

STAT512 – Exam 1

Spring 2018

Honor Pledge: I have not given, received, or used any unauthorized assistance on this exam.

Signature: _____

Printed Name: _____

Instructions:

- **Open book, open notes, calculator required. No computers or cell phones.**
- **Time limit is 1 hour 50 minutes - strictly enforced!**
- If an answer is in the computer output, use it; don't calculate it by hand.
- Show your work where appropriate. Put your final answer in the box (if provided).
- Make explanations brief and legible.
- All questions are worth 4 points except where noted. Maximum score is 100.
- Computer input/output is provided at the end of the exam.
- The exam contains a total of 13 pages (including computer input/output).
- If you run out of space, you may use the blank area on page 6.

Questions 1 through 5 (Model Selection): Suppose that we have a response variable (Y) and 5 predictor variables (X1 through X5). We are interested in model selection with main effects only (no interactions or polynomial terms). Circle one answer; no need to justify your response. **2 pts** per question.

1. (Multiple) R^2 can be used for model selection.

TRUE FALSE

2. The model with the lowest AIC will satisfy all model assumptions (equal variance, normality, etc).

TRUE FALSE

3. The model with the lowest AIC could include predictors that are not statistically significant.

TRUE FALSE

4. Forward and backwards selection will always arrive at the same model.

TRUE FALSE

5. Using AIC all subsets selection, suppose the model including X1 and X2 is selected. Using hypothesis testing forward selection, the model including X1 and X2 is also selected. The estimated coefficients ("betas") for these models will be the same.

TRUE FALSE

Questions 6 through 17 (MS Activity): In the article “Functional factors that are important correlates to physical activity in people with multiple sclerosis: a pilot study” (Ketelhut et al., 2017), the researchers were interested in identifying variables that are associated with physical activity in people with multiple sclerosis. A total of $n = 34$ subjects were included. Two response variables were considered: **MVPA** (moderate vigorous physical activity) and **Total.Activity**. Several potential predictor variables were considered:

Walk.Speed: Walking speed (m/sec)

Avg.Peg.Test: Average time to complete peg test (sec)

Chair.Rise: Time to rise from a seated position (sec)

TUG.Avg: Average “timed up and go” (sec)

LA.TotStr: Strength of the “less affected” leg (N/kg)

MA.TotStr: Strength of the “more affected” leg (N/kg)

Questions 6 through 10 (MS Activity 1): In this group of questions, we use **MVPA** as the response.

6. Which variable would be added first using “traditional” (hypothesis testing) forward selection?

Variable:

Brief Justification:

7. Considering the full model (**MVPAFull**), several of the VIF values are greater than 4 suggesting collinearity. However for the “final” model (**MVPAFinal**) all the VIF values are close to 1 suggesting no problems with collinearity. Explain why the collinearity is greatly reduced for the “final” model. Be specific.
8. Explain how the “final” model (**MVPAFinal**) was selected. Be sure to state both the method and the criteria.
9. Using **MVPAFinal**, predict MVPA for a subject with Avg.Peg.Test = 20 and LA.TotStr = 8. Give your final answer to one decimal place.
-
10. In the summary() output for **MVPAFinal**, an F test statistic ($F = 6.975$, $p\text{-value} = 0.003$) are shown. What is being tested? State the null hypothesis.

Questions 11 through 15 (MS Activity 2): In this group of questions, we use **Total.Activity** as the response.

11. Explain how the “final” model (**TotActFinal**) was selected. Be sure to state both the method and the criteria.

12. Is the assumption of equal variances satisfied? Name the plot you are considering (be specific) and briefly discuss whether the assumption is satisfied.

Plot:

Brief discussion:

13. Is the assumption of normality satisfied? Name the plot you are considering (be specific) and briefly discuss whether the assumption is satisfied.

