

Statistics project to use R evaluation tools

Predictors in Tipping Amount

December 14

Contents

Introduction	1
Methods	3
Results	18
Discussion	18
Appendix	19

[Back to Home Page](#)

Introduction

This is a project I worked on with 2 other students, so I'm not claiming that all the code is mine, although I certainly contributed my share. I've removed my colleagues full names, since I didn't ask for their permission to share this.

The goal of this assignment was just to demonstrate proficiency in R's regression analysis tools, such as whether the data conforms to our assumptions, and the ability to experiment with any combination of features. Therefore, we weren't looking to find anything uninituitive, since the focus was to be on the use of the tools.

Team Members

Jerry Finn
Greg ...
Rich ...

Personal Interest

We've all worked in the service sector in our youth or have kids that have, so we wanted to know what really determines service-based compensation. The nature of the data set also allows us to explore some interesting statistical regression, such as multiple variables: continuous, discrete, and factor variables and their interaction.

How We Found The Data Set

From the RStudio Console, we invoked "`??datasets`" to see all 90+ datasets available in R. The `tips` data set was the one we found most interesting for the reasons given above.

- `reshape2::tips` Tipping data

Data Set and Package Overview

We have chosen a straight forward dataset of a waiter that recorded his tips over the period of several months to look at what predictors determine the overall tip amount. The data comes with the **reshape2** R package that is used to manipulate and clean data, although that is not our primary goal here. As stated previously, it was collected by a waiter tracking his own job performance.

Detailed Description of Data and Cited Reference

As we can note when clicking on the link above, there are a number of interesting variables to explore here, such as the numeric variables of total bill amount and size of party. We also have a fair number of factor variables that we can explore in terms of additive or interaction properties, such as Smoker or Non-Smoker, Male or Female, Meal Type (Dinner or Lunch) and Day of the Week. Intuitively, we suspect that **total_bill** is the most significant variable. We would also suspect that **size** of the party seems highly correlated to **total_bill**.

Quick Look at the Observation Data

```
## 'data.frame':   244 obs. of  7 variables:
## $ total_bill: num  17 10.3 21 23.7 24.6 ...
## $ tip       : num  1.01 1.66 3.5 3.31 3.61 4.71 2 3.12 1.96 3.23 ...
## $ sex       : Factor w/ 2 levels "Female","Male": 1 2 2 2 1 2 2 2 2 2 ...
## $ smoker    : Factor w/ 2 levels "No","Yes": 1 1 1 1 1 1 1 1 1 1 ...
## $ day       : Factor w/ 4 levels "Fri","Sat","Sun",...: 3 3 3 3 3 3 3 3 3 3 ...
## $ time      : Factor w/ 2 levels "Dinner","Lunch": 1 1 1 1 1 1 1 1 1 1 ...
## $ size      : int   2 3 3 2 4 4 2 4 2 2 ...
```

We were instructed to use any dataset of our choice, as long as it contained a minimum 200 observations (ours has 244), a numeric response variable (**tip**), at least one categorical predictor (**sex**, **smoker**, **day**, **time**), and at least two numeric predictors (**total_bill** and **size**).

Goal of This Analysis and Proposed Model

We are creating a model for this tipping data to determine whether **total_bill** is the only predictor relevant for determining the response, **tip**. Our assumption, from life experience, is that it is. The alternative theory is to assume there are other predictors that are significant to determining **tip**, other than **total_bill**.

We want to answer the following questions about our data and proposed model:

- What are the key predictors to receiving the highest tip? Does it matter if the person paying the bill is a smoker or not? Does gender matter? Do people tip better over a meal at dinner time rather than lunch? Are there days in the week that result in the highest tips? We have preconceived ideas on the answers to these questions but want to use the statistical methods learned in class to make informed decisions.
- Per the questions given above, we will take a closer look at the factor variables in the data set. For example, we will find out if the mean **tip** of a group of smokers differs than the mean **tip** of female bill payers? We will use box plots, interaction plots, and perform TukeyHSD tests to confirm whether there is any interaction between factors.
- The primary goal of our analysis is to recommend a model that is best for Prediction, and then a model best for Explanation. We will comment on whether these proposed models violate assumptions of linearity.

- Finally, we will use our proposed model to predict tip amounts for 3 types of bills: 1) a very small bill amount (near zero), 2) for a bill near the mean, and 3) for a bill at the very high end of our observed values.

General Approach

- Team collaboration via Zoom video conferencing and email.
- Use GitHub as our web-based Git repository hosting service for all source code.
- Use Git as our distributed version control and source code management (SCM) tool.

We utilized the following data analysis and statistical methods, learned in STAT 420, to complete our analysis:

- Read about the source of our data, its attributes, and the original study.
- Inspect our data (244 observations)
- Fix any N/A or missing fields and look for factors that need to be coerced.
- Look for response or predictor variables that may be always positive and orders of magnitude greater than the other variables (indicating possible log transformation required).
- Report on description statistics: e.g., mean and variance of response variable.
- Before trying to fit a good model, perform significance of regression to confirm that one or more predictors has a linear relationship with the response variable.
- Use a fitted versus residual plot and a Q-Q plot to visually check for violations of linearity.
- Perform Shapiro-Wilk test to validate normality.
- Perform Breusch-Pagan Test to validate constant variance.
- Confirm if transformation of response variables or predictors is necessary (log or polynomial).
- Determine if simple linear or multi-linear regression.
- For an MLR model, determine if an additive model only, or will there also be interaction terms, polynomial terms.
- Comment on dummy variables, factor variables, binary variables. Test for interaction.
- Perform 2-way ANOVA to test the main effects of more than one factor variable.
- Use interaction plots to show possible interaction between factors.
- Use TukeyHSD (Honest Significance Difference) function to determine significance of interaction between factors.
- Test data observations for outliers, high influence (Cook's Distance), and high leverage.
- Calc VIFs and partial correlation to check for collinearity of predictors.
- Use ANOVA (Analysis of Variance) and inference to compare models. Test for significance.
- Find a good model using step() function using forward, backward, and both to calculate lowest AIC and BIC
- Determine an exhaustive model using regsubsets() function in library(leaps).
- For all models built, calculate Adjusted R^2 (larger is better), RMSE (smaller is better), LOOCV_RMSE (smaller is better), AIC (smaller is better), and BIC (smaller is better) when comparing different models.
- Summarize all model results into a table for easy comparison.
- Randomly split observation data into Train data (66%) and Test data (34%) to validate the chosen model. Select model with lowest Test RMSE to ensure fewest predictors (avoid overfitting).
- If applicable, comment on causation vs. correlation.
- Recommend best model for prediction vs. best model for explanation.

Methods

Here we'll take a quick look at our data.

```
summary(tips)
```

```
##      total_bill      tip      sex      smoker      day      time
##  Min.   : 3.07   Min.   : 1.00  Female: 87   No :151   Fri :19   Dinner:176
##  1st Qu.:13.35   1st Qu.: 2.00   Male  :157   Yes: 93   Sat :87   Lunch : 68
##  Median :17.80   Median : 2.90                                     Sun :76
##  Mean   :19.79   Mean    : 3.00                                     Thur:62
##  3rd Qu.:24.13   3rd Qu.: 3.56
##  Max.   :50.81   Max.    :10.00
##      size
##  Min.    :1.00
##  1st Qu.:2.00
##  Median :2.00
##  Mean    :2.57
##  3rd Qu.:3.00
##  Max.    :6.00
```

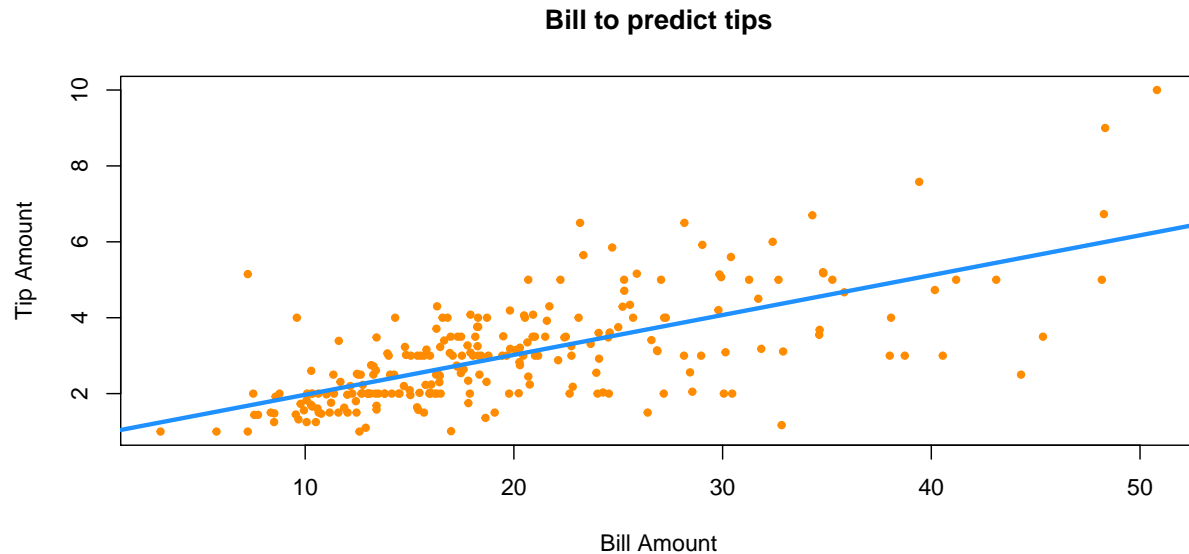
A tip is usually calculated as a percentage of a the total bill. Therefore, we could hypothesize that a simple linear model with one predictor variable, `total_bill`, would be the best model. Let's take look at a summary of the model and at a graph.

```
# Here we build a simple linear model with an obvious variable
tips_slr = lm(tip ~ total_bill, data = tips)
summary(tips_slr)
```

```
##
## Call:
## lm(formula = tip ~ total_bill, data = tips)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.198 -0.565 -0.097  0.486  3.743
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.92027    0.15973     5.76 2.5e-08 ***
## total_bill   0.10502    0.00736    14.26 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.02 on 242 degrees of freedom
## Multiple R-squared:  0.457, Adjusted R-squared:  0.454
## F-statistic: 203 on 1 and 242 DF, p-value: <2e-16
```

```
# The plot shows that the data gets sparcer as the bill grows.
par(mfrow = c(1, 1))
plot(tip ~ total_bill, data = tips,
     xlab = "Bill Amount",
     ylab = "Tip Amount",
     main = "Bill to predict tips",
     pch = 20,
     cex = 1,
```

```
col = "darkorange")
abline(tips_slr, lwd = 3, col = "dodgerblue")
```



With a p-value of 6.6925×10^{-34} this model would be considered linearly significant at any α .

We'll accept this, `tips_slr = lm(tip ~ total_bill, data = tips)`, as our baseline model to beat.

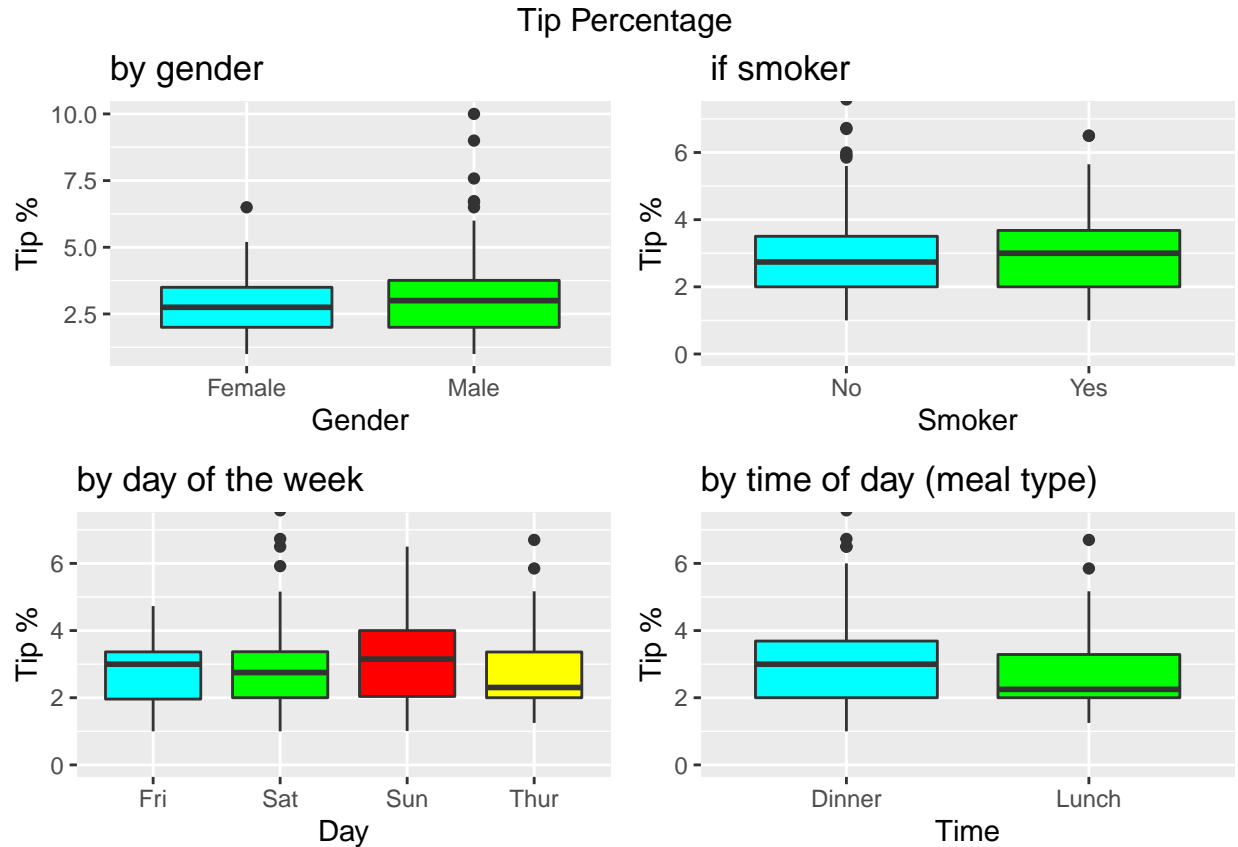
First let's explore the data a little more and check for outliers. For now, we'll just define an outlier as an extreme value. We'll return to the definition given in the text later. Here we are going to check for a value more than 3 standard deviations above the mean. Since the lower bound of a tip is zero, we won't check the lower bound.

```
# Check for tip outliers
tips$tip[ tips$tip > mean(tips$tip) + 3*sd(tips$tip)]
```

```
## [1] 7.58 10.00 9.00
```

Here we see that we have 3 outliers. The reason, for now, that we want to explore this is, to do further exploration we are going to put the tips on a y-axis of some box plots, but we don't want the outliers to dominate the scale and obscure other information.

Recall that factor variables could affect the outcome, if there is a significant difference between the means and variation between the levels within a factor variable, and a box plot is a good way to visualize this.



