

STAT 542 / CS 598: Project 2

Fall 2019, by Prathamesh(satpute3), Vivek(vivekg3) and Athul(as81)

Due: Monday, Dec 16 by 11:59 PM Pacific Time

[10 Points, half a page] Project description and summary.

Goal

Goal of this project is to identify malignant moles from benign moles. Provided image dataset has 150 benign and 150 malignant moles with varying resolutions. Some of the images have moles with hairs, markings on the microscope for measurement. This data has to be processed and classified to identify malignant moles from the benign ones.

Approach

Approach for Q1 - we resized the images to 100x100 resolution and used that all pixel information to perform classification.

For the first model we choose Lasso with best lambda which provided minimum misclassification error by doing 10 fold cross-validation and using that model for predicting the results.

For the second model we choose to do the PCA as it's an excellent choice when it comes to images. Using the selected principal components, we trained a SVM model. SVM Model is sensitive to parameters chosen. We used `tune.svm()` for trying different values of gamma and cost. We picked the gamma and cost which gives minimum classification error.

For the third model we use PCA features for training Random Forest model. We tried different values of no. of trees and no. of nodes to get the best parameters for the random forest model.

Approach for Q2 - as a starting point, we began from the directions provided by the American Academy of Dermatology association, which targeted the ABCD of melanoma detection. Darrell Rigel (n.d. [Link here](#)). This provided us the intuition needed to generate the features. ABCDE is an acronym for Asymmetry, Border, Color, Diameter and Evolving respectively. We extracted features which correlate to the attributes like Area, Perimeter, Circularity index, Color in the segmented region. We used Ridge Regression and Random Forest Model on this transformed dataset to see whether it improved the accuracy.

Results

Table 1: Q1 Results

Model	Accuracy
Lasso	64.00%
SVM	75.00%
Random Forest	93.00%

Table 2: Q2 Results

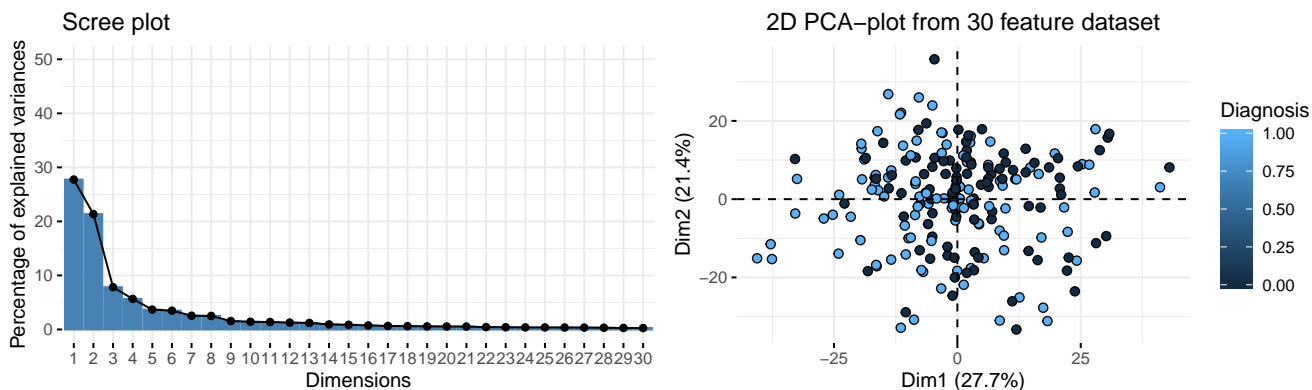
Model	Accuracy
Lasso	74.00%
Random Forest	91.00%

We can see an improvement in classification accuracy for Lasso Model after we performed the transformation mentioned in second approach. We chose Lasso Model for both approaches as the model interpretability is high.

[5 Points, half a page] Data processing for Question 1.

We have used the EBImage package for processing and getting image data. We have created the 2 matrix of 150x30000 each for benign and malignant files. We resize each image to 100x100 pixels and then used this resize image data as vector. As this data consists of 3 channels we got matrices with 100x100x3 cols and 150 rows. We used all the pixel information for the lasso model.

We then choose to do the PCA as PCA is an excellent choice when it comes to images because inherently due to its nature, there is spatial correlation among pixel. Instead of using all the pixels, we can significantly reduce the number of features which encompass most of the variation. After plotting some information as shown below we choose 27 components which accounts for over 90% of the variation. We split our data in 70/30 train and test data split and performed SVM and RF on it.

