", "false"), main = "edibile or not",col = rainbow(2))
402. pie( c(NumberOf(df$population, "abundant"),
403.        NumberOf(df$population, "clustered"),
404.        NumberOf(df$population, "numerous"),
405.        NumberOf(df$population, "scattered"),
406.        NumberOf(df$population, "several"),
407.        NumberOf(df$population, "solitary")) ,
408.     population_variables, main = "population",
409.     col = rainbow(length(population_variables)))
410. pie( c(NumberOf(df$habitat, "grasses"),
411.        NumberOf(df$habitat, "leaves"),
412.        NumberOf(df$habitat, "meadows"),
413.        NumberOf(df$habitat, "paths"),
414.        NumberOf(df$habitat, "urban"),
415.        NumberOf(df$habitat, "waste"),
416.        NumberOf(df$habitat, "woods")) ,
417.     habitat_variables, main = "habitat",
418.     col=rainbow(length(habitat_variables)))
419. par(mfrow = c(1, 1))
420.   #makes it so that only one graph is one each screen for later plots
421. #one thing to notice from these charts is that there are some values that are
422. #incredibly one sided and should be investigated more
```

```
423.  NumberOf(df$veilType, "universal") #checks how many "universal" types mushrooms
424.                          #there are
425.  #as it turns out there are only partial type so this whole column can
426.  #safely be removed with no consequence
427.  df <- subset(df, select = -veilType )
428.
429.
430.  ######probability tables##########
431.
432.  table(df$capShape, df$edible,dnn =
433.      c("capShape", "edible"))
434.  #creates table comparing the cap shape on whether it is edible
435.  prop.table(table(df$capShape,
436.              df$edible,
437.              dnn = c("capShape", "edible")),1)
438.  #probability table of above
439.  #most of the varaibles are not to far from 50/50 so there is unlikely to be
440.  #much we can get from just these specific values
441.
442.  table(df$capSurface, df$edible,dnn =
443.      c("cap-Surface", "edible"))
444.  #creates table comparing edible with cap surgfce
445.  #with wiether they are edible or not
446.  prop.table(table(df$capSurface,
447.              df$edible,
448.              dnn = c("cap-Surface", "edible")),1)#probability table of above
449.  #similar to the table above there isn't much to show  a direct corelation
450.  #for these characteristics  other than groves of which all are poisonous
451.
452.  table(df$capColor,
453.      df$edible,
454.      dnn = c("capColor", "edible"))
455.  #creates table comparing if the mushroom is edible and with its cap colour
456.  prop.table(table(df$capColor,
457.              df$edible,
458.              dnn = c("capColor", "edible")),1) #probability table of above
459.  #probability table of above
460.  #this shows some more correlation such as pink and goldenrod4 so these might be
461.  #more defining when doing later algorithms
462.
463.
464.
465.  #######Cross-Validation and the Bootstrap###########
466.
467.  set.seed(1)
468.  train <- sample(length(df$edible), length(df$edible)*0.5)
469.  lm.fit <- lm(df$edible ~ df$capColor , data = df , subset = train)
470.  mean((df$edible - predict(lm.fit, df))[-train]^2)
471.  #only 0.1825958 very low error
472.
473.  #the following checks with a diffrent random set to  make sure that
474.  #it is really that low
475.  set.seed(2)
476.  train <- sample(length(df$edible), length(df$edible)*0.5)
477.  lm.fit <- lm(df$edible ~ df$capColor , data = df , subset = train)
478.  mean((df$edible - predict(lm.fit, df))[-train]^2)
479.  #only 0.1819106 which is very similar and still very low
480.  #however it would be better to use logistic regression as the data is entirely
481.  #categorical
482.
483.
484.  ######glm functions########
485.
