

Assignment

Hans

April 9, 2018

Load the data

Data is in quite a weird format. Let us apply some changes to the data such that observations correspond to rows and features correspond to columns (i.e. translate it to a tidy format)

```
train <- read_delim("Train_call.txt", delim="\t") %>%
  mutate(chrom_loc = as.factor(paste0("C", Chromosome, "_", Start, "_", End)))

## Parsed with column specification:
## cols(
##   .default = col_integer()
## )

## See spec(...) for full column specifications.
metadata <- train %>% subset(select=c(Chromosome, Start, End, Nclone, chrom_loc))

train <- train %>% gather(Sample, "measurements", which(grepl("Array", names(train)))) %>%
  subset(select=c(Sample, measurements, chrom_loc)) %>% spread(chrom_loc, measurements)

train <- read_delim("Train_clinical.txt", delim="\t") %>%
  mutate(Subgroup = as.factor(Subgroup)) %>%
  merge(train) %>% as.data.frame()

## Parsed with column specification:
## cols(
##   Sample = col_character(),
##   Subgroup = col_character()
## )

set.seed(130910)
is.train <- sample(c(TRUE, FALSE), nrow(train), replace = TRUE, prob = c(0.8, 0.2))
test <- train[!is.train,] %>% as.data.frame()
train <- train[is.train,] %>% as.data.frame()

train$Sample <- NULL
test$Sample <- NULL
```

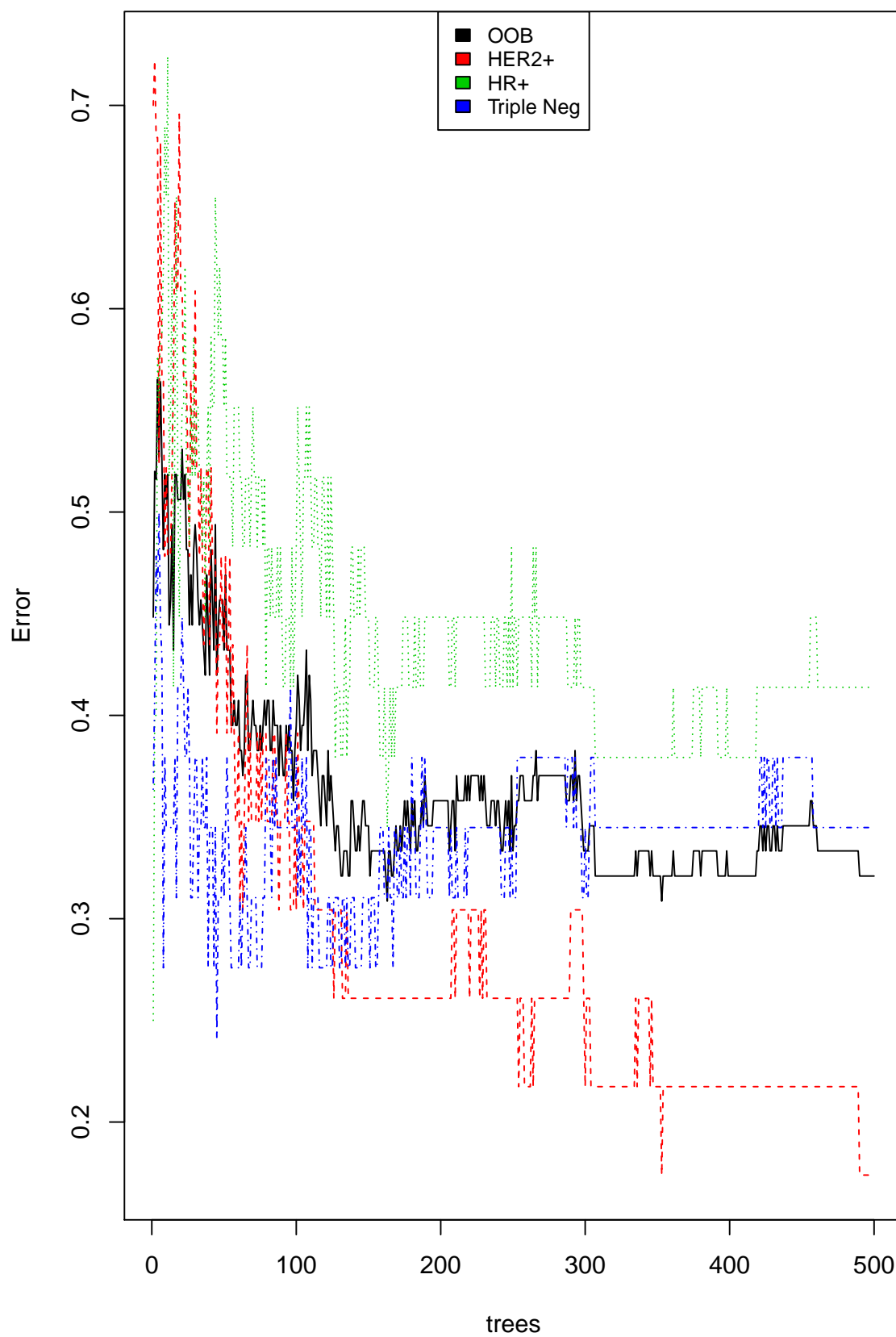
Fit a random forest, k-NN and support vector classifier.