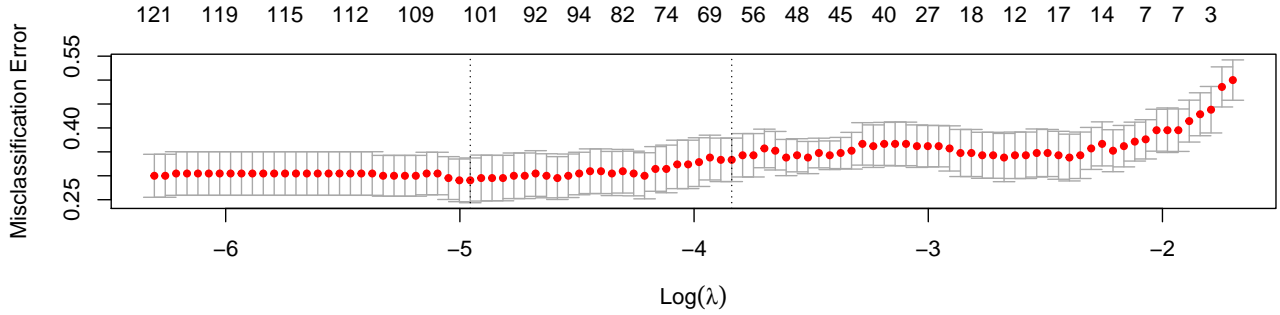


[30 Points, within 5 pages] Classification models based on pixels.

Classification Model 1. Lasso

For performing classification with the lasso we are selecting the family as binomial. We are using all the pixels which we got after resizing the image to 100x100 and we are using 3 channels. As we have split the data in 70/30 ratio, we got xtrain data as 210x30000 matrix and ytrain as 210 corresponding responses. xtest is 90x30000 matrix and ytest is corresponding 90 responses.

```
# use parallel for performance
registerDoMC(cores = 4)
set.seed(1)
cv_glmnet_model <- cv.glmnet(xtrain, ytrain, parallel = TRUE, alpha=1, nfolds = 10,
                             family="binomial", type.measure = "class")
# selecting the lambda which provided minimum misclassification error as best lambda
best_lambda = cv_glmnet_model$lambda.min
# training with best lambda selected with 10 fold cross validation
glmnet_model <- glmnet(xtrain, ytrain, lambda = best_lambda, alpha=1,
                       family="binomial", type.measure = "class")
# predicting on test data
pred <- predict(glmnet_model, s = best_lambda, newx = xtest, type = "class")
accuracy = mean(ytest == pred) * 100
lasso_results = data.frame("Best lambda" = best_lambda, "Accuracy" = as.numeric(accuracy,3))
lasso_confusionMatrix <- confusionMatrix(as.factor(pred), as.factor(ytest))
```



After performing the 10 fold validation and selecting the best lambda we retrained the model with the selected best lambda and family binomial and type.measure class and used this model for predicting class on the test data. We got around 68% accuracy for the test data.

Table 3: Lasso Results

Best.lambda	Accuracy
0.0070434	64.44444

Table 4: Lasso Confusion Matrix

	0	1
0	35	20
1	12	23

Classification Model 2. SVM with PCA

When we select only few features from the original dataset and build a linear SVM, we noticed an 80.47% training accuracy and a 65% validation set accuracy. After increasing the number of features to 50, 70, 90, and 100, the training accuracy increased to 98% and flattened. However the test set accuracy dropped from 65% to 50%. We performed PCA and based on the graphs and the scores, we retained the 27 features PCA as the final choice.

```
set.seed(1)
tuned_svm <- tune(
  svm,
  train.x = train.data,
  train.y = train.data[, 1],
  kernel = "linear",
  range = list(
    cost = 10 ^ (-2:2),
    gamma = c(0.1, 0.25, 0.5, 0.75, 1, 2)
  )
)
```

We noticed SVM is very sensitive to the choice of parameters. We trained a linear SVM model with the best parameters we got from tuning, $\gamma = 0.1$ and $\text{cost} = 1$. We are also performing 10-fold cross validation to assess the quality of the model.

```
best_svm <- svm(
  as.factor(y) ~ .,
  data = train.data,
  kernel = "linear",
  cost = tuned_svm$best.parameters$cost,
  gamma = tuned_svm$best.parameters$gamma,
  cross = 10
)
```

We tried different hyperplanes to separate the dataset. We noticed that a simple linear split is consistently performing well. 133 support vectors have been identified by the best svm model.

