

Read Data

11/23/2021

```
library(foreign)
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --
## v ggplot2 3.3.5      v purrr  0.3.4
## v tibble  3.1.4      v dplyr  1.0.7
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   2.0.1      v forcats 0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()

library(e1071)
library(tree)

## Registered S3 method overwritten by 'tree':
##   method      from
##   print.tree cli

library(randomForest)

## randomForest 4.6-14
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:dplyr':
##
##   combine
## The following object is masked from 'package:ggplot2':
##
##   margin

library(caret)

## Loading required package: lattice
##
## Attaching package: 'caret'
## The following object is masked from 'package:purrr':
##
##   lift

library(ggplot2)
library(dplyr)
library(tidyr)
```

```

library(tidyverse)
library(patchwork)
sesame <- read.dta("sesame.dta")
sesame <- sesame %>%
  mutate(site=factor(site)) %>%
  mutate(bodyDiff = postbody - prebody,
         letDiff = postlet - prelet,
         formDiff = postform - preform,
         numbDiff = postnumb - prenumb,
         relatDiff = postrelat - prerelat,
         clasfDiff = postclasf - preclasf)
sesame.sd <- sesame%>%
  mutate(sd_pBod = scale(prebody, center = TRUE, scale = TRUE),
         sd_plet = scale(prelet, center = TRUE, scale = TRUE),
         sd_pform = scale(preform, center = TRUE, scale = TRUE),
         sd_pnumb = scale(prenumb, center = TRUE, scale = TRUE),
         sd_prelat = scale(prerelat, center = TRUE, scale = TRUE),
         sd_pclasf = scale(preclasf, center = TRUE, scale = TRUE),
         sd_peabody = scale(peabody, center = TRUE, scale = TRUE),
         sd_age = scale(age, center = TRUE, scale = TRUE),
         male=if_else(sex==1, 1, 0),
         female=if_else(sex==2, 1, 0))

```

Exploratory Data Analysis

```
head(sesame)
```

```

##   rownames id site sex age viewcat setting viewenc prebody prelet preform
## 1      1  1  1  1  66      1      2      1     16     23     12
## 2      2  2  1  2  67      3      2      1     30     26      9
## 3      3  3  1  1  56      3      2      2     22     14      9
## 4      4  4  1  1  49      1      2      2     23     11     10
## 5      5  5  1  1  69      4      2      2     32     47     15
## 6      6  6  1  2  54      3      2      2     29     26     10
##   prenumb prerelat preclasf postbody postlet postform postnumb postrelat
## 1      40      14      20      18      30      14      44      14
## 2      39      16      22      30      37      17      39      14
## 3       9       9       8      21      46      15      40       9
## 4      14       9      13      21      14      13      19       8
## 5      51      17      22      32      63      18      54      14
## 6      33      14      14      27      36      14      39      16
##   postclasf peabody agecat encour _Isite_2 _Isite_3 _Isite_4 _Isite_5 regular
## 1      23      62      1      1      0      0      0      0      0
## 2      22       8      1      1      0      0      0      0      1
## 3      19      32      1      0      0      0      0      0      1
## 4      15      27      0      0      0      0      0      0      0
## 5      21      71      1      0      0      0      0      0      1
## 6      24      32      1      0      0      0      0      0      1
##   bodyDiff letDiff formDiff numbDiff relatDiff clasfDiff
## 1       2       7       2       4       0       3
## 2       0      11       8       0      -2       0
## 3      -1      32       6      31       0      11
## 4      -2       3       3       5      -1       2

```

## 5	0	16	3	3	-3	-1
## 6	-2	10	4	6	2	10

Variables:

The ID refers to a subject's identification number. The site refers to the age and background information of the child. A site value of 1 indicates a 3-5 year old disadvantaged child from the inner city. A site value of 2 represents a 4 year old advantaged child from the suburbs. A value of 3 represents an advantaged rural child. A site value of 4 indicates a disadvantaged rural child. Lastly, a value of 5 represents a disadvantaged Spanish speaking child. For the sex, a value of 1 indicates male, and a value of 2 indicates female. The age category is the child's age in months. The viewcat column is the frequency of viewing Sesame Street (1 = rarely, 2 = once/twice per week, 3 = 3-5 times a week, 4 = more than 5 times per week). The setting is where Sesame Street was viewed; a value of 1 indicates home and a value of 2 indicates school. The viewenc column refers to if the child was encouraged to watch or not (1 = child not encouraged, 2 = child encouraged). Encour is the same variable but with values 0 and 1, respectively. Regular is an indicator variable representing if a child is a regular viewer (0 = rarely watched, 1 = watched once per week or greater).