Plot:

Brief discussion:

14. From the results of outlierTest(), we see that obs #29 has the largest magnitude Rstudent residual. Notice the large difference between the unadjusted p-value ($p = 0.008$) versus the Bonferoni adjusted p-value ($p = 0.2809$). Which is appropriate here? Briefly justify your response:

Which is appropriate: Unadjusted Bonferonni

Brief Justification:

15. Considering the diagnostic plots, we see that obs #24 appears to be the most influential. Give the approximate Cook’s distance for this observation and compare to the rule of thumb from the notes. (6 pts)

Approximate Cook’s D for Obs #24:

Compare to rule of thumb:

High influence: Yes No

Questions 16 and 17 (MS Activity): For this question only, we look at both models (**MVPAFinal** and **TotActFinal**).

16. A colleague looks at the AIC values for the two models (**MVPAFinal** AIC = 198.99, **TotActFinal** AIC = 269.19). He concludes that “since the AIC is smaller for MVPAFinal (as compared to TotActFinal) that model fits the data better.” Discuss whether it is appropriate to compare AIC values in this way.
17. Thinking about the question above, provide an alternative way to compare the “fit” of the two models.

Questions 18 through 26 (Firing Range): A study was done to examine noise exposure at police firing ranges. The primary question was whether the response variable **Noise** (measured in dB) differed based on the shot **Weight** (measured in grams) and **Location** (“Indoor” firing range, “Outdoor” firing range, sound proof “Control” box). Information was recorded for a total of **n=36** shots across all 3 Locations. The R input and output are labeled **Firing Range**. Use $\alpha=0.05$.

18. A colleague looks at the ANOVA table and sees that Location ($F = 1.15$, $p\text{-value} = 0.33$) is not statistically significant. He concludes that there is “no difference between Locations”. Do you agree? Briefly discuss.
19. Again, a colleague looks at the ANOVA table and sees Weight ($F = 0.96$, $p\text{-value} = 0.33$) is not statistically significant. He suggests dropping Weight from the model. Do you agree? Briefly discuss.
20. Calculate AIC for the model. Also give the value of p (# parameters).

AIC:

p:

21. Test the null hypothesis that the intercepts are the same for the three Locations. Give a test statistic and p-value.

Test Statistic:

p-value:

22. Test the null hypothesis that the slopes are the same for the three Locations. Give a test statistic and p-value.

Test Statistic:

p-value:

23. Identify the estimated intercept and slope for Indoor Location. (8 pts) Give your answers to two decimal places.

Intercept:

Slope:

24. One goal of the study is to examine the relationship between Noise and Weight at each Location. Hence the investigators are interested in testing whether the slope at each Location is different from zero. Using the output provided, which Locations have slopes that are significantly different from zero. Hint: Think about confidence intervals.

Which Locations have slopes significantly different from zero? Circle all that apply.

Control Indoor Outdoor

Brief justification:

25. Another goal is to compare the slopes for the three Locations. Identify pairs of Locations that have slopes that are significantly different from each other.

26. Yet another goal is to test for differences between the mean Noise comparing the three Locations (1) at low weight (say 5 grams) and (2) at high weight (say 200 grams). Explain how you would do this. You can either respond in words or provide (brief) R code if that is easier.