SVM for Training data : First two Principle Components

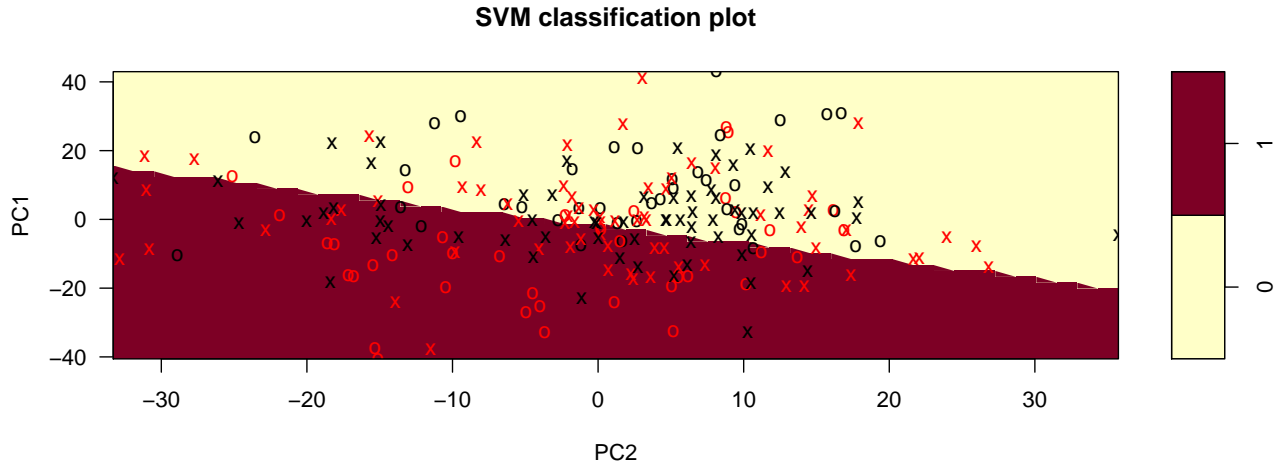


Table 5: SVM with PCA Results

Complete.Data.Accuracy
75

Table 6: SVM Confusion Matrix

	0	1
0	123	48
1	27	102

Linear SVM with 133 support vectors, trained over 27 Principal Components performs consistently with training accuracy of 76% and testing accuracy of 71%.

Classification Model 3. Random Forest with PCA

We used the dataset with 27 Principal Components to train a tree-based classification model. Random Forest is a popular decision tree model as it performs well for high-dimensional data and bins the outliers. Choice of no. of trees and no. of nodes is very important, so we did a grid search.

```
set.seed(1)
num_trees <- c(1,10,50,100,500,1000,1250,2000)
num_nodes <- c(1)
all_data_x <- rbind(test.data, train.data[, -1])
all_data_y <- c(pca.test[, 1], train.data[, 1])
best_i <- 1
best_j <- 1
best_all_error <- 100

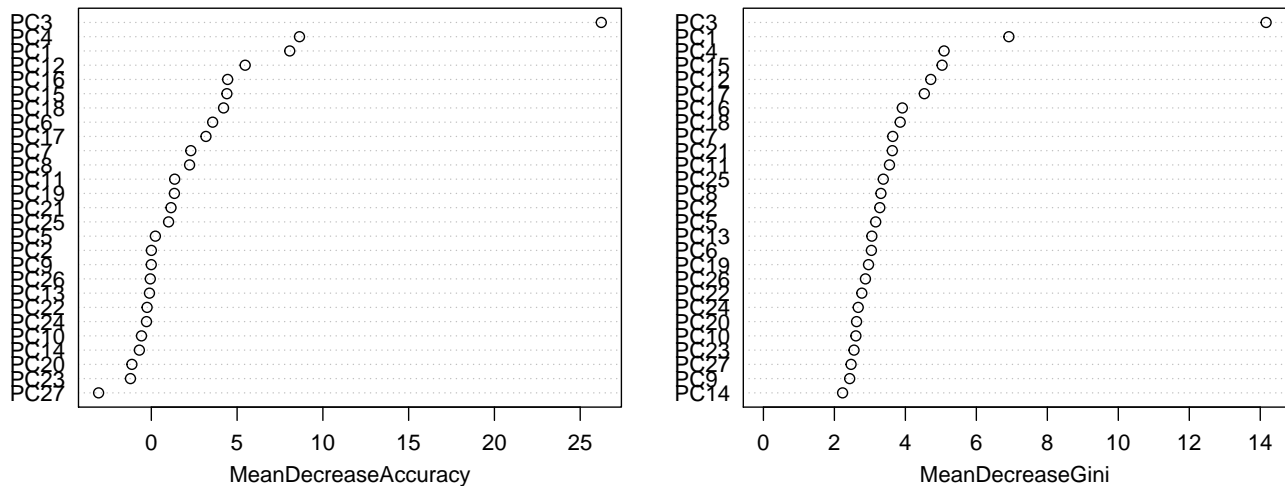
for(i in num_trees) {
  for (j in num_nodes) {
    rfModel = randomForest(
      formula = as.factor(y) ~ .,
      data = train.data,
      importance = T,
      ntree = i,
      nodesize = j
```

```

)
temp <- cbind(as.numeric(as.character(rfModel$predicted)), (train.data$y))
yhat.test = predict(rfModel, test.data)
training_error <- length(which(temp[, 1] != temp[, 2])) * 100 / nrow(train.data)
temp <- cbind(as.numeric(as.character(yhat.test)), (pca.test[, 1]))
test_error <- length(which(temp[, 1] != temp[, 2])) * 100 / nrow(test.data)
yhat.all <- predict(rfModel, all_data_x)
temp <- cbind(as.numeric(as.character(yhat.all)), as.numeric(all_data_y))
all_error <- length(which(temp[, 1] != temp[, 2])) * 100 / nrow(all_data_x)
if (all_error < best_all_error) {
  best_i <- i
  best_j <- j
  best_all_error <- all_error
  best_rfModel <- rfModel
}
}
}

