# The impact of double blind reviewing at EvoLang 12: statistics

## Contents

## Introduction

## Data

This script uses the data file `EvoLang_Scores_8_to_12.csv`:

- conference: Which conference the paper was submitted to
- gender: Gender of first author
- Score.Mean: Mean raw score given by reviewers (scaled between 0 and 1, higher = better paper)
- student: The student status of the first author at submission.

All variables with an underscore are measures of readability. Below we calculate a variable `review`, which represents the type of review (Single / Double blind).

## Loading data for first analysis

Load libraries.

```
# Load data
library(lattice)
library(ggplot2)
library(gplots)
library(lme4)
library(car)
library(caret)
library(dplyr)
library(party)
library(lmerTest)
```

```r
# read data
allData = read.csv("../data/EvoLang_Scores_8_to_12.csv",stringsAsFactors = F)
# relabel factor
allData$FirstAuthorGender = factor(allData$FirstAuthorGender,labels=c("F","M"))
allData$review = factor(c("Single","Double")[(allData$conference %in% c("E11","E12"))+1])
allData$conference = factor(allData$conference,levels = c("E8","E9","E10","E11","E12"))
allData$format = factor(allData$format)

allData$student[!is.na(allData$student) &
                  allData$student=="Faculty"] = "Non-Student"
allData$student[!is.na(allData$student) &
                  allData$student=="EC"] = "Non-Student"
allData$student = factor(allData$student)

#allData$Score.mean = scale(allData$Score.mean)

for(conf in levels(allData$conference)){
  allData$Score.mean[allData$conference==conf] = scale(allData$Score.mean[allData$conference==conf])
}
```

Look at the distribution of submissions:

```r
table(allData$FirstAuthorGender,allData$conference)
```

```
##
##      E8  E9 E10 E11 E12
##   F  58  52  67  76  84
##   M  95 130 124 119 122
```

```r
prop.table(table(allData$FirstAuthorGender,allData$conference),2)
```

```
##
##            E8        E9       E10       E11       E12
##   F 0.3790850 0.2857143 0.3507853 0.3897436 0.4077670
##   M 0.6209150 0.7142857 0.6492147 0.6102564 0.5922330
```

```r
gtable = table(allData$FirstAuthorGender,allData$conference,allData$student)
write.csv(cbind(t(gtable[,,1]),t(gtable[,,2])),
          "../results/CountTable.csv")
gtable
```

```
## , ,  = Non-Student
##
##
##      E8 E9 E10 E11 E12
##   F   0 34  55  41  54
##   M   0 85  94  77  93
##
## , ,  = Student
##
##
##      E8 E9 E10 E11 E12
##   F   0 18  12  35  30
##   M   0 45  30  42  29
```
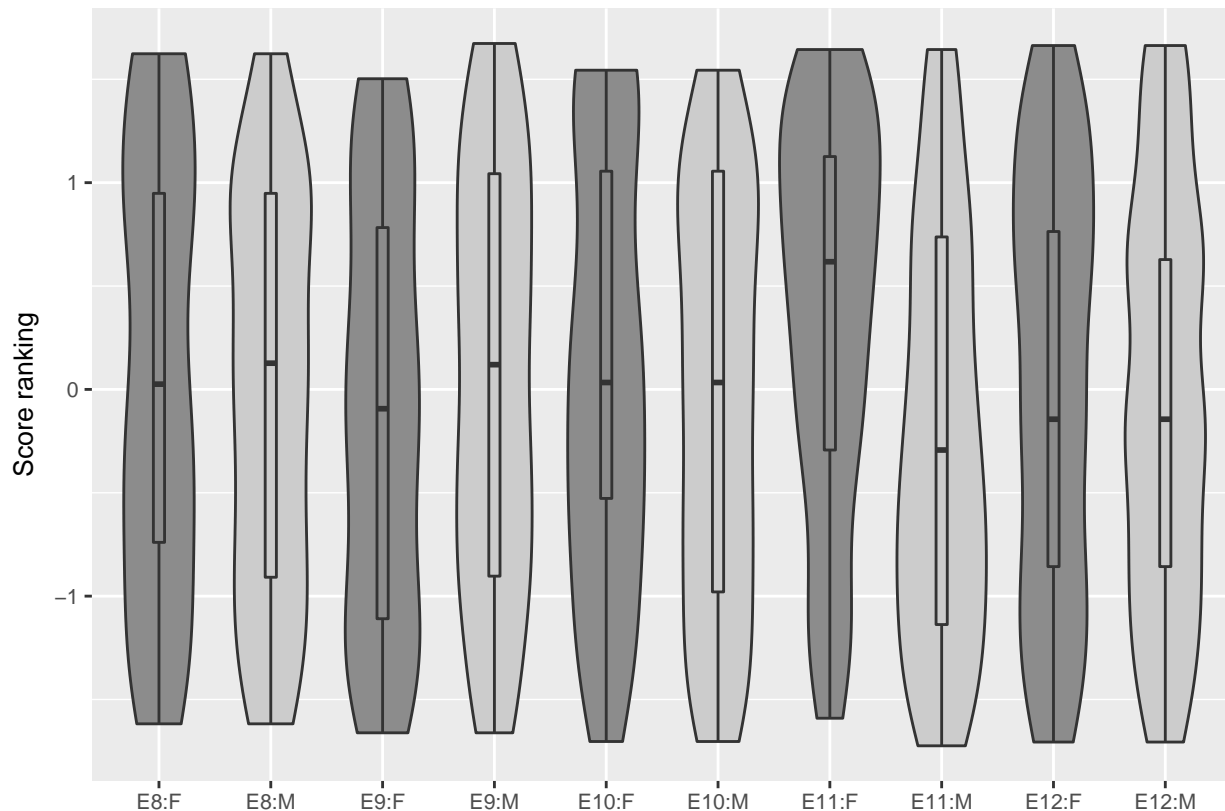
# Plots

Rank by gender. It seems that the difference in E11 is not replicated in E12.

```r
source("../analysis/summarySE.r")
p2 <- ggplot(allData,
             aes((conference):(FirstAuthorGender), Score.mean,
                 fill=FirstAuthorGender))

p2 <- p2 + geom_violin() + geom_boxplot(width=0.1) +
  theme(text=element_text(size=20), legend.position="none") +
  scale_y_continuous(name="Score ranking")+
  scale_x_discrete(name="")+
  scale_fill_grey(start = 0.55, end=0.8) +
  theme(text = element_text(size=10))

p2
```
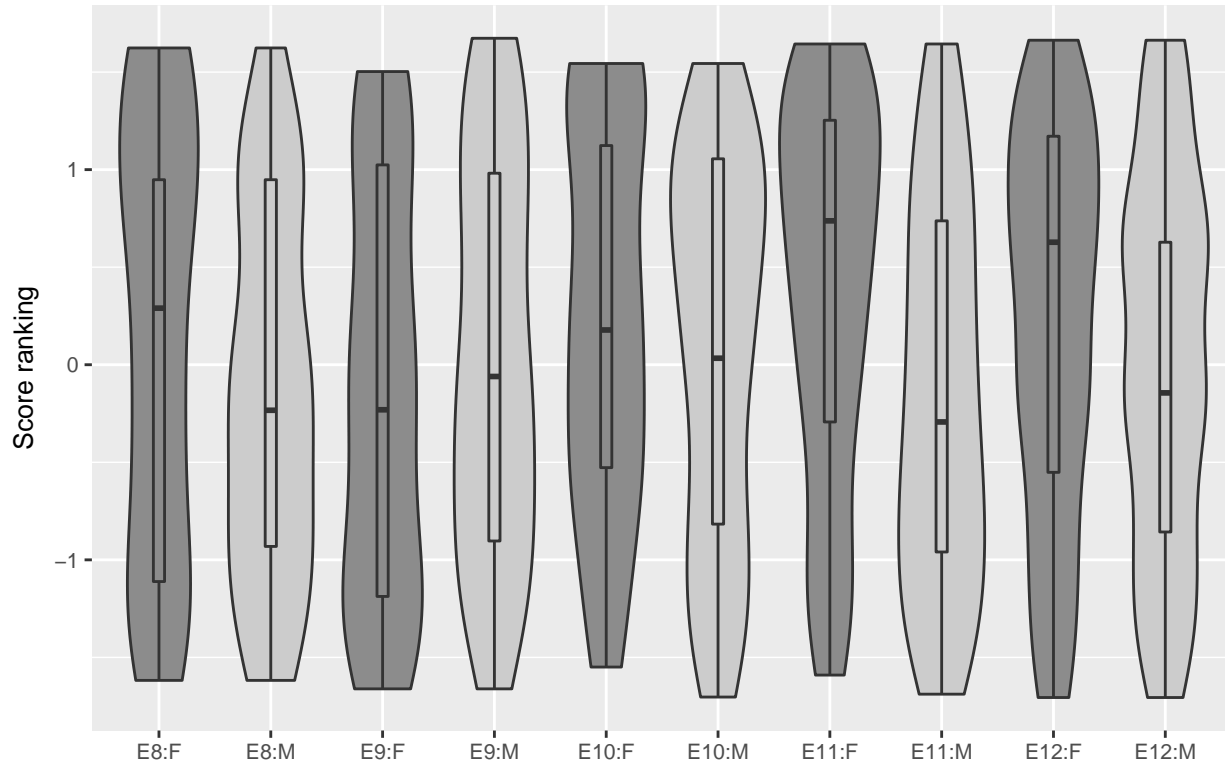


```r
pdf("../results/Results_Gender_3conf.pdf", width = 12, height= 6)
p2
dev.off()
```

```
## pdf
##   2
```

```r
p2Abstract <- ggplot(allData[allData$format=="Abstract",],
             aes((conference):(FirstAuthorGender), Score.mean,
                 fill=FirstAuthorGender))
```

```
p2Abstract <- p2Abstract + geom_violin() + geom_boxplot(width=0.1) +
  theme(text=element_text(size=20), legend.position="none") +
  scale_y_continuous(name="Score ranking")+
  scale_x_discrete(name="")+
  scale_fill_grey(start = 0.55, end=0.8) +
  theme(text = element_text(size=10)) +
  ggtitle("Scores for abstracts only")
p2Abstract
```
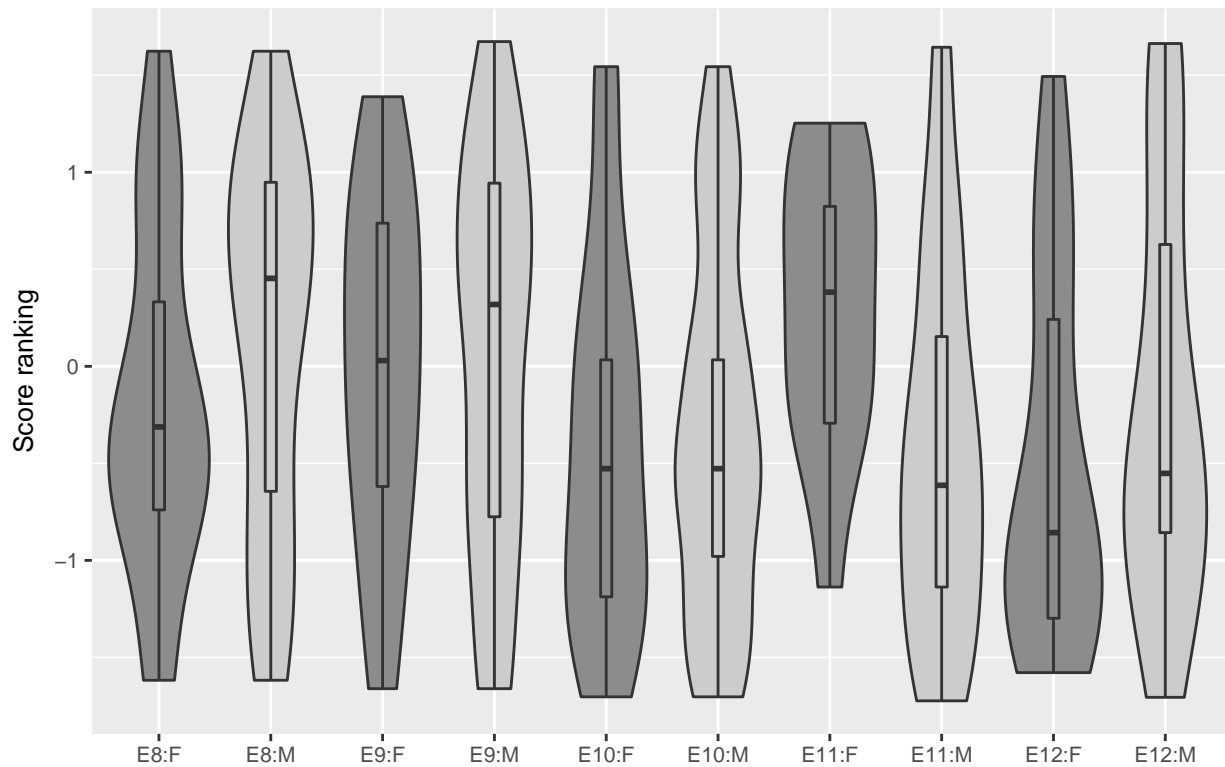
## Scores for abstracts only



```
p2Paper <- ggplot(allData[allData$format=="Paper",],
          aes((conference):(FirstAuthorGender), Score.mean,
              fill=FirstAuthorGender))

p2Paper <- p2Paper + geom_violin() + geom_boxplot(width=0.1) +
  theme(text=element_text(size=20), legend.position="none") +
  scale_y_continuous(name="Score ranking")+
  scale_x_discrete(name="")+
  scale_fill_grey(start = 0.55, end=0.8) +
  theme(text = element_text(size=10)) +
  ggtitle("Scores for full papers only")
p2Paper
```