The prebody, prelet, preform, prenumb, prerelat, and preclasf columns all describe pretest scores on varying types of assessments (body parts, letters, forms, numbers, relational terms, and classification skills, respectively). The columns labelled postbody, postlet, postform, postnumb, postrelat, and postclasf are the children's respective posttest scores. Above, we created the following variables - bodyDiff, letDiff, formDiff, numbDiff, relatDiff, clasfDiff - to represent the difference in posttest scores and pretest scores for each child. Lastly, peabody represents a score of "mental age" for vocabulary maturity from the Peabody Picture Vocabulary Test.

Our main focus will be on the new variables we created (bodyDiff, letDiff, formDiff, numbDiff, relatDiff, clasfDiff) and variables related to how often the children watch Sesame Street (namely, viewcat and regular). Lastly, we will look into the backgrounds of the children, including site, sex, and age.

Distributions:

For the purposes of our analysis, we will first look at the distributions of bodyDiff, letDiff, formDiff, numbDiff, relatDiff, and clasfDiff.

#want to visualize distributions of bodyDiff, letDiff, formDiff, numbDiff, relatDiff, clasfDiff

```
bodyDiffplot <- ggplot(sesame, aes(x = bodyDiff)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of bodyDiff", x = "Post - Pre on Body Parts", y = "Count") +
  theme_minimal()

letDiffplot <- ggplot(sesame, aes(x = letDiff)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of letDiff", x = "Post - Pre on Letters", y = "Count") +
  theme_minimal()

formDiffplot <- ggplot(sesame, aes(x = formDiff)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of formDiff", x = "Post - Pre on Forms", y = "Count") +
  theme_minimal()

numbDiffplot <- ggplot(sesame, aes(x = numbDiff)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of numbDiff", x = "Post - Pre on Numbers", y = "Count") +
  theme_minimal()
```

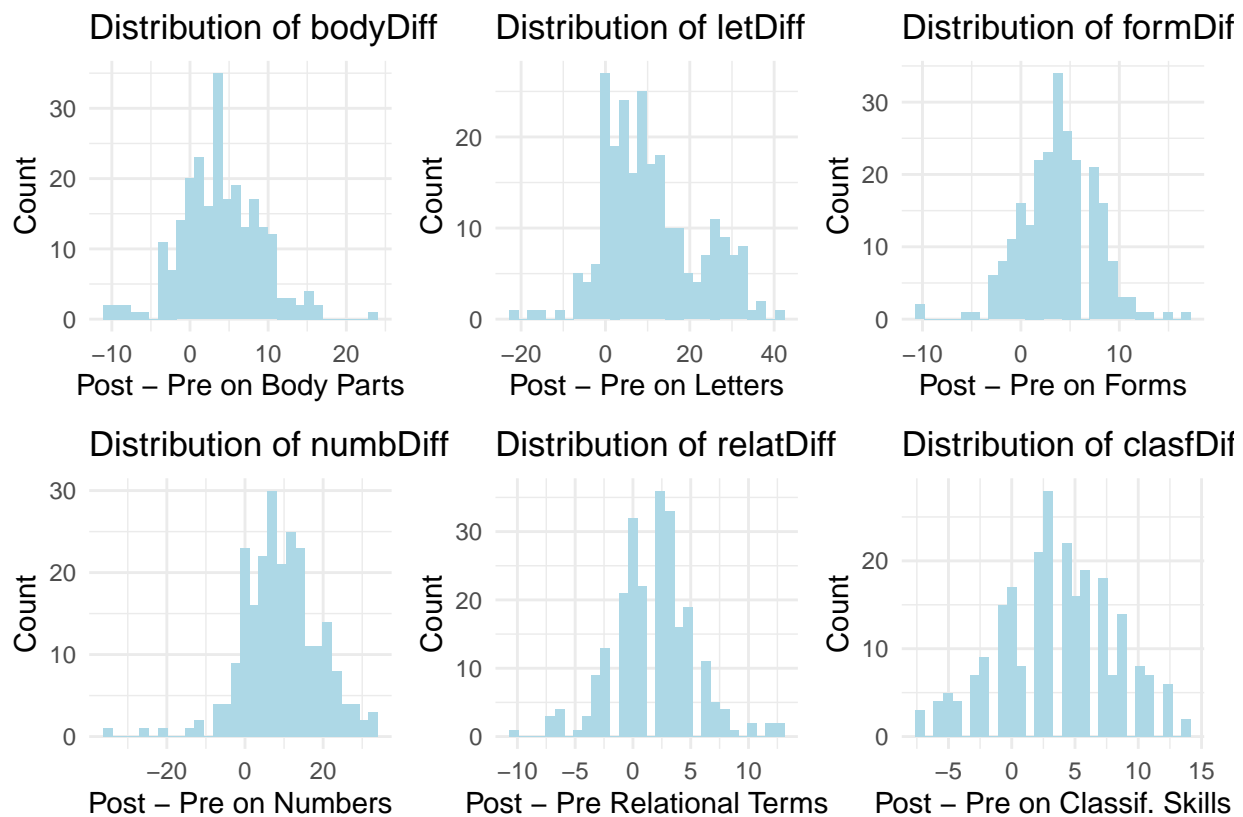
```

relatDiffplot <- ggplot(sesame, aes(x = relatDiff)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of relatDiff", x = "Post - Pre Relational Terms", y = "Count") +
  theme_minimal()

clasfDiffplot <- ggplot(sesame, aes(x = clasfDiff)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of clasfDiff", x = "Post - Pre on Classif. Skills", y = "Count") +
  theme_minimal()

bodyDiffplot + letDiffplot + formDiffplot + numbDiffplot + relatDiffplot + clasfDiffplot