```

Variable Importance Plot



This plot helps us in understanding the important variables chosen by the best random forest model. MeanDecreaseAccuracy helps us in understanding the importance of a variable's inclusion in reducing the classification error. MeanDecreaseGini helps us in understanding the relevance of variables in classification using node impurity (Gini).

Best Random Forest Model

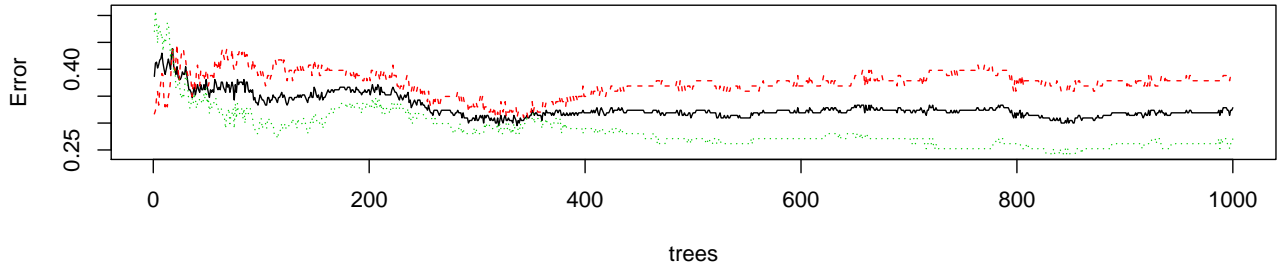


Table 7: RandomForest with PCA Results

Accuracy
93.33333

Table 8: RandomForest Confusion Matrix

	0	1
0	139	9
1	11	141

We run a random forest through a grid search through of number of trees (1 to 5000) and nodesize(1 to 50). We find that our best model is with 500 trees and a nodesize of 20. Our accuracy on the entire set is 93.33%