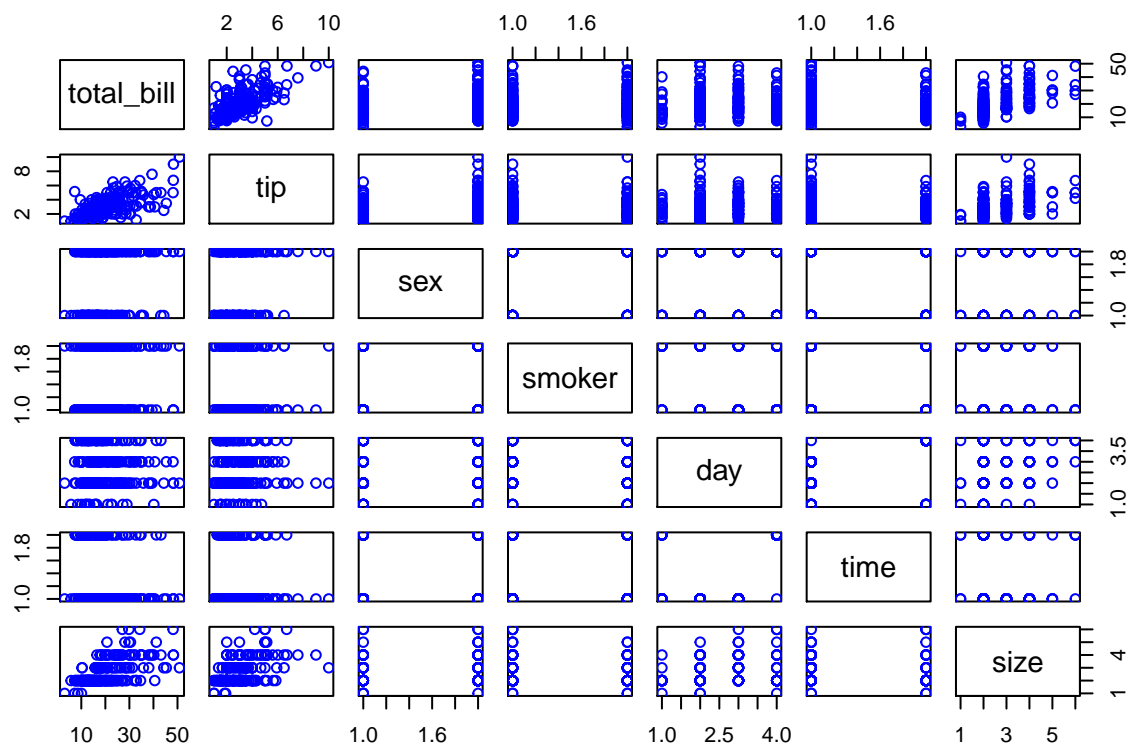
The first box plot has all points on the scale, while we adjusted the scale on the following 3 plots.

In any case, there **does not** seem to be a large variation in mean or variance of the factor variables in the data set. From looking at the interaction plots and TukeyHSD tests results (found in the Appendix), there doesn't seem to be any significant interaction between factors. This is consistent with what we see in the box plots above.

Nonetheless, we'll create full additive and interaction models and see which predictors are eliminated by the techniques taught in this course.

We'll take a quick look at correlation between all variables.

```
# cor(tips)           # Error: must be all numeric
pairs(tips, col="blue")
```



If we look at the correlations. The only obvious linear relation is between tip and bill total. It seems our original intuition should hold.

But we'll build larger models and try to find a good smaller model. See the appendix for the summary output of the models.

```
# An additive model
tip_add = lm(tip ~ ., data = tips)

# A full interaction model.
tip_int = lm(tip ~ total_bill * sex * smoker * day * time * size, data = tips)
```

```
## Warning in v1 * v2: longer object length is not a multiple of shorter object
## length
```

Now we can search for a better model by taking our larger models and eliminating less significant coefficients to reduce them to a manageable size by using AIC and BIC step function.

```
# Backward AIC and BIC models
# additive models
aic_back_add = step(tip_add, direction = "backward", trace = 0)
n_a = length(resid(tip_add))
mod_list = list(aic_back_add)
bic_back_add = step(tip_add, direction = "backward", k = log(n_a), trace = 0)
append_model_list(bic_back_add)
```

```

# interaction models
aic_back_int = step(tip_int, direction = "backward", trace = 0)
n_i = length(resid(tip_int))
append_model_list(aic_back_int)
bic_back_int = step(tip_int, direction = "backward", k = log(n_i), trace = 0)
append_model_list(bic_back_int)

# Forward AIC and BIC models
# additive models
tip_start = lm(tip ~ 1, data = tips)
aic_for_add = step(tip_start, tip ~ total_bill + sex + smoker + day + time + size, direction = "forward")
append_model_list(aic_for_add)
bic_for_add = step(tip_start, tip ~ total_bill + sex + smoker + day + time + size, direction = "forward")
append_model_list(bic_for_add)

# interaction models
aic_for_int = step(tip_start, tip ~ total_bill * sex * smoker * day * time * size, direction = "forward")
append_model_list(aic_for_int)
bic_for_int = step(tip_start, tip ~ total_bill * sex * smoker * day * time * size, direction = "forward")
append_model_list(bic_for_int)

# Both AIC and BIC models
# additive models
tip_start = lm(tip ~ 1, data = tips)
aic_both_add = step(tip_start, tip ~ total_bill + sex + smoker + day + time + size, direction = "both",
append_model_list(aic_both_add)
bic_both_add = step(tip_start, tip ~ total_bill + sex + smoker + day + time + size, direction = "both",
append_model_list(bic_both_add)

# interaction models
aic_both_int = step(tip_start, tip ~ total_bill * sex * smoker * day * time * size, direction = "both",
append_model_list(aic_both_int)
bic_both_int = step(tip_start, tip ~ total_bill * sex * smoker * day * time * size, direction = "both",
append_model_list(bic_both_int)

```

We eliminated duplicate models as we went along invoking the `append_model_list()` function which only appends unique models to the official list. We end up with 4 models, one of them the simple linear model, `tips_slr`, with `total_bill` as the only predictor, that we originally postulated as the best.

Final 4 Models

```

# look at the model coefficients to see how many unique models we have
for(i in 1:length(mod_list)) {
  cat("For model ", i , " coefficients are ", names(coef(mod_list[[i]])), "\n")
}

```

```

## For model 1 coefficients are (Intercept) total_bill size
## For model 2 coefficients are (Intercept) total_bill
## For model 3 coefficients are (Intercept) total_bill sexMale smokerYes daySat daySun dayThur size
## For model 4 coefficients are (Intercept) total_bill smokerYes total_bill:smokerYes

```

Comments on Final 4 Models

It seems we have gotten 3 manageable models out of the four unique models chosen by our exploration. Model 3 is not much more manageable than our largest models and would not be a comprehensible model even with a low Adjusted R^2 value.

The second model is the same as our baseline, `tips_slr`.

Assumptions of linearity

Now let's look at the assumptions of linearity for our models. We will do so visually and via statistical tests. We will create a Fitted versus Residuals plot and a Q-Q plot for each model.

For the Fitted versus Residuals plots, we want to see 1) at any fitted value, the mean of the residuals should be roughly 0, and 2) at every fitted value, the spread of the residuals should be roughly the same. If 1) is true, the **linearity** assumption is valid. If 2) is true, the **constant variance** assumption is valid.

For each Q-Q plot, we desire the points of the plot to closely follow a straight line. If so, this would suggest that the data comes from a **normal distribution**.

From the plots below, we will show the results from running the Shapiro-Wilk and Breusch-Pagan tests.

Testing the assumption of normality

The **Shapiro-Wilk Test** is a test for **normality**. The null hypothesis assumes the data follows a normal distribution. The alternative hypothesis assumes a non-normal distribution. When we see a low p-value for a reasonable *alpha*, it suggests the normality assumption is violated. A large p-value indicates a normal distribution.

Testing the assumption of constant variance

The **Breusch-Pagan Test** is a test for **constant variance**. The null and alternative hypotheses can be considered to be:

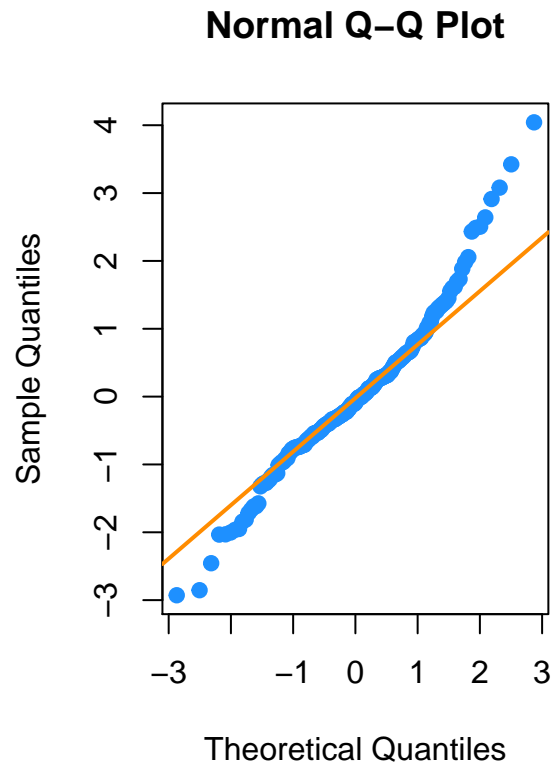
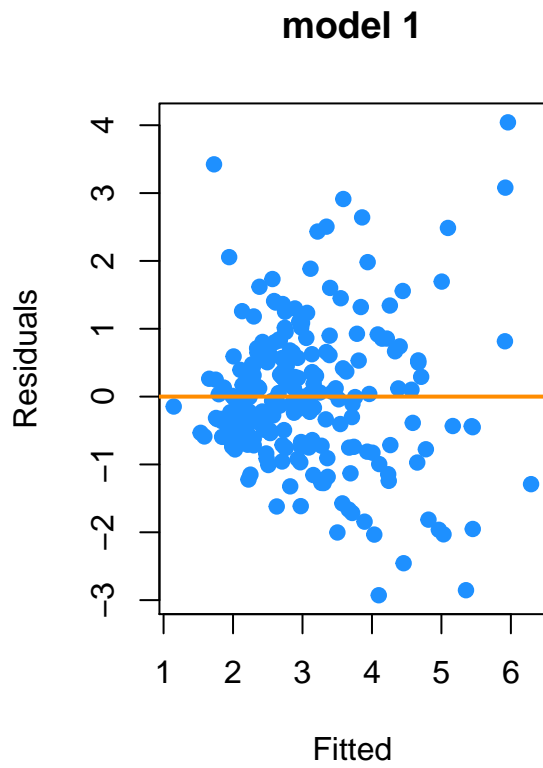
H_0 : Homoscedasticity. The errors have constant variance about the true model.

H_1 : Heteroscedasticity. The errors have non-constant variance about the true model.

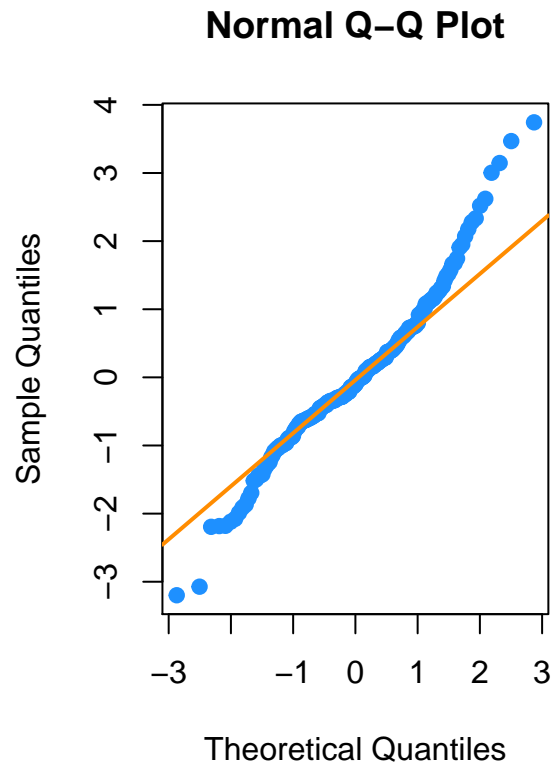
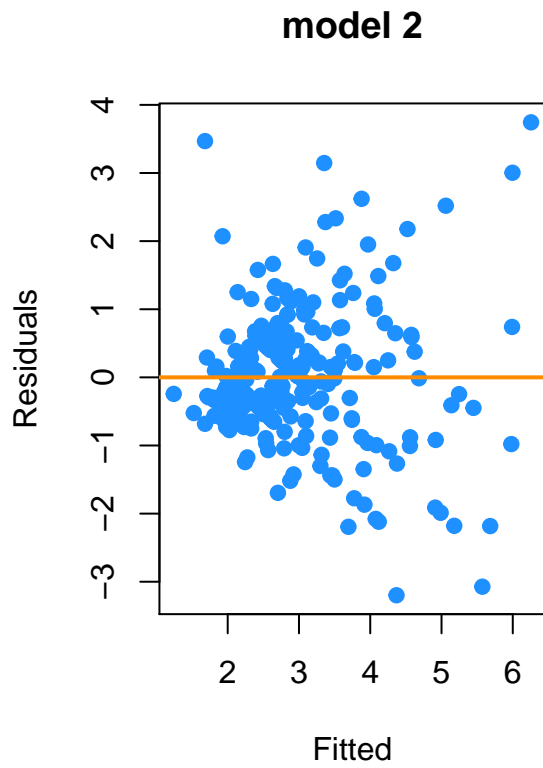
For a large p-value, so we fail to reject the null of homoscedasticity. Therefore, the constant variance assumption would not be violated.

For a small p-value (for a reasonable *alpha*), so we reject the null of homoscedasticity. In this case, the constant variance assumption would be violated.