```



The six variables above were calculated by subtracting pre-test scores from post-test scores, so they are all numerical. The distributions of these six variables (bodyDiff, letDiff, formDiff, numbDiff, relatDiff, and clasfDiff) all appear to be roughly normal and unimodal. BodyDiff, letDiff, formDiff, relatDiff, and classDiff do not appear to have any obvious extreme outliers. Numbdiff, however, seems to be slightly left-skewed with outliers to the left -20. All of the six variables appear to have centers between 2 and 4.

We will now examine the distributions of the variables related to how often children watch Sesame Street (namely, viewcat and regular).

```

# want to visualize distributions of viewcat and regular

viewcatplot <- ggplot(sesame, aes(x = factor(viewcat))) +
  geom_bar(fill = "lightblue") +
  labs(title = "Distribution of Viewcat", x = "Frequency of viewing Sesame Street", y = "Count") +
  scale_x_discrete("Frequency of Viewing Sesame Street", labels=c("rarely", "1-2 pw", "3-5 pw", ">5 pw")) +
  theme_minimal() +

```

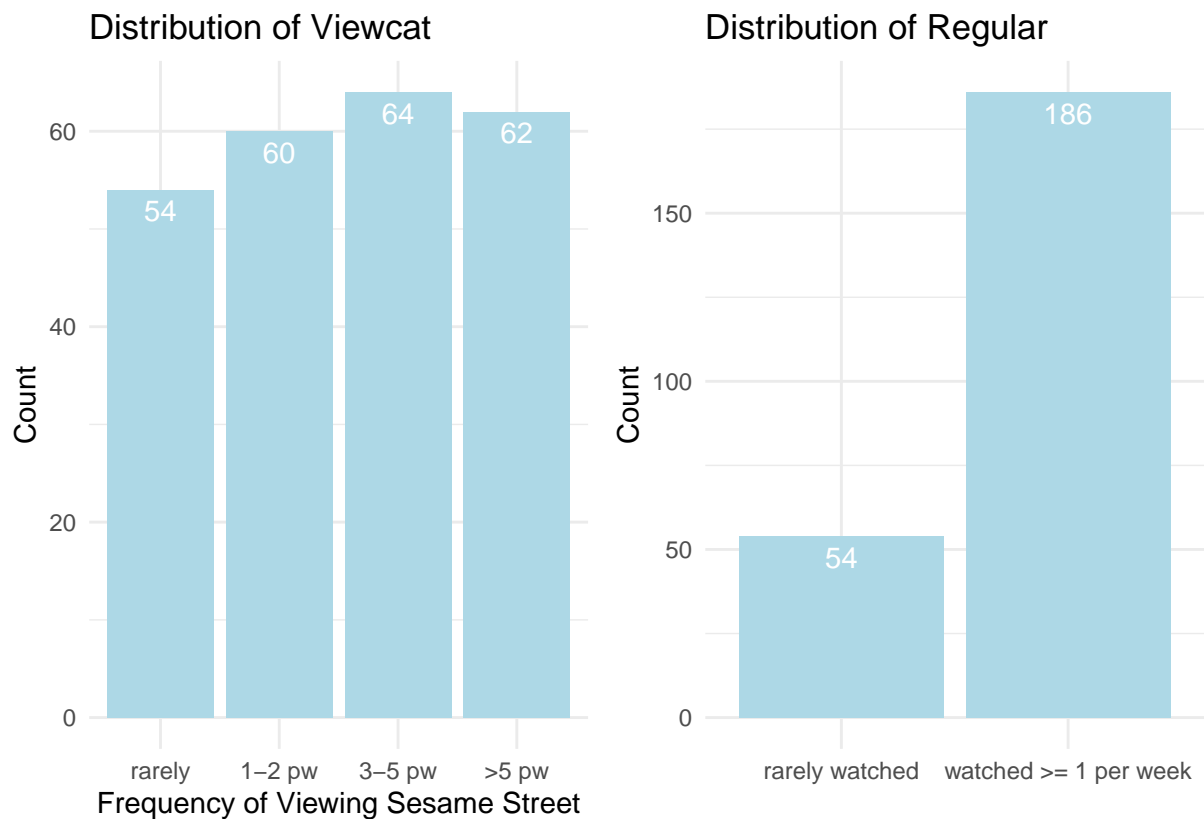
```

geom_text(aes(label = ..count..), stat = "count", vjust = 1.5, colour = "white")

regularplot <- ggplot(sesame, aes(x = factor(regular))) +
  geom_bar(fill = "lightblue") +
  labs(title = "Distribution of Regular", y = "Count") +
  scale_x_discrete(labels=c("rarely watched", "watched >= 1 per week")) +
  theme_minimal() +
  theme(axis.title.x = element_blank()) +
  geom_text(aes(label = ..count..), stat = "count", vjust = 1.5, colour = "white")

viewcatplot + regularplot

```



Both of these variables are categorical. On the left, viewcat appears to have a roughly uniform distribution, with “rarely” having the least amount of children and 3-5 times per week having the most (the range is only 10 children, so all of the bars are relatively close in height). For the variable regular, the category “watched once per week or greater” has far more observations than “rarely watched.” The former category has more than triple the amount of the latter. We will be aware of this disparity in our analysis and continue with caution towards potential bias.

Lastly, we want to examine the distributions of site, sex, and age, all variables that relate to a child’s background.

```

# want to visualize distributions of site, sex, and age

siteplot <- ggplot(sesame, aes(x = factor(site))) +
  geom_bar(fill = "lightblue") +
  labs(title = "Distribution of Site", y = "Count") +
  scale_x_discrete(labels=c("3-5, disadv., inner city", "4, adv., suburb", "adv., rural", "disadv., rural"))