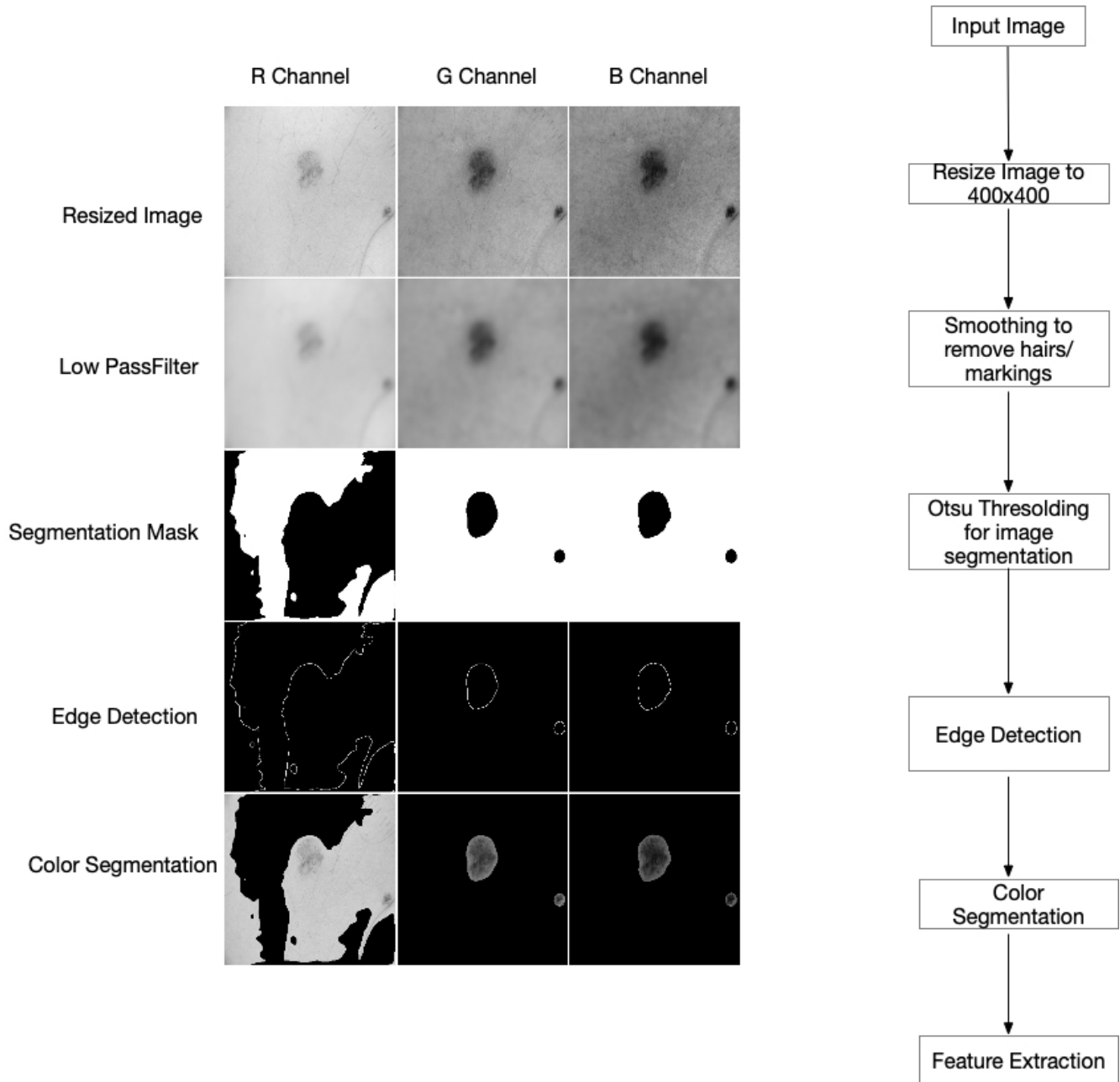
[10 Points, 1 page] Literature review. You should search and read existing literature and summarize clinically relevant characteristics that could be used for skin cancer image diagnosis. There is no limitation on what type of literature you could use. However, the goal should be motivating your feature engineering approaches from a clinical and analytic point of view. Please give appropriate citations to the literature you read.

The American Academy of Dermatology Association (Darrell Rigel (n.d.)) lists the ABCDE of detecting Melanoma. ABCDE is an acronym for Asymmetry, Border, Color, Diameter and Evolving respectively. Since in this project statement, are given a series of pictures over a lesion over a period of time, we will not be able to create a specific feature for Evolving. Hence, we concentrate on the remaining.

A great tools for us is EBImage which is package developed by Pau et al. (2010) which provides us very usable functions for image processing tasks. We use various features of this package for our feature engineering. Amelard, Wong, and Clausi (2012) provides high level features one could use to help engineer our features. Their main work talks about the different way to quantize irregularities. They use a mix of both coarse and fine grain methods to achieve this.

Jain, Jagtap, and Pise (2015) Shivangi et al, also talks about a good pipeline to this. Their work differs in that they introduce quantitative metrics to mark irregularities in shape. They also take into account the size of the lesion. However to this accurately they must have ensured that the each image is taken from the same distance and focus. We do not have that information regarding our dataset, so we must rely on ration. One of the interesting features proposed in using Circularity Index, it is the ration of $(4\pi \text{Area})/(\text{Perimeter} \times \text{Perimeter})$. This is a great metric because this is scale invariant. Of course is Jain, Jagtap, and Pise (2015) and (???) (???), they do concentrate a fair amount on preprocessing steps like illumination and segmentation.

[10 Points, 1 page] Feature engineering. Motivated by what you have read (or your understanding), process the data in a reasonable way such that the new variables are more intuitive to your collaborator/clinicians. You need to describe clearly what is your data processing criteria and how your variables are calculated.