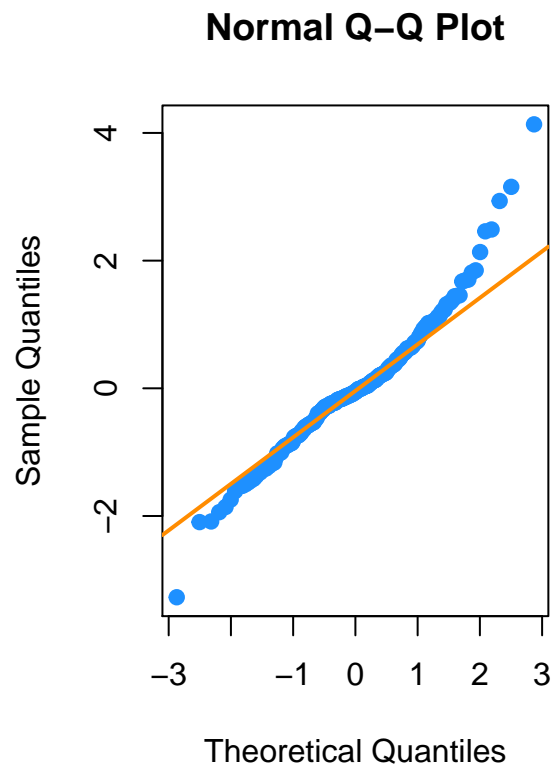
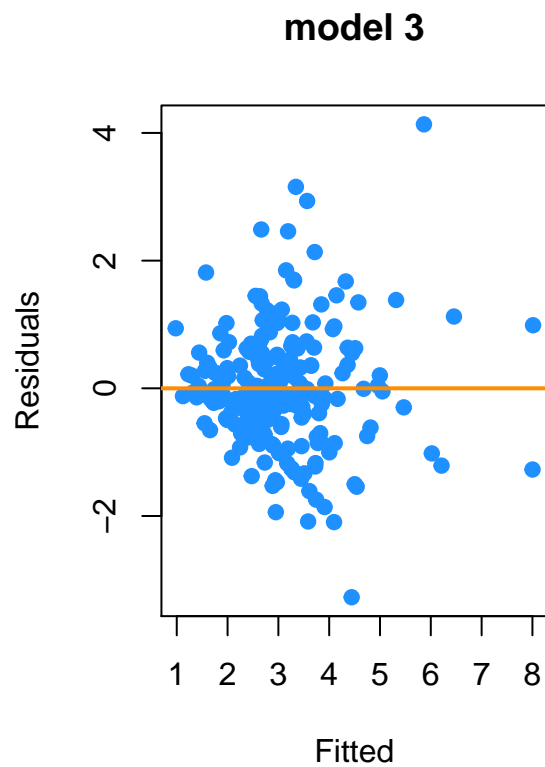
```
library(lmtest)
for(i in 1:length(mod_list)) {
  plot_title = sprintf("model %i", i)
  var_lin_plots(mod_list[[i]], plot_title)
  # BP test of constant variance. $H_0$ is constant variance
  # Shapiro test of normal dist. $H_0$ is normal dist
  cat("The Breusch-Pagan p-value is ", bptest(mod_list[[i]])$p.value, " with ", bptest(mod_list[[i]])$p)
}
```



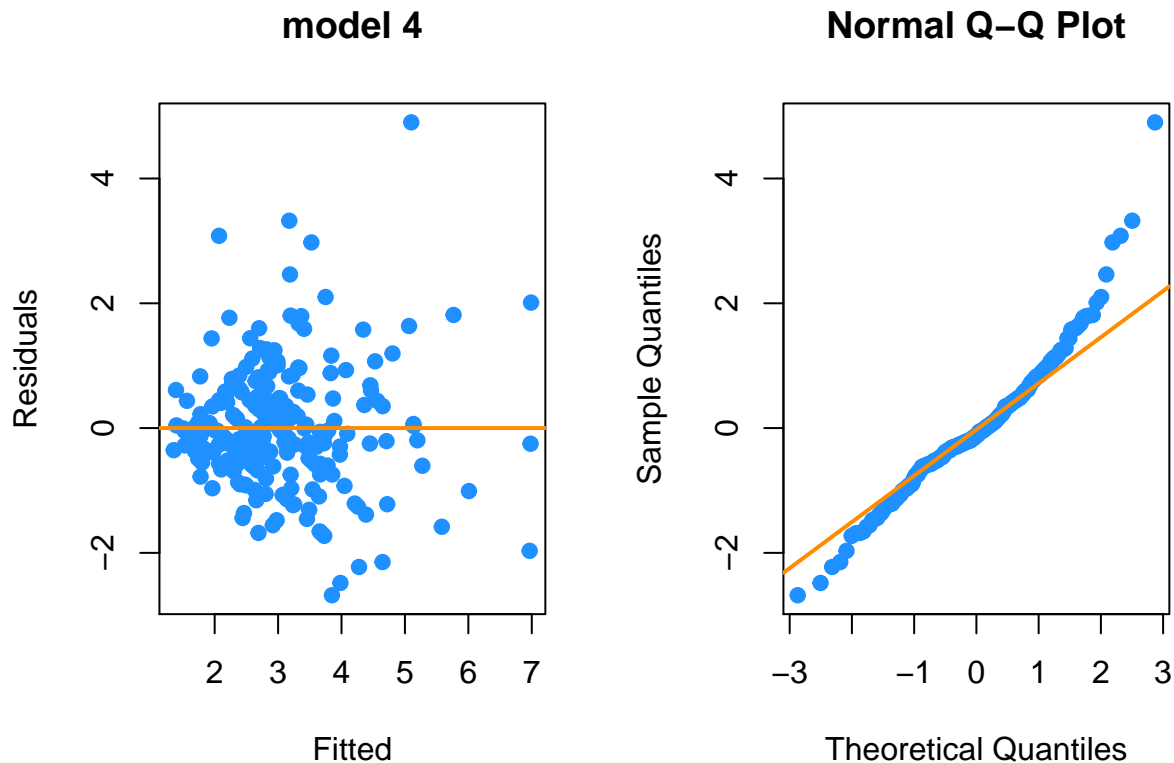
```
## The Breusch-Pagan p-value is 3.499e-11 with 2 degrees of freedom
## The Shapiro-Wilk test has a p-value of 0.00001777
```



```
## The Breusch-Pagan p-value is 4.539e-12 with 1 degrees of freedom
## The Shapiro-Wilk test has a p-value of 0.00002171
```



```
## The Breusch-Pagan p-value is 1.115e-06 with 21 degrees of freedom
## The Shapiro-Wilk test has a p-value of 6.431e-06
```



```
## The Breusch-Pagan p-value is 7.848e-09 with 3 degrees of freedom
## The Shapiro-Wilk test has a p-value of 7.591e-07
```

None of our models have data that has perfectly constant variance nor a perfect normal distribution, as shown by the charts and tests. Variance seems to get wider to the right of graphs, and linearity does not hold on the edges of our Q-Q graph.

We'll continue, but since some assumptions are weak, we have to be careful interpreting our results, especially if we try to extrapolate and do predictions at or beyond the extreme points of the data.

Let's further explore the nature of our data and look for high influence points by calculating Cook's Distance.

```
for (i in 1:length(mod_list)) {
  # high leverage points
  lev = which.max(hatvalues(mod_list[[i]]))
  # outliers i.e. point with high residuals
  temp = abs(rstandard(mod_list[[i]])) > 2
  outl = as.integer(names(temp[temp == TRUE]))
  # Both high "residual" and "leverage" is "influential"
  cooki = intersect(lev, outl)
  if (length(cooki) > 0) {
    cat("The following data point(s) are influential and therefore their omission could change intercept")
    for (j in 1:length(cooki)) {
      if (cooks.distance(mod_list[[i]])[unname(cooki[j])] > 4 / length(cooks.distance(mod_list[[i]]))) {
        print(tips[cooki[j],])
      }
    }
  }
}
```

```

    }
  }
}

```

```

## The following data point(s) are influential and therefore their omission could change intercept and
##      total_bill tip  sex smoker day   time size
## 171      50.81  10 Male   Yes Sat Dinner    3
## The following data point(s) are influential and therefore their omission could change intercept and
##      total_bill tip  sex smoker day   time size
## 171      50.81  10 Male   Yes Sat Dinner    3
## The following data point(s) are influential and therefore their omission could change intercept and
##      total_bill tip  sex smoker day   time size
## 171      50.81  10 Male   Yes Sat Dinner    3

```

Now let's see how well our models do in avoiding collinearity.

```

for(i in 1:length(mod_list)) {
  print(vif(mod_list[[i]]))
}

```

```

## total_bill      size
##      1.558      1.558
## total_bill
##      1
##
##      total_bill      sexMale
##      124.51      99.31
##      smokerYes      daySat
##      86.97      81.71
##      daySun      dayThur
##      91.39      70.46
##      size      total_bill:sexMale
##      13.80      251.39
##      total_bill:smokerYes      total_bill:daySat
##      93.90      187.90
##      total_bill:daySun      total_bill:dayThur
##      200.11      116.20
##      sexMale:daySat      sexMale:daySun
##      94.21      102.81
##      sexMale:dayThur      total_bill:size
##      58.34      42.25
##      smokerYes:size      total_bill:sexMale:daySat
##      101.41      183.61
##      total_bill:sexMale:daySun      total_bill:sexMale:dayThur
##      204.30      88.37
##      total_bill:smokerYes:size
##      108.02
##      total_bill      smokerYes      total_bill:smokerYes
##      1.884      5.861      7.083

```

Our 3rd and 4th models have collinearity problems, and 3 is already an unwieldy model and unnecessarily large.

According to the VIF test, in model 1, the parameters total bill and size don't have a collinearity problem. Would it be reasonable to assume that the size of the bill and the size of the party eating are not related? We don't think so.

Let's check if any of the parameters beyond our baseline are significant.

```
(anova24 = anova(mod_list[[2]], mod_list[[4]]))

## Analysis of Variance Table
##
## Model 1: tip ~ total_bill
## Model 2: tip ~ total_bill + smoker + total_bill:smoker
##   Res.Df RSS Df Sum of Sq  F    Pr(>F)
## 1      242 253
## 2      240 230  2          23 12 0.000011 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(anova21 = anova(mod_list[[2]], mod_list[[1]]))

## Analysis of Variance Table
##
## Model 1: tip ~ total_bill
## Model 2: tip ~ total_bill + size
##   Res.Df RSS Df Sum of Sq  F Pr(>F)
## 1      242 253
## 2      241 248  1      5.23 5.1 0.025 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

With a p-value of 1.0795×10^{-5} we would reject the null at any reasonable α that addition β s are 0 in our model 4. But our larger model here has some collinearity problems and is not a good model by that standard.

Comparing our model 1 and 2 (our baseline) we have a p-value of 0.0249 and we would reject the null hypothesis at an α of 0.05, but not 0.01.

Now let's look at our models' fit.

```
# initiaize vectors
loocv_rmse = rep(0, length(mod_list))
adr2_v = rep(0, length(mod_list))

# What is our LOOCV_RMSE? Smaller is better.
for(i in 1:length(mod_list)) {
  loocv_rmse[i] = calc_loocv_rmse(mod_list[[i]])
}

# What is the Adjusted R^2? Larger is better.
for(i in 1:length(mod_list)) {
  adr2_v[i] = summary(mod_list[[i]])$adj.r.squared
}
```

	Model_description	LOOCV_RMSE	Adjusted R ²
model 1	tip ~ total_bill + size	1.030	0.4635
model 2	tip ~ total_bill	1.033	0.4544
model 3	tip ~ total_bill + I(sex == Male) + I(smoker == Yes) + I(day == Sat) + I(day == Sun) + I(day == Thur) + ...	1.044	0.5230
model 4	tip ~ total_bill + I(smoker == Yes)	1.003	0.4998

Determining Best Model for Explanation

Model 1 seems to be the best at explanation, if we reject models 3 and 4 for other reasons.

We see that model 3 has the highest Adjusted R^2 and model 4 is next. This would usually indicate that these are the best to explain the model, but we have noted several other problems with these models.

Determining Best Model for Prediction

Model 4 with the lowest LOOCV_RMSE would seem to be the best for prediction, but we will put this to the test.

We'll split the data to train and test and see if the predictions are good or not.

```
# randomly split observation data into train (66%) and test data (34%)
```

```
train_index = sample(1:nrow(tips), .66*nrow(tips))
tip_train_data = tips[train_index,]
tip_test_data = tips[-train_index,]
```

```
# re-create the 4 final models from analysis above
```

```
model_1_train = lm(tip ~ total_bill + size, data = tip_train_data)
model_2_train = lm(tip ~ total_bill, data = tip_train_data)
model_3_train = lm(tip ~ total_bill + I(sex == "Male") + I(smoker == "Yes") + I(day == "Sat") + I(day == "Sun") + I(day == "Thur"), data = tip_train_data)
model_4_train = lm(tip ~ total_bill + I(smoker == "Yes"), data = tip_train_data)
```

```
# forget about the unwieldy model!
```

```
# = lm(tip ~ total_bill + sex + smoker + day + size + total_bill*sex + total_bill*smoker + total_bill*day + total_bill*size, data = tips)
```

```
#tip_add_train = lm(tip ~ ., data = tip_train_data)
```

```
mod_list_test_train = list(model_1_train)
RMSE_1_train = sqrt(mean(resid(model_1_train) ^ 2))
RMSE_2_train = sqrt(mean(resid(model_2_train) ^ 2))
#append_model_list_test_train(tip_int_train)
RMSE_3_train = sqrt(mean(resid(model_3_train) ^ 2))
RMSE_4_train = sqrt(mean(resid(model_4_train) ^ 2))
```

```
test_1 = predict(model_1_train, newdata = tip_test_data)
RMSE_1_test = sqrt(mean(test_1-tip_test_data$tip) ^ 2)
```

```
test_2 = predict(model_2_train, newdata = tip_test_data)
RMSE_2_test = sqrt(mean(test_2-tip_test_data$tip) ^ 2)
```

```
test_3 = predict(model_3_train, newdata = tip_test_data)
```



```

RMSE_3_test = sqrt(mean(test_3-tip_test_data$tip) ^ 2)

test_4 = predict(model_4_train, newdata = tip_test_data)
RMSE_4_test = sqrt(mean(test_4-tip_test_data$tip) ^ 2)

```

	Model_description	RMSE_train	RMSE_test
model 1	tip ~ total_bill + size	0.9880	0.2841
model 2	tip ~ total_bill	1.0030	0.2941
model 3	tip ~ total_bill + I(sex == Male) + I(smoker == Yes) + I(day == Sat) + I(day == Sun) + I(day == Thur) + ...	0.8786	0.3622
model 4	tip ~ total_bill + I(smoker == Yes)	0.9917	0.3330

Using the 4 trained models, models 2 and 4 have roughly the lowest RMSEs for the trained and test data indicating they probably do about the same in predicting model outcomes. Model 3 has a myriad of variables and is probably overfit, and adding `size` as a predictor to the `total_bill` doesn't seem to help matters.