```

```

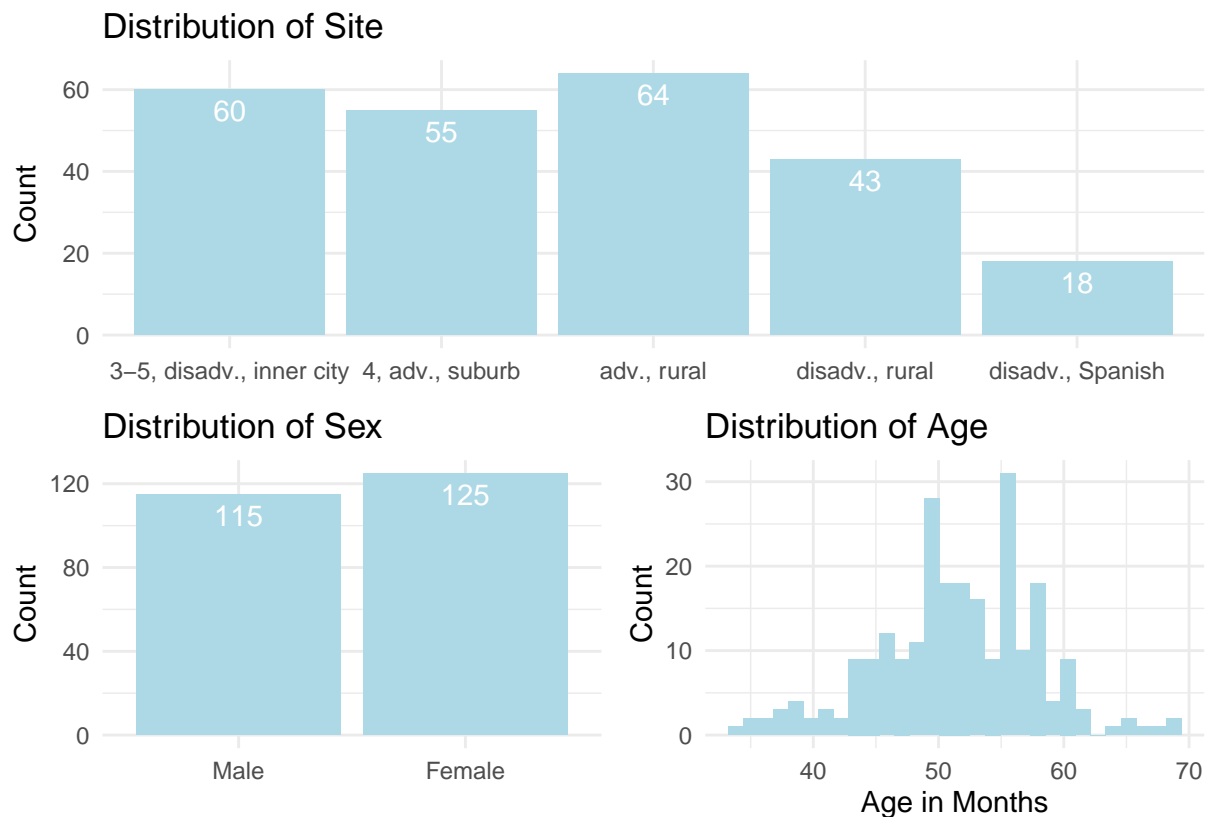
theme_minimal() +
  theme(axis.title.x = element_blank()) +
  geom_text(aes(label = ..count..), stat = "count", vjust = 1.5, colour = "white")

sexplot <- ggplot(sesame, aes(x = factor(sex))) +
  geom_bar(fill = "lightblue") +
  labs(title = "Distribution of Sex", y = "Count") +
  scale_x_discrete(labels=c("Male", "Female")) +
  theme_minimal() +
  theme(axis.title.x = element_blank()) +
  geom_text(aes(label = ..count..), stat = "count", vjust = 1.5, colour = "white")

ageplot <- ggplot(sesame, aes(x = age)) +
  geom_histogram(fill = "lightblue") +
  labs(title = "Distribution of Age", x = "Age in Months", y = "Count") +
  theme_minimal()

siteplot / (sexplot + ageplot)

```



Distribution of Site and Sex are both categorical variables. Distribution of Site has four categories with roughly the same amount of children (ranging from 43 to 64), but one category with far fewer observations (disadvantaged Spanish-speaking). This category has less than half of the observations as the next smallest category, which is a relatively large disparity. We will continue our analysis with caution towards this bias in the data. The distribution of sex is very even - the male category has 115 observations, while the female category has 125 observations. Age is a numerical variable that appears to be normal and bimodal, with two peaks around 50 and 56. There do not appear to be any extreme outliers in the distribution of age.

Q.2 Classification Question: Can we use the pre-test scores and other demographic variables to predict which region the children came from?