Our pipeline takes an input image and prepares it for feature extraction. Once, we have our features, we then run it through various classifiers to observe the results. All steps in the preprocessing steps are done in 3 planes. We keep results from all these planes, as the classifier should be able to choose from them. Preprocessing steps :

1. Resize all images to 400 x400
2. Convert Color to Grayscale but also preserve the color image.
3. We run it through a low pass filter to remove hairs and scale markings.

Preprocessing steps : Do notice from the image above that the red channel does not yield very good preprocessing results while the green and blue channels are very successful in isolation the regions.

We generate 5 features for each of the channels:

1. Area: We count the number pixels in the mask of segmentation which are 0.
2. Perimeter: In the edge detected image, we count the number of pixels which are 1.
3. After applying the segmentation mask on the color image, we sum the total pixels that are in the segmented area.
4. Regularity Parameter: The ration of area / perimeter
5. Circularity Index: The ratio of $(4\pi \text{area})/(\text{perimeter}^2)$

[20 Points, 2 page] Classification models based on new features. Fit two different classification models to identify malignant moles. You can either use the ones from Question 1 or use some new models if you believe they may perform better on the new features. Same requirements of Question 1 apply to this part. Besides, you should focus more on variable selection and interpretation.

```
extract_features <- function(img_in) {
  resize_w <- 400
  resize_h <- 400
  img_resize <- resize(img_in, resize_w, resize_h)
  img_gray <- img_resize
  colorMode(img_gray) = Grayscale

  #Low Pass filter
  w = makeBrush(size = 31, shape = 'gaussian', sigma = 5)
  img_lp = filter2(img_gray, w)
  img_in_th <- img_lp
  threshold <- otsu(img_in_th)
  img_th = EBImage::combine(mapply(function(frame, th) frame > th,
                                    getFrames(img_in_th), threshold, SIMPLIFY = FALSE))

  img_th_val <- img_th
  img_th_val[which(img_th) == TRUE] <- 1
  img_th_val[which(img_th) == FALSE] <- 0
  fhi = matrix(1, nrow = 3, ncol = 3)
  fhi[2, 2] = -8
  img_fhi = filter2(img_th_val, fhi)
  img_fhi_col <- filter2(img_gray, fhi)

  #Features based on preprocessing
  area_f <- c(0, 0, 0)
  perimeter_f <- c(0, 0, 0)
  rgb_f <- c(0, 0, 0)
  img_col_thresh <- img_resize * img_th_val

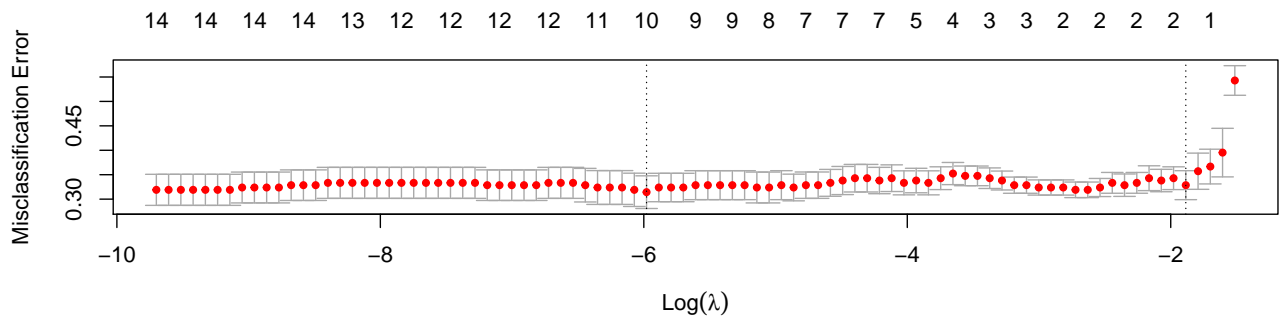
  for (i in 1:3) {
    area_f[i] <- length(which(img_th_val[, , i] == 0)) / (resize_h * resize_w)
    perimeter_f[i] <- length(which(img_fhi[, , i] == 1)) / (resize_h * resize_w)
    rgb_f[i] <- sum(img_col_thresh[which(img_th_val[, , i] == 0)]) / (resize_h * resize_w)
  }

  reg_f <- area_f / perimeter_f
  return(c(area_f, perimeter_f, rgb_f, (area_f/perimeter_f), (4*pi*area_f/perimeter_f^2)))
}
```