```

# now look at the model coefficients to how many unique models we have
for(i in 1:length(mod_list_test_train)) {
  cat("For model ", i , " coefficients are ", names(coef(mod_list_test_train[[i]])), "\n")
}

# What is our LOOCV_RMSE? Smaller LOOCV_RMSE better
for(i in 1:length(mod_list_test_train)) {
  print(calc_loocv_rmse(mod_list_test_train[[i]]))
}

# What is the Adjusted R^2? Larger is better.
for(i in 1:length(mod_list_test_train)) {
  print(summary(mod_list_test_train[[i]])$adj.r.squared)
}

for(i in 1:length(mod_list_test_train)) {
  print(vif(mod_list_test_train[[i]]))
}

```

	Model_description	LOOCV_RMSE	Adjusted R ² _train
model 1	tip ~ total_bill + size	1.018	0.4728
model 2	tip ~ total_bill	1.023	0.4602
model 3	tip ~ total_bill + I(sex == Male) + I(smoker == Yes) + I(day == Sat) + I(day == Sun) + ...	1.085	0.5261
model 4	tip ~ total_bill + I(smoker == Yes)	1.018	0.4690

No surprise that the Adjusted R^2 value for model 3, given all the variables used to fit the model. However,

lowest LOOCV_RMSEs appear for the first two models indicating **size** (of party) and its closely related predictor **total_bill**.

Results

Here are the final 4 models that we built and are under consideration:

model 1: `lm(tip ~ total_bill + size, data = tips)`

model 2: `lm(tip ~ total_bill, data = tips)`

model 3: `lm(tip ~ total_bill + sex + smoker + day + size + total_bill:sex + total_bill:smoker + total_bill:day + sex:day + total_bill:size + smoker:size + total_bill:sex:day + total_bill:smoker:size, data = tips)`

model 4: `lm(tip ~ total_bill + smoker + total_bill:smoker, data = tips)`

For all models built, we assumed the following: Adjusted R^2 (larger is better), RMSE (smaller is better), LOOCV_RMSE (smaller is better), AIC (smaller is better), and BIC (smaller is better) when comparing models.

Numerical Summary of Model Results

	LOOCV_RMSE	LOOCV_RMSE_train	Adj_R2	Adj_R2_train	RMSE_train	RMSE_test
model 1	1.030	1.018	0.4635	0.4728	0.9880	0.2841
model 2	1.033	1.023	0.4544	0.4602	1.0030	0.2941
model 3	1.044	1.085	0.5230	0.5261	0.8786	0.3622
model 4	1.003	1.018	0.4998	0.4690	0.9917	0.3330

Recommended Model for Prediction

model 2,

Recommended Model for Explanation

model 1, though we would suspect collinearity with size and total bill, the VIF test does not indicate this.

Rejected Models

model 3 and **model 4**

Discussion

At the onset, we suspected that **total_bill** would be the most significant variable in our data set. We also surmised that **size** of the party would be highly correlated to **total_bill**. After a thorough analysis and building 4 unique models, we conclude that a linear model with predictors variable, **total_bill** and **size**, would be the best model for explanation according to our tests. Furthermore, we didn't find any significant interaction among our 4 factor variables: **time** (Lunch or Dinner), **day** (Thur, Fri, Sat, or Sun), **sex** (Male or Female), and **smoker** (Yes or No).

Our final model is useful because it is easy to understand and confirms most people's intuition on what determines tip amount. It also confirms that there are no "hidden" or unusually complex predictors in the data set. That said, it would have been ideal to have at least one more predictor called **service_rating**. By having diners rate the service experience of their waiter (e.g., Terrible, Poor, Fair, Good, Excellent), would have given us much greater accuracy in predicting the response variable, **tip**.

Finally, we will use our proposed model to predict tip amounts for 3 types of bills: 1) a very small bill amount (near zero), 2) for a bill near the mean, and 3) for a bill at the very high end of our observed values.

```
# find lowest bill amount in data set
min(tips$total_bill)
```

```
## [1] 3.07
```

```
low_bill = 2
```

```
# find average bill amount
mean(tips$total_bill)
```

```
## [1] 19.79
```

```
avg_bill = 20
```

```
# find largest bill amount
max(tips$total_bill)
```

```
## [1] 50.81
```

```
high_bill = 50
```

```
# create data fram for these 3 sample bill amounts
new_tips = data.frame(total_bill = c(low_bill, avg_bill, high_bill))
```

```
# feed the 3 sammples bills into our prediction model
model1 = lm(tip ~ total_bill + size, data = tips)
model2 = lm(tip ~ total_bill, data = tips)
predict(model2, newdata = new_tips, interval = "prediction",
        level = 0.99)
```

```
##      fit      lwr    upr
## 1 1.130 -1.5503 3.811
## 2 3.021  0.3618 5.680
## 3 6.171  3.4505 8.893
```

Observations on Predicted Tips

We see from the results above that for a small bill (2), the tip would be between x and y dollars.
For an average bill (20), the tip would be between a and b dollars.
And for an unusually large bill (50), the tip would be c and d dollars.

Appendix

The **appendix** section contains code and analysis that is used, but that may clutter the report, or is not directly related to the model choice.

Code for helper functions hidden above

```

# Helper functions

# Function to add new model to list of best models
append_model_list = function(mod_new) {
  for(i in 1:length(mod_list)) {
    if (identical(coef(mod_list[[i]]), coef(mod_new))) {break}
    if (i == length(mod_list)) {
      .GlobalEnv$mod_list[[length(.GlobalEnv$mod_list)+1]] <- mod_new
    }
  }
}

# Function to add new model to list of best train and best test models
set.seed(42)
append_model_list_test_train = function(mod_new) {
  for(i in 1:length(mod_list_test_train)) {
    if (identical(coef(mod_list_test_train[[i]]), coef(mod_new))) {break}
    if (i == length(mod_list_test_train)) {
      .GlobalEnv$mod_list_test_train[[length(.GlobalEnv$mod_list_test_train)+1]] <- mod_new
    }
  }
}

# Function to run Leave-One-Out Cross-Validation RMSE from text
calc_loocv_rmse = function(model) {
  sqrt(mean((resid(model) / (1 - hatvalues(model))) ^ 2))
}

# Helper functions from text
plot_fitted_resid = function(model, model_name, pointcol = "dodgerblue", linecol = "darkorange") {
  plot(fitted(model), resid(model),
       col = pointcol, pch = 20, cex = 1.5,
       main=model_name, xlab = "Fitted", ylab = "Residuals")
  abline(h = 0, col = linecol, lwd = 2)
}

plot_qq = function(model, pointcol = "dodgerblue", linecol = "darkorange") {
  qqnorm(resid(model), col = pointcol, pch = 20, cex = 1.5)
  qqline(resid(model), col = linecol, lwd = 2)
}

var_lin_plots <- function(model, model_name){
  par(mfrow = c(1, 2))
  plot_fitted_resid(model, model_name)
  plot_qq(model)
}

```

Code to produce Box Plots

```
library(ggplot2)
```

```

# Here we do 1 box plot with outliers and then 3 without to show how
# the scale with outliers make the plot harder to see.
g <- ggplot(tips, aes(sex, tip)) +
  geom_boxplot(aes(fill=sex))
# guide=FALSE to suppress legend
g <- g + scale_fill_manual(values=c("cyan","green"), guide=FALSE)
g1 <- g + labs(title="by gender", x="Gender", y="Tip %")

g <- ggplot(tips, aes(smoker, tip)) +
  geom_boxplot(aes(fill=smoker))
# guide=FALSE to suppress legend
g <- g + scale_fill_manual(values=c("cyan","green"), guide=FALSE)
g <- g + coord_cartesian(ylim = c(0,mean(tips$tip) + 3*sd(tips$tip)))
g2 <- g + labs(title="if smoker", x="Smoker", y="Tip %")

g <- ggplot(tips, aes(day, tip)) +
  geom_boxplot(aes(fill=day))
# guide=FALSE to suppress legend
g <- g + scale_fill_manual(values=c("cyan","green", "red", "yellow"), guide=FALSE)
g <- g + labs(title="by day of week", x="Day", y="Tip %")
g3 <- g + coord_cartesian(ylim = c(0,mean(tips$tip) + 3*sd(tips$tip)))

g <- ggplot(tips, aes(time, tip)) +
  geom_boxplot(aes(fill=time))
# guide=FALSE to suppress legend
g <- g + scale_fill_manual(values=c("cyan","green", "red"), guide=FALSE)
g <- g + coord_cartesian(ylim = c(0,mean(tips$tip) + 3*sd(tips$tip)))
g4 <- g + labs(title="by time of day (meal type)", x="Time", y="Tip %")

library(gridExtra)
grid.arrange(g1, g2, g3, g4, nrow=2, ncol = 2, top="Tip Percentage")

```

Summary of full additive and interactive models run through the AIC and BIC process.

```

summary(tip_add)

##
## Call:
## lm(formula = tip ~ ., data = tips)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.848 -0.573 -0.103  0.476  4.108
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.8038     0.3527   2.28   0.024 *
## total_bill    0.0945     0.0096   9.84 <2e-16 ***
## sexMale      -0.0324     0.1416  -0.23   0.819
## smokerYes    -0.0864     0.1466  -0.59   0.556
## daySat       -0.1215     0.3097  -0.39   0.695

```

```
## daySun      -0.0255      0.3213     -0.08      0.937
## dayThur     -0.1623      0.3934     -0.41      0.680
## timeLunch    0.0681      0.4446      0.15      0.878
## size         0.1760      0.0895      1.97      0.051 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.02 on 235 degrees of freedom
## Multiple R-squared:  0.47,    Adjusted R-squared:  0.452
## F-statistic: 26.1 on 8 and 235 DF,  p-value: <2e-16
```