Random forest may be interesting as it's not been used by the authors of paper Iwan found. k-NN is a very naive but sometimes powerful feature. It suffers from high variance, however. The support vector classifier isn't tried in that same paper with a non-linear kernel. May be a good idea to try this too.

Finally, we will try a boosted decision tree classifier from the gbm package. Such boosting can be very powerful in reducing the bias of the classifier (but not necessarily the variance). It is prone to overfit but still it may be good to use it here as boosting allows for heterogeneous instances which is clearly something we are dealing with in the case of cancer cells; that was also the motivation for using the knowledge-based network as we saw in the lecture.

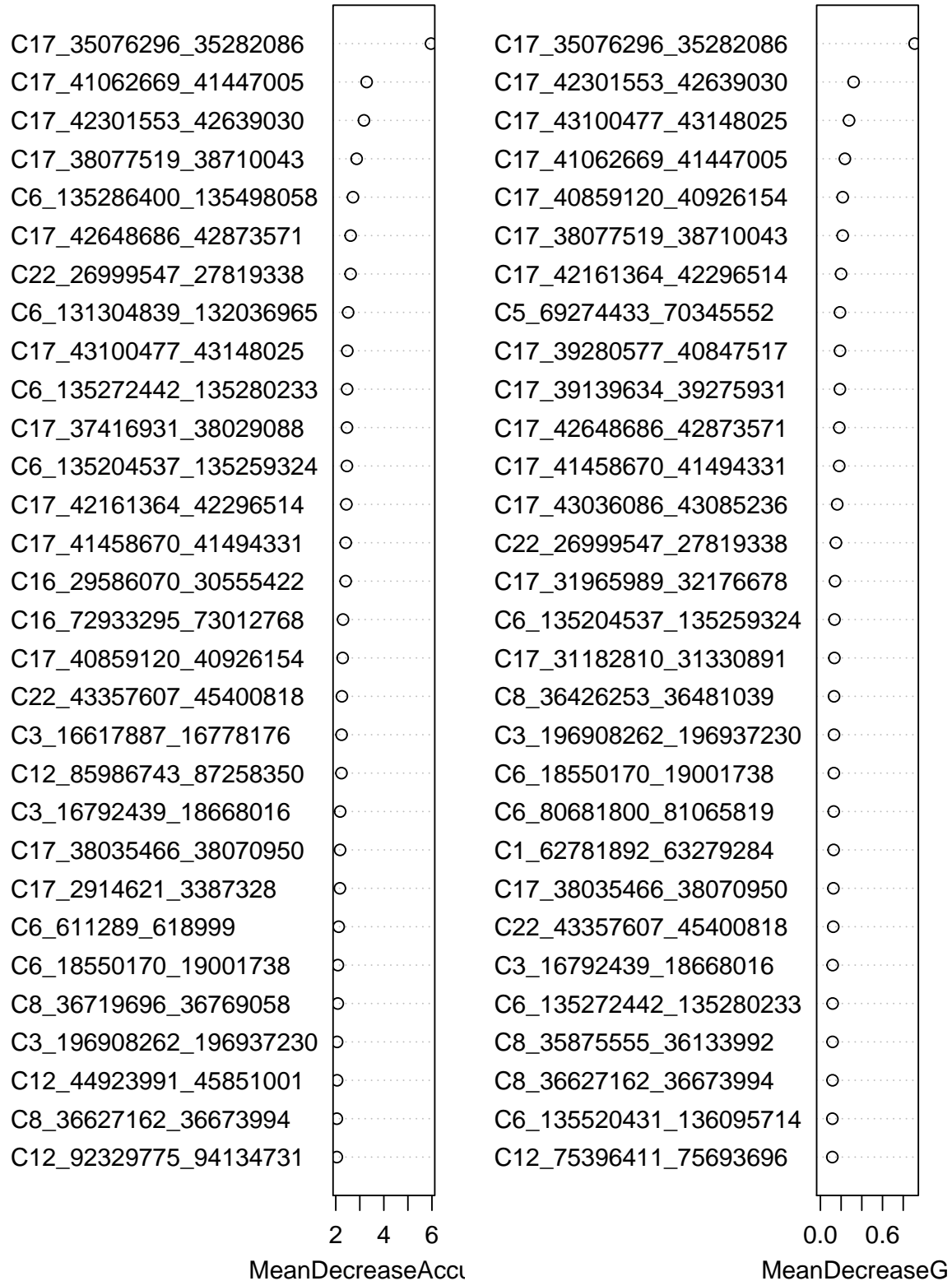
```
set.seed(12)
rf.model <- randomForest(Subgroup ~ ., data=train, ntrees = 1000, importance=TRUE)
plot(rf.model)
legend("top", colnames(rf.model$err.rate),col=1:4,cex=0.8,fill=1:4)
```

rf.model



```
varImpPlot(rf.model)
```

rf.model



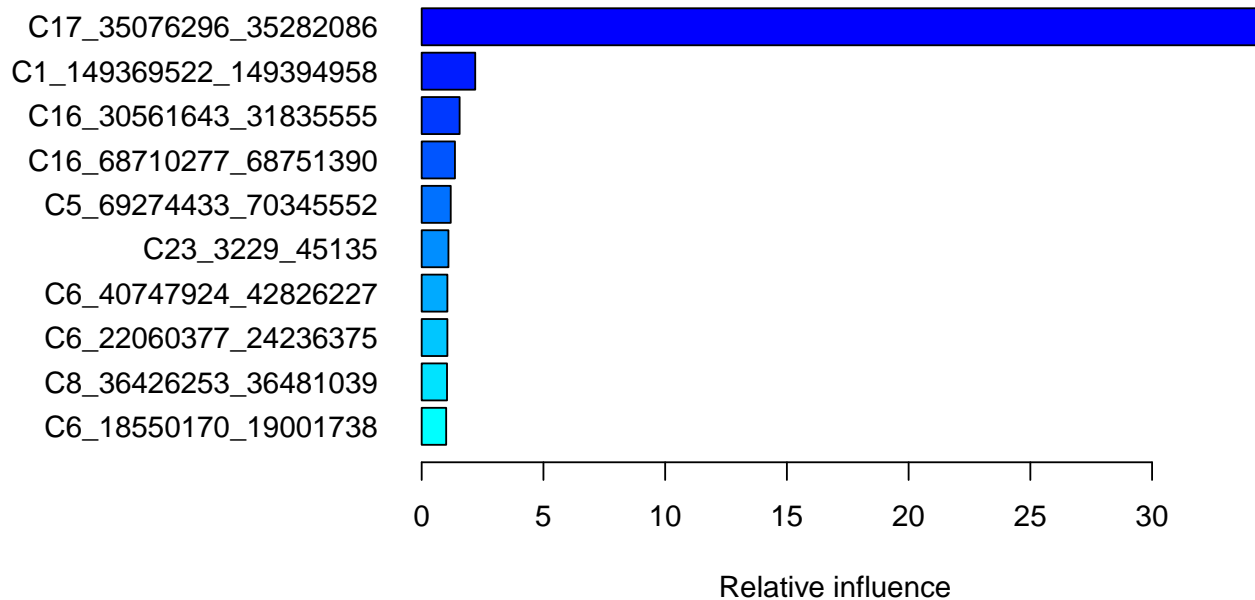
```
rf.predict <- predict(rf.model, test)
mean(test$Subgroup == rf.predict)
```

```
## [1] 0.6315789
```