486.  ggplot(df, aes(y = edible, x = capColor)) +
487.    geom_jitter(width = 0.05, height = 0.05, alpha = 0.5) +
488.    theme_classic(base_size = 15)
489.  #creates a jittered plot of cap colour and if the mushroom is poisonous
490.  #the tables above indicated that there are more likely to be a correlation for
491.  #capColors
492.  capColor_edible.glm <- glm(edible ~ capColor,
493.                      data = df,
494.                      family = binomial)
495.  summary(capColor_edible.glm)
496.    #creates a binomial model of the same data as above
497.  plogis(coef(capColor_edible.glm)[1]) #looks at the intercept
498.  plogis(coef(capColor_edible.glm)[2]) #looks at the slope
499.  predicted <- predict(capColor_edible.glm, type = "response", se.fit = TRUE)
500.  df$fit <- predicted$fit
501.  df$se <- predicted$se.fit #creates predictions for what has been plotted above
502.  ggplot(df, aes( y = edible, x = capColor)) +
503.    geom_jitter(width = 0.05, height = 0.05, alpha = 0.5) +
504.    geom_smooth(method = "glm", method.args = list(family = "binomial"),
505.          se = TRUE, color = "blue") +
506.    theme_classic(base_size = 15)
507.    #puts the plot on the same graph as above to compare it visually
508.    #we can see that there is not a specific correlation between these two type
```

```
509.   #most likely because there are purely discrete data (although inteersting
510.   #to note there are no goldenrod4 edible mushrooms)
511.  df$se <- NULL
512.  df$fit <- NULL  #removes the extra columns that where created since they are
513.          #no longer needed
514.
515.  edible_model_capColor <- glm(edible ~ capColor, data = df)
516.   #fits a logistic regression model with edible as our response variable and
517.   #and capColor as the explanatory variable
518.   #I chose these values since they seemed the most likely from above to give
519.   #us useful values
520.  summary(edible_model_capColor) # creates a summary of above
521.   #the summary shows that all reduce the odds of edibility to varying degrees
522.   #except for red which slightly increases the odds of edibility
523.
524.  edible_model_all <- glm(edible ~ .,data = df,family = binomial)
525.
526.                      #fits a logistic regression model with edible as our
527.  #response variable and all other variables as covariates
528.  #we are unfortunately told that the algorithm didn't converge
529.  summary(edible_model_all) #creates a summary of above
530.  #this shows us the chances of each mushroom being edible for each feature
531.  #compared to the first variable of each feature
532.
533.  predict(edible_model_all,
534.      newdata = data.frame(capShape = "bell",
535.                  capSurface = "fibrous",
536.                  capColor = "brown",
537.                  bruises = TRUE,
538.                  odor = "almond",
539.                  gillAttachment = "attached",
540.                  gillSpacing = "close",
541.                  gillSize = "broad",
542.                  gillColor = "black",
543.                  stalkShape = "enlarging",
544.                  stalkRoot = "club",
545.                  stalkSurfaceAboveRing = "smooth",
546.                  stalkSurfaceBelowRing = "smooth",
547.                  stalkColorAboveRing = "brown",
548.                  stalkColorBelowRing = "brown",
549.                  veilColor = "white",
550.                  ringNumber = 0,
551.                  ringType = "evanescent",
552.                  sporePrintColor = "white",
553.                  population = "clustered",
554.                  habitat = "woods"),
555.      type = "response")
556.  # Compare the estimated probabilities of a mushroom with the above qualities
557.  #it is near zero that this mushroom is edible according to this predictor
558.  #it does give a warning which is most likely because some of the features have
559.  #very little variability however this is just the nature of the data and
560.  #mushrooms themselves.
561.  #This is very useful for when you want to check for a specific mushroom since
562.  #you could enter observed characteristics of the mushroom in the above algorthm
563.  #and it would be able to tell you the probability of it being edible or not.
564.
565.
566.  model_stalkSurfaceBelowRing <- glm(edible ~ stalkSurfaceBelowRing,
567.                      data = df,
568.                      family = binomial)
569.  hnp(model_stalkSurfaceBelowRing)
570.  model_stalkColorBelowRing <- glm(edible ~ stalkColorBelowRing,
571.              data = df,
572.              family = binomial)
573.  hnp(model_stalkColorBelowRing)
574.