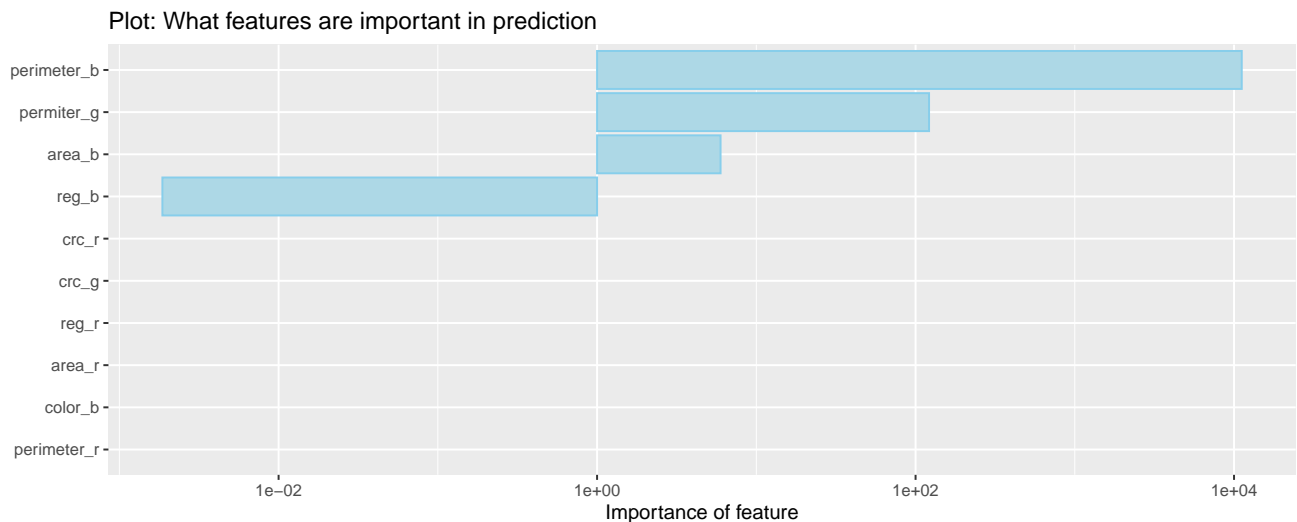

Classification Model 1. Lasso

For performing classification with the lasso we are selecting the family as binomial. We are using selected 15 features as xtrain and ytrain is corresponding response.

```
set.seed(1)
cv_glmnet_model <- cv.glmnet(xtrain, ytrain, parallel = TRUE, alpha=1, nfolds = 10,
                             family="binomial", type.measure = "class")
# selecting the lambda which provides the minimum clasification error as best lambda
best_lambda = cv_glmnet_model$lambda.min
# training with best lambda selected with 10 fold cross validation
glmnet_model <- glmnet(xtrain, ytrain, lambda = best_lambda, alpha=1,
                       family="binomial", type.measure = "class")
# predicting on test data
pred <- predict(glmnet_model, s = best_lambda, newx = xtest, type = "class")
accuracy = mean(ytest == pred) * 100
lasso_results = data.frame("Best lambda" = best_lambda, "Accuracy" = as.numeric(accuracy,3))
lasso_confusionMatrix <- confusionMatrix(as.factor(pred), as.factor(ytest))
```



After performing the 10 fold validation and selecting the best lambda we retrained the model with the selected best lambda and family binomial and type.measure class and used that model for predicting class on the test data. Accuracy we got is around 72%



By looking at this plot we can say that perimeter is the important fearure especially with color b and g after that we can see area also plays and important role.

Table 9: Lasso Results	
Best.lambda	Accuracy
0.0025312	74.44444

Table 10: Lasso Confusion Matrix		
	0	1
0	36	12
1	11	31

Classification Model 2. Random Forest

```

set.seed(1)
num_trees <- c(1,10,50,100,500,1000,1250,1500,2000)
num_nodes <- c(1)
all_data_x <- rbind(test.data[, -1], train.data[, -1])
all_data_y <- c(test.data[, 1], train.data[, 1])
best_i <- 1;
best_j <- 1;
best_all_error <- 100
for(i in num_trees){
  for(j in num_nodes){
    rfModel = randomForest(formula = as.factor(y) ~ ., data = train.data, importance = T, ntree=i, nodesize=j)
    temp <- cbind(as.numeric(as.character(rfModel$predicted)), (train.data[, 1]))
    yhat.test = predict(rfModel, test.data)
    training_error <- length(which(temp[, 1] != temp[, 2]))*100/nrow(train.data)
    temp <- cbind(as.numeric(as.character(yhat.test)), (test.data[, 1]))
    test_error <- length(which(temp[, 1] != temp[, 2]))*100/nrow(test.data)
    yhat.all <- predict(rfModel, all_data_x)
    temp <- cbind(as.numeric(as.character(yhat.all)), as.numeric(all_data_y))
    all_error <- length(which(temp[, 1] != temp[, 2]))*100/nrow(all_data_x)
    #cat("Num of Trees:", i, "Num of nodes:", j, " Training Error:", training_error, " Test Error:", test_error)
    if(all_error < best_all_error){
      best_i <- i
      best_j <- j
      best_all_error <- all_error
      best_rfModel <- rfModel
    }
  }
}