```
summary(tip_int)
```

```
##
## Call:
## lm(formula = tip ~ total_bill * sex * smoker * day * time * size,
##     data = tips)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.230 -0.392  0.000   0.425   3.933
##
## Coefficients: (65 not defined because of singularities)
##                                     Estimate Std. Error t value
## (Intercept)                      -12.7826    25.6129   -0.50
## total_bill                        0.1004     0.2849    0.35
## sexMale                          -2.6962     6.4129   -0.42
## smokerYes                        -6.0623     8.2666   -0.73
## daySat                           13.2863    25.7245    0.52
## daySun                           16.2647    25.8793    0.63
## dayThur                          -8.2200    18.5927   -0.44
## timeLunch                        26.9932    48.5299    0.56
## size                             6.1532    12.3493    0.50
## total_bill:sexMale                0.1062     0.2626    0.40
## total_bill:smokerYes             0.4060     0.3810    1.07
## sexMale:smokerYes                -8.8010    26.1692   -0.34
## total_bill:daySat                -0.0103     0.3173   -0.03
## total_bill:daySun                -0.1405     0.3455   -0.41
## total_bill:dayThur               0.1331     0.3252    0.41
## sexMale:daySat                   2.9020     7.0690    0.41
## sexMale:daySun                   -0.8668     7.4842   -0.12
## sexMale:dayThur                  -5.0996    15.1521   -0.34
## smokerYes:daySat                 3.4705     8.8373    0.39
## smokerYes:daySun                11.2212    23.0859    0.49
## smokerYes:dayThur                2.7630     9.1537    0.30
## total_bill:timeLunch             -0.1781     0.2549   -0.70
## sexMale:timeLunch                5.8528    14.6952    0.40
## smokerYes:timeLunch              -4.9917     9.0688   -0.55
## daySat:timeLunch                 NA          NA      NA
## daySun:timeLunch                 NA          NA      NA
## dayThur:timeLunch               -3.9866     9.0105   -0.44
## total_bill:size                   0.0317     0.0395    0.80
## sexMale:size                     1.0939     1.6705    0.65
## smokerYes:size                   2.9776     2.2563    1.32
```

## daySat:size	-5.8233	12.3937	-0.47
## daySun:size	-6.9876	12.4117	-0.56
## dayThur:size	3.0605	8.6389	0.35
## timeLunch:size	-9.9825	20.0256	-0.50
## total_bill:sexMale:smokerYes	0.2708	1.2824	0.21
## total_bill:sexMale:daySat	-0.0262	0.3149	-0.08
## total_bill:sexMale:daySun	0.0676	0.3317	0.20
## total_bill:sexMale:dayThur	0.3904	1.1431	0.34
## total_bill:smokerYes:daySat	-0.3130	0.4146	-0.75
## total_bill:smokerYes:daySun	-0.7452	1.2226	-0.61
## total_bill:smokerYes:dayThur	0.0701	0.1712	0.41
## sexMale:smokerYes:daySat	10.0205	26.4947	0.38
## sexMale:smokerYes:daySun	1.3277	34.0754	0.04
## sexMale:smokerYes:dayThur	1.1279	2.5332	0.45
## total_bill:sexMale:timeLunch	-0.4296	1.1394	-0.38
## total_bill:smokerYes:timeLunch	NA	NA	NA
## sexMale:smokerYes:timeLunch	NA	NA	NA
## total_bill:daySat:timeLunch	NA	NA	NA
## total_bill:daySun:timeLunch	NA	NA	NA
## total_bill:dayThur:timeLunch	NA	NA	NA
## sexMale:daySat:timeLunch	NA	NA	NA
## sexMale:daySun:timeLunch	NA	NA	NA
## sexMale:dayThur:timeLunch	NA	NA	NA
## smokerYes:daySat:timeLunch	NA	NA	NA
## smokerYes:daySun:timeLunch	NA	NA	NA
## smokerYes:dayThur:timeLunch	NA	NA	NA
## total_bill:sexMale:size	-0.0353	0.0544	-0.65
## total_bill:smokerYes:size	-0.1503	0.0950	-1.58
## sexMale:smokerYes:size	5.0935	12.7065	0.40
## total_bill:daySat:size	-0.0372	0.0679	-0.55
## total_bill:daySun:size	0.0146	0.0681	0.21
## total_bill:dayThur:size	NA	NA	NA
## sexMale:daySat:size	-1.8724	2.0522	-0.91
## sexMale:daySun:size	0.1683	2.1136	0.08
## sexMale:dayThur:size	NA	NA	NA
## smokerYes:daySat:size	-1.2955	2.7054	-0.48
## smokerYes:daySun:size	-3.8568	10.5640	-0.37
## smokerYes:dayThur:size	NA	NA	NA
## total_bill:timeLunch:size	NA	NA	NA
## sexMale:timeLunch:size	NA	NA	NA
## smokerYes:timeLunch:size	NA	NA	NA
## daySat:timeLunch:size	NA	NA	NA
## daySun:timeLunch:size	NA	NA	NA
## dayThur:timeLunch:size	NA	NA	NA
## total_bill:sexMale:smokerYes:daySat	-0.3642	1.2999	-0.28
## total_bill:sexMale:smokerYes:daySun	0.1227	1.7349	0.07
## total_bill:sexMale:smokerYes:dayThur	NA	NA	NA
## total_bill:sexMale:smokerYes:timeLunch	NA	NA	NA
## total_bill:sexMale:daySat:timeLunch	NA	NA	NA
## total_bill:sexMale:daySun:timeLunch	NA	NA	NA
## total_bill:sexMale:dayThur:timeLunch	NA	NA	NA
## total_bill:smokerYes:daySat:timeLunch	NA	NA	NA
## total_bill:smokerYes:daySun:timeLunch	NA	NA	NA
## total_bill:smokerYes:dayThur:timeLunch	NA	NA	NA

## sexMale:smokerYes:daySat:timeLunch	NA	NA	NA
## sexMale:smokerYes:daySun:timeLunch	NA	NA	NA
## sexMale:smokerYes:dayThur:timeLunch	NA	NA	NA
## total_bill:sexMale:smokerYes:size	-0.2300	0.6231	-0.37
## total_bill:sexMale:daySat:size	0.0468	0.0824	0.57
## total_bill:sexMale:daySun:size	-0.0229	0.0789	-0.29
## total_bill:sexMale:dayThur:size	NA	NA	NA
## total_bill:smokerYes:daySat:size	0.0911	0.1158	0.79
## total_bill:smokerYes:daySun:size	0.2304	0.5487	0.42
## total_bill:smokerYes:dayThur:size	NA	NA	NA
## sexMale:smokerYes:daySat:size	-5.8956	12.8327	-0.46
## sexMale:smokerYes:daySun:size	-1.8397	16.4353	-0.11
## sexMale:smokerYes:dayThur:size	NA	NA	NA
## total_bill:sexMale:timeLunch:size	NA	NA	NA
## total_bill:smokerYes:timeLunch:size	NA	NA	NA
## sexMale:smokerYes:timeLunch:size	NA	NA	NA
## total_bill:daySat:timeLunch:size	NA	NA	NA
## total_bill:daySun:timeLunch:size	NA	NA	NA
## total_bill:dayThur:timeLunch:size	NA	NA	NA
## sexMale:daySat:timeLunch:size	NA	NA	NA
## sexMale:daySun:timeLunch:size	NA	NA	NA
## sexMale:dayThur:timeLunch:size	NA	NA	NA
## smokerYes:daySat:timeLunch:size	NA	NA	NA
## smokerYes:daySun:timeLunch:size	NA	NA	NA
## smokerYes:dayThur:timeLunch:size	NA	NA	NA
## total_bill:sexMale:smokerYes:daySat:timeLunch	NA	NA	NA
## total_bill:sexMale:smokerYes:daySun:timeLunch	NA	NA	NA
## total_bill:sexMale:smokerYes:dayThur:timeLunch	NA	NA	NA
## total_bill:sexMale:smokerYes:daySat:size	0.2622	0.6280	0.42
## total_bill:sexMale:smokerYes:daySun:size	0.0728	0.8263	0.09
## total_bill:sexMale:smokerYes:dayThur:size	NA	NA	NA
## total_bill:sexMale:smokerYes:timeLunch:size	NA	NA	NA
## total_bill:sexMale:daySat:timeLunch:size	NA	NA	NA
## total_bill:sexMale:daySun:timeLunch:size	NA	NA	NA
## total_bill:sexMale:dayThur:timeLunch:size	NA	NA	NA
## total_bill:smokerYes:daySat:timeLunch:size	NA	NA	NA
## total_bill:smokerYes:daySun:timeLunch:size	NA	NA	NA
## total_bill:smokerYes:dayThur:timeLunch:size	NA	NA	NA
## sexMale:smokerYes:daySat:timeLunch:size	NA	NA	NA
## sexMale:smokerYes:daySun:timeLunch:size	NA	NA	NA
## sexMale:smokerYes:dayThur:timeLunch:size	NA	NA	NA
## total_bill:sexMale:smokerYes:daySat:timeLunch:size	NA	NA	NA
## total_bill:sexMale:smokerYes:daySun:timeLunch:size	NA	NA	NA
## total_bill:sexMale:smokerYes:dayThur:timeLunch:size	NA	NA	NA
##	Pr(> t)		
## (Intercept)	0.62		
## total_bill	0.72		
## sexMale	0.67		
## smokerYes	0.46		
## daySat	0.61		
## daySun	0.53		
## dayThur	0.66		
## timeLunch	0.58		
## size	0.62		

## total_bill:sexMale	0.69
## total_bill:smokerYes	0.29
## sexMale:smokerYes	0.74
## total_bill:daySat	0.97
## total_bill:daySun	0.68
## total_bill:dayThur	0.68
## sexMale:daySat	0.68
## sexMale:daySun	0.91
## sexMale:dayThur	0.74
## smokerYes:daySat	0.69
## smokerYes:daySun	0.63
## smokerYes:dayThur	0.76
## total_bill:timeLunch	0.49
## sexMale:timeLunch	0.69
## smokerYes:timeLunch	0.58
## daySat:timeLunch	NA
## daySun:timeLunch	NA
## dayThur:timeLunch	0.66
## total_bill:size	0.42
## sexMale:size	0.51
## smokerYes:size	0.19
## daySat:size	0.64
## daySun:size	0.57
## dayThur:size	0.72
## timeLunch:size	0.62
## total_bill:sexMale:smokerYes	0.83
## total_bill:sexMale:daySat	0.93
## total_bill:sexMale:daySun	0.84
## total_bill:sexMale:dayThur	0.73
## total_bill:smokerYes:daySat	0.45
## total_bill:smokerYes:daySun	0.54
## total_bill:smokerYes:dayThur	0.68
## sexMale:smokerYes:daySat	0.71
## sexMale:smokerYes:daySun	0.97
## sexMale:smokerYes:dayThur	0.66
## total_bill:sexMale:timeLunch	0.71
## total_bill:smokerYes:timeLunch	NA
## sexMale:smokerYes:timeLunch	NA
## total_bill:daySat:timeLunch	NA
## total_bill:daySun:timeLunch	NA
## total_bill:dayThur:timeLunch	NA
## sexMale:daySat:timeLunch	NA
## sexMale:daySun:timeLunch	NA
## sexMale:dayThur:timeLunch	NA
## smokerYes:daySat:timeLunch	NA
## smokerYes:daySun:timeLunch	NA
## smokerYes:dayThur:timeLunch	NA
## total_bill:sexMale:size	0.52
## total_bill:smokerYes:size	0.12
## sexMale:smokerYes:size	0.69
## total_bill:daySat:size	0.58
## total_bill:daySun:size	0.83
## total_bill:dayThur:size	NA
## sexMale:daySat:size	0.36

## sexMale:daySun:size	0.94
## sexMale:dayThur:size	NA
## smokerYes:daySat:size	0.63
## smokerYes:daySun:size	0.72
## smokerYes:dayThur:size	NA
## total_bill:timeLunch:size	NA
## sexMale:timeLunch:size	NA
## smokerYes:timeLunch:size	NA
## daySat:timeLunch:size	NA
## daySun:timeLunch:size	NA
## dayThur:timeLunch:size	NA
## total_bill:sexMale:smokerYes:daySat	0.78
## total_bill:sexMale:smokerYes:daySun	0.94
## total_bill:sexMale:smokerYes:dayThur	NA
## total_bill:sexMale:smokerYes:timeLunch	NA
## total_bill:sexMale:daySat:timeLunch	NA
## total_bill:sexMale:daySun:timeLunch	NA
## total_bill:sexMale:dayThur:timeLunch	NA
## total_bill:smokerYes:daySat:timeLunch	NA
## total_bill:smokerYes:daySun:timeLunch	NA
## total_bill:smokerYes:dayThur:timeLunch	NA
## sexMale:smokerYes:daySat:timeLunch	NA
## sexMale:smokerYes:daySun:timeLunch	NA
## sexMale:smokerYes:dayThur:timeLunch	NA
## total_bill:sexMale:smokerYes:size	0.71
## total_bill:sexMale:daySat:size	0.57
## total_bill:sexMale:daySun:size	0.77
## total_bill:sexMale:dayThur:size	NA
## total_bill:smokerYes:daySat:size	0.43
## total_bill:smokerYes:daySun:size	0.67
## total_bill:smokerYes:dayThur:size	NA
## sexMale:smokerYes:daySat:size	0.65
## sexMale:smokerYes:daySun:size	0.91
## sexMale:smokerYes:dayThur:size	NA
## total_bill:sexMale:timeLunch:size	NA
## total_bill:smokerYes:timeLunch:size	NA
## sexMale:smokerYes:timeLunch:size	NA
## total_bill:daySat:timeLunch:size	NA
## total_bill:daySun:timeLunch:size	NA
## total_bill:dayThur:timeLunch:size	NA
## sexMale:daySat:timeLunch:size	NA
## sexMale:daySun:timeLunch:size	NA
## sexMale:dayThur:timeLunch:size	NA
## smokerYes:daySat:timeLunch:size	NA
## smokerYes:daySun:timeLunch:size	NA
## smokerYes:dayThur:timeLunch:size	NA
## total_bill:sexMale:smokerYes:daySat:timeLunch	NA
## total_bill:sexMale:smokerYes:daySun:timeLunch	NA
## total_bill:sexMale:smokerYes:dayThur:timeLunch	NA
## total_bill:sexMale:smokerYes:daySat:size	0.68
## total_bill:sexMale:smokerYes:daySun:size	0.93
## total_bill:sexMale:smokerYes:dayThur:size	NA
## total_bill:sexMale:smokerYes:timeLunch:size	NA
## total_bill:sexMale:daySat:timeLunch:size	NA

```
## total_bill:sexMale:daySun:timeLunch:size NA
## total_bill:sexMale:dayThur:timeLunch:size NA
## total_bill:smokerYes:daySat:timeLunch:size NA
## total_bill:smokerYes:daySun:timeLunch:size NA
## total_bill:smokerYes:dayThur:timeLunch:size NA
## sexMale:smokerYes:daySat:timeLunch:size NA
## sexMale:smokerYes:daySun:timeLunch:size NA
## sexMale:smokerYes:dayThur:timeLunch:size NA
## total_bill:sexMale:smokerYes:daySat:timeLunch:size NA
## total_bill:sexMale:smokerYes:daySun:timeLunch:size NA
## total_bill:sexMale:smokerYes:dayThur:timeLunch:size NA
##
## Residual standard error: 0.987 on 181 degrees of freedom
## Multiple R-squared: 0.621, Adjusted R-squared: 0.491
## F-statistic: 4.78 on 62 and 181 DF, p-value: <2e-16
```

Code to produce Interaction Plots and perform TukeyHSD tests on factors

There are the 4 factor variables in our data set. They will be tested for interaction by viewing interaction plots and running Tukey's Honest Significance Difference test. We will first confirm `time`, `smoker`, `sex`, and `day` are factor variables, then look at their levels.

```
# confirm time, smoker, sex, and day are factor variables
is.factor(tips$time)
```

```
## [1] TRUE
```

```
is.factor(tips$smoker)
```

```
## [1] TRUE
```

```
is.factor(tips$sex)
```

```
## [1] TRUE
```

```
is.factor(tips$day)
```

```
## [1] TRUE
```

```
# what meal was consumed, or time of day?
levels(tips$time)
```

```
## [1] "Dinner" "Lunch"
```

```
# was the bill payer a smoker?
levels(tips$smoker)
```

```
## [1] "No" "Yes"
```

```
# gender of the bill payer
levels(tips$sex)
```

```
## [1] "Female" "Male"
```

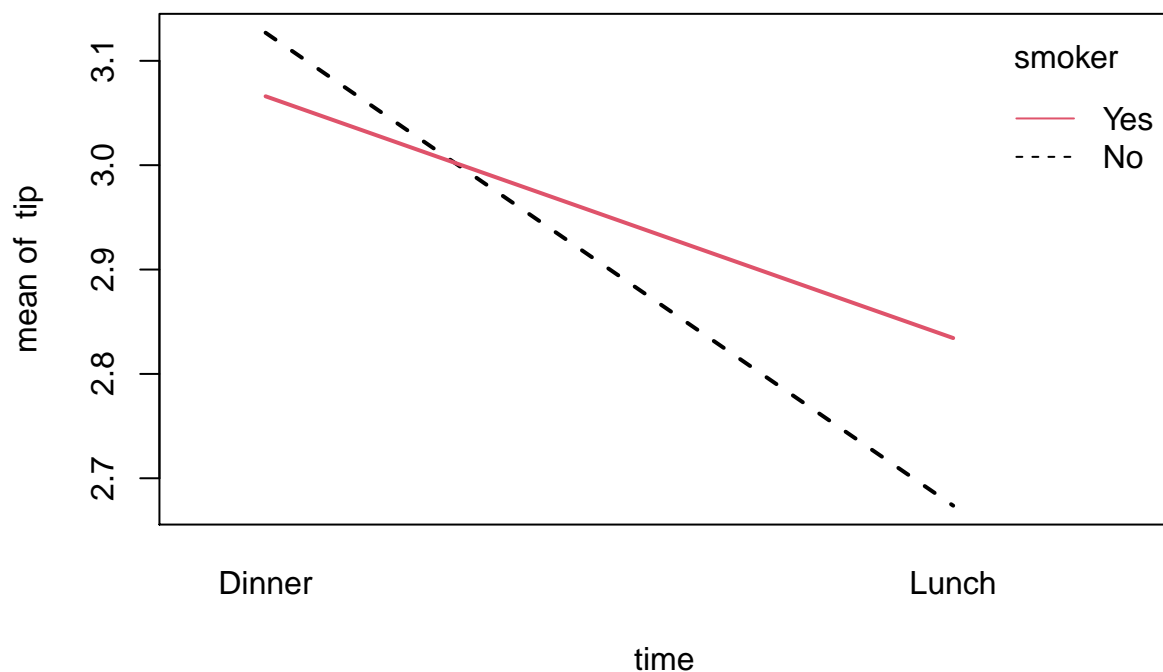
```
# day of the week (that the dining experience occurred)
levels(tips$day)
```

```
## [1] "Fri" "Sat" "Sun" "Thur"
```

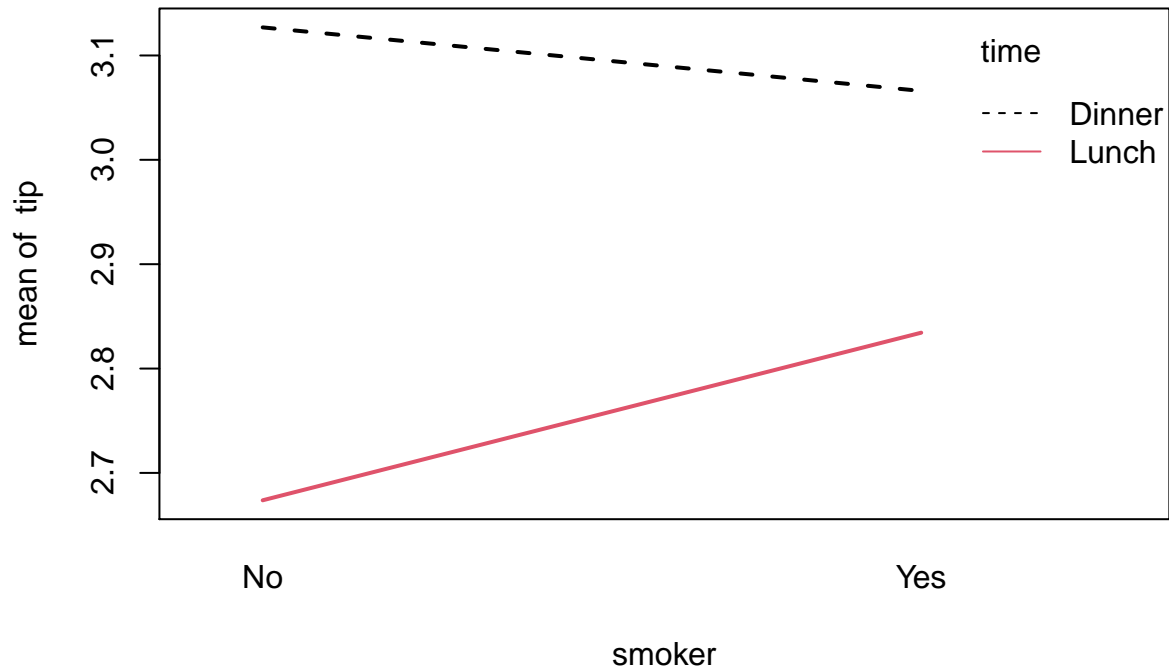
```
# interaction plots
# par(mfrow = c(1, 2))
# Before running any tests, we should first look at the data
# We will create interaction plots, which will help us visualize the effect
# of one factor, as we move through the levels of another factor.
```

```
# par(mfrow = c(1, 2))
# tip ~ total_bill * size * time * smoker * sex * day

# time, smoker factors
# tip ~ total_bill * size * time * smoker * sex * day
# interaction plots for time, smoker
with(tips, interaction.plot(time, smoker, tip, lwd = 2, col = 1:4))
```