Scores for full papers only

Rank by student status in each conference.

```
p <- ggplot(allData[complete.cases(allData),], aes(conference:student, Score.mean, fill=student))

p <- p + geom_violin() + geom_boxplot(width=0.1) +
  theme(text=element_text(size=20), legend.position="none") +
  scale_y_continuous(name="Score ranking")+
  scale_x_discrete(name="")+
  scale_fill_grey(start = 0.55, end=0.8)+
  theme(text = element_text(size=10))
p
```

```
pdf("../results/Results_Student_3conf.pdf", width = 12, height= 6)
p
dev.off()
```
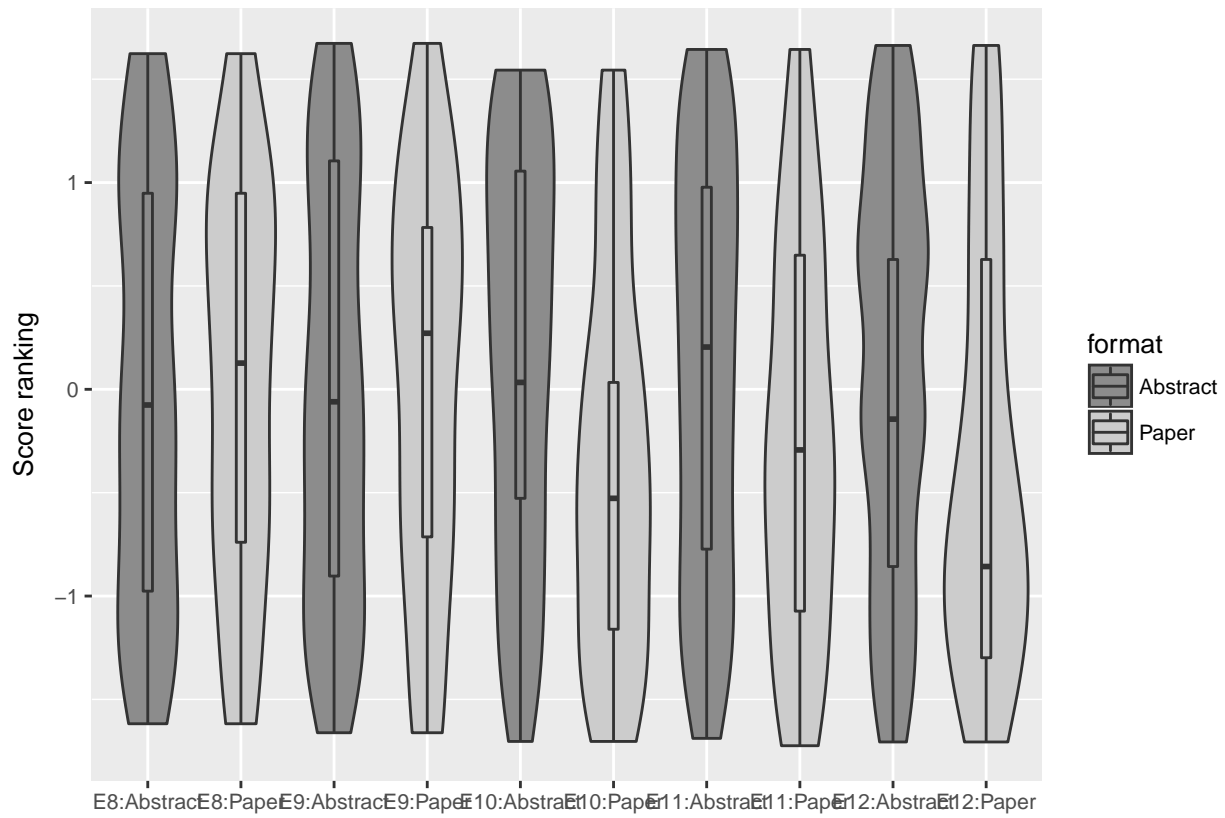
```
## pdf
##   2
```

Format:

```
p <- ggplot(allData, aes(conference:format, Score.mean, fill=format))

p <- p + geom_violin() + geom_boxplot(width=0.1) +
  theme(text=element_text(size=10)) +
  scale_y_continuous(name="Score ranking")+
  scale_x_discrete(name="")+
  scale_fill_grey(start = 0.55, end=0.8)
p
```

Combined student and gender:

```
ggplot(allData[allData$conference!="E8",],
       aes(y=Score.mean,x=paste(student,FirstAuthorGender),colour=conference))+ geom_boxplot(varwidth =
```

```
allData$stuGen = factor(paste(allData$conference,
                 allData$FirstAuthorGender),
                levels=c("E8 F","E8 M","E9 F","E9 M","E10 F","E10 M","E11 F","E11 M","E12 F",'E12 M'))

ad2 = allData[allData$conference!="E8",]

ggplot(ad2, mapping = aes(y=Score.mean,
         x=stuGen,
         colour=student))+
  geom_boxplot(varwidth = 0.5)
```

# Review ranks by gender and student status

Are papers with female first authors ranked higher than those with male first authors under double-blind review?

Using a simple anova, there's a significant interaction between gender and review type:

```
summary(aov(Score.mean ~ FirstAuthorGender*student*review*format,
            data=allData[allData$conference!="E8",]))
```

```
##                                      Df Sum Sq Mean Sq F value
## FirstAuthorGender                     1    5.4   5.366   5.551
## student                               1    0.4   0.423   0.438
## review                                1    0.1   0.054   0.056
## format                                1   11.7  11.747  12.151
## FirstAuthorGender:student             1    0.8   0.758   0.784
## FirstAuthorGender:review              1    4.3   4.278   4.425
## student:review                        1    0.3   0.302   0.313
## FirstAuthorGender:format              1    0.9   0.946   0.979
## student:format                        1   10.1  10.079  10.426
## review:format                         1    0.7   0.701   0.725
## FirstAuthorGender:student:review      1    0.0   0.005   0.005
## FirstAuthorGender:student:format      1    0.0   0.037   0.038
## FirstAuthorGender:review:format       1    0.3   0.270   0.279
## student:review:format                 1    2.1   2.124   2.197
## FirstAuthorGender:student:review:format 1  0.1   0.080   0.082
## Residuals                           758  732.8   0.967
##                                      Pr(>F)
## FirstAuthorGender                    0.018726 *
## student                              0.508378
## review                               0.813575
## format                               0.000519 ***
## FirstAuthorGender:student            0.376058
## FirstAuthorGender:review             0.035743 *
## student:review                       0.576264
## FirstAuthorGender:format             0.322788
## student:format                       0.001296 **
## review:format                        0.394665
## FirstAuthorGender:student:review     0.943998
## FirstAuthorGender:student:format     0.844520
## FirstAuthorGender:review:format      0.597387
## student:review:format                0.138730
## FirstAuthorGender:student:review:format 0.774242
## Residuals
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

However, it looks like this is driven just by EvoLang11:

```
t.test.string = function(tx){
  t = signif(tx$statistic,2)
  df = tx$parameter['df']
  p = signif(tx$p.value,3)
  est = signif(diff(tx$estimate),2)

  paste("(difference in means = ",est,", t = ",t,", p = ",p,")",sep = "")
```

```
}
for(conf in levels(allData$conference)){
  print(conf)
  print(t.test.string(t.test(Score.mean~FirstAuthorGender, data=allData[allData$conference==conf,])))
}
```

```
## [1] "E8"
## [1] "(difference in means = -0.1, t = 0.6, p = 0.55)"
## [1] "E9"
## [1] "(difference in means = 0.14, t = -0.87, p = 0.386)"
## [1] "E10"
## [1] "(difference in means = -0.12, t = 0.75, p = 0.454)"
## [1] "E11"
## [1] "(difference in means = -0.61, t = 4.4, p = 1.93e-05)"
## [1] "E12"
## [1] "(difference in means = -0.058, t = 0.4, p = 0.687)"
```

There is also a significant main effect of first author gender.

The model above mots EvoLang 8 because it has no data for student status. We get the same results if we omit student status and run the test for all conferences:

```
summary(aov(Score.mean ~ FirstAuthorGender*review*format,
            data=allData))
```

```
##                                 Df Sum Sq Mean Sq F value   Pr(>F)
## FirstAuthorGender                1    5.6   5.594   5.719  0.01699 *
## review                           1    0.0   0.023   0.024  0.87744
## format                           1    8.5   8.497   8.687  0.00329 **
## FirstAuthorGender:review         1    4.8   4.756   4.862  0.02770 *
## FirstAuthorGender:format         1    2.5   2.500   2.556  0.11024
## review:format                    1    1.7   1.695   1.732  0.18843
## FirstAuthorGender:review:format  1    0.0   0.023   0.023  0.87872
## Residuals                      919  898.9   0.978
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## Mixed effects model

Alternatively, we can use a mixed effects model, with random slopes for conference and test whether the interaction between gender and review type is a significant fixed predictor. A random intercept is not necessary, because the data is scaled to be centered around 0 within each conference. A random slope for the interaction between gender and review is also not permissable, since review type does not vary by conference.

```
contrasts(allData$FirstAuthorGender) <- contr.sum(2)/2
contrasts(allData$review) <- contr.sum(2)/2
contrasts(allData$student) <- contr.sum(2)/2
contrasts(allData$format) <- contr.sum(2)/2

m0 <- lmer(
      Score.mean ~
        1 + (FirstAuthorGender*review*student*format) +
        (0+FirstAuthorGender+student+format|conference),
      allData[allData$conference!="E8",],
  control=lmerControl(optimizer="bobyqa",optCtrl = list(maxfun=10000000)),
  REML = T
)

summary(m0)
```

```
## Linear mixed model fit by REML t-tests use Satterthwaite approximations
##   to degrees of freedom [lmerMod]
## Formula:
## Score.mean ~ 1 + (FirstAuthorGender * review * student * format) +
##     (0 + FirstAuthorGender + student + format | conference)
##    Data: allData[allData$conference != "E8", ]
## Control: lmerControl(optimizer = "bobyqa", optCtrl = list(maxfun = 1e+07))
##
## REML criterion at convergence: 2175.5
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -2.0469 -0.8318 -0.0649  0.8731  2.1003
##
## Random effects:
##  Groups     Name              Variance Std.Dev. Corr
##  conference FirstAuthorGenderF 0.049878 0.22333
##             FirstAuthorGenderM 0.002765 0.05258  -0.97
##             student1           0.045642 0.21364  -0.87  0.73
##             format1            0.023844 0.15441   0.37 -0.14 -0.77
##  Residual                      0.950489 0.97493
## Number of obs: 774, groups:  conference, 4
##
## Fixed effects:
##                                       Estimate Std. Error
## (Intercept)                          -0.005526   0.063844
## FirstAuthorGender1                    0.146719   0.166438
## review1                              -0.094290   0.127687
## student1                             -0.203825   0.142736
## format1                               0.154509   0.121783
## FirstAuthorGender1:review1            0.256651   0.332875
## FirstAuthorGender1:student1          -0.208867   0.189766
```

```
## review1:student1                                0.217541   0.285473
## FirstAuthorGender1:format1                       0.088045   0.188464
## review1:format1                                  0.286881   0.243566
## student1:format1                                 0.620946   0.189427
## FirstAuthorGender1:review1:student1              0.070548   0.379532
## FirstAuthorGender1:review1:format1               0.178654   0.376927
## FirstAuthorGender1:student1:format1              0.250443   0.377860
## review1:student1:format1                        -0.543252   0.378853
## FirstAuthorGender1:review1:student1:format1      0.151257   0.755720
##                                                         df t value Pr(>|t|)
## (Intercept)                                       2.900000  -0.087   0.9367
## FirstAuthorGender1                                2.600000   0.882   0.4519
## review1                                           2.900000  -0.738   0.5163
## student1                                          2.900000  -1.428   0.2528
## format1                                           3.400000   1.269   0.2845
## FirstAuthorGender1:review1                        2.600000   0.771   0.5046
## FirstAuthorGender1:student1                     674.800000  -1.101   0.2714
## review1:student1                                  2.900000   0.762   0.5040
## FirstAuthorGender1:format1                      749.300000   0.467   0.6405
## review1:format1                                   3.400000   1.178   0.3148
## student1:format1                                615.300000   3.278   0.0011 **
## FirstAuthorGender1:review1:student1             674.800000   0.186   0.8526
## FirstAuthorGender1:review1:format1              749.300000   0.474   0.6357
## FirstAuthorGender1:student1:format1             719.900000   0.663   0.5077
## review1:student1:format1                        615.300000  -1.434   0.1521
## FirstAuthorGender1:review1:student1:format1 719.900000   0.200   0.8414
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation matrix not shown by default, as p = 16 > 12.
## Use print(x, correlation=TRUE)  or
##     vcov(x)       if you need it
```

The results above suggest that there's no overall interaction between gender and review type. The tendency is there, but from the plots it's probably just driven by EvoLang 11.