Boosted decision trees. I took this implementation from some Rbloggers blog. Maybe we should look at the parameters too (but it works already quite well.) Summary shows one mutation that is really important for this cell (there is a one-to-one correspondence between this chromosome having a value 2 and it being a HER-2 cell if we inspect the validation and training data).

Train the boosted decision tree classifier (gbm).

```
rf.boost = gbm(Subgroup ~ . ,data = train, distribution = "multinomial", n.trees = 10000, shrinkage = 0.1)
par(oma=c(0,10,0,0))
output <- summary(rf.boost, cBars=10, order=TRUE, las=1)
```



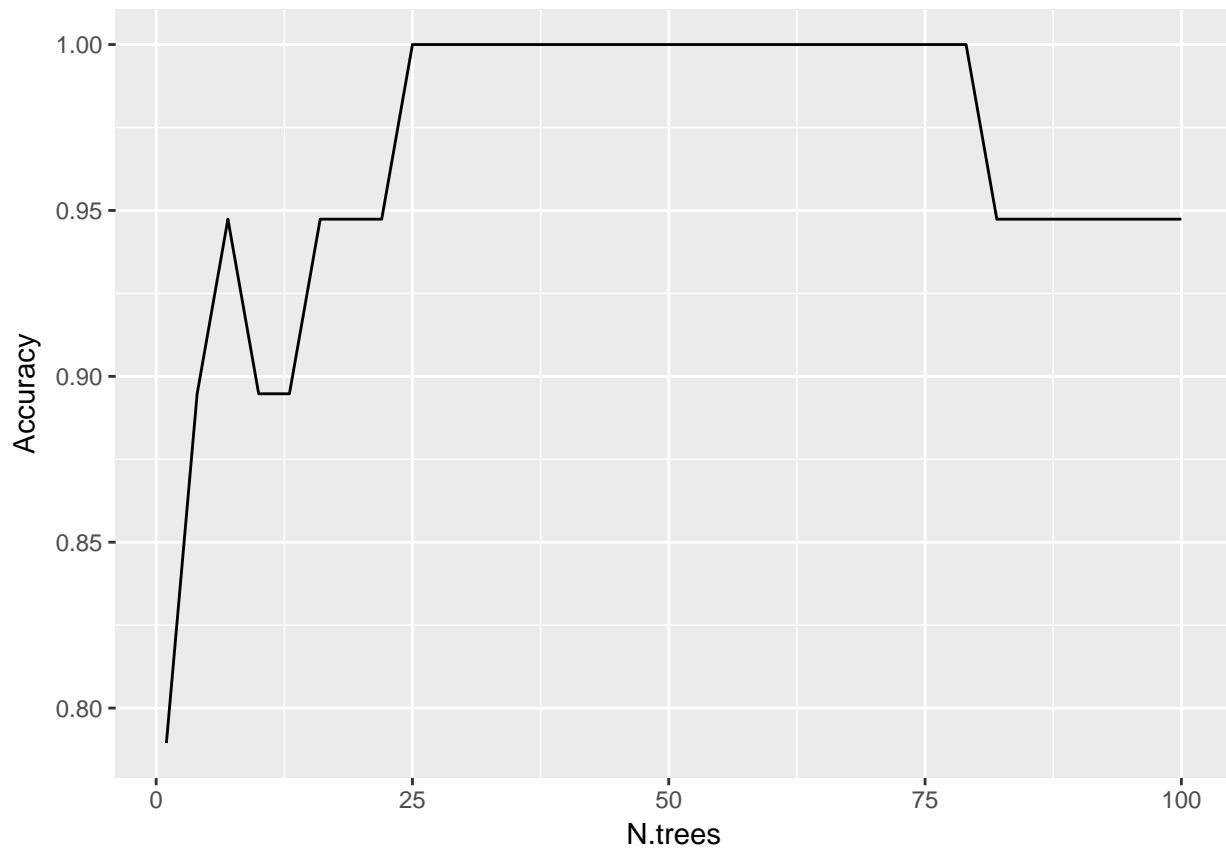
Evaluate accuracy of this classifier on the test set. We see that it has extremely high accuracy. May be reflecting overfitting.

```
n.trees = seq(from=1 ,to=100, by=3) #no of trees-a vector of 100 values

#Generating a dataframe of predictions for the number of trees.
prob_df <- predict(rf.boost,test,n.trees = n.trees, type="response") %>% melt()
test_labels <- test$Subgroup %>% as.data.frame() %>% mutate(Sample = row_number())
names(test_labels) <- c("True_type","Sample")
names(prob_df) <- c("Sample","Type", "N.trees", "p_hat")
predictions <- prob_df %>%
  group_by(Sample, N.trees) %>%
  filter(p_hat == max(p_hat)) %>%
  arrange(N.trees, Sample, Type) %>%
  merge(test_labels, by=c("Sample")) %>%
  group_by(N.trees) %>%
  summarize(Accuracy = sum(Type == True_type)/n())

#Calculating The Mean squared Test Error
predictions %>% ggplot(aes(x=N.trees, y=Accuracy)) +
```

```
geom_line()
```