```
with(tips, interaction.plot(smoker, time, tip, lwd = 2, col = 1:4))
```



```
## time + smoker
# within the additive model, we do further testing about the main effects
TukeyHSD(aov(tip ~ time + smoker, data = tips))
```

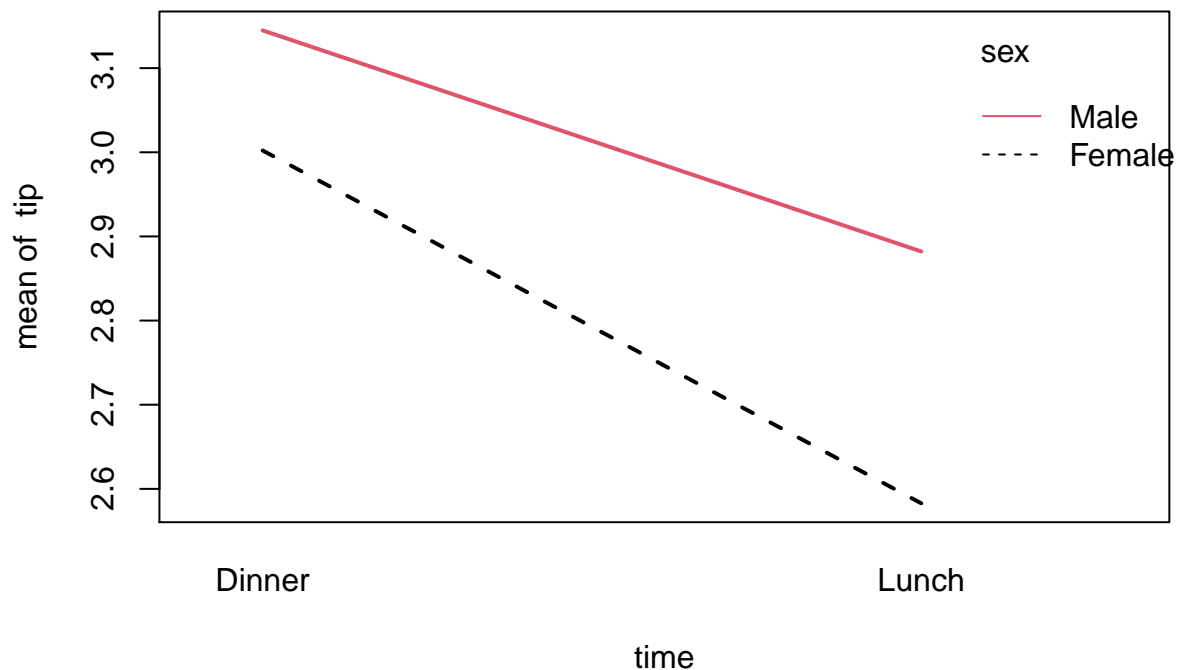
```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ time + smoker, data = tips)
##
## $time
##           diff      lwr      upr p adj
## Lunch-Dinner -0.3746 -0.7625 0.0133 0.0583
##
## $smoker
##           diff      lwr      upr p adj
## Yes-No -0.002136 -0.3602 0.3559 0.9906
```

```
## time * smoker
# within the interaction model, we do further testing about the main effects
TukeyHSD(aov(tip ~ time * smoker, data = tips))
```

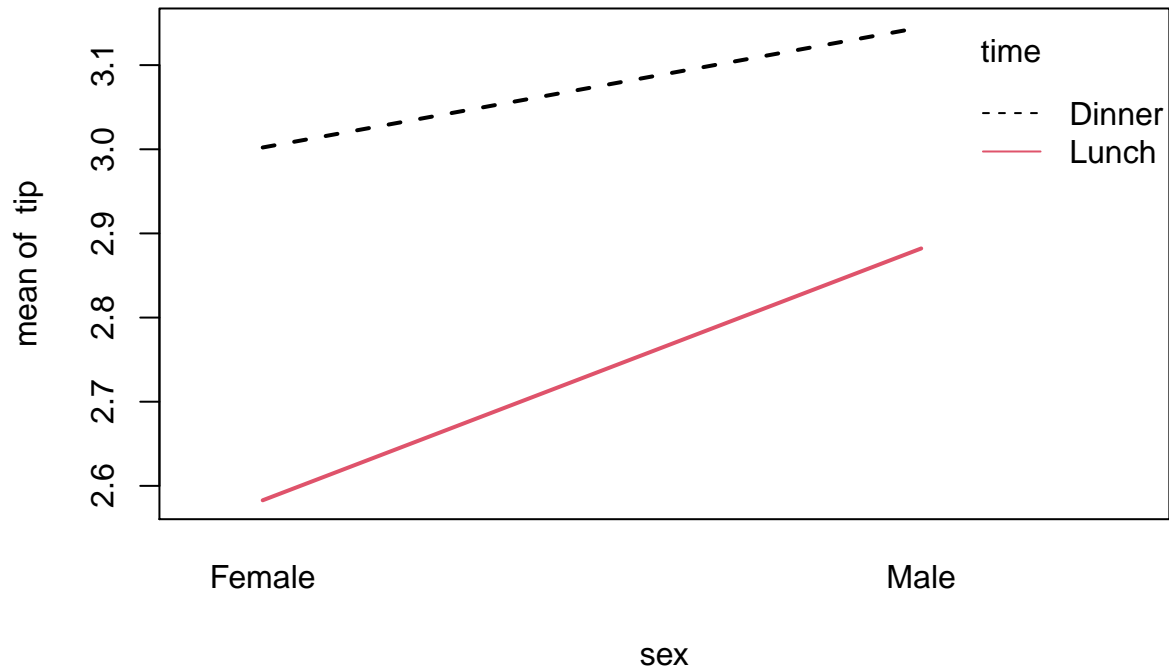
```
## Tukey multiple comparisons of means
```

```
##      95% family-wise confidence level
##
## Fit: aov(formula = tip ~ time * smoker, data = tips)
##
## $time
##              diff      lwr      upr  p adj
## Lunch-Dinner -0.3746 -0.763 0.01388 0.0587
##
## $smoker
##              diff      lwr      upr  p adj
## Yes-No -0.002136 -0.3608 0.3565 0.9906
##
## $'time:smoker'
##              diff      lwr      upr  p adj
## Lunch:No-Dinner:No -0.45311 -1.0888 0.1826 0.2554
## Dinner:Yes-Dinner:No -0.06089 -0.6112 0.4894 0.9918
## Lunch:Yes-Dinner:No -0.29254 -1.1144 0.5294 0.7937
## Dinner:Yes-Lunch:No 0.39222 -0.2905 1.0749 0.4473
## Lunch:Yes-Lunch:No 0.16057 -0.7553 1.0764 0.9689
## Lunch:Yes-Dinner:Yes -0.23165 -1.0904 0.6271 0.8978
```

```
# time, sex factors
# tip ~ total_bill * size * time * smoker * sex * day
# interaction plots for time, sex
with(tips, interaction.plot(time, sex, tip, lwd = 2, col = 1:4))
```



```
with(tips, interaction.plot(sex, time, tip, lwd = 2, col = 1:4))
```



```
## time + sex
# within the additive model, we do further testing about the main effects
TukeyHSD(aov(tip ~ time + sex, data = tips))
```

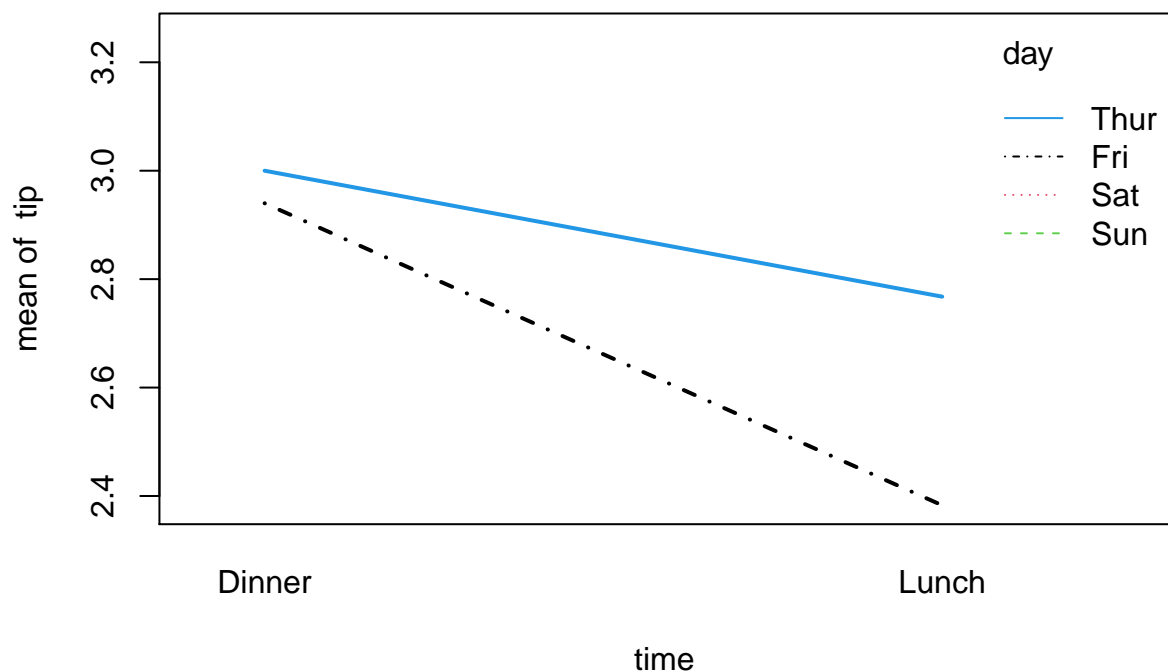
```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ time + sex, data = tips)
##
## $time
##           diff      lwr      upr p adj
## Lunch-Dinner -0.3746 -0.7616 0.01246 0.0578
##
## $sex
##           diff      lwr      upr p adj
## Male-Female 0.1842 -0.1781 0.5465 0.3175
```

```
## time * sex
# within the interaction model, we do further testing about the main effects
TukeyHSD(aov(tip ~ time * sex, data = tips))
```

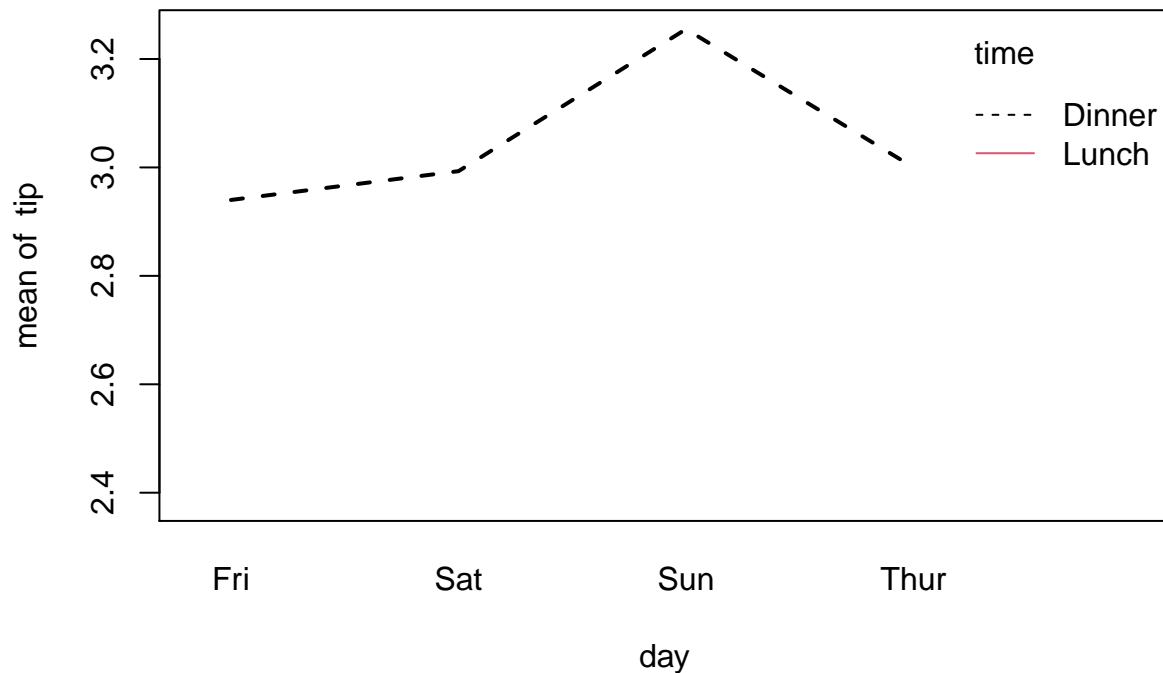
```
## Tukey multiple comparisons of means
```

```
##      95% family-wise confidence level
##
## Fit: aov(formula = tip ~ time * sex, data = tips)
##
## $time
##              diff      lwr      upr  p adj
## Lunch-Dinner -0.3746 -0.7623 0.01315 0.0582
##
## $sex
##              diff      lwr      upr  p adj
## Male-Female 0.1842 -0.1787 0.5472 0.3184
##
## $'time:sex'
##              diff      lwr      upr  p adj
## Lunch:Female-Dinner:Female -0.4193 -1.1990 0.3605 0.5062
## Dinner:Male-Dinner:Female  0.1427 -0.4465 0.7319 0.9234
## Lunch:Male-Dinner:Female  -0.1200 -0.9137 0.6737 0.9797
## Dinner:Male-Lunch:Female   0.5620 -0.1206 1.2446 0.1466
## Lunch:Male-Lunch:Female    0.2993 -0.5661 1.1646 0.8076
## Lunch:Male-Dinner:Male    -0.2627 -0.9613 0.4358 0.7650

# time, day factors
# tip ~ total_bill * size * time * smoker * sex * day
# interaction plots for time, day
with(tips, interaction.plot(time, day, tip, lwd = 2, col = 1:4))
```