We can run the same model without student status to include data from EvoLang 8:

```r
m0 <- lmer(
      Score.mean ~
        1 + (FirstAuthorGender*review*format) +
        (0+FirstAuthorGender+format|conference),
      allData,
  control=lmerControl(optimizer="bobyqa",optCtrl = list(maxfun=10000000)),
  REML = T
)

summary(m0)
```

```
## Linear mixed model fit by REML t-tests use Satterthwaite approximations
##   to degrees of freedom [lmerMod]
## Formula: Score.mean ~ 1 + (FirstAuthorGender * review * format) + (0 +
##     FirstAuthorGender + format | conference)
##    Data: allData
## Control: lmerControl(optimizer = "bobyqa", optCtrl = list(maxfun = 1e+07))
```

13

```
##
## REML criterion at convergence: 2619.3
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -2.00117 -0.86420 -0.03787  0.89640  2.01372
##
## Random effects:
##  Groups     Name               Variance Std.Dev. Corr
##  conference FirstAuthorGenderF 0.018981 0.1378
##             FirstAuthorGenderM 0.005492 0.0741   -0.62
##             format1            0.052243 0.2286   -0.40 -0.46
##  Residual                      0.966309 0.9830
## Number of obs: 927, groups:  conference, 5
##
## Fixed effects:
##                                  Estimate Std. Error        df t value
## (Intercept)                      -0.05026    0.04642   6.10000  -1.083
## FirstAuthorGender1                0.11759    0.11825   3.90000   0.994
## review1                          -0.02712    0.09284   6.10000  -0.292
## format1                           0.26284    0.13066   3.40000   2.012
## FirstAuthorGender1:review1        0.28928    0.23649   3.90000   1.223
## FirstAuthorGender1:format1        0.21307    0.15752 905.20000   1.353
## review1:format1                   0.17634    0.26131   3.40000   0.675
## FirstAuthorGender1:review1:format1 -0.05977   0.31505 905.20000  -0.190
##                                  Pr(>|t|)
## (Intercept)                         0.320
## FirstAuthorGender1                  0.377
## review1                             0.780
## format1                             0.127
## FirstAuthorGender1:review1          0.290
## FirstAuthorGender1:format1          0.177
## review1:format1                     0.543
## FirstAuthorGender1:review1:format1  0.850
##
## Correlation of Fixed Effects:
##                 (Intr) FrsAG1 reviw1 formt1 FrstAthrGndr1:r1
## FrstAthrGn1      0.448
## review1          0.202  0.068
## format1         -0.605 -0.159 -0.164
## FrstAthrGndr1:r1 0.068  0.204  0.448 -0.035
## FrstAthrGndr1:f1 -0.197 -0.333 -0.044  0.205 -0.125
## revw1:frmt1     -0.164 -0.035 -0.605  0.202 -0.159
## FrstAG1:1:1     -0.044 -0.125 -0.197  0.019 -0.333
##                 FrstAthrGndr1:f1 rvw1:1
## FrstAthrGn1
## review1
## format1
## FrstAthrGndr1:r1
## FrstAthrGndr1:f1
## revw1:frmt1      0.019
## FrstAG1:1:1      0.206            0.205
```

Again, there's no interaction between gender and review type.

## Permutation test

The distributions of score means are not very normal within conferences. We run a permutation test to address this. We calculate the average difference between single blind and double blind scores for males (dM) and for females (dF). Then we calculate dF - dM. A value > 0 means females scores increase more than male scores under double blind review. This 'true difference' is compared to a 'permuted difference'. The association between review scores and review type is randomly permuted, and dF - dM is calculated again. This is done 10,000 times to compare the true difference to a distribution of random differences.

```
meanDifferenceBetweenGenders = function(d){
  # difference in means between review types
  # for males
  # (change from single to double)
  diffMales = diff(rev(tapply(d[d$FirstAuthorGender=="M",]$Score.mean,
              d[d$FirstAuthorGender=="M",]$review,
              mean)))
  # for females
  diffFemales = diff(rev(tapply(d[d$FirstAuthorGender=="F",]$Score.mean,
              d[d$FirstAuthorGender=="F",]$review,
              mean)))
  # difference in differences
  # value > 0 means female scores increase
  # more under double-blind review than male scores
  return(diffFemales-diffMales)
}

perm = function(d){
  d$review = sample(d$review)
  meanDifferenceBetweenGenders(d)
}

perm.test = function(d,title){
  n = 10000
  trueDiff = meanDifferenceBetweenGenders(d)
  permDiff = replicate(n, perm(d))

  p = sum(permDiff>trueDiff) / n
  z = (trueDiff-mean(permDiff)) / sd(permDiff)
  print(paste("p=",p,", z=",z))
  hist(permDiff,xlab="Female advantage in double-blind",main=title)
  abline(v=trueDiff,col=2)
}
```
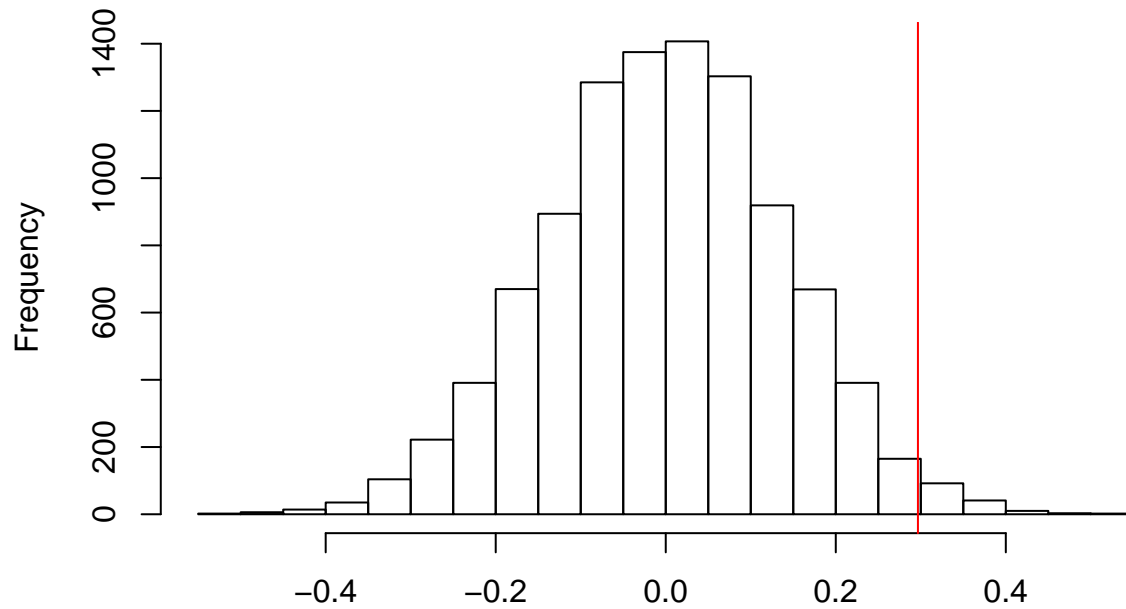
Permutation test for all data:

```
perm.test(allData,
          "All conferences")
```

```
## [1] "p= 0.0155 , z= 2.15609184767127"
```
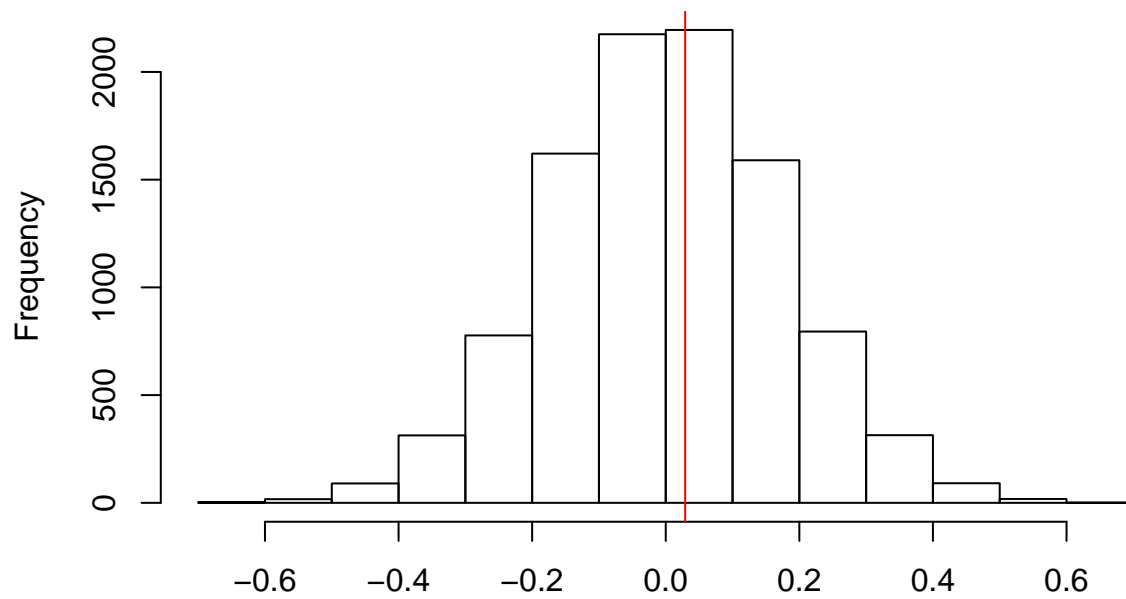
## All conferences



Permutation test without E11 data:

```
perm.test(allData[allData$conference!="E11",],
          "Without E11")
```

```
## [1] "p= 0.4331 , z= 0.168484592717635"
```
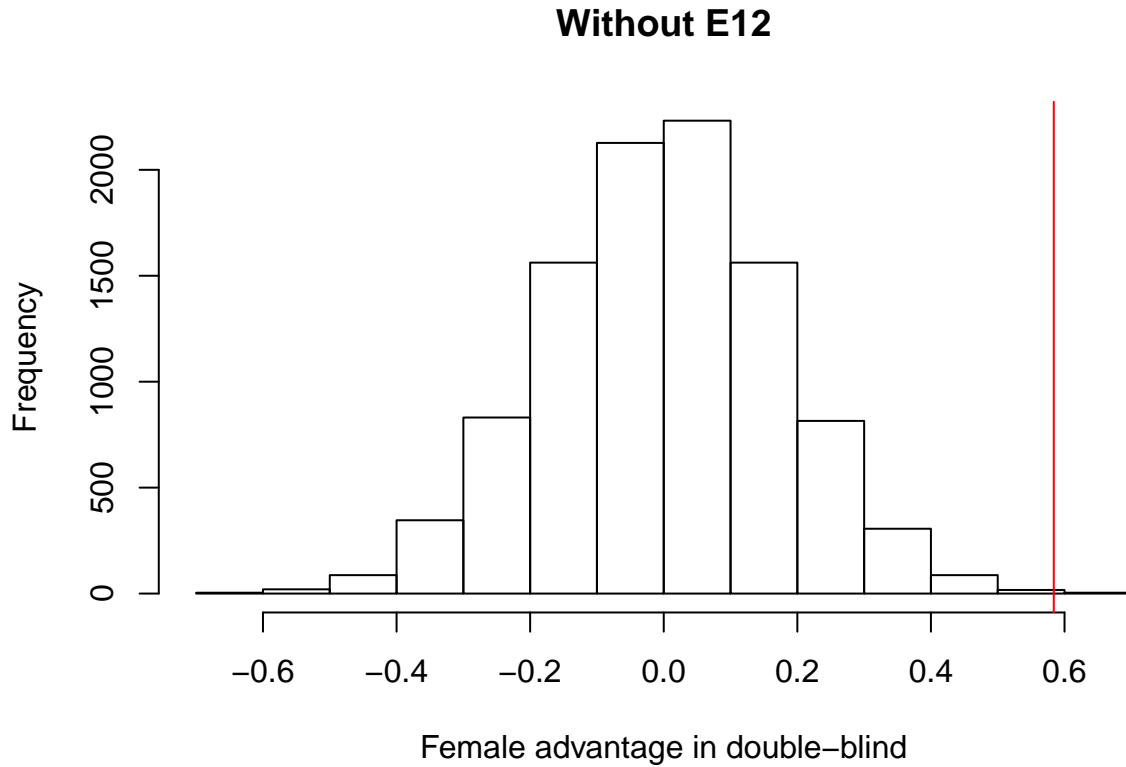
## Without E11

Permutation test without E12 data:

```
perm.test(allData[allData$conference!="E12",],
          "Without E12")
```

```
## [1] "p= 4e-04 , z= 3.34657009220246"
```

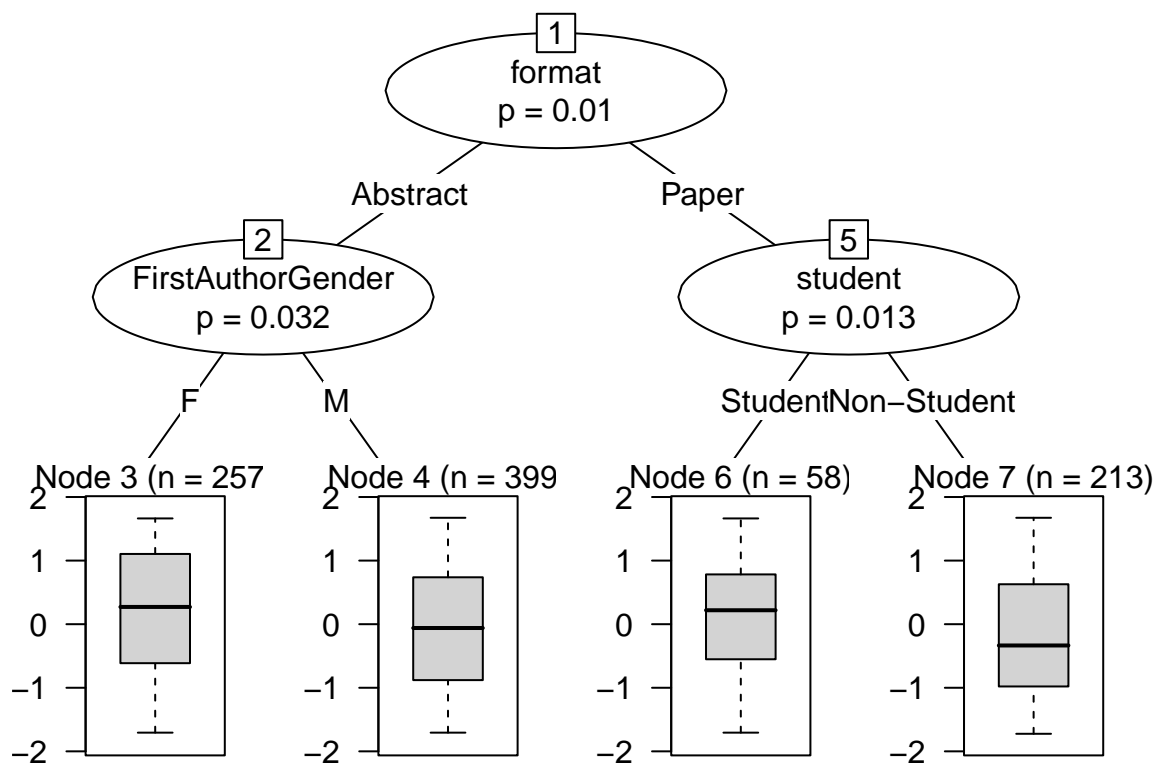## Without E12



Female advantage in double−blind

The results are in line with the test above. Across the whole data, females are given higher scores in double-blind, but this is driven by E11 alone.