```

As with question1, we run through a random forest model with a grid search on number of trees and nodesize. The best model was with 100 trees and 5 nodesize with an accuracy of 92%.

best_rfModel

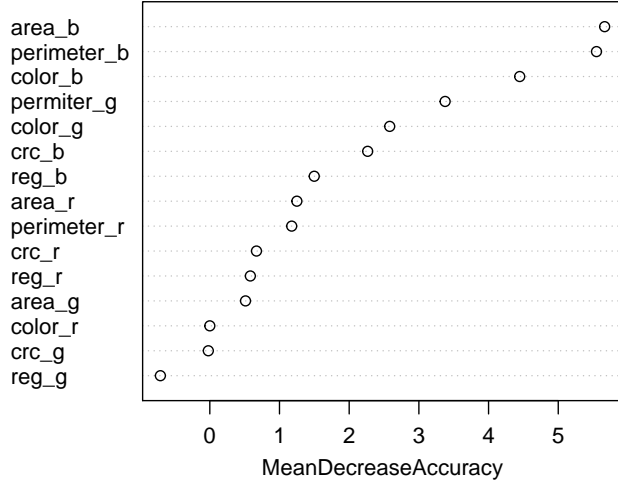


Table 11: RandomForest Results

Accuracy
91.66667

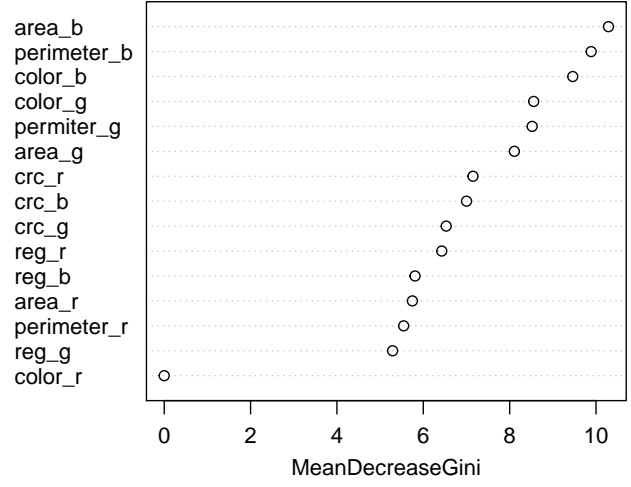


Table 12: RandomForest Confusion Matrix

	0	1
0	139	14
1	11	136

By Making understandable features, we see from above that: - Area, color and perimeter are key features. If the lesion contains multiple colors, then as Darrell Rigel (n.d.) explained, it is a high chance of cancer. - The next set of feature are the regularity and circularity index respectively. - Red is the least important channel in the model. We saw during preprocessing as well, the Red plane didn not yield good results as well.

References

- Amelard, R, A Wong, and D A. Clausi. 2012. “Extracting Morphological High-Level Intuitive Features (Hlif) for Enhancing Skin Lesion Classification.” In *34th Annual International Conference of the Ieee Engineering in Medicine and Biology Society*. San Diego.
- Darrell Rigel, American Academy of Dermatologist Association. n.d. “ABCDE of Melanoma.” <https://www.aad.org/diseases/skin-cancer/abcde-of-melanoma>.
- Jain, Shivangi, Vandana Jagtap, and Nitin Pise. 2015. “Computer Aided Melanoma Skin Cancer Detection Using Image Processing.” *Procedia Computer Science* 48 (December): 736–41. <https://doi.org/10.1016/j.procs.2015.04.209>.
- Pau, Grégoire, Florian Fuchs, Oleg Sklyar, Michael Boutros, and Wolfgang Huber. 2010. “EBImage - an R Package for Image Processing with Applications to Cellular Phenotypes.” *Bioinformatics* 26 (7): 979–81. <http://dblp.uni-trier.de/db/journals/bioinformatics/bioinformatics26.html#PauFSBH10>.