***** This page intentionally left blank *****

MS Activity 1 (Questions 6 through 10)

```
library(car)
library(MuMIn)
#Drop Total.Activity since not used in this group of questions
ActivityData <- ActivityData[,-2]
str(ActivityData)
round(cor(ActivityData),2)
#Full Model
MVPAFull <- lm(MVPA ~ Walk.Speed + Avg.Peg.Test + Chair.Rise + TUG.AVG + L
A.TotStr + MA.TotStr, data = ActivityData)
summary(MVPAFull)
vif(MVPAFull)
options(na.action = "na.fail")
MVPAcompare <- dredge(MVPAFull)
head(MVPAcompare)
#Final Model
MVPAFinal <- lm(MVPA ~ Avg.Peg.Test + LA.TotStr, data = ActivityData)
summary(MVPAFinal)
vif(MVPAFinal)
extractAIC(MVPAFinal)
```

```
> library(car)
> library(MuMIn)
> #Drop Total.Activity since not used in this group of questions
> ActivityData <- ActivityData[,-2]
> str(ActivityData)
'data.frame':      34 obs. of  8 variables:
 $ MVPA      : int  13 32 33 18 60 39 16 25 22 21 ...
 $ Walk.Speed : num  1.28 1.49 1.54 1.46 1.68 2.1 1.63 2.11 1.25 ..
 $ Avg.Peg.Test : num  23.1 17.7 17.9 23.1 18 ...
 $ Chair.Rise  : num  20.75 9.85 11.43 12.22 11.59 ...
 $ TUG.AVG    : num  8.55 6.88 6.5 7.41 7.36 ...
 $ LA.TotStr  : num  4.64 7.25 10.57 8.58 11.14 ...
 $ MA.TotStr  : num  4.62 6.1 9.42 7.76 9.8 ...

> round(cor(ActivityData),2)
      MVPA Walk.Speed Avg.Peg.Test Chair.Rise TUG.AVG LA.TotStr MA.TotStr
MVPA      1.00      0.44      -0.35      -0.47      -0.44      0.49      0.44
Walk.Speed 0.44      1.00      -0.62      -0.80      -0.84      0.50      0.60
Avg.Peg.Test -0.35    -0.62      1.00      0.64      0.72     -0.18     -0.31
Chair.Rise  -0.47    -0.80      0.64      1.00      0.91     -0.50     -0.61
TUG.AVG     -0.44    -0.84      0.72      0.91      1.00     -0.41     -0.55
LA.TotStr   0.49      0.50     -0.18     -0.50     -0.41      1.00      0.90
MA.TotStr   0.44      0.60     -0.31     -0.61     -0.55      0.90      1.00
```

MS Activity 1 continued (Questions 6 through 10)

```
> #Full Model
> MVPAFull <- lm(MVPA ~ Walk.Speed + Avg.Peg.Test + Chair.Rise + TUG.AVG +
LA.TotStr + MA.TotStr, data = ActivityData)
```

```
> summary(MVPAFull)
```

Call:

```
lm(formula = MVPA ~ Walk.Speed + Avg.Peg.Test + Chair.Rise +
    TUG.AVG + LA.TotStr + MA.TotStr, data = ActivityData)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	29.9360	36.5421	0.819	0.4198
Walk.Speed	0.9285	15.6403	0.059	0.9531
Avg.Peg.Test	-0.5942	0.8501	-0.699	0.4905
Chair.Rise	-0.4872	1.3135	-0.371	0.7136
TUG.AVG	-0.4667	2.4549	-0.190	0.8506
LA.TotStr	5.3292	3.0788	1.731	0.0949 .
MA.TotStr	-2.7294	3.1855	-0.857	0.3991

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

Residual standard error: 18.75 on 27 degrees of freedom

Multiple R-squared: 0.3406, Adjusted R-squared: 0.1941

F-statistic: 2.324 on 6 and 27 DF, p-value: 0.06133

```
> vif(MVPAFull)
```

Walk.Speed	Avg.Peg.Test	Chair.Rise	TUG.AVG	LA.TotStr	MA.TotStr
3.874243	2.178964	6.851501	9.123130	5.791174	6.680443

```
> options(na.action = "na.fail")
```

```
> MVPAcompare <- dredge(MVPAFull)
```

Fixed term is "(Intercept)"

```
> head(MVPAcompare)
```

Global model call: lm(formula = MVPA ~ Walk.Speed + Avg.Peg.Test + Chair.Rise +

TUG.AVG + LA.TotStr + MA.TotStr, data = ActivityData)

Model selection table

	(Int)	Avg.Peg.Tst	Chr.Ris	LA.TtS	TUG.AVG	Wlk.Spd	df	logLik	AICc	delta	weight
6	25.250	-0.9874		3.616			4	-144.739	298.9	0.00	0.219
21	18.510			3.080	-1.471		4	-144.802	299.0	0.13	0.206
7	19.990		-0.9449	2.808			4	-144.854	299.1	0.23	0.196
5	-1.853			4.019			3	-146.369	299.5	0.68	0.156
37	-12.030			2.976		12.83	4	-145.271	299.9	1.06	0.129
3	48.920		-1.4990				3	-146.880	300.6	1.70	0.094

Models ranked by AICc(x)

MS Activity 1 continued (Questions 6 through 10)