## Decision tree exploration

Construct a decision tree, attempting to predict review socres by format, student status, gender, review model and conference.

```
set.seed(2389)
for(f in c("conference","format",'student','FirstAuthorGender','review')){
  allData[,f] = as.factor(allData[,f])
}
ct = ctree(Score.mean ~ format + student  +
             FirstAuthorGender + review + conference, data=allData)
plot(ct)
```



Work out differences between leaves of the tree:

```
paperVabstract = tapply(allData$Score.mean,allData$format,mean)
paperVabstract
```

```
##    Abstract       Paper
##  0.06519752 -0.15782129
```

```
pStudentVpNonStuent = tapply(allData[
  allData$format=="Paper",]$Score.mean,
  allData[allData$format=="Paper",]$student,mean)
pStudentVpNonStuent
```

```
## Non-Student     Student
##  -0.3300312   0.1235369
```

The tree suggests that full papers are given lower ratings than abstracts on average (about 6.6% difference). For full papers, students are given higher ratings than non-students (about 13.4% difference).

18

# Readability scores

This section uses the file `EvoLang_ReadingScores_E8_to_E12.csv`. It includes the following variables:

- `conference`: Conference
- `gender`: Gender of first author
- `student`: Student status
- `format`: Full paper or short abstract
- `char_count`, `word_count`, `sent_count`, `sybl_count`: Number of characters, words, sentences and syllables. These distributions have been scaled and centrered.
- `*_score`: Various measures of readability, calculated using the tools from Hengel (2016).
- Score.mean: Mean raw score given by reviewers (scaled between 0 and 1, higher = better paper)

Read the data:

```
readScores = read.csv("../data/EvoLang_ReadingScores_E8_to_E12.csv",stringsAsFactors = F)
```

We'll focus on the Flesch-Kinkaid score (since most other measures are highly correlated with it and it's easy to interpret) and the Dale-Chall score (which is not highly correlated with the other measures):

```
round(cor(readScores[,c("flesch_score","fleschkincaid_score",
                        "gunningfog_score" ,"smog_score","dalechall_score"
                        )]),2)
```

```
##                      flesch_score fleschkincaid_score gunningfog_score
## flesch_score                 1.00               -0.93            -0.92
## fleschkincaid_score         -0.93                1.00             0.99
## gunningfog_score            -0.92                0.99             1.00
## smog_score                  -0.94                0.97             0.99
## dalechall_score             -0.64                0.55             0.55
##                      smog_score dalechall_score
## flesch_score              -0.94           -0.64
## fleschkincaid_score        0.97            0.55
## gunningfog_score           0.99            0.55
## smog_score                 1.00            0.56
## dalechall_score            0.56            1.00
```

Scale the variables:

```
readScores$fleschkincaid_score_scaled = scale(readScores$fleschkincaid_score)
readScores$dalechall_score_scaled = scale(readScores$dalechall_score)
readScores$student[readScores$student=="EC"] = "Non-Student"
readScores$student[readScores$student=="Faculty"] = "Non-Student"
# Remove an outlier
readScores = readScores[readScores$fleschkincaid_score_scaled<6,]
readScores$gender = factor(readScores$gender)

readScores$conference = factor(readScores$conference,
                               levels = c("E8","E9","E10","E11","E12"))

# Box-Cox scaling
pp = preProcess(readScores[,
        c('fleschkincaid_score',"dalechall_score")],
        method="BoxCox")
lambda.fk = pp$bc$fleschkincaid_score$lambda
lambda.dc = pp$bc$dalechall_score$lambda
readScores$fleschkincaid_score_norm =
```

```
    bcPower(readScores$fleschkincaid_score, lambda = lambda.fk)
readScores$dalechall_score_norm =
    bcPower(readScores$dalechall_score, lambda = lambda.dc)
readScores$Score.mean.norm = scale(readScores$Score.mean)

readScores$review = factor(c("Single","Double")[(readScores$conference %in% c("E11","E12"))+1])
readScores$student = factor(readScores$student)
readScores$format = factor(readScores$format)
```

Create `time` variable: a continuous variable increasing with each conference.

```
readScores$time = as.numeric(readScores$conference)-3
```

Number of available datapoints (less than the total because some papers could not be automatically converted to text):

```
table(readScores$conference,readScores$gender)
```

```
##
##          F   M
##   E8    56  94
##   E9    52 130
##   E10   67 120
##   E11   68 111
##   E12   84 121
```

```
gtable2 = table(readScores$gender,readScores$conference,readScores$student)
write.csv(cbind(t(gtable2[,,1]),t(gtable2[,,2])),
          "../results/CountTable_Readability.csv")
gtable2
```
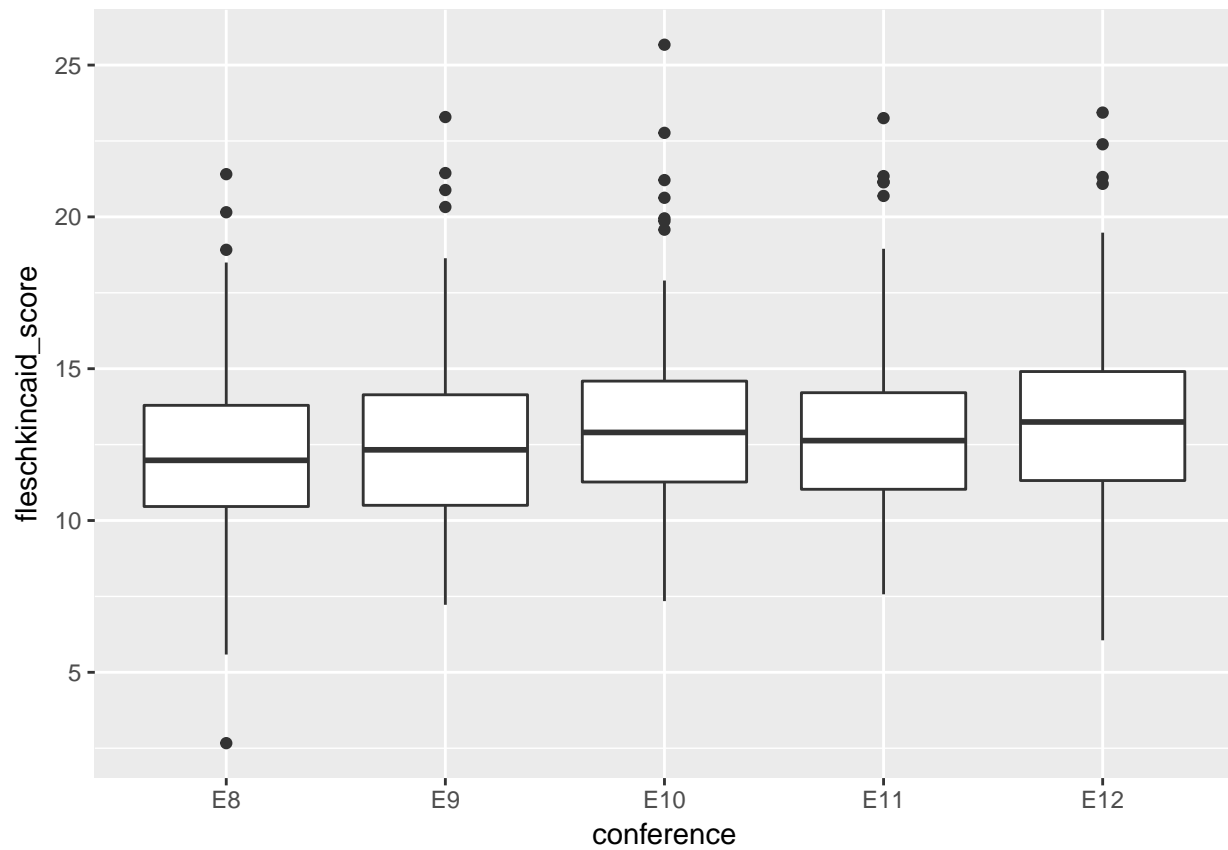
```
## , ,  = Non-Student
##
##
##      E8 E9 E10 E11 E12
##   F   0 34  55  38  54
##   M   0 85  90  72  92
##
## , ,  = Student
##
##
##      E8 E9 E10 E11 E12
##   F   0 18  12  30  30
##   M   0 45  30  39  29
```
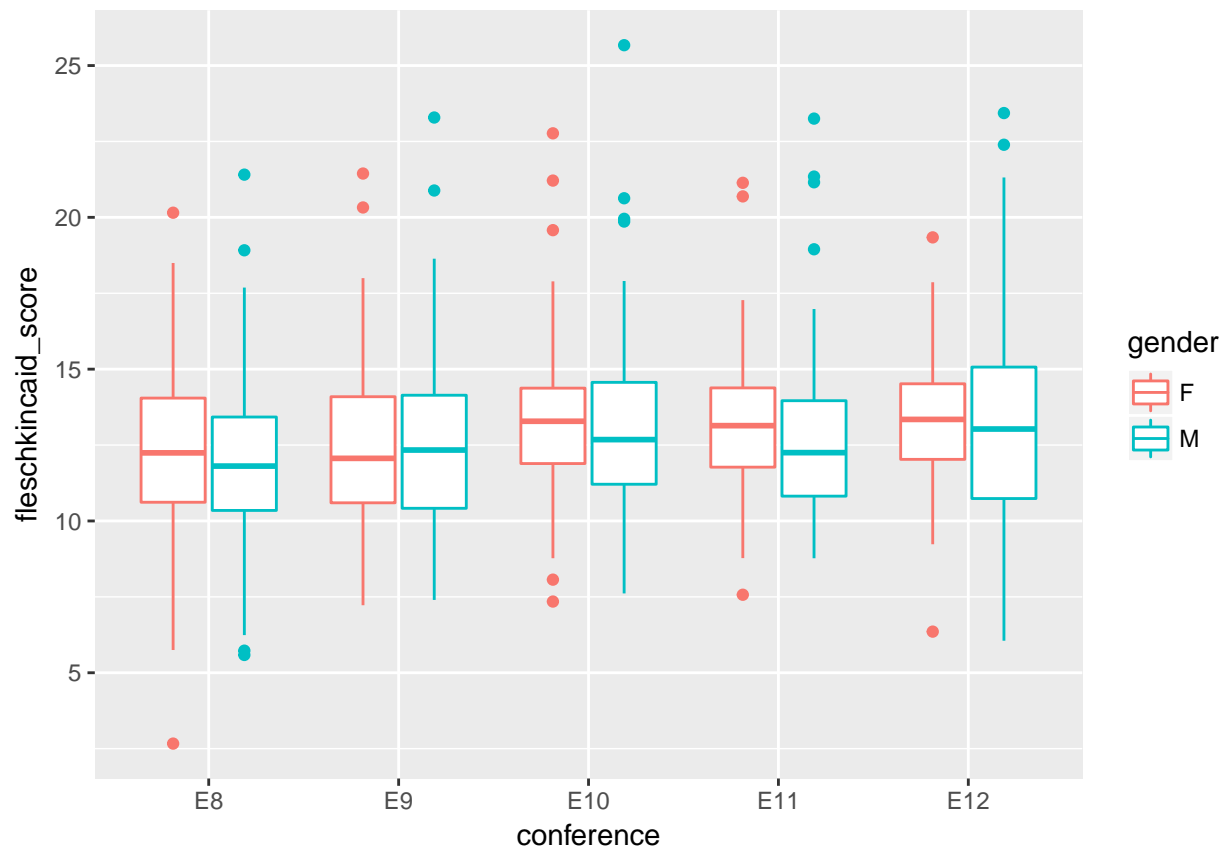
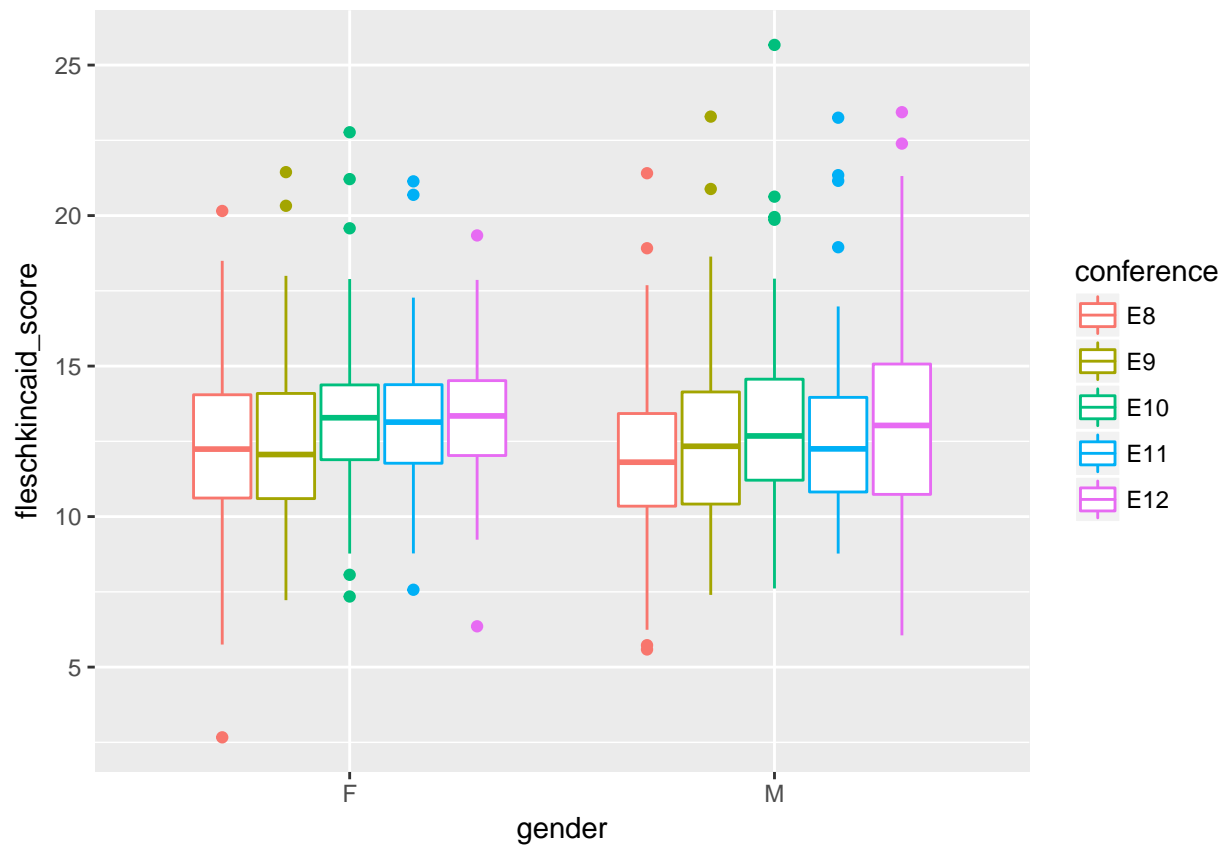**Flesch-Kinkaid score**