SVM

```
set.seed(3241)

n <- nrow(sesame)
train.index <- sample(1:n, size = floor(0.7*n), replace=FALSE)
train.data <- sesame.sd[train.index,]
test.data <- sesame.sd[-train.index,]

train.data %>%
  count(site)

##   site  n
## 1     1 40
## 2     2 42
## 3     3 48
## 4     4 25
## 5     5 13

# Response: site (categorical)
set.seed(315)
costs <- c(0.001, 0.01, 0.1, 1, 5, 10, 100)
gammas <- c(0.1, 0.5, 1, 2, 3, 4)

linear.tune <- tune(svm, site~female+ male + sd_age+sd_pBod+sd_plet+sd_pform + sd_pnumb+sd_prelat+sd_pclat,
  data=train.data, kernel="linear",
  ranges=list(cost=costs))

radial.tune <- tune(svm, site~female + male + sd_age+sd_pBod+sd_plet+sd_pform + sd_pnumb+sd_prelat+sd_pclat,
  data=train.data, kernel="radial",
  ranges=list(cost=costs,
              gamma=gammas))

#radial.tune <- tune(svm, site~sex+age+prebody+prelet+preform+prenumb+prerelat+preclasf,
#  data=train.data, kernel="radial",
#  ranges=list(cost=costs,
#              gamma=gammas))

linear.conMatrix <- table(true=test.data[, "site"],
  pred=predict(linear.tune$best.model, newdata=test.data))

radial.conMatrix <- table(true=test.data[, "site"],
  pred=predict(radial.tune$best.model, newdata=test.data))

confusionMatrix(linear.conMatrix)

## Confusion Matrix and Statistics
##
##      pred
## true  1  2  3  4  5
##    1  2  5 13  0  0
##    2  0  8  5  0  0
##    3  1  1 14  0  0
##    4  0  4 14  0  0
```

```
##      5  0  1  4  0  0
##
## Overall Statistics
##
##           Accuracy : 0.3333
##           95% CI : (0.2266, 0.4543)
##       No Information Rate : 0.6944
##       P-Value [Acc > NIR] : 1
##
##           Kappa : 0.1523
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3 Class: 4 Class: 5
## Sensitivity      0.66667  0.4211  0.2800      NA      NA
## Specificity      0.73913  0.9057  0.9091    0.75  0.93056
## Pos Pred Value   0.10000  0.6154  0.8750      NA      NA
## Neg Pred Value   0.98077  0.8136  0.3571      NA      NA
## Prevalence       0.04167  0.2639  0.6944    0.00  0.00000
## Detection Rate   0.02778  0.1111  0.1944    0.00  0.00000
## Detection Prevalence 0.27778  0.1806  0.2222    0.25  0.06944
## Balanced Accuracy 0.70290  0.6634  0.5945      NA      NA
```

```
confusionMatrix(radial.conMatrix)
```

```
## Confusion Matrix and Statistics
##
##      pred
## true  1  2  3  4  5
##      1  7  3 10  0  0
##      2  1  8  4  0  0
##      3  1  1 14  0  0
##      4  1  3 14  0  0
##      5  0  1  4  0  0
##
## Overall Statistics
##
##           Accuracy : 0.4028
##           95% CI : (0.2888, 0.525)
##       No Information Rate : 0.6389
##       P-Value [Acc > NIR] : 1
##
##           Kappa : 0.2337
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3 Class: 4 Class: 5
## Sensitivity      0.70000  0.5000  0.3043      NA      NA
## Specificity      0.79032  0.9107  0.9231    0.75  0.93056
## Pos Pred Value   0.35000  0.6154  0.8750      NA      NA
## Neg Pred Value   0.94231  0.8644  0.4286      NA      NA
```


## Prevalence	0.13889	0.2222	0.6389	0.00	0.00000
## Detection Rate	0.09722	0.1111	0.1944	0.00	0.00000
## Detection Prevalence	0.27778	0.1806	0.2222	0.25	0.06944
## Balanced Accuracy	0.74516	0.7054	0.6137	NA	NA

Radial kernel improves prediction on class 1.

RBF slightly improved after standardizing? (it seems slightly more likely to predict on class 1.) thought, simpler models still retain the same performance (arguably better) `sd_age+sd_pBod+sd_plet`. But we are still not getting any prediction on class 4 & 5.

Tree

```
set.seed(3215)
#tree.data <- sesame %>%
# select(site, sex, age, viewcat, setting, viewenc, prebody, prelet, preform,
#        prenumb, prerelat, preclasf)

n <- nrow(sesame)
train.index <- sample(1:n, size = floor(0.7*n), replace=FALSE)
#train.tree <- tree.data[train.index,]
#test.tree <- tree.data[-train.index,]

tree.features <- c("site", "age", "viewcat", "setting", "viewenc", "prebody", "prelet",
                  "preform", "prenumb", "prerelat", "preclasf", "postbody", "postlet",
                  "postform", "postnumb", "postrelat", "postclasf", "peabody")

tree.data <- sesame[, tree.features]
train.data <- tree.data[train.index,]
test.data <- tree.data[-train.index,]

rf.tree<- randomForest(site~., data=tree.data, subset=train.index,
                       mtry=4, importance=TRUE)

importance(rf.tree)
```