575.  model_both <- glm(edible ~ stalkSurfaceBelowRing + stalkColorBelowRing,
576.              data = df,
577.              family = binomial)
578.  hnp(model_both)
579.  #the vast majority of all three of these plots lie within the envelope
580.  #so they are likely a good fit
581.
582.
583.
584.
585.  #####Leave-One-Out Cross-Validation##########
586.
587.  glm.fit <- glm(df$edible ~ df$capShape , data = df)
588.  coef (glm.fit)
589.  #the reference here is bell shaped and almost all are less likely to be edible
590.  #than the reference
591.
592.  glm.fit <- glm(df$edible ~ df$capShape , data = df)
593.  cv.err <- cv.glm (df , glm.fit)
594.  cv.err$delta
```

```
595.
596.
597.
598.
599.
600.
601.
602.
603.
604.
605.
606.
607.
608.
609.
610.   ######fitting trees#########
611.
612.   #through all the previous analysis I have decided that since almost all the
613.   #data is categorical that decision trees will be a very good way of
614.   #classifying this data
615.
616.   #the classifications only work with factors so all mut be converted to them
617.   df_class <- df
618.   df_class$edible <- as.factor(df$edible)
619.   df_class$capShape <- as.factor(df$capShape)
620.   df_class$capSurface <- as.factor(df$capSurface)
621.   df_class$capColor <- as.factor(df$capColor)
622.   df_class$bruises <- as.factor(df$bruises)
623.   df_class$odor <- as.factor(df$odor)
624.   df_class$gillAttachment <- as.factor(df$gillAttachment)
625.   df_class$gillSpacing <- as.factor(df$gillSpacing)
626.   df_class$gillSize <- as.factor(df$gillSize)
627.   df_class$gillColor <- as.factor(df$gillColor)
628.   df_class$stalkShape <- as.factor(df$stalkShape)
629.   df_class$stalkRoot <- as.factor(df$stalkRoot)
630.   df_class$stalkSurfaceAboveRing <- as.factor(df$stalkSurfaceAboveRing)
631.   df_class$stalkSurfaceBelowRing <- as.factor(df$stalkSurfaceBelowRing)
632.   df_class$stalkColorAboveRing <- as.factor(df$stalkColorAboveRing)
633.   df_class$stalkColorBelowRing <- as.factor(df$stalkColorBelowRing)
634.   df_class$veilColor <- as.factor(df$veilColor)
635.   df_class$ringNumber <- as.factor(df$ringNumber)
636.   df_class$ringType <- as.factor(df$ringType)
637.   df_class$sporePrintColor <- as.factor(df$sporePrintColor)
638.   df_class$population <- as.factor(df$population)
639.   df_class$habitat <- as.factor(df$habitat)
640.
641.   set.seed(45)
642.   train <- sample(length(df$edible), length(df$edible)*0.5)
643.
644.   tree.df_class=tree(edible~. , data = df_class)
645.   summary(tree.df_class)
646.   plot(tree.df_class)
647.   text(tree.df_class, pretty = 0)
648.   tree.df_class
649.   #this has created a decision tree to help identify edible mushrooms which is
650.   #exactly what I wanted, however I should do more checking to make sure it is
651.   #accurate
652.
653.   df_class.test = df_class[-train,]
654.   edible.test=df_class$edible[-train]
655.   tree.df_class=tree(edible~. , data = df_class,subset=train)
656.   tree.pred=predict(tree.df_class,df_class.test,type = "class")
657.   table(tree.pred,edible.test)
658.   #this shows that the accuracy is 100%, there are no false positives or false
659.   #negatives, initial this looks to be a good thing however this could very
660.   #possibly be a result of over fitting
661.
662.   set.seed(3)
663.   cv.df_class=cv.tree(tree.df_class,FUN=prune.misclass)
664.   names(cv.df_class)
665.