Various Plots:

```
ggplot(readScores, aes(y=fleschkincaid_score,x=conference)) + geom_boxplot()
```



```
ggplot(readScores, aes(y=fleschkincaid_score,x=conference,colour=gender)) + geom_boxplot()
```

```
ggplot(readScores, aes(y=fleschkincaid_score,x=gender,colour=conference)) + geom_boxplot()
```
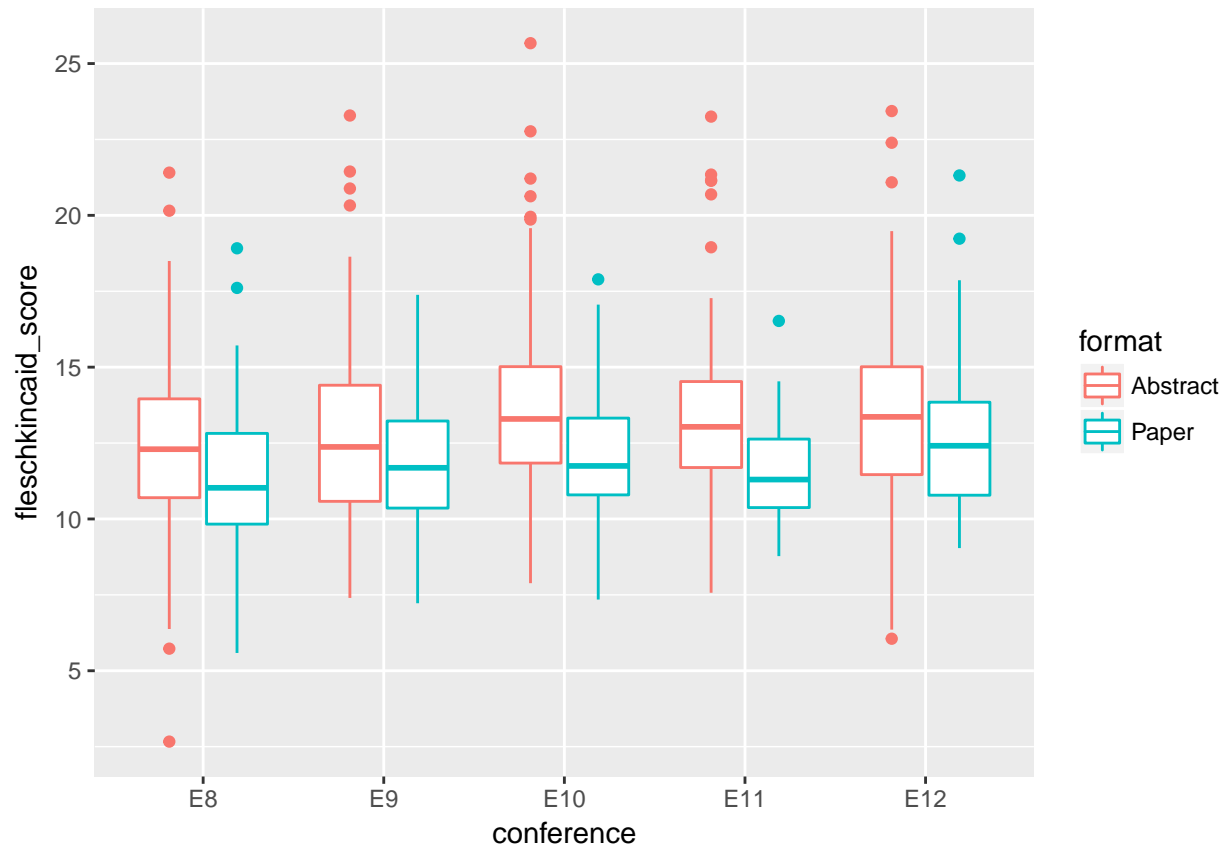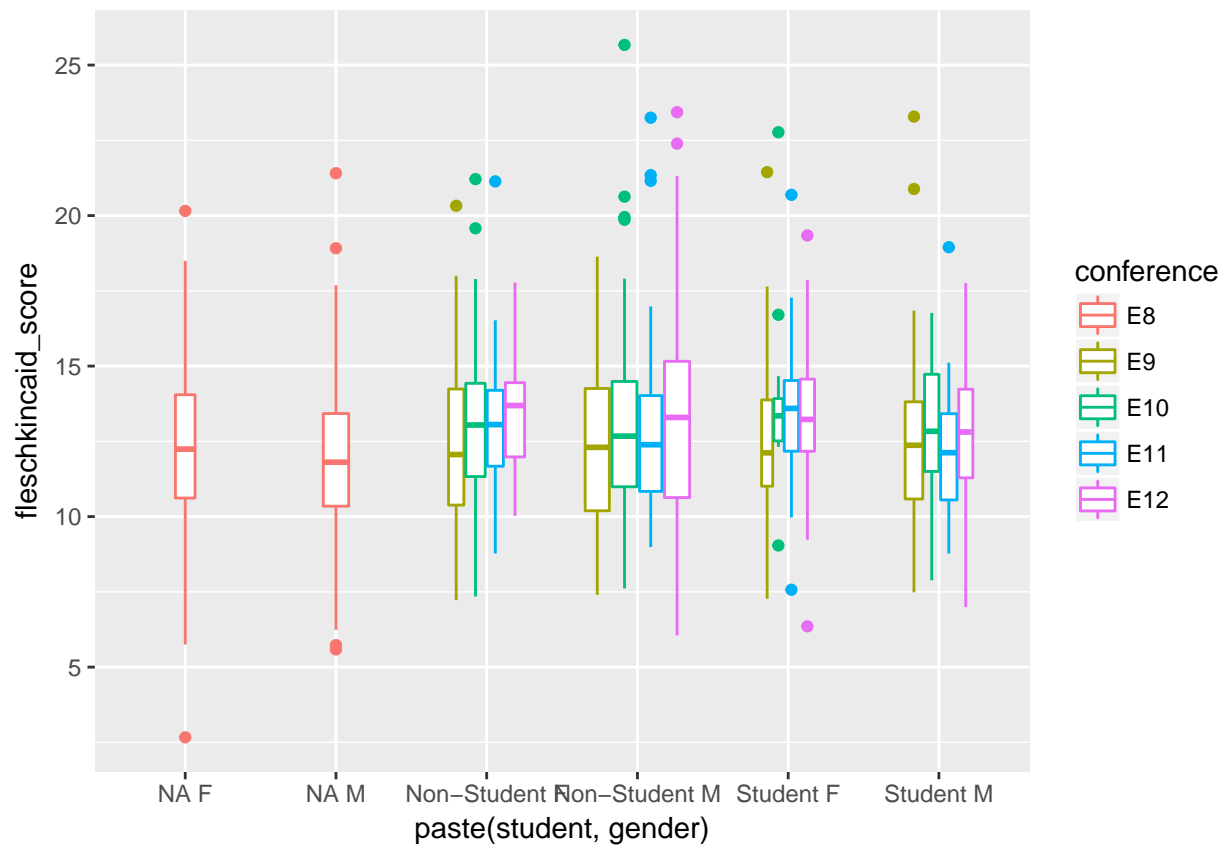
```
ggplot(readScores, aes(y=fleschkincaid_score,x=conference,colour=format)) + geom_boxplot()
```

```
ggplot(readScores, aes(y=fleschkincaid_score,x=paste(student,gender),colour=conference))+ geom_boxplot(v
```

```
x = readScores %>% group_by(conference,gender,student) %>%
  summarise(dalechall_score=mean(dalechall_score),
            fleschkincaid_score=mean(fleschkincaid_score))
ggplot(x,aes(x=(conference),y=fleschkincaid_score,
             group=paste(gender,student),
             colour=paste(gender,student))) +
  geom_line() + geom_point()
```

```
ggplot(readScores,
       aes(x=fleschkincaid_score,
           y=dalechall_score,
           colour=format)) +
  geom_point()
```
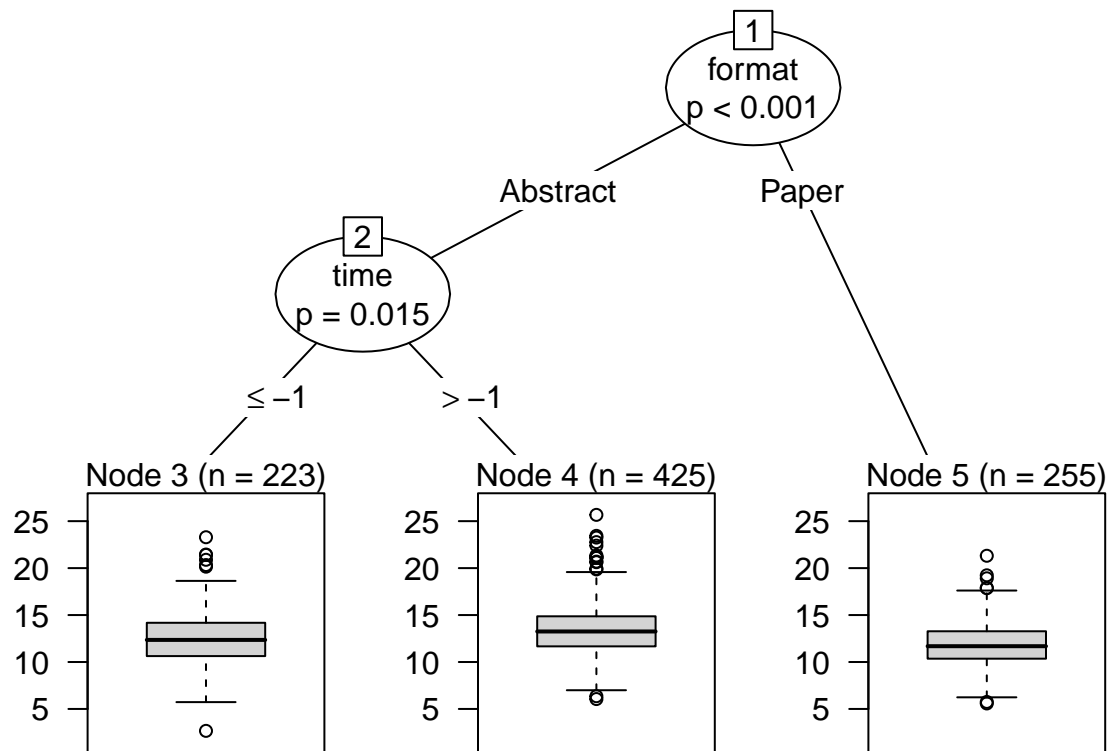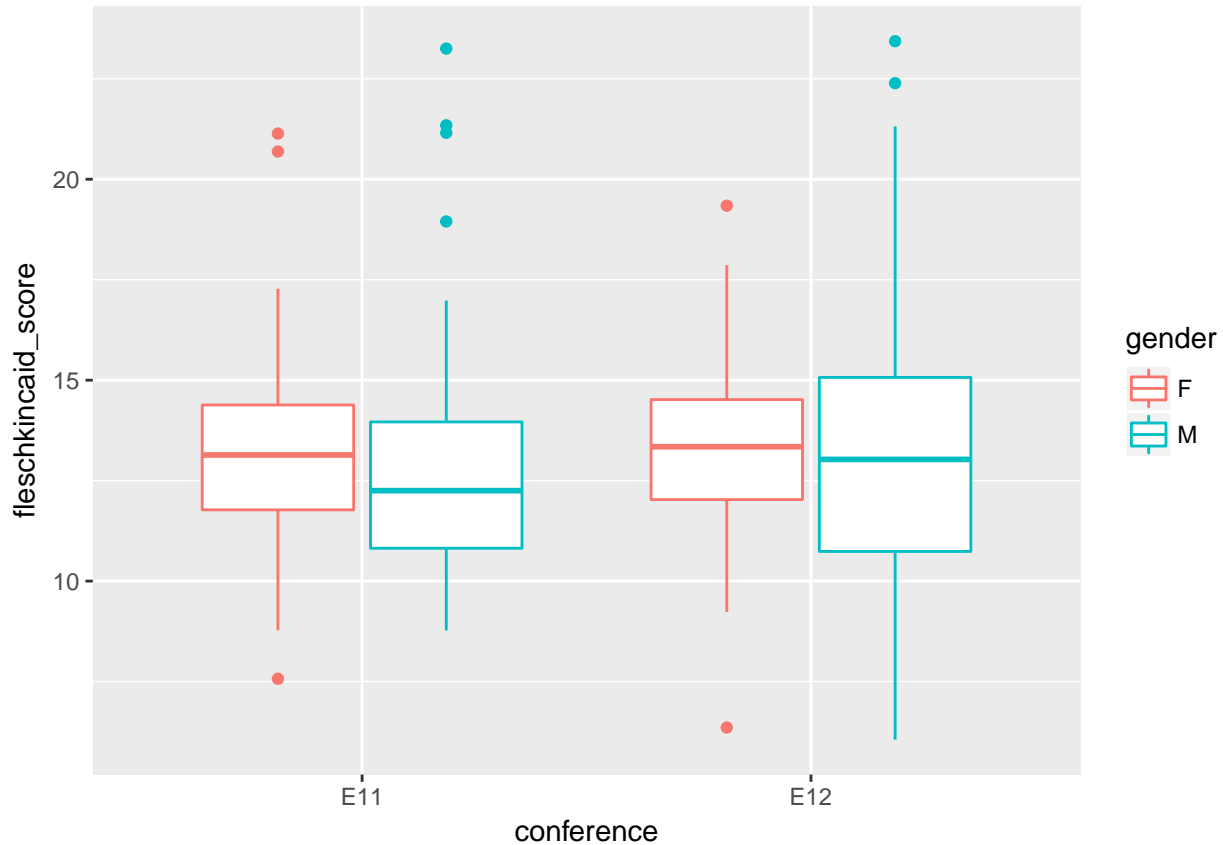
Decision tree

```
plot(ctree(fleschkincaid_score~
              review+gender+time+format,
          data=readScores))
```