```
> #Final Model
> MVPAFinal <- lm(MVPA ~ Avg.Peg.Test + LA.TotStr, data = ActivityData)
> summary(MVPAFinal)
```

```
Call:
lm(formula = MVPA ~ Avg.Peg.Test + LA.TotStr, data = ActivityData)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	25.2476	18.2530	1.383	0.17649
Avg.Peg.Test	-0.9874	0.5592	-1.766	0.08728 .
LA.TotStr	3.6158	1.2422	2.911	0.00662 **

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

Residual standard error: 17.89 on 31 degrees of freedom
Multiple R-squared: 0.3104, Adjusted R-squared: 0.2659
F-statistic: 6.975 on 2 and 31 DF, p-value: 0.003153

```
> vif(MVPAFinal)
Avg.Peg.Test    LA.TotStr
    1.034945    1.034945
> extractAIC(MVPAFinal)
[1] 3.000 198.991
```

MS Activity 2 (Questions 11 through 15)

```
library(car)
str(ActivityData)
#Full Model
TotActFull <- lm(Total.Activity ~ Walk.Speed + Avg.Peg.Test + Chair.Rise +
TUG.AVG + LA.TotStr + MA.TotStr, data = ActivityData)
#Final Model
TotActFinal <- step(TotActFull, direction = "backward", trace = 0)
summary(TotActFinal)
extractAIC(TotActFinal)
outlierTest(TotActFinal)
par(mfrow=c(2,2))
plot(TotActFinal, which = c(1:2,4:5))
```

```
> library(car)
> str(ActivityData)
'data.frame':      34 obs. of  8 variables:
 $ MVPA      : int  13 32 33 18 60 39 16 25 22 21 ...
 $ Total.Activity: int 179 272 237 201 286 276 232 267 219 204 ...
 $ Walk.Speed   : num  1.28 1.49 1.54 1.46 1.68 2.1 1.63 2.11 1.25 ..
 $ Avg.Peg.Test : num  23.1 17.7 17.9 23.1 18 ...
 $ Chair.Rise   : num  20.75 9.85 11.43 12.22 11.59 ...
 $ TUG.AVG      : num  8.55 6.88 6.5 7.41 7.36 ...
 $ LA.TotStr    : num  4.64 7.25 10.57 8.58 11.14 ...
 $ MA.TotStr    : num  4.62 6.1 9.42 7.76 9.8 ...
> #Full Model
> TotActFull <- lm(Total.Activity ~ Walk.Speed + Avg.Peg.Test + Chair.Rise
+ TUG.AVG + LA.TotStr + MA.TotStr, data = ActivityData)
> #Final Model
> TotActFinal <- step(TotActFull, direction = "backward", trace = 0)
> summary(TotActFinal)
```

```
Call:
lm(formula = Total.Activity ~ Walk.Speed + Chair.Rise + TUG.AVG,
    data = ActivityData)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	148.171	88.500	1.674	0.10448
Walk.Speed	69.767	39.310	1.775	0.08608 .
Chair.Rise	-9.057	3.283	-2.759	0.00979 **
TUG.AVG	12.287	5.849	2.101	0.04417 *

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