##		1	2	3	4	5
## age	10.6367151	3.5481566	1.6383512	4.24992100	-1.7747638	
## viewcat	1.4256856	3.8720209	2.5376319	7.54692254	1.2853072	
## setting	2.8373514	0.3753826	2.4853549	5.63595749	4.6148439	
## viewenc	2.3532392	0.3486145	3.8298659	1.90553085	-1.6264562	
## prebody	-2.8186130	10.8165553	4.2727110	2.49478287	-0.7357324	
## prelet	2.1177636	-0.6934365	0.6900640	-2.94473567	-1.6952551	
## preform	3.3471057	3.9742915	8.3220425	-2.72484260	0.9514853	
## prenumb	5.3244531	0.5629957	4.1266902	0.46613752	-1.1757433	
## prerelat	5.7847029	4.3568496	2.3557526	0.41517054	-0.7713901	
## preclasf	4.2934252	3.4337476	4.4184872	-2.41447115	1.6455558	
## postbody	-0.4435060	8.1660541	7.3159714	-2.45506814	-3.9035158	
## postlet	1.6542452	11.2009875	5.6507274	0.33103211	-1.3059478	
## postform	-0.7233607	2.6516348	5.6151521	-1.25130634	1.5584843	
## postnumb	2.8756649	1.6981872	0.4828499	0.05621188	0.7314776	
## postrelat	1.4839647	1.9238753	1.4612181	-1.15455355	-0.2312538	
## postclasf	4.5310373	11.0219926	3.2089732	0.37348213	2.9720521	
## peabody	-2.2077087	21.4186474	11.6438247	-0.41364479	-0.1274023	
##		MeanDecreaseAccuracy	MeanDecreaseGini			

```
## age          10.4592525      11.163988
## viewcat      7.4350785      5.559341
## setting      6.1976799      3.318730
## viewenc      4.0313230      2.348905
## prebody      7.5362285      8.417371
## prelet       -0.3182436      6.801044
## preform      7.5511763      8.203401
## prenumb      5.0459069      8.472669
## prerelat     6.4707132      6.588968
## preclasf     5.7883091      7.710234
## postbody     7.7126091      7.617352
## postlet     10.8432090      9.543848
## postform     3.8259496      6.077053
## postnumb     2.8137769      8.046909
## postrelat    1.9491940      5.578544
## postclasf    10.8962374      8.841455
## peabody     17.4895159     15.594451
```

```
rf.pred <- predict(rf.tree, newdata=test.data)
```

```
tree.conMatrix <- table(true=test.data[, "site"],
                        pred=rf.pred)
confusionMatrix(tree.conMatrix)
```

```
## Confusion Matrix and Statistics
```

```
##
```

```
##      pred
## true  1  2  3  4  5
##      1  7  1  6  4  0
##      2  4 10  2  2  0
##      3  3  3 12  1  0
##      4  0  2  6  4  0
##      5  0  2  1  2  0
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##              Accuracy : 0.4583
##              95% CI : (0.3402, 0.58)
##      No Information Rate : 0.375
##      P-Value [Acc > NIR] : 0.09136
```

```
##
```

```
##              Kappa : 0.2871
```

```
##
```

```
## McNemar's Test P-Value : NA
```

```
##
```

```
## Statistics by Class:
```

```
##
```

```
##              Class: 1 Class: 2 Class: 3 Class: 4 Class: 5
## Sensitivity      0.50000    0.5556    0.4444    0.30769      NA
## Specificity      0.81034    0.8519    0.8444    0.86441    0.93056
## Pos Pred Value    0.38889    0.5556    0.6316    0.33333      NA
## Neg Pred Value    0.87037    0.8519    0.7170    0.85000      NA
## Prevalence        0.19444    0.2500    0.3750    0.18056    0.00000
## Detection Rate    0.09722    0.1389    0.1667    0.05556    0.00000
## Detection Prevalence 0.25000    0.2500    0.2639    0.16667    0.06944
```

Balanced Accuracy 0.65517 0.7037 0.6444 0.58605 NA

0.42 – 0.4861. (Add more features, not including the standardized variables)

Questions for OH:

anything else needed in EDA?

Both linear and radial kernels never output predictions for 4 & 5?

polynomial kernel? Which variables to give polynomial terms

use PCA to perform feature selection?

feature selections for SVM in general?

how to interpret the confusion matrix tables for SVM & Trees

How to interpret the importance variance for multiclass classification

How to write for SVM interpretations