```
with(tips, interaction.plot(day, time, tip, lwd = 2, col = 1:4))
```



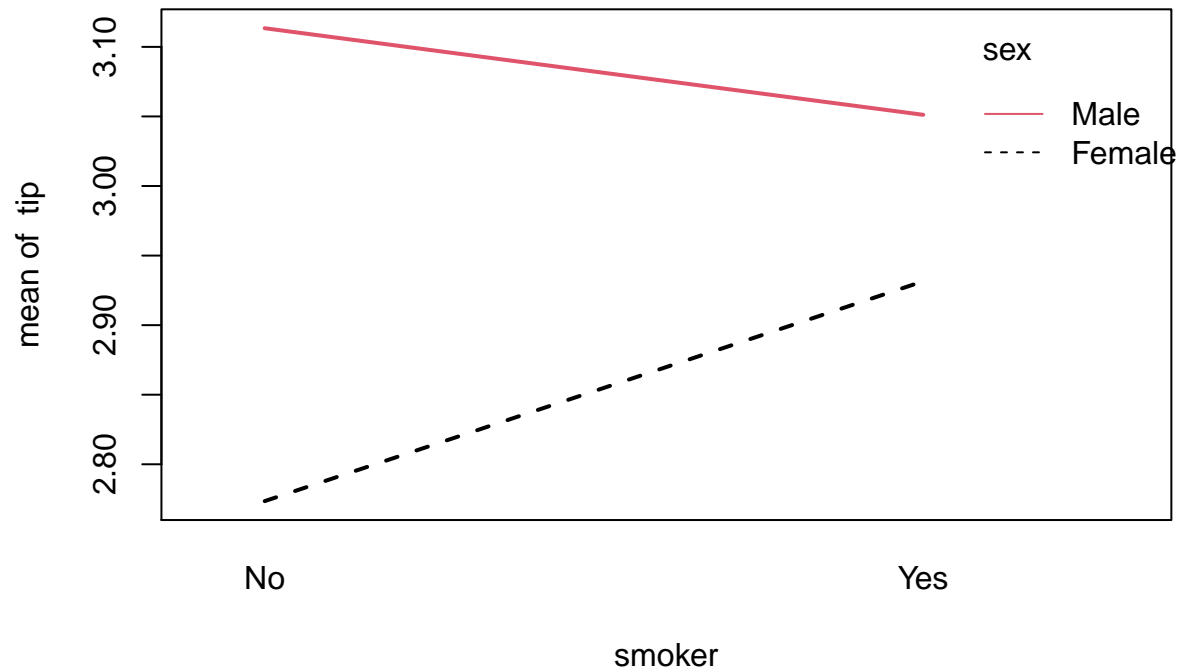
```
## time + day
# within the additive model, we do further testing about the main effects
TukeyHSD(aov(tip ~ time + day, data = tips))
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ time + day, data = tips)
##
## $time
##          diff      lwr      upr p adj
## Lunch-Dinner -0.3746 -0.7624 0.01324 0.0583
##
## $day
##          diff      lwr      upr p adj
## Sat-Fri    0.1204 -0.7830 1.0237 0.9859
## Sun-Fri    0.3824 -0.5326 1.2974 0.7013
## Thur-Fri   0.2673 -0.6681 1.2026 0.8812
## Sun-Sat    0.2620 -0.2981 0.8221 0.6209
## Thur-Sat   0.1469 -0.4460 0.7398 0.9186
## Thur-Sun  -0.1151 -0.7256 0.4953 0.9617
```

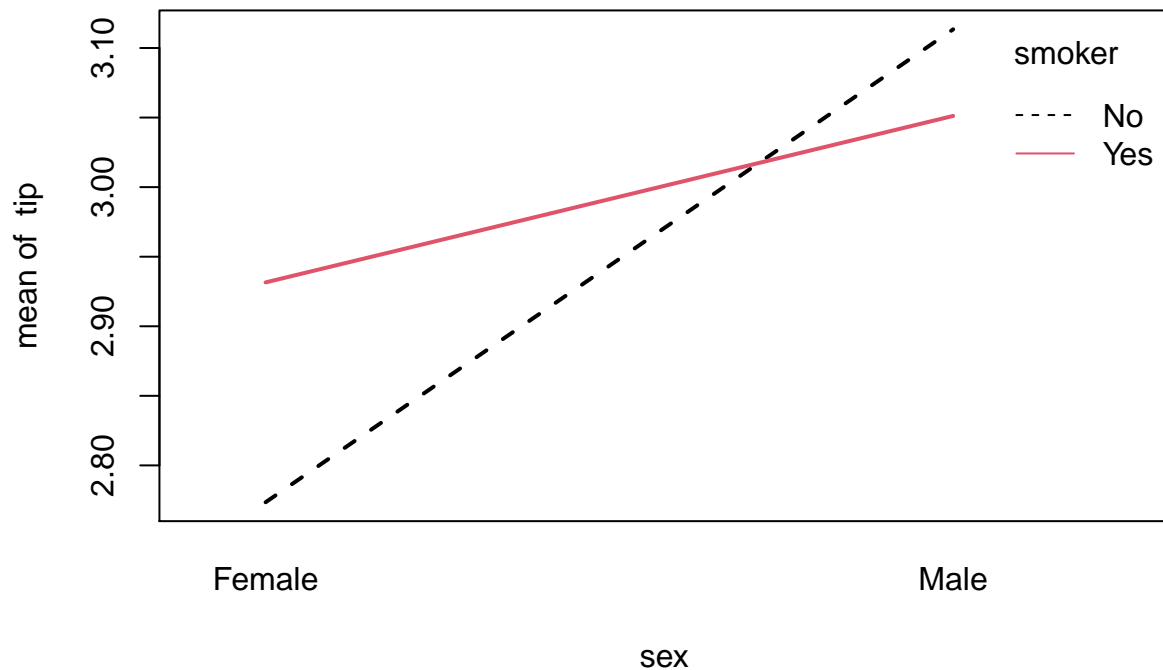
```
## time * day
# within the interaction model, we do further testing about the main effects
TukeyHSD(aov(tip ~ time * day, data = tips))
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ time * day, data = tips)
##
## $time
##          diff      lwr      upr p adj
## Lunch-Dinner -0.3746 -0.7632 0.01402 0.0588
##
## $day
##          diff      lwr      upr p adj
## Sat-Fri    0.1204 -0.7848 1.0255 0.9860
## Sun-Fri    0.3824 -0.5344 1.2992 0.7026
## Thur-Fri   0.2673 -0.6700 1.2045 0.8818
## Sun-Sat    0.2620 -0.2992 0.8232 0.6224
## Thur-Sat   0.1469 -0.4472 0.7410 0.9190
## Thur-Sun  -0.1151 -0.7268 0.4966 0.9619
##
## $'time:day'
##          diff      lwr      upr p adj
## Lunch:Fri-Dinner:Fri -0.557143 -2.5665 1.4522 0.9901
## Dinner:Sat-Dinner:Fri  0.053103 -1.2479 1.3541 1.0000
## Lunch:Sat-Dinner:Fri      NA      NA      NA      NA
## Dinner:Sun-Dinner:Fri  0.315132 -0.9973 1.6275 0.9959
## Lunch:Sun-Dinner:Fri      NA      NA      NA      NA
## Dinner:Thur-Dinner:Fri  0.060000 -4.3374 4.4574 1.0000
## Lunch:Thur-Dinner:Fri -0.172295 -1.5065 1.1619 0.9999
## Dinner:Sat-Lunch:Fri    0.610246 -1.0496 2.2701 0.9510
## Lunch:Sat-Lunch:Fri      NA      NA      NA      NA
## Dinner:Sun-Lunch:Fri    0.872274 -0.7965 2.5411 0.7509
## Lunch:Sun-Lunch:Fri      NA      NA      NA      NA
## Dinner:Thur-Lunch:Fri  0.617143 -3.8995 5.1338 0.9999
## Lunch:Thur-Lunch:Fri   0.384848 -1.3012 2.0708 0.9970
## Lunch:Sat-Dinner:Sat      NA      NA      NA      NA
## Dinner:Sun-Dinner:Sat   0.262028 -0.4013 0.9254 0.9290
## Lunch:Sun-Dinner:Sat      NA      NA      NA      NA
## Dinner:Thur-Dinner:Sat  0.006897 -4.2422 4.2560 1.0000
## Lunch:Thur-Dinner:Sat -0.225399 -0.9309 0.4801 0.9773
## Dinner:Sun-Lunch:Sat      NA      NA      NA      NA
## Lunch:Sun-Lunch:Sat      NA      NA      NA      NA
## Dinner:Thur-Lunch:Sat      NA      NA      NA      NA
## Lunch:Thur-Lunch:Sat      NA      NA      NA      NA
## Lunch:Sun-Dinner:Sun      NA      NA      NA      NA
## Dinner:Thur-Dinner:Sun -0.255132 -4.5077 3.9975 1.0000
## Lunch:Thur-Dinner:Sun -0.487427 -1.2137 0.2389 0.4491
## Dinner:Thur-Lunch:Sun      NA      NA      NA      NA
## Lunch:Thur-Lunch:Sun      NA      NA      NA      NA
## Lunch:Thur-Dinner:Thur -0.232295 -4.4917 4.0271 1.0000
```

```
# smoker, sex factors
# tip ~ total_bill * size * time * smoker * sex * day
# interaction plots for smoker, sex
with(tips, interaction.plot(smoker, sex, tip, lwd = 2, col = 1:4))
```



```
with(tips, interaction.plot(sex, smoker, tip, lwd = 2, col = 1:4))
```



```
## smoker + sex
# within the additive model, we do further testing about the main effects
TukeyHSD(aov(tip ~ smoker + sex, data = tips))
```

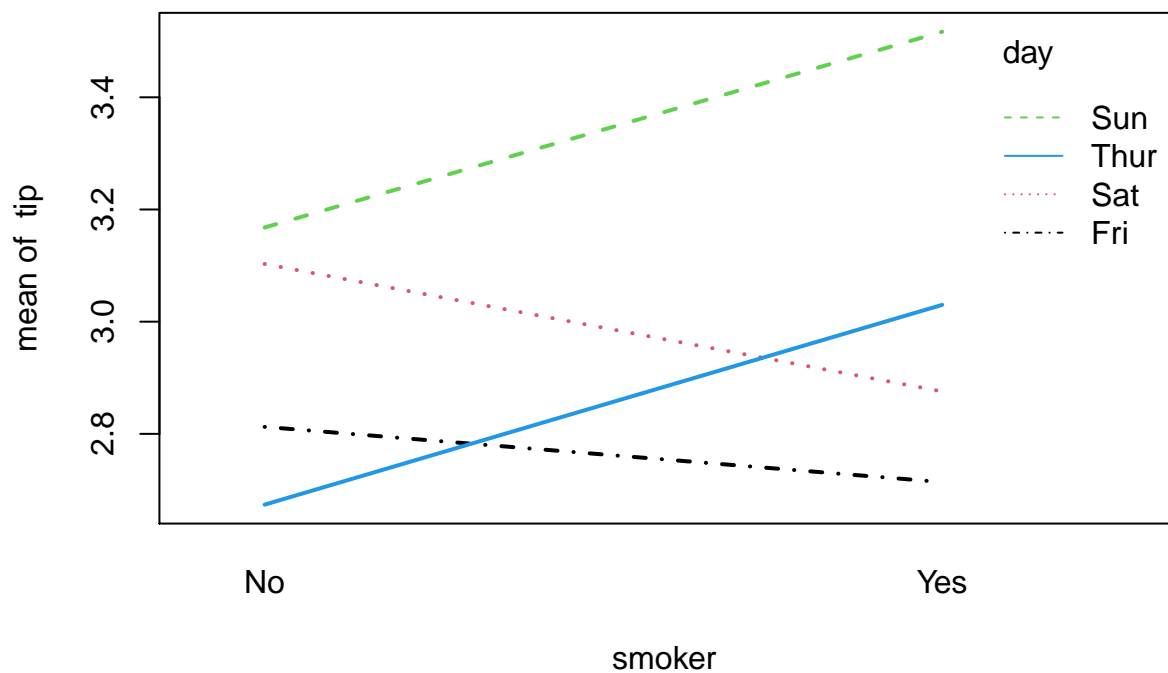
```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ smoker + sex, data = tips)
##
## $smoker
##          diff      lwr      upr p adj
## Yes-No 0.01686 -0.3425 0.3762 0.9265
##
## $sex
##          diff      lwr      upr p adj
## Male-Female 0.2561 -0.1082 0.6205 0.1674
```

```
## smoker * sex
# within the interaction model, we do further testing about the main effects
TukeyHSD(aov(tip ~ smoker * sex, data = tips))
```

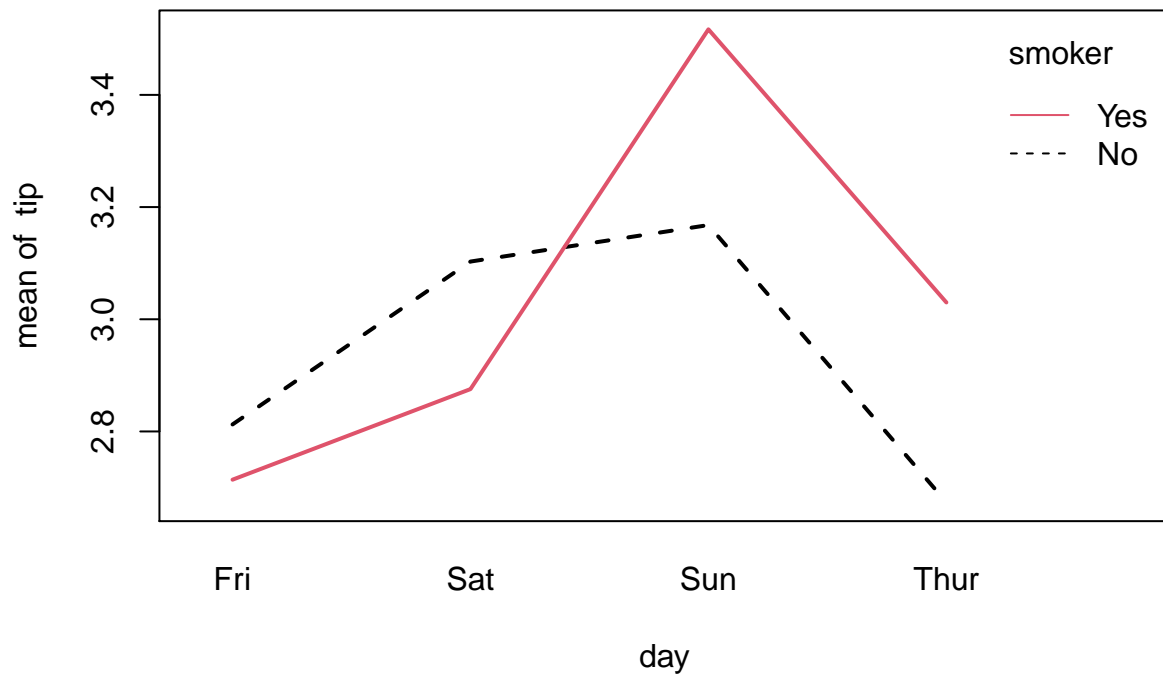
```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ smoker * sex, data = tips)
```

```
##
## $smoker
##           diff      lwr      upr  p adj
## Yes-No 0.01686 -0.343 0.3767 0.9266
##
## $sex
##           diff      lwr      upr  p adj
## Male-Female 0.2561 -0.1087 0.621 0.168
##
## $'smoker:sex'
##           diff      lwr      upr  p adj
## Yes:Female-No:Female 0.15800 -0.6342 0.9502 0.9552
## No:Male-No:Female 0.33988 -0.2688 0.9486 0.4729
## Yes:Male-No:Female 0.27765 -0.3948 0.9501 0.7093
## No:Male-Yes:Female 0.18189 -0.5406 0.9044 0.9150
## Yes:Male-Yes:Female 0.11965 -0.6573 0.8966 0.9785
## Yes:Male-No:Male -0.06224 -0.6511 0.5266 0.9928
```

```
# smoker, day factors
# tip ~ total_bill * size * time * smoker * sex * day
# interaction plots for smoker, day
with(tips, interaction.plot(smoker, day, tip, lwd = 2, col = 1:4))
```