```
Residual standard error: 49.58 on 30 degrees of freedom
Multiple R-squared:  0.44,    Adjusted R-squared:  0.384
F-statistic: 7.858 on 3 and 30 DF,  p-value: 0.0005148
```

MS Activity 2 continued (Questions 11 through 15)

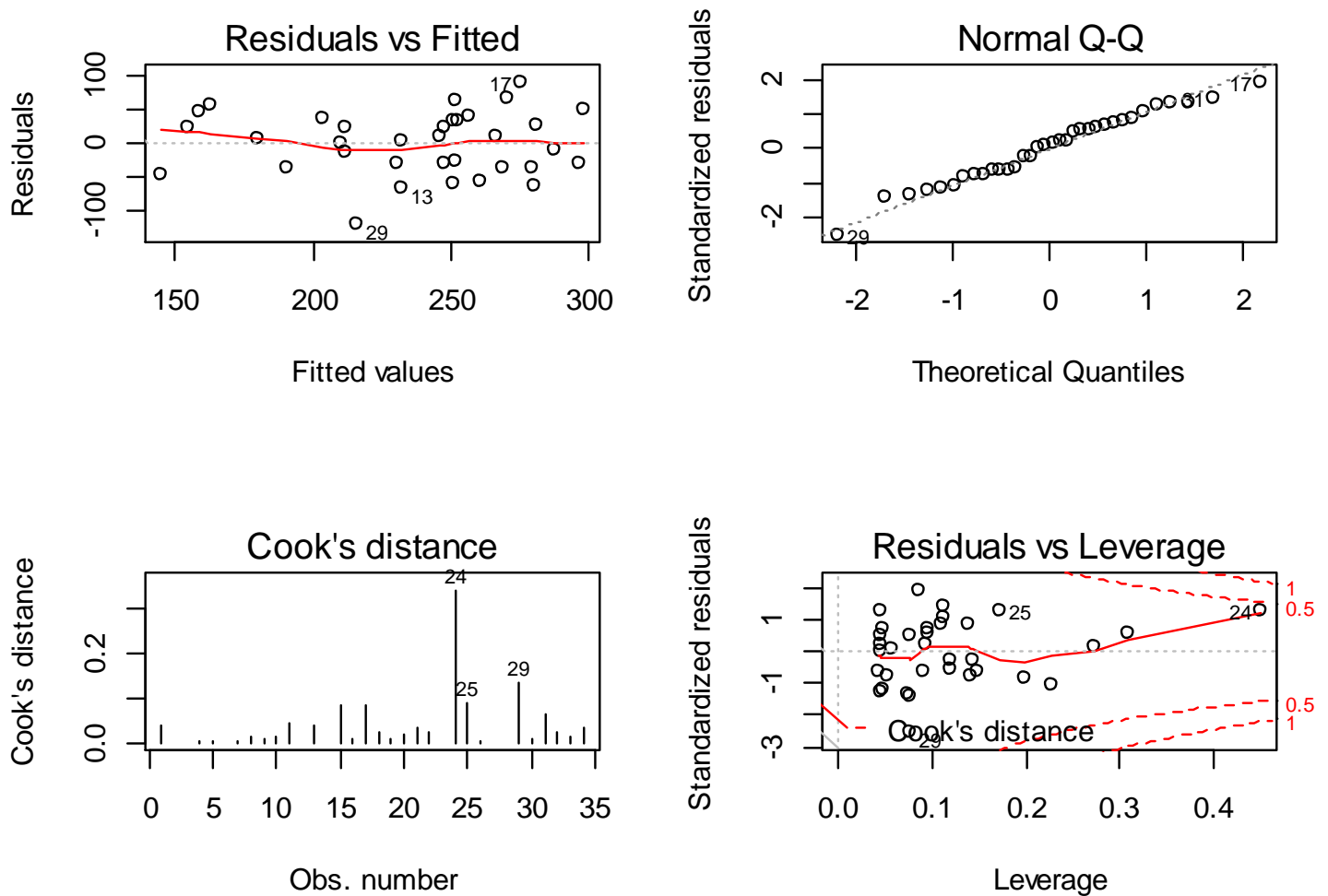
```
> extractAIC(TotActFinal)
[1] 4.0000 269.1942
> outlierTest(TotActFinal)
```

No Studentized residuals with Bonferonni $p < 0.05$

Largest $|r_{\text{student}}|$:

	r_{student}	unadjusted p-value	Bonferonni p
29	-2.835026	0.008263	0.28094

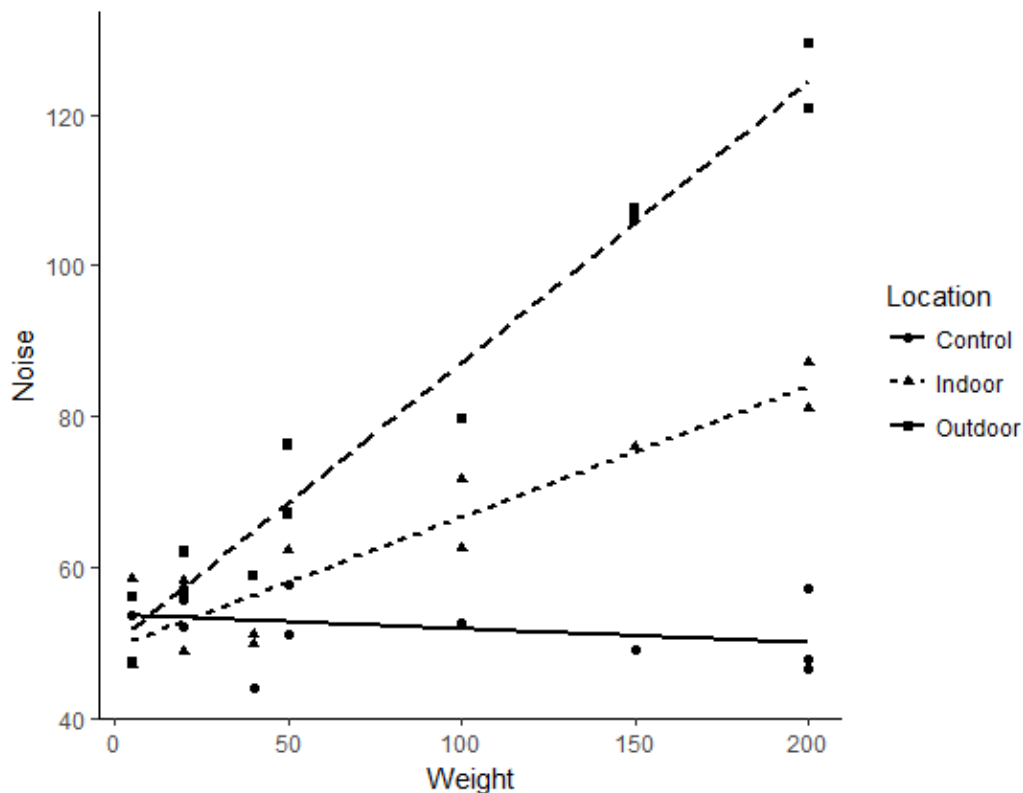
```
> par(mfrow=c(2,2))
> plot(TotActFinal, which = c(1:2,4:5))
```



Firing Range (Questions 18 through 26)