666.   par(mfrow=c(1,2))
667.   plot(cv.df_class$size,cv.df_class$dev,type="b")
668.   plot(cv.df_class$k,cv.df_class$dev,type="b")
669.   #the above graphs show that it is unlikely that pruning the tree will help at
670.   #all as it is very accurate and compact as it is
671.   par(mfrow=c(1,1))
672.
673.   prune.df_class=prune.misclass(tree.df_class,best=5)
674.   plot(prune.df_class)
675.   text(prune.df_class,pretty=0)
676.   #the above shows that it the original number of branches is in fact correct
677.
678.   #########random forest##############
679.
680.   set.seed(1)
```

```
681.  bag.df_class=randomForest(df_class$edible~., data=df_class, subset=train,
682.                    mtry=21, importance=TRUE)
683.  bag.df_class
684.  # the above creates random forests which shows very similar results to the
685.  #trees above with no false posatives or false negatives
686.
687.  df_class.train = df_class[train,]
688.  for (col in names(df_class.train)) {
689.    if (is.factor(df_class.train[[col]])) {
690.      df_class.train[[col]] <- factor(df_class.train[[col]],
691.                          levels = levels(df_class[[col]]))
692.    }
693.  }
694.
695.  df_class.test = df_class[-train,] #creates the test set from the training set
696.  bag.df_class = randomForest(edible ~ ., data = df_class.train)
697.  yhat.bag = predict(bag.df_class, df_class.test)
698.
699.  varImpPlot(bag.df_class)
700.  #this shows how important the odor is specifically for checking the edibility of
701.  #the mushrooms in the data set with sporePrintColor and a few others coming
702.  #just after it
703.
704.  #########classification forest##########
705.
706.  #splits data into training and testing
707.  set.seed(150)
708.  split<-sample.split(df$edible, SplitRatio = 0.5)
709.  training_set<-subset(df_class,split==TRUE)
710.  test_set<-subset(df_class,split==FALSE)
711.  topredict_set <-  dplyr::select(test_set,-edible)
712.
713.  #fits a decision tree
714.  model_dtree<- rpart(edible ~ ., data=training_set)
715.  preds_dtree<-predict(model_dtree,newdata=topredict_set,type="class")
716.
717.  conf_matrix_dtree <- table(preds_dtree, test_set$edible)
718.  confusionMatrix(conf_matrix_dtree)
719.  #this tree misclassifies 11 mushrooms as poisonous when they where really
720.  #edible this is very accurate (and for someone who where to try to use this
721.  #algorithm its good that's its being safe)
722.  #the Kappa value and accuracy show how well this model does, and from a quick
723.  #observation doesn't look quite as susceptible to overfitting
724.
725.  #creates forest
726.  set.seed(1234)
727.  model_rf<-randomForest(edible~.,data=training_set,importance=TRUE, ntree=1000)
728.  preds_rf <- predict(model_rf, topredict_set)
729.  (conf_matrix_forest <- table(preds_rf, test_set$edible))
730.  confusionMatrix(conf_matrix_forest)
731.  #above creates a forest which has a 100% accuracy no there are no
732.  #misclassifcations unfortunately even though this looks like a good thing it
733.  #it could be a result of over fitting, it could be that the data is very
734.  #reliant on a few data points as seen in some of the graphs above
735.
736.  train_control <- trainControl(method = "cv", number = 10)
737.  cv_model <- train(edible ~ ., data = df_class,
738.                method = "rf",
739.                trControl = train_control)
740.  print(cv_model)
741.  #the above creates and compares different subsets to compare the model
742.  #it shows that that appose to what I'd initially thought the data is unlikely
743.  #to be over fitting instead there are just a few elements that can describe
744.  #edibility of the mushrooms in the data set very well
745.  #However there is a chance that this would be specific to the mushrooms within
746.  #this data set and could not be used on others.
```

# Bibliography

Mushroom [Dataset]. (1981). UCI Machine Learning Repository. https://doi.org/10.24432/C5959T.