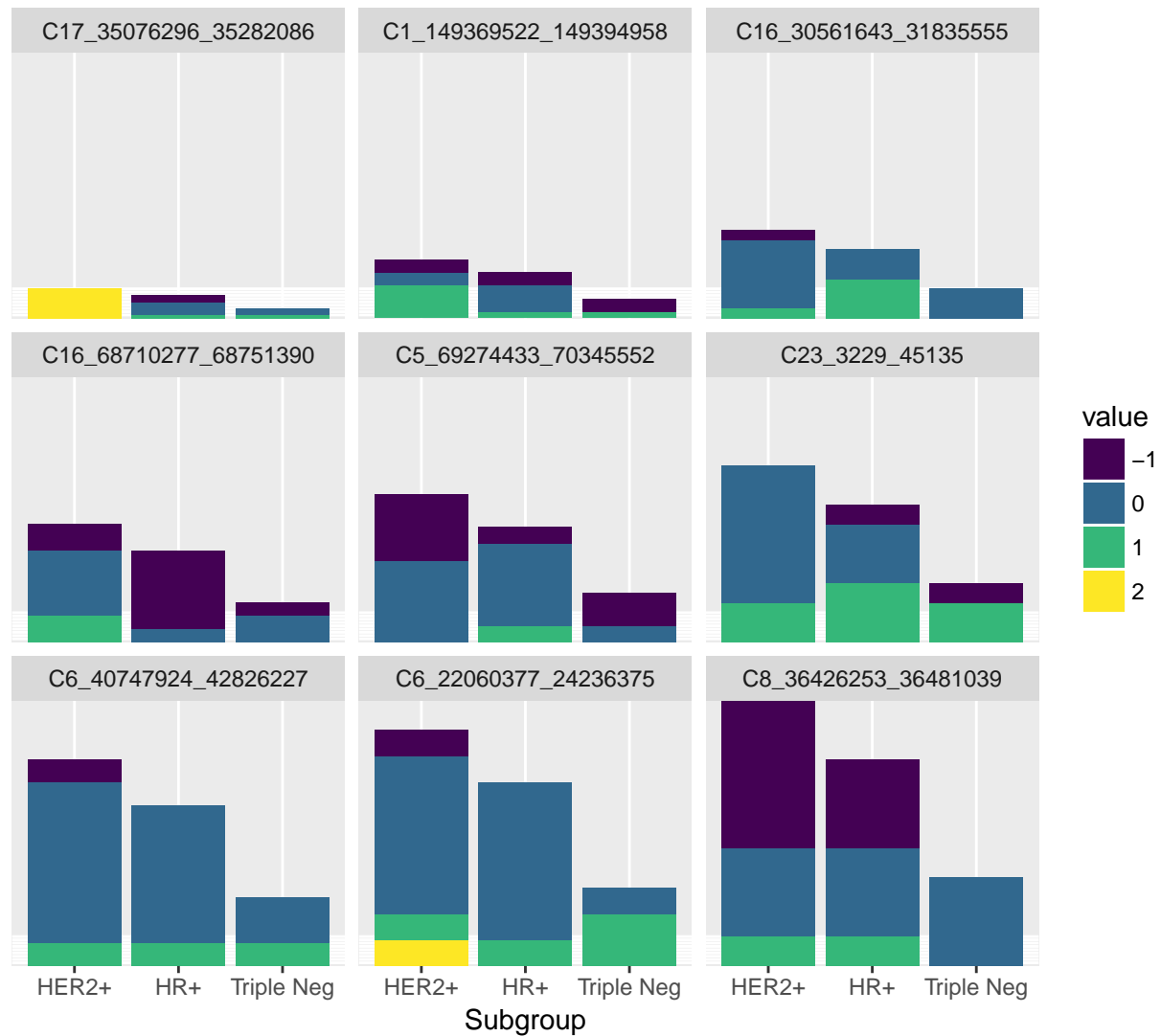


gbm() also reports on the importances of different features (i.e. chromosome regions). Let us plot the associations between the k most important features for test and training data. Note that there does not seem to be another magic marker that corresponds one-to-one to a cell-type.

```
k = 9
numRows = floor(k/3)
selected_cols <- output$var[1:k] %>% as.character()
test[c(selected_cols, "Subgroup")] %>% melt() %>%
  mutate(value = as.factor(value)) %>%
  ggplot(aes(x=Subgroup, y=variable, fill = value)) +
  geom_col() + facet_wrap( ~ variable, nrow = numRows) +
  theme(axis.ticks.y = element_blank(), axis.text.y=element_blank(), axis.title.y = element_blank()) +
  ggtitle("Distributions of gains/losses over most important chromosomes as identified by the GBM class")
  scale_fill_viridis(discrete=T)
```

```
## Using Subgroup as id variables
```

Distributions of gains/losses over most important chromosomes as identified by 1



```
train[c(selected_cols, "Subgroup")] %>% melt() %>%
  mutate(value = as.factor(value)) %>%
  ggplot(aes(x=Subgroup, y=variable, fill = value)) +
  geom_col() + facet_wrap( ~ variable, nrow = numRows) +
  theme(axis.ticks.y = element_blank(), axis.text.y=element_blank(), axis.title.y = element_blank()) +
  ggtitle("Distributions of gains/losses over most important chromosomes as identified by the GBM class")
  scale_fill_viridis(discrete=T)
```

```
## Using Subgroup as id variables
```


