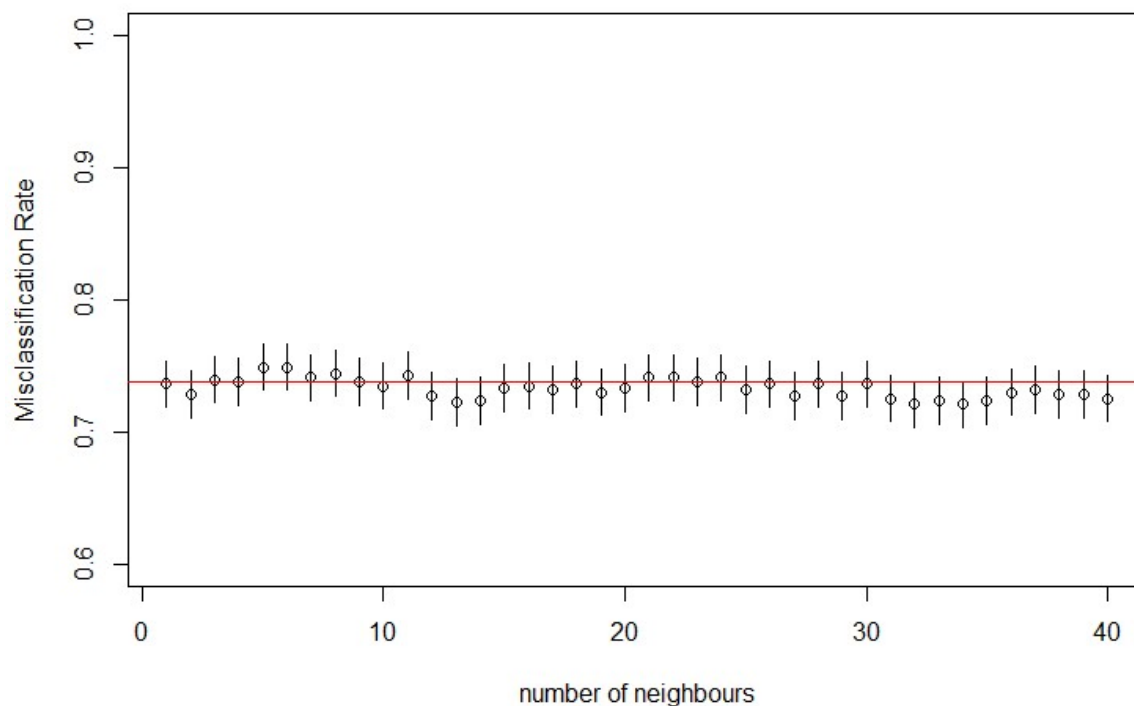4. Reset the seed to set.seed(9910314, kind="Mersenne-Twister"). Tune the KNN using CV error with knn.cv(). Use a grid from $m = 1, . . . , 40$. This may take a few seconds or minutes.

(a) Plot the validation error with standard errors against the number of neighbours. **Show the plot and comment: is there a clear best $m$ or is there a broad range of similar values, according to the SE?**

```
88   ##### 4. Reset the seed to set.seed(9910314, kind="Mersenne-Twister"). Tune the KNN
89   # using CV error with knn.cv(). Use a grid from m = 1, . . . , 40. This may take a few
90   # seconds or minutes.
91   set.seed(9910314, kind="Mersenne-Twister")
92
93
94   # (a) Plot the validation error with standard errors against the number of neighbours.
95   # Show the plot and comment: is there a clear best m or is there a broad
96   # range of similar values, according to the SE?
97
98
99   K.max = 40 # Maximum number of neighbours
00
01   ### Container to store CV misclassification rates
02   mis.CV = rep(0, times = K.max)
03
04 - for(i in 1:K.max){
05     ### Progress update
06     print(paste0(i, " of ", K.max))
07
08     ### Fit leave-one-out CV
09     this.knn = knn(X.train, X.valid, cl = set1[,19], k=i)
10
11     ### Get and store CV misclassification rate
12     this.mis.CV = mean(this.knn != Y.train)
13     mis.CV[i] = this.mis.CV
14 - }
15
16   ### Get SEs
17 - SE.mis.CV = sapply(mis.CV, function(r){
18     sqrt(r*(1-r)/nrow(X.train))
19 - })
20
21   plot(1:K.max, mis.CV, xlab = "number of neighbours", ylab = "Misclassification Rate",
22       ylim = c(0.6, 1))
124 - for(i in 1:K.max){
125     lower = mis.CV[i] - SE.mis.CV[i]
126     upper = mis.CV[i] + SE.mis.CV[i]
127
128     lines(x = c(i, i), y = c(lower, upper))
129 - }
130
131   ### Get CV min value for K
132   k.min = which.min(mis.CV)
133
134   thresh = mis.CV[k.min] + SE.mis.CV[k.min]
135   abline(h = thresh, col = "red")
```

➔ The error is the lowest at k=32 but there is no clear difference.

(b) Report the value of *m* with lowest error, as well as the one selected by the 1SE rule.

M=32

```
### Get CV 1SE value for K
k.1se = max(which(mis.CV <= thresh))

> k.1se
[1] 40
```

(c) Compute the test misclassification rate with both of these parameter values. **Report both error rates, using only one more digit than the first digit in their standard errors.** (For example, if the SE is 0.00375, the first digit in the SE 3 after the decimal, so report error to 4 after the decimal.)
Keep these error rates handy. You will use them for comparing all classification methods.
It would be better to use a more rigourous method like multiple reps of CV for computing error rates to make an "arena" for comparing methods. Classification is no different from regression in this way. However, this process is time consuming, and you have already practiced doing this on regression. I want to keep the classification assignments lighter, because they will accumulate rather quickly.

```
### Finally, let's see how our tuned KNN models do
knn.min = knn(X.train, X.valid, Y.train, k.min)
knn.1se = knn(X.train, X.valid, Y.train, k.1se)

(mis.min = mean(Y.valid != knn.min))
(mis.1se = mean(Y.valid != knn.1se))
I

> (mis.min = mean(Y.valid != knn.min))
[1] 0.3584906
> (mis.1se = mean(Y.valid != knn.1se))
[1] 0.3867925
```

➔ 2 after the decimal