Is there a gender difference between E11 and E12?

```
ggplot(readScores[readScores$conference %in% c("E11","E12"),],
       aes(x = conference, y=fleschkincaid_score, colour=gender)) +
  geom_boxplot()
```



```
summary(aov(fleschkincaid_score_norm~
            format*conference*student*gender,
            data = readScores[readScores$conference %in% c("E11","E12"),]))
```

```
##                                Df Sum Sq Mean Sq F value   Pr(>F)
## format                          1   4.00   4.003  13.269 0.000309 ***
## conference                      1   0.43   0.434   1.440 0.230956
## student                         1   0.40   0.396   1.314 0.252439
## gender                          1   0.44   0.440   1.458 0.227976
## format:conference               1   1.15   1.150   3.813 0.051614 .
## format:student                  1   0.45   0.447   1.481 0.224434
## conference:student              1   0.00   0.003   0.009 0.926404
## format:gender                   1   0.27   0.270   0.896 0.344608
## conference:gender               1   0.02   0.019   0.064 0.801013
## student:gender                  1   0.56   0.556   1.842 0.175582
## format:conference:student       1   0.23   0.234   0.776 0.378962
## format:conference:gender        1   0.01   0.012   0.038 0.845177
## format:student:gender           1   0.08   0.081   0.270 0.603717
## conference:student:gender       1   0.11   0.113   0.374 0.541115
## format:conference:student:gender 1  0.42   0.424   1.406 0.236534
## Residuals                     368 111.02   0.302
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

There is an effect for format, but nothing else.

Mixed effects model across the whole readability data. The model was not converging with a random slope for student, so:

```
contrasts(readScores$gender) <- contr.sum(2)/2
contrasts(readScores$student) <- contr.sum(2)/2
contrasts(readScores$format) <- contr.sum(2)/2

m0 = lmer(fleschkincaid_score_scaled~ 1 +
          (format*student*gender*review) + time +
          (1 + format + student + gender | conference),
       data = readScores[readScores$conference!="E8",])
summary(m0)
```

```
## Linear mixed model fit by REML t-tests use Satterthwaite approximations
##   to degrees of freedom [lmerMod]
## Formula: fleschkincaid_score_scaled ~ 1 + (format * student * gender *
##     review) + time + (1 + format + student + gender | conference)
##    Data: readScores[readScores$conference != "E8", ]
##
## REML criterion at convergence: 2047.4
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -2.8332 -0.6348 -0.0696  0.5286  4.5830
##
## Random effects:
##  Groups      Name        Variance  Std.Dev. Corr
##  conference (Intercept) 8.026e-05 0.008959
##             format1     2.903e-02 0.170371 1.00
##             student1    1.297e-03 0.036010 1.00 1.00
##             gender1     4.107e-03 0.064089 1.00 1.00 1.00
##  Residual               8.667e-01 0.930984
## Number of obs: 753, groups:  conference, 4
##
## Fixed effects:
##                                      Estimate Std. Error         df
## (Intercept)                          -0.21747    0.13356   21.80000
## format1                               0.25639    0.18453    3.50000
## student1                             -0.01655    0.14217   45.00000
## gender1                               0.25464    0.14762   18.80000
## reviewSingle                          0.24196    0.18251   20.10000
## time                                  0.16781    0.07526   13.10000
## format1:student1                      0.37429    0.27997  734.50000
## format1:gender1                      -0.37183    0.28122  701.00000
## student1:gender1                     -0.44971    0.28067  714.40000
## format1:reviewSingle                  0.11120    0.25634    3.20000
## student1:reviewSingle                -0.11476    0.19531   39.00000
## gender1:reviewSingle                 -0.16032    0.20217   16.40000
## format1:student1:gender1              0.39529    0.56039  720.10000
## format1:student1:reviewSingle        -0.39261    0.38496  689.30000
## format1:gender1:reviewSingle          0.34493    0.38415  724.60000
## student1:gender1:reviewSingle         0.22435    0.38472  719.70000
## format1:student1:gender1:reviewSingle -0.27203    0.76923  707.60000
##                                      t value Pr(>|t|)
## (Intercept)                           -1.628   0.1178
```

```
## format1                                     1.389   0.2473
## student1                                    -0.116   0.9078
## gender1                                      1.725   0.1009
## reviewSingle                                 1.326   0.1998
## time                                         2.230   0.0438 *
## format1:student1                             1.337   0.1817
## format1:gender1                             -1.322   0.1865
## student1:gender1                            -1.602   0.1095
## format1:reviewSingle                         0.434   0.6919
## student1:reviewSingle                       -0.588   0.5602
## gender1:reviewSingle                        -0.793   0.4391
## format1:student1:gender1                     0.705   0.4808
## format1:student1:reviewSingle               -1.020   0.3082
## format1:gender1:reviewSingle                 0.898   0.3695
## student1:gender1:reviewSingle                0.583   0.5600
## format1:student1:gender1:reviewSingle  -0.354   0.7237
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation matrix not shown by default, as p = 17 > 12.
## Use print(x, correlation=TRUE)   or
##   vcov(x)      if you need it
```
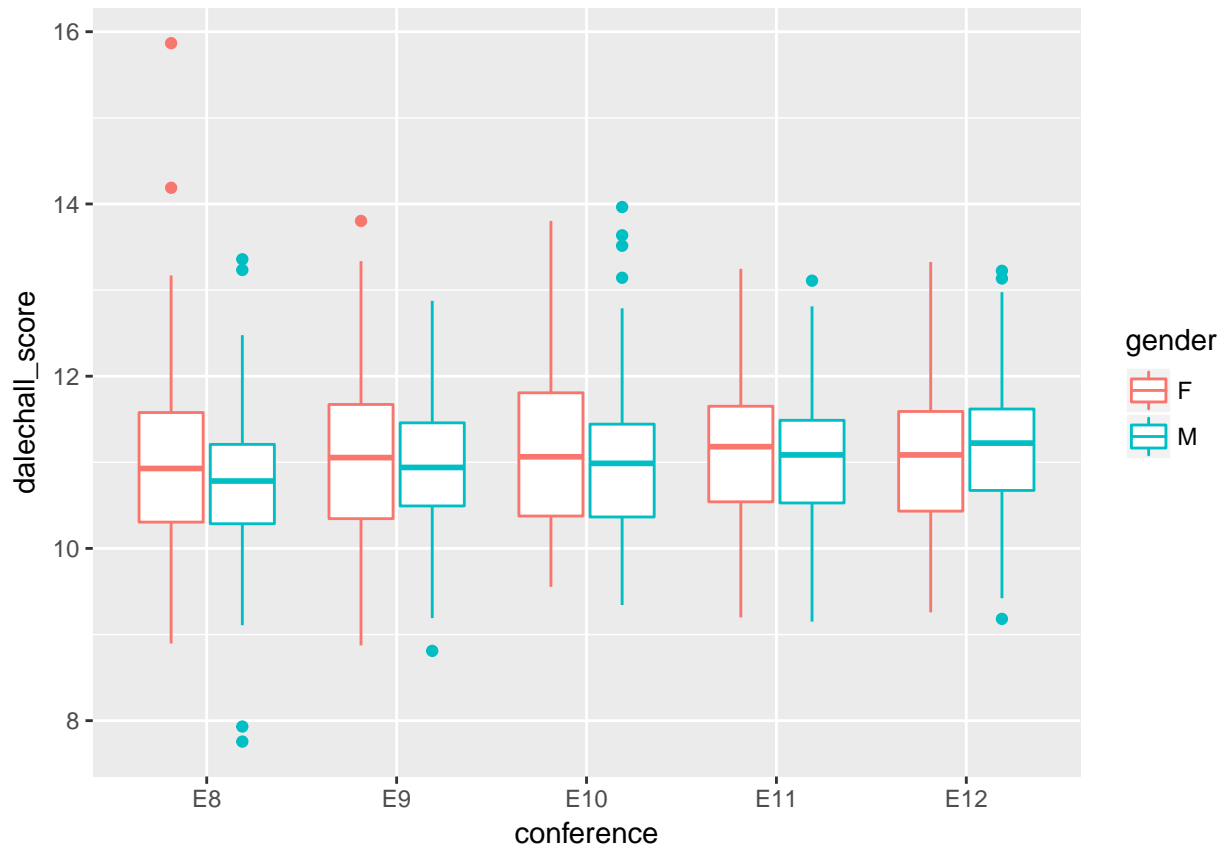
Abstracts have higher reading scores than papers, and socres are increasing over time, but there are no other significant effects.
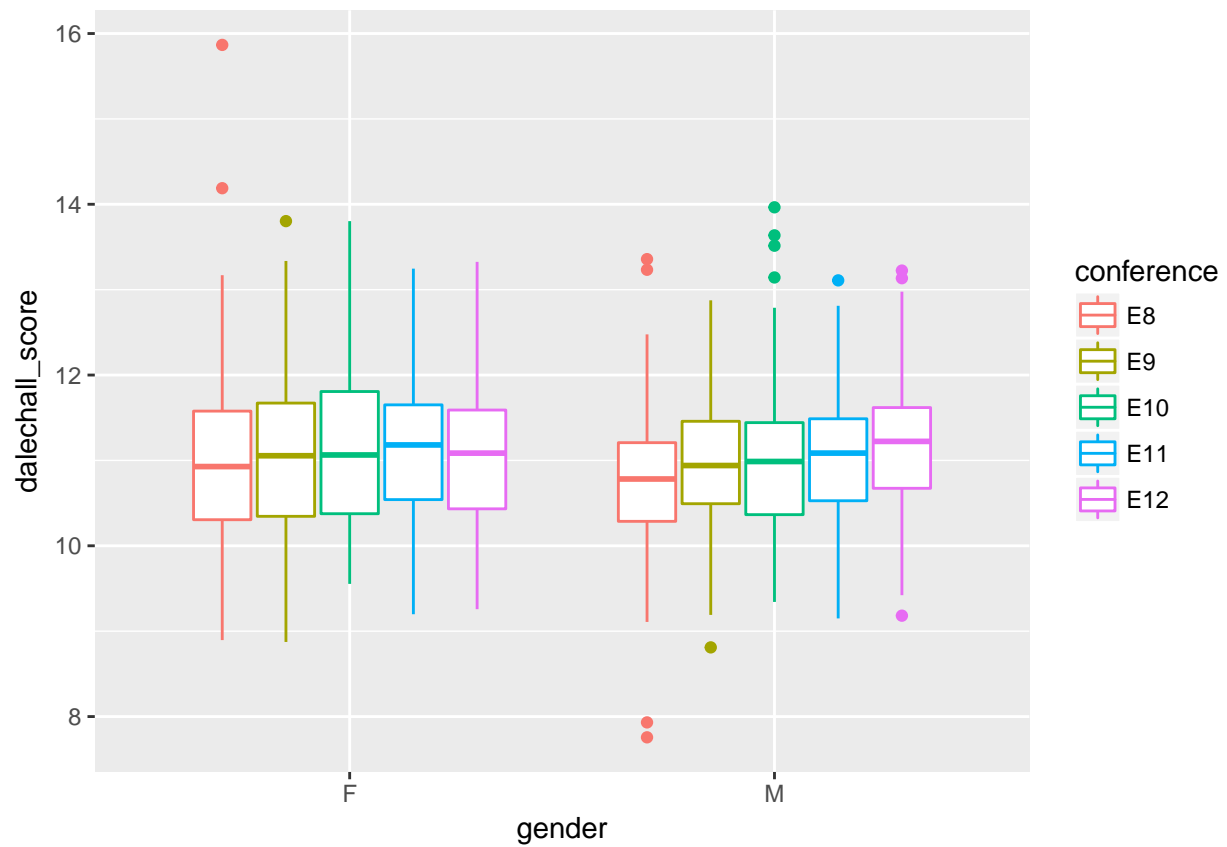

**Dale-Chall scale**

Plots

```
ggplot(readScores, aes(y=dalechall_score,x=conference,colour=gender)) + geom_boxplot()
```