Distributions of gains/losses over most important chromosomes as identified by 1



The other classifiers

Very simple k-nearest neighbor classifier. Doesn't seem to work very well. May suffer from noise. I hypothesized a correlation-based measure would maybe be more appropriate. However, problem is that the matrix of observations isn't invertible so it does not work... We could try some dimension reduction e.g. through PCA and then do it but not a fan.

```
nn3 <- knn(train[,-1], test[,-1], unlist(train$Subgroup), k=5, prob=TRUE)
mean(nn3 == test$Subgroup)
```

```
## [1] 0.4736842
```

Fit an SVM classifier We will try both the linear and RBF kernel.

```
tuned_svm<-tune(svm, train.x=train[,-1], train.y = unlist(train$Subgroup),kernel="linear", range=list(c
print(tuned_svm)
```

```
##
```

```
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##   cost
##   0.01
##
## - best performance: 0.3083333
svm_good_model<-svm(Subgroup~., data=train, kernel="linear",cost=tuned_svm$best.parameters$cost)
svm.predict <- predict(svm_good_model, test)
mean(svm.predict == test$Subgroup)

## [1] 0.7368421

Tuned SVM seems to be overfitting; performance on test set declines.

library(e1071)
tuned_svm<-tune(svm, train.x=train[, -1], train.y = unlist(train$Subgroup),kernel="radial", range=list(c
print(tuned_svm)

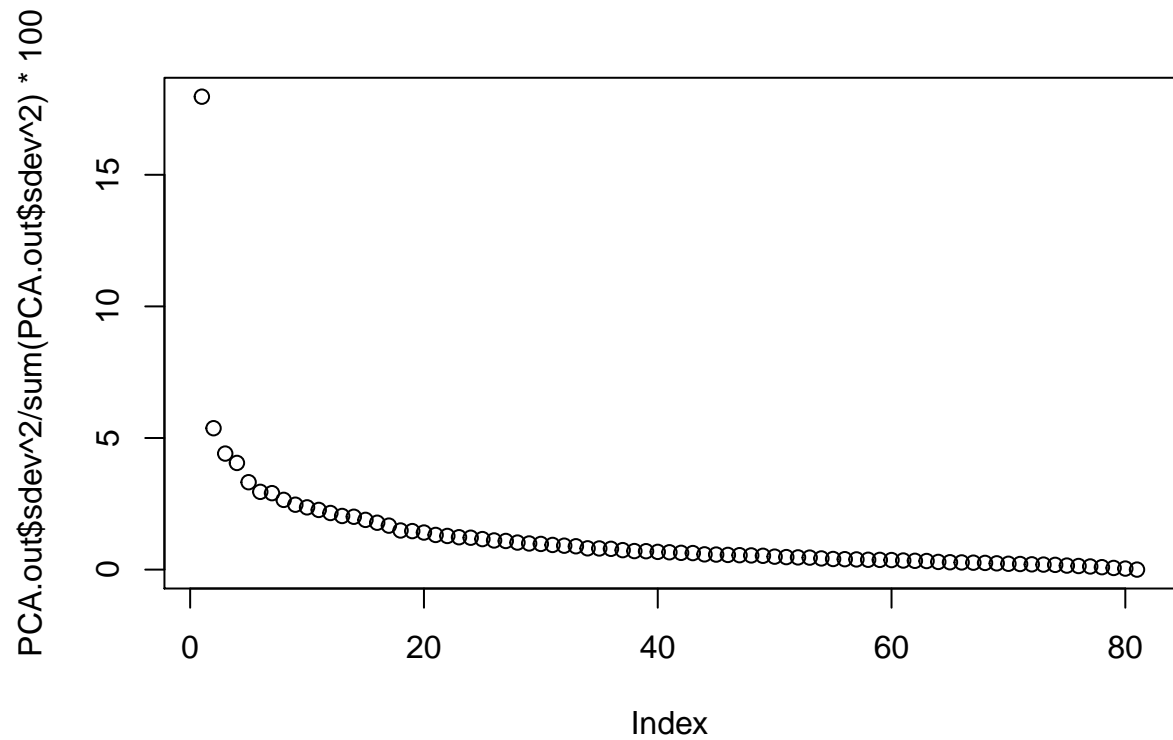
##
## Parameter tuning of 'svm':
##
## - sampling method: 10-fold cross validation
##
## - best parameters:
##   cost      gamma
##   8 3.051758e-05
##
## - best performance: 0.2819444
svm_good_model<-svm(Subgroup~., data=train, kernel="radial",cost=tuned_svm$best.parameters$cost)
svm.predict <- predict(svm_good_model, test)
mean(svm.predict == test$Subgroup)

## [1] 0.6842105
```