```
with(tips, interaction.plot(day, smoker, tip, lwd = 2, col = 1:4))
```



```
## smoker + day
# within the additive model, we do further testing about the main effects
TukeyHSD(aov(tip ~ smoker + day, data = tips))
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ smoker + day, data = tips)
##
## $smoker
##      diff      lwr      upr p adj
## Yes-No 0.01686 -0.3416 0.3753 0.9263
##
## $day
##      diff      lwr      upr p adj
## Sat-Fri 0.2635 -0.6408 1.1679 0.8749
## Sun-Fri 0.5295 -0.3865 1.4455 0.4418
## Thur-Fri 0.0454 -0.8910 0.9818 0.9993
## Sun-Sat 0.2660 -0.2948 0.8267 0.6102
## Thur-Sat -0.2181 -0.8117 0.3754 0.7773
## Thur-Sun -0.4841 -1.0952 0.1271 0.1731
```

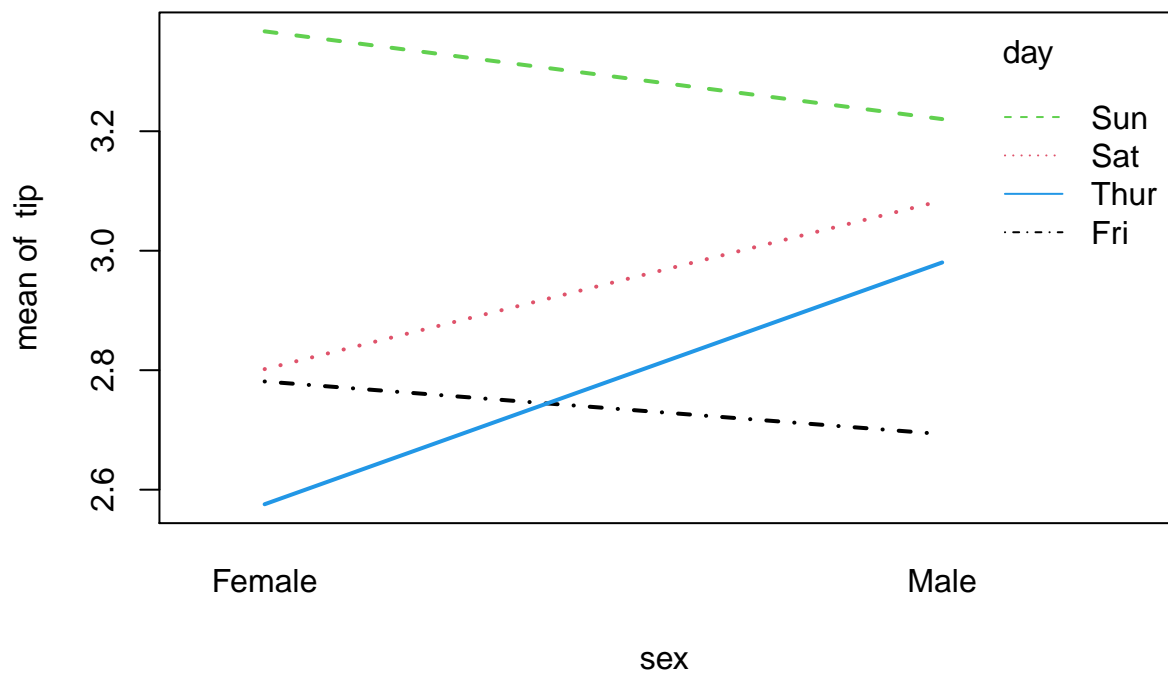
```
## smoker * day
# within the interaction model, we do further testing about the main effects
TukeyHSD(aov(tip ~ smoker * day, data = tips))
```

```

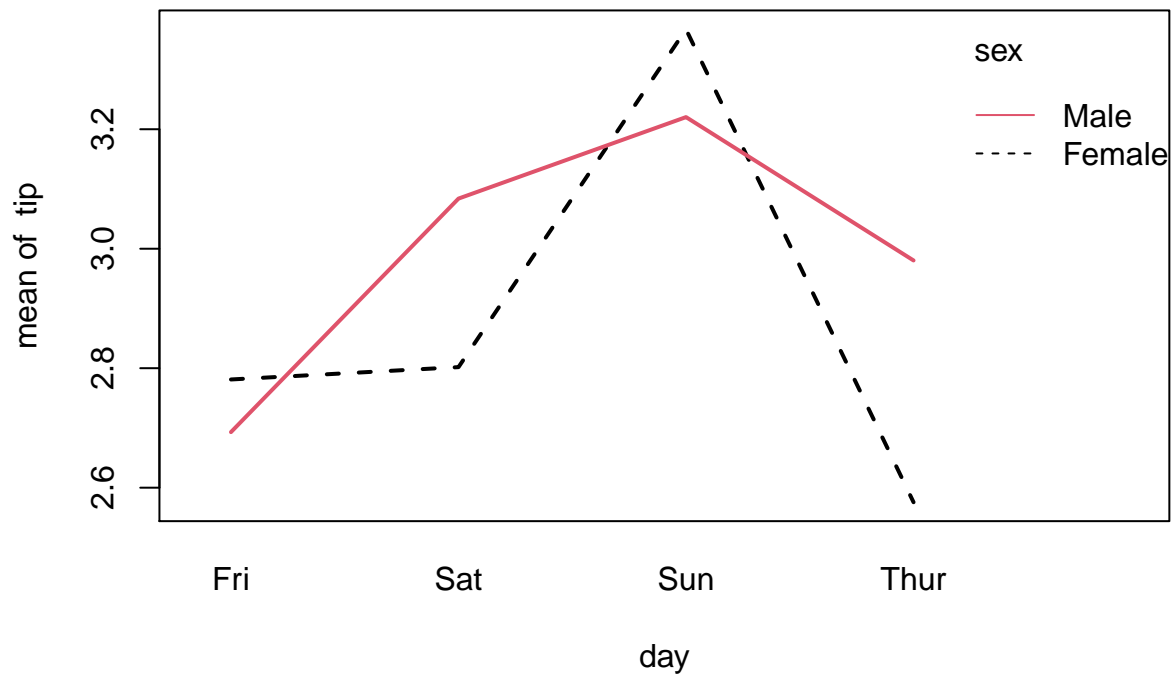
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ smoker * day, data = tips)
##
## $smoker
##      diff      lwr      upr p adj
## Yes-No 0.01686 -0.3422 0.3759 0.9264
##
## $day
##      diff      lwr      upr p adj
## Sat-Fri 0.2635 -0.6425 1.1695 0.8755
## Sun-Fri 0.5295 -0.3882 1.4472 0.4434
## Thur-Fri 0.0454 -0.8928 0.9836 0.9993
## Sun-Sat 0.2660 -0.2958 0.8277 0.6116
## Thur-Sat -0.2181 -0.8128 0.3765 0.7783
## Thur-Sun -0.4841 -1.0964 0.1282 0.1744
##
## $'smoker:day'
##      diff      lwr      upr p adj
## Yes:Fri-No:Fri -0.09850 -2.4782 2.2812 1.0000
## No:Sat-No:Fri 0.29039 -1.9160 2.4968 0.9999
## Yes:Sat-No:Fri 0.06298 -2.1499 2.2758 1.0000
## No:Sun-No:Fri 0.35539 -1.8320 2.5428 0.9997
## Yes:Sun-No:Fri 0.70434 -1.6221 3.0308 0.9833
## No:Thur-No:Fri -0.13872 -2.3452 2.0677 1.0000
## Yes:Thur-No:Fri 0.21750 -2.1326 2.5676 1.0000
## No:Sat-Yes:Fri 0.38889 -0.8719 1.6497 0.9814
## Yes:Sat-Yes:Fri 0.16148 -1.1106 1.4335 0.9999
## No:Sun-Yes:Fri 0.45389 -0.7733 1.6811 0.9494
## Yes:Sun-Yes:Fri 0.80284 -0.6578 2.2635 0.6996
## No:Thur-Yes:Fri -0.04022 -1.3010 1.2206 1.0000
## Yes:Thur-Yes:Fri 0.31600 -1.1821 1.8141 0.9982
## Yes:Sat-No:Sat -0.22741 -1.1347 0.6799 0.9946
## No:Sun-No:Sat 0.06501 -0.7783 0.9083 1.0000
## Yes:Sun-No:Sat 0.41395 -0.7431 1.5710 0.9576
## No:Thur-No:Sat -0.42911 -1.3206 0.4624 0.8216
## Yes:Thur-No:Sat -0.07289 -1.2768 1.1310 1.0000
## No:Sun-Yes:Sat 0.29242 -0.5676 1.1524 0.9678
## Yes:Sun-Yes:Sat 0.64137 -0.5278 1.8106 0.7017
## No:Thur-Yes:Sat -0.20170 -1.1090 0.7056 0.9975
## Yes:Thur-Yes:Sat 0.15452 -1.0611 1.3702 0.9999
## Yes:Sun-No:Sun 0.34895 -0.7713 1.4692 0.9803
## No:Thur-No:Sun -0.49412 -1.3374 0.3492 0.6262
## Yes:Thur-No:Sun -0.13789 -1.3065 1.0308 1.0000
## No:Thur-Yes:Sun -0.84306 -2.0001 0.3139 0.3386
## Yes:Thur-Yes:Sun -0.48684 -1.8987 0.9250 0.9652
## Yes:Thur-No:Thur 0.35622 -0.8477 1.5601 0.9854

# sex, day factors
# tip ~ total_bill * size * time * smoker * sex * day
# interaction plots for sex, day
with(tips, interaction.plot(sex, day, tip, lwd = 2, col = 1:4))

```



```
with(tips, interaction.plot(day, sex, tip, lwd = 2, col = 1:4))
```

```
## sex + day
# within the additive model, we do further testing about the main effects
TukeyHSD(aov(tip ~ sex + day, data = tips))
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ sex + day, data = tips)
##
## $sex
##          diff      lwr      upr p adj
## Male-Female 0.2562 -0.1068 0.6191 0.1657
##
## $day
##          diff      lwr      upr p adj
## Sat-Fri    0.21947 -0.6836 1.1225 0.9227
## Sun-Fri    0.45972 -0.4550 1.3744 0.5636
## Thur-Fri   0.04759 -0.8875 0.9827 0.9992
## Sun-Sat    0.24025 -0.3197 0.8002 0.6837
## Thur-Sat  -0.17188 -0.7646 0.4208 0.8765
## Thur-Sun  -0.41214 -1.0224 0.1982 0.3017
```

```
## sex * day
# within the interaction model, we do further testing about the main effects
TukeyHSD(aov(tip ~ sex * day, data = tips))
```

```

## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = tip ~ sex * day, data = tips)
##
## $sex
##           diff      lwr      upr  p adj
## Male-Female 0.2562 -0.108 0.6203 0.1671
##
## $day
##           diff      lwr      upr  p adj
## Sat-Fri    0.21947 -0.6866 1.1255 0.9234
## Sun-Fri    0.45972 -0.4580 1.3775 0.5663
## Thur-Fri   0.04759 -0.8907 0.9858 0.9992
## Sun-Sat    0.24025 -0.3215 0.8020 0.6859
## Thur-Sat  -0.17188 -0.7666 0.4228 0.8775
## Thur-Sun  -0.41214 -1.0245 0.2002 0.3046
##
## $'sex:day'
##           diff      lwr      upr  p adj
## Male:Fri-Female:Fri -0.08811 -2.0313 1.8551 1.0000
## Female:Sat-Female:Fri 0.02067 -1.5999 1.6412 1.0000
## Male:Sat-Female:Fri 0.30279 -1.2107 1.8163 0.9987
## Female:Sun-Female:Fri 0.58611 -1.1405 2.3127 0.9681
## Male:Sun-Female:Fri 0.43923 -1.0760 1.9544 0.9870
## Female:Thur-Female:Fri -0.20549 -1.8012 1.3902 0.9999
## Male:Thur-Female:Fri 0.19922 -1.4081 1.8066 0.9999
## Female:Sat-Male:Fri 0.10879 -1.4493 1.6668 1.0000
## Male:Sat-Male:Fri 0.39090 -1.0554 1.8372 0.9915
## Female:Sun-Male:Fri 0.67422 -0.9938 2.3423 0.9203
## Male:Sun-Male:Fri 0.52734 -0.9208 1.9755 0.9534
## Female:Thur-Male:Fri -0.11738 -1.6496 1.4148 1.0000
## Male:Thur-Male:Fri 0.28733 -1.2570 1.8316 0.9992
## Male:Sat-Female:Sat 0.28211 -0.6884 1.2527 0.9868
## Female:Sun-Female:Sat 0.56544 -0.7123 1.8431 0.8770
## Male:Sun-Female:Sat 0.41856 -0.5547 1.3918 0.8924
## Female:Thur-Female:Sat -0.22616 -1.3206 0.8683 0.9984
## Male:Thur-Female:Sat 0.17855 -0.9328 1.2899 0.9997
## Female:Sun-Male:Sat 0.28332 -0.8555 1.4221 0.9948
## Male:Sun-Male:Sat 0.13645 -0.6456 0.9185 0.9995
## Female:Thur-Male:Sat -0.50827 -1.4368 0.4202 0.7039
## Male:Thur-Male:Sat -0.10356 -1.0519 0.8448 1.0000
## Male:Sun-Female:Sun -0.14688 -1.2880 0.9942 0.9999
## Female:Thur-Female:Sun -0.79160 -2.0377 0.4545 0.5229
## Male:Thur-Female:Sun -0.38689 -1.6478 0.8740 0.9820
## Female:Thur-Male:Sun -0.64472 -1.5760 0.2866 0.4069
## Male:Thur-Male:Sun -0.24001 -1.1911 0.7111 0.9944
## Male:Thur-Female:Thur 0.40471 -0.6701 1.4795 0.9444

```

Final 4 Models Considered

```
# model 1
mod_list[1]
```

```
## [[1]]
##
## Call:
## lm(formula = tip ~ total_bill + size, data = tips)
##
## Coefficients:
## (Intercept)    total_bill         size
##      0.6689         0.0927         0.1926
```

```
# model 2
mod_list[2]
```

```
## [[1]]
##
## Call:
## lm(formula = tip ~ total_bill, data = tips)
##
## Coefficients:
## (Intercept)    total_bill
##      0.920         0.105
```

```
# model 3
mod_list[3]
```

```
## [[1]]
##
## Call:
## lm(formula = tip ~ total_bill + sex + smoker + day + size + total_bill:sex +
##      total_bill:smoker + total_bill:day + sex:day + total_bill:size +
##      smoker:size + total_bill:sex:day + total_bill:smoker:size,
##      data = tips)
##
## Coefficients:
##              (Intercept)              total_bill
##              -0.05804              0.17910
##              sexMale              smokerYes
##              -0.19094              -0.87667
##              daySat              daySun
##              0.88319              1.27165
##              dayThur              size
##              0.15325              -0.02690
##      total_bill:sexMale      total_bill:smokerYes
##              -0.02547              0.01674
##      total_bill:daySat      total_bill:daySun
##              -0.07895              -0.08015
##      total_bill:dayThur      sexMale:daySat
##              -0.03458              -0.95776
##              sexMale:daySun      sexMale:dayThur
```

```
##              0.45640              0.59733
##      total_bill:size              smokerYes:size
##              0.00321              0.93764
## total_bill:sexMale:daySat total_bill:sexMale:daySun
##              0.08719              0.00487
## total_bill:sexMale:dayThur total_bill:smokerYes:size
##              0.00678              -0.03432
```

```
# model 4
mod_list[4]
```

```
## [[1]]
##
## Call:
## lm(formula = tip ~ total_bill + smoker + total_bill:smoker, data = tips)
##
## Coefficients:
##      (Intercept)      total_bill      smokerYes
##           0.3601           0.1372           1.2042
## total_bill:smokerYes
##           -0.0676
```