```
ggplot(readScores, aes(y=dalechall_score,x=gender,colour=conference)) + geom_boxplot()
```
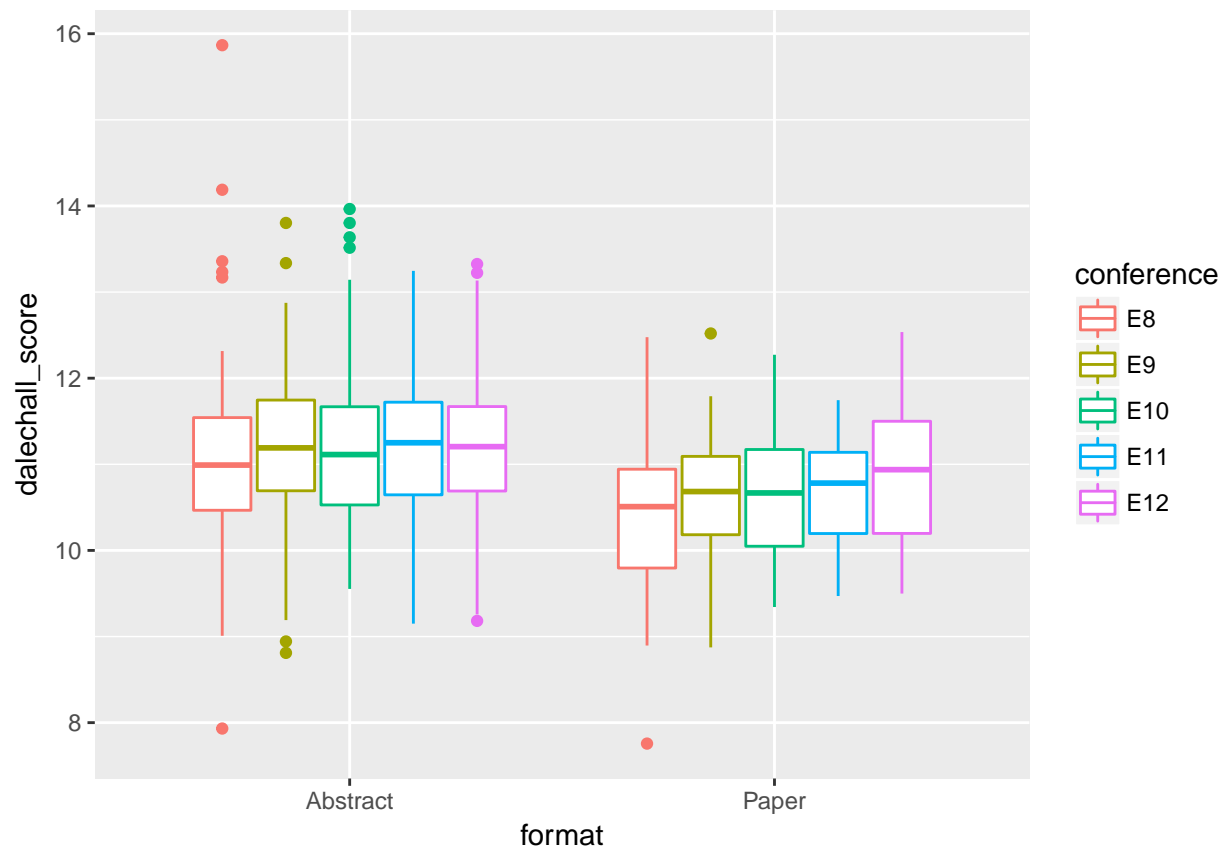
```
ggplot(x,aes(x=(conference),y=dalechall_score,group=paste(gender,student),colour=paste(gender,student))]
```
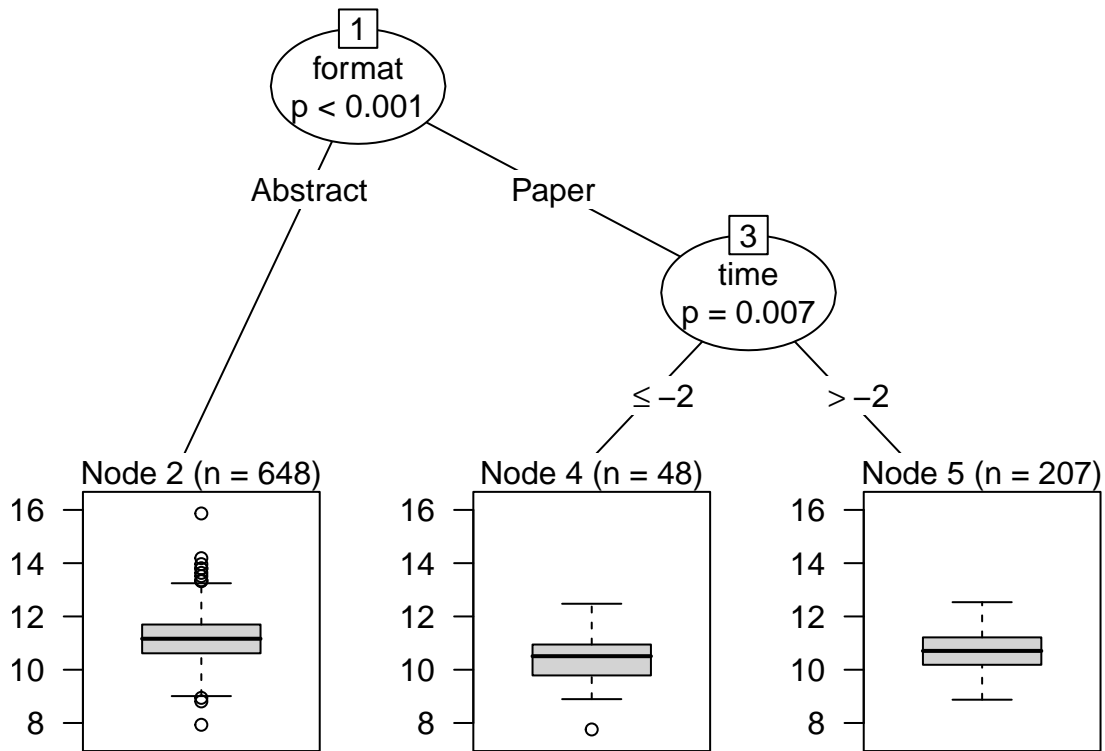
```
ggplot(readScores, aes(y=dalechall_score,x=format,colour=conference)) + geom_boxplot()
```
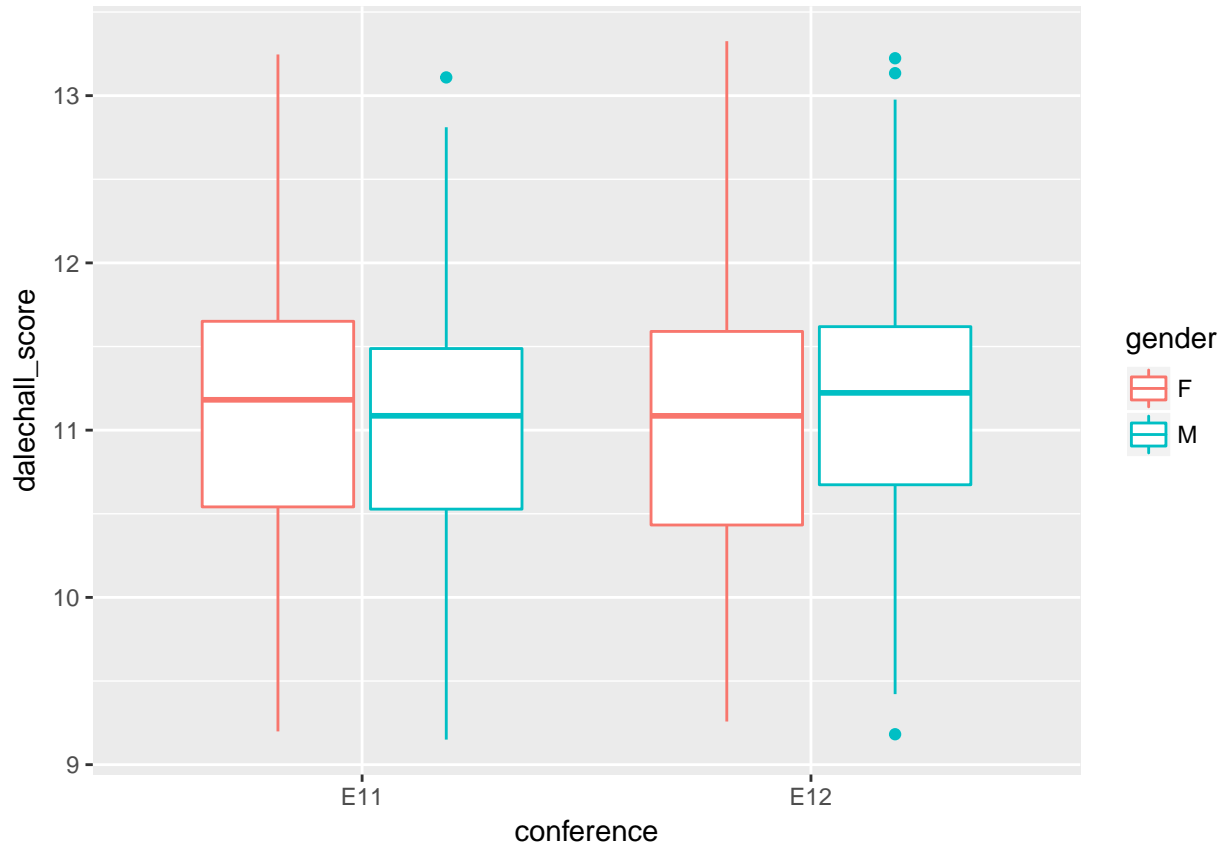
Decision tree:

```
plot(ctree(dalechall_score~review+gender+
           time+format,data=readScores))
```

Is there a gender difference between E11 and E12?

```
ggplot(readScores[readScores$conference %in% c("E11","E12"),],
       aes(x = conference, y=dalechall_score, colour=gender)) +
  geom_boxplot()
```



```
summary(aov(dalechall_score_norm~
            format*conference*student*gender,
            data = readScores[readScores$conference %in% c("E11","E12"),]))
```

```
##                                Df Sum Sq Mean Sq F value   Pr(>F)
## format                          1 0.1645 0.16449  19.627 1.24e-05 ***
## conference                      1 0.0005 0.00052   0.062   0.8035
## student                         1 0.0032 0.00318   0.379   0.5385
## gender                          1 0.0036 0.00360   0.430   0.5123
## format:conference               1 0.0224 0.02242   2.675   0.1028
## format:student                  1 0.0254 0.02539   3.029   0.0826 .
## conference:student              1 0.0003 0.00032   0.039   0.8443
## format:gender                   1 0.0007 0.00073   0.087   0.7687
## conference:gender               1 0.0035 0.00346   0.412   0.5212
## student:gender                  1 0.0032 0.00324   0.387   0.5345
## format:conference:student       1 0.0100 0.01000   1.193   0.2755
## format:conference:gender        1 0.0002 0.00018   0.021   0.8847
## format:student:gender           1 0.0049 0.00489   0.584   0.4454
## conference:student:gender       1 0.0032 0.00321   0.383   0.5363
## format:conference:student:gender 1 0.0036 0.00361   0.431   0.5119
## Residuals                     368 3.0841 0.00838
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

There's an effect for format, but nothing else.

Mixed effects model across whole data:

Scale and center the distribution, removing some outliers:

```
#readScores = readScores[readScores$student!="Student",]
sdx = 1.96 * sd(readScores$dalechall_score_norm)
mx = mean(readScores$dalechall_score_norm)
readScoresDC = readScores[
  readScores$dalechall_score_norm < (mx +sdx) &
  readScores$dalechall_score_norm > (mx -sdx)
,]
readScoresDC$dalechall_score_norm = scale(readScoresDC$dalechall_score_norm)

contrasts(readScoresDC$gender) <- contr.sum(2)/2
contrasts(readScoresDC$format) <- contr.sum(2)/2
contrasts(readScoresDC$student) <- contr.sum(2)/2
contrasts(readScoresDC$review) <- contr.sum(2)/2
```

Run mixed effects model:

```
m0 = lmer(dalechall_score_norm~ 1 +
            (format*student*gender*review) + time +
          (1 + format + student + gender | conference),
        data = readScoresDC[readScoresDC$conference!="E8",])
summary(m0)
```

```
## Linear mixed model fit by REML t-tests use Satterthwaite approximations
##   to degrees of freedom [lmerMod]
## Formula:
## dalechall_score_norm ~ 1 + (format * student * gender * review) +
##     time + (1 + format + student + gender | conference)
##    Data: readScoresDC[readScoresDC$conference != "E8", ]
##
## REML criterion at convergence: 2023.2
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -2.49923 -0.75103  0.04842  0.68345  2.38857
##
## Random effects:
##  Groups     Name        Variance Std.Dev. Corr
##  conference (Intercept) 0.008579 0.09262
##             format1     0.047145 0.21713  -1.00
##             student1    0.001599 0.03998  -1.00  1.00
##             gender1     0.001048 0.03238   1.00 -1.00 -1.00
##  Residual               0.936430 0.96769
## Number of obs: 724, groups:  conference, 4
##
## Fixed effects:
##                                 Estimate Std. Error        df t value
## (Intercept)                     -0.02927    0.07877   0.90000  -0.372
## format1                          0.48043    0.14897   2.70000   3.225
## student1                        -0.09727    0.10512  40.10000  -0.925
## gender1                         -0.01229    0.10376  53.30000  -0.118
## review1                          0.24252    0.20606   1.80000   1.177
## time                            -0.05479    0.07552   5.00000  -0.726
```

```
## format1:student1                 0.32725    0.20604 620.70000  1.588
## format1:gender1                  -0.04245    0.20464 689.40000 -0.207
## student1:gender1                 -0.07918    0.20523 645.30000 -0.386
## format1:review1                  -0.07322    0.29789   2.70000 -0.246
## student1:review1                 -0.10208    0.20965  39.60000 -0.487
## gender1:review1                  -0.08197    0.20727  54.00000 -0.395
## format1:student1:gender1          0.41973    0.41076 623.80000  1.022
## format1:student1:review1          0.02998    0.41246 633.10000  0.073
## format1:gender1:review1          -0.06977    0.40945 698.10000 -0.170
## student1:gender1:review1         -0.04201    0.41080 668.60000 -0.102
## format1:student1:gender1:review1 -1.54455    0.82220 653.80000 -1.879
##                                  Pr(>|t|)
## (Intercept)                       0.7792
## format1                           0.0564 .
## student1                          0.3603
## gender1                           0.9061
## review1                           0.3712
## time                              0.5009
## format1:student1                  0.1127
## format1:gender1                   0.8357
## student1:gender1                  0.6998
## format1:review1                   0.8234
## student1:review1                  0.6290
## gender1:review1                   0.6940
## format1:student1:gender1          0.3073
## format1:student1:review1          0.9421
## format1:gender1:review1           0.8647
## student1:gender1:review1          0.9186
## format1:student1:gender1:review1  0.0607 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation matrix not shown by default, as p = 17 > 12.
## Use print(x, correlation=TRUE)  or
##   vcov(x)      if you need it
```

Differences by format, but no other effects.