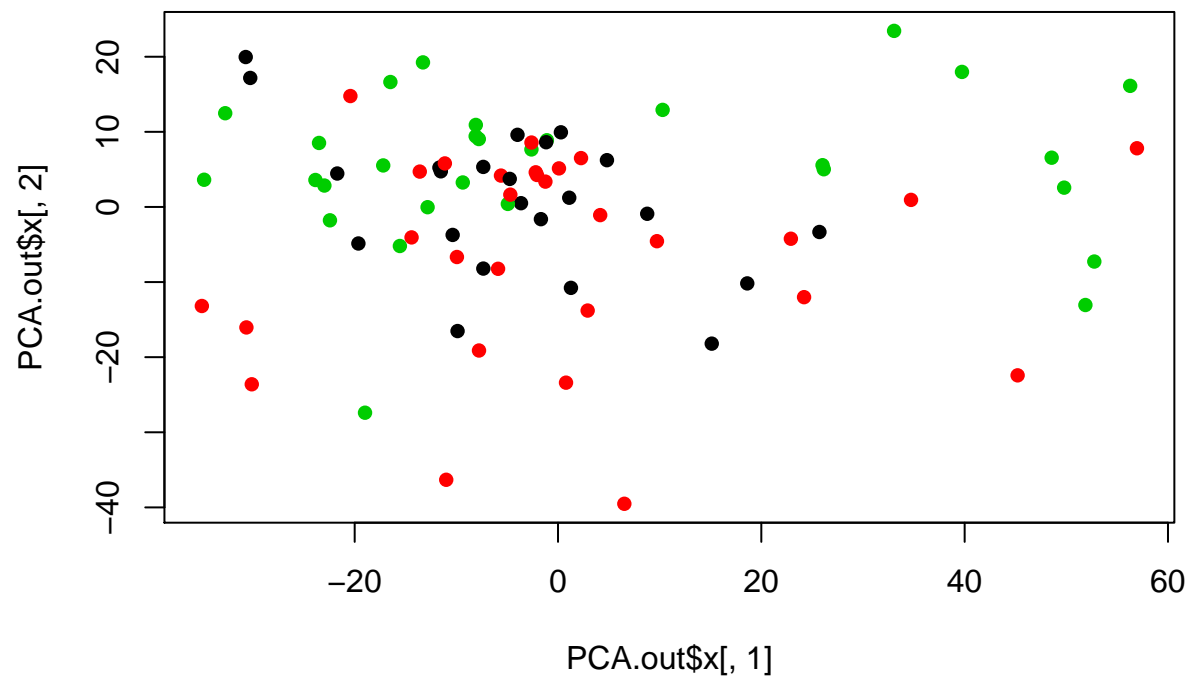
Scores on the first two principal components.

PCA doesn't yield clear patterns in the scores. Also, it's discussed in literature that it's only really appropriate when you think some underlying unobservable process generates the labels. Here, that seems a bit far-fetched.

```
PCA.out <- scale(train[, -1], center = TRUE, scale = TRUE) %>% prcomp()
plot(PCA.out$sdev^2/sum(PCA.out$sdev^2)*100)
```



```
plot(PCA.out$x[,1], PCA.out$x[,2], col = train$Subgroup, pch=16)
```



Other ideas:

- Investigate if perfect association stems from clinical test. If so, discuss that it may be cheating.
- Try binary classification between the triple negative and HR+ cells to select features that distinguish the two.

•