## Reading scores and review scores

The simple correlations between reading score and review scores are weak, but suggest that higher scores are given to submissions with higher reading grades:

```
cor.test(readScores$Score.mean, readScores$fleschkincaid_score)


##
##  Pearson's product-moment correlation
##
## data:  readScores$Score.mean and readScores$fleschkincaid_score
## t = 3.2308, df = 901, p-value = 0.001279
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.04207075 0.17106141
## sample estimates:
```
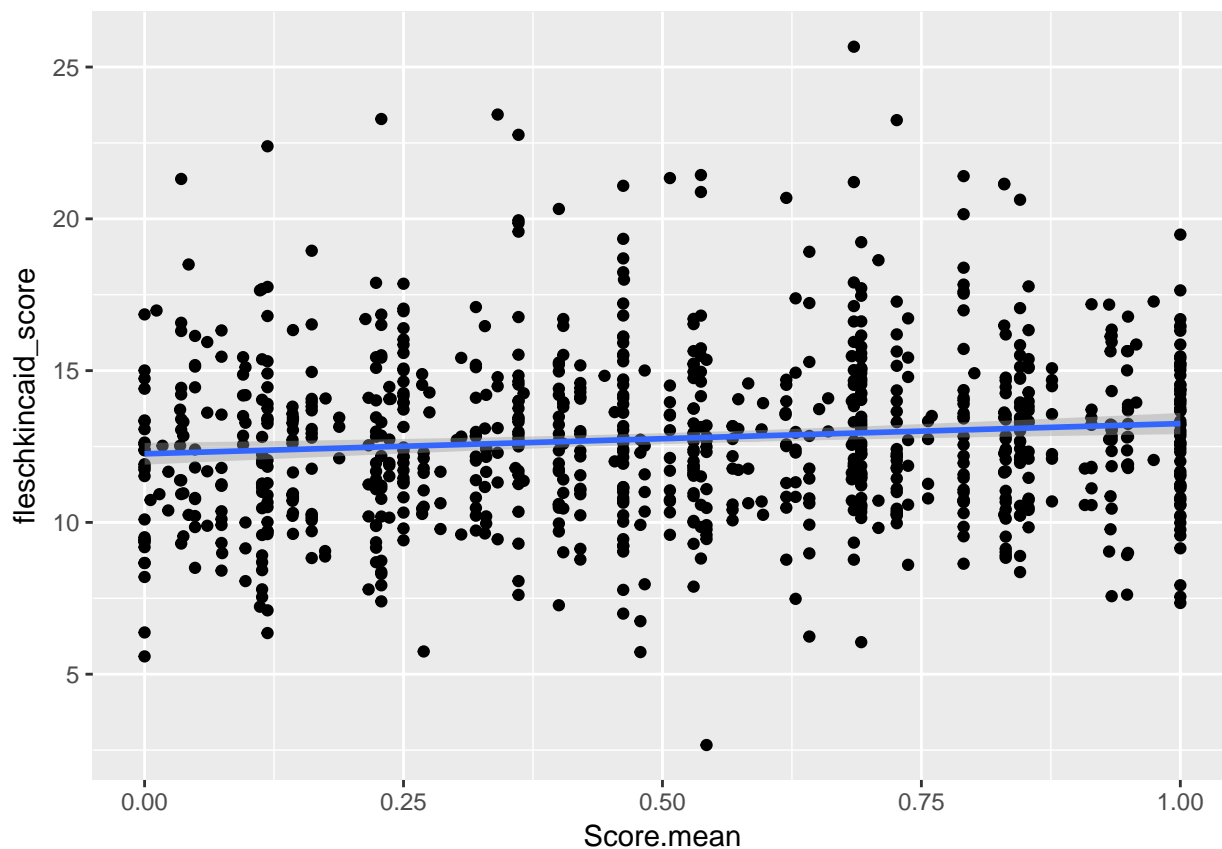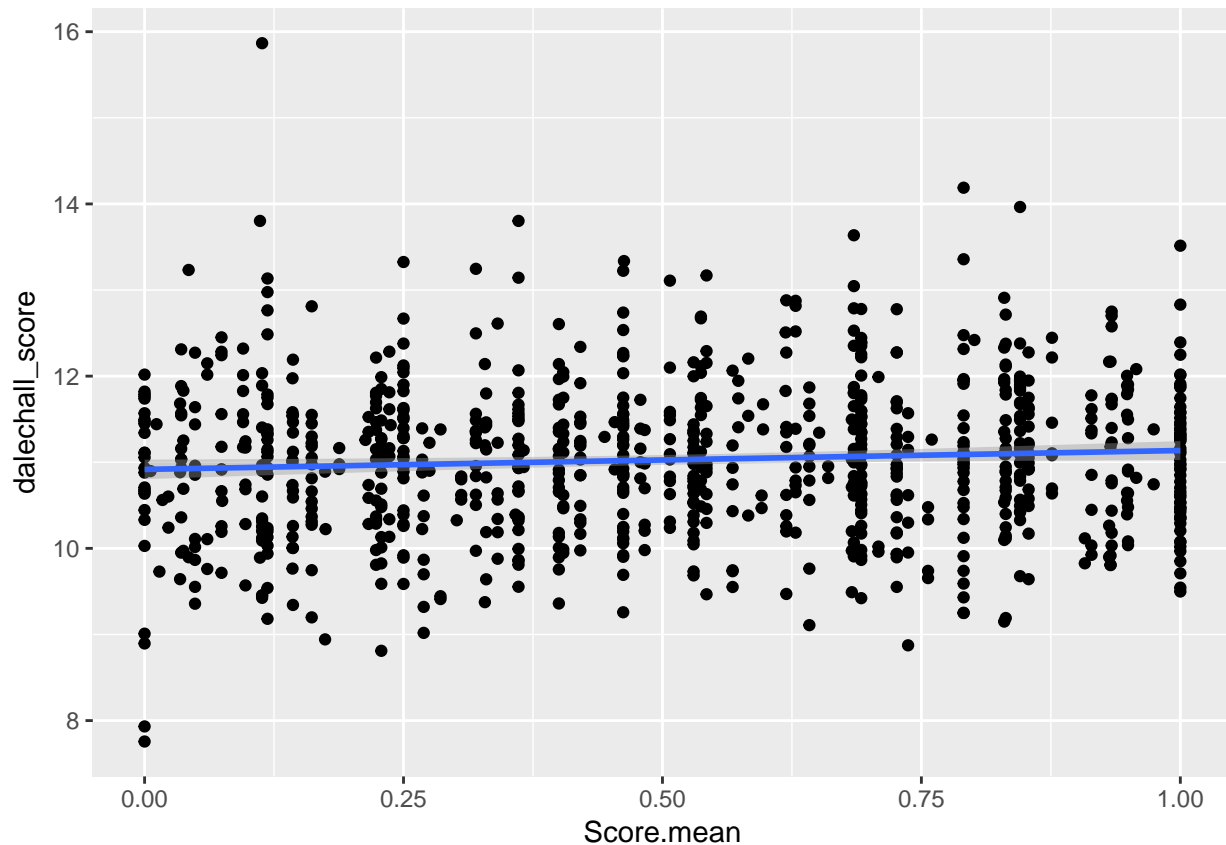
```
##        cor
## 0.1070164
```

```
cor.test(readScores$Score.mean, readScores$dalechall_score)
```

```
##
##  Pearson's product-moment correlation
##
## data:  readScores$Score.mean and readScores$dalechall_score
## t = 2.2498, df = 901, p-value = 0.02471
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.009547784 0.139300674
## sample estimates:
##        cor
## 0.07474057
```

```
ggplot(readScores,
       aes(y=fleschkincaid_score,
           x=Score.mean)) +
  geom_point() +
  stat_smooth(method = 'lm')
```



```
ggplot(readScores,
       aes(y=dalechall_score,
           x=Score.mean)) +
  geom_point() +
  stat_smooth(method = 'lm')
```

Are there interactions between reading scores and gender?

```
m0 = lmer(Score.mean.norm~ 1 +
            format + student + gender +
          (1 | conference),
        data = readScores,
        control = lmerControl(optimizer = 'Nelder_Mead'),
        REML = F)
m1 = update(m0,~.+fleschkincaid_score_scaled)
m2 = update(m1,~.+fleschkincaid_score_scaled:gender)
anova(m0,m1,m2)
```

```
## Data: readScores
## Models:
## object: Score.mean.norm ~ 1 + format + student + gender + (1 | conference)
## ..1: Score.mean.norm ~ format + student + gender + (1 | conference) +
## ..1:     fleschkincaid_score_scaled
## ..2: Score.mean.norm ~ format + student + gender + (1 | conference) +
## ..2:     fleschkincaid_score_scaled + gender:fleschkincaid_score_scaled
##        Df    AIC    BIC  logLik deviance  Chisq Chi Df Pr(>Chisq)
## object  6 2126.3 2154.1 -1057.2   2114.3
## ..1     7 2126.4 2158.8 -1056.2   2112.4 1.8815      1     0.1702
## ..2     8 2128.3 2165.3 -1056.2   2112.3 0.1260      1     0.7226
```

```
summary(m2)
```

```
## Linear mixed model fit by maximum likelihood  ['lmerMod']
## Formula: Score.mean.norm ~ format + student + gender + (1 | conference) +
```

```
##      fleschkincaid_score_scaled + gender:fleschkincaid_score_scaled
##    Data: readScores
## Control: lmerControl(optimizer = "Nelder_Mead")
##
##      AIC      BIC   logLik deviance df.resid
##   2128.3   2165.3  -1056.2   2112.3      745
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -1.93160 -0.90703 -0.01343  0.89065  1.93833
##
## Random effects:
##  Groups      Name         Variance Std.Dev.
##  conference (Intercept) 0.0000   0.0000
##  Residual               0.9678   0.9838
## Number of obs: 753, groups:  conference, 4
##
## Fixed effects:
##                                  Estimate Std. Error t value
## (Intercept)                      -0.03209    0.04536  -0.707
## format1                           0.24489    0.08256   2.966
## student1                         -0.01764    0.07803  -0.226
## gender1                           0.12210    0.07553   1.617
## fleschkincaid_score_scaled        0.05795    0.04112   1.410
## gender1:fleschkincaid_score_scaled  0.02890  0.08139   0.355
##
## Correlation of Fixed Effects:
##            (Intr) formt1 stdnt1 gendr1 flsc__
## format1    -0.460
## student1   -0.362  0.091
## gender1     0.274 -0.106  0.023
## flschkncd__ 0.001 -0.155  0.003 -0.065
## gndr1:fls__ -0.103  0.039  0.054 -0.082  0.354
```

Dale-Chall scores:

```
m0 = lmer(Score.mean.norm~ 1 +
          format + student + gender +
          (1 | conference),
      data = readScoresDC,
      REML = F)
m1 = update(m0,~.+dalechall_score_scaled)
m2 = update(m1,~.+dalechall_score_scaled:gender)
anova(m0,m1,m2)
```

```
## Data: readScoresDC
## Models:
## object: Score.mean.norm ~ 1 + format + student + gender + (1 | conference)
## ..1: Score.mean.norm ~ format + student + gender + (1 | conference) +
## ..1:     dalechall_score_scaled
## ..2: Score.mean.norm ~ format + student + gender + (1 | conference) +
## ..2:     dalechall_score_scaled + gender:dalechall_score_scaled
##        Df    AIC    BIC  logLik deviance Chisq Chi Df Pr(>Chisq)
## object  6 2047.1 2074.6 -1017.5   2035.1
## ..1     7 2048.4 2080.5 -1017.2   2034.5  0.65      1     0.4201
```

```
## ..2      8 2050.4 2087.0 -1017.2   2034.4  0.09       1      0.7642
```

**summary**(m2)

```
## Linear mixed model fit by maximum likelihood  ['lmerMod']
## Formula: Score.mean.norm ~ format + student + gender + (1 | conference) +
##     dalechall_score_scaled + gender:dalechall_score_scaled
##    Data: readScoresDC
##
##      AIC      BIC   logLik deviance df.resid
##   2050.4   2087.0  -1017.2   2034.4      716
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -1.93536 -0.89658 -0.00423  0.87978  1.92779
##
## Random effects:
##  Groups      Name        Variance Std.Dev.
##  conference (Intercept) 0.0000   0.0000
##  Residual               0.9724   0.9861
## Number of obs: 724, groups:  conference, 4
##
## Fixed effects:
##                               Estimate Std. Error t value
## (Intercept)                   -0.01180    0.04622  -0.255
## format1                        0.24773    0.08525   2.906
## student1                      -0.04666    0.07963  -0.586
## gender1                        0.14041    0.07712   1.821
## dalechall_score_scaled         0.03876    0.04565   0.849
## gender1:dalechall_score_scaled 0.02664    0.08880   0.300
##
## Correlation of Fixed Effects:
##             (Intr) formt1 stdnt1 gendr1 dlch__
## format1     -0.457
## student1    -0.361  0.085
## gender1      0.284 -0.115  0.009
## dlchll_scr_  0.069 -0.235  0.010 -0.006
## gndr1:dlc__ -0.049  0.046  0.003 -0.040  0.192
```

No interactions.