```
library(ggplot2)
library(car)
library(emmeans)
str(RangeData)
p <- qplot(Weight, Noise, shape = Location, group = Location, data =
RangeData)
p + geom_smooth(method = "lm", se = FALSE, aes(linetype = Location), color
= "black") + theme_classic()
FRModel <- lm(Noise ~ Location*Weight, data = RangeData)
Anova(FRModel, type = 3)
summary(FRModel)
emtrends(FRModel, pairwise ~ Location, var = "Weight")
```

```
> library(car)
> library(emmeans)
> str(RangeData)
'data.frame':      36 obs. of  3 variables:
 $ Location: Factor w/ 3 levels "Control","Indoor",...: 2 2 2 2 2 2 2...
 $ Weight  : int   5 5 20 20 40 40 50 100 100 150 ...
 $ Noise   : num   58.4 47 48.9 58.1 49.7 51 62.3 71.6 62.5 76 ...
> p <- qplot(Weight, Noise, shape = Location, group = Location, data = Ran
geData)
> p + geom_smooth(method = "lm", se = FALSE, aes(linetype = Location), col
or = "black") + theme_classic()
```



```
> FRModel <- lm(Noise ~ Location*Weight, data = RangeData)
```

Firing Range continued (Questions 18 through 26)

```
> Anova(FRModel, type = 3)
Anova Table (Type III tests)
```

Response: Noise

	Sum Sq	Df	F value	Pr(>F)
(Intercept)	14562.7	1	605.3154	< 2.2e-16 ***
Location	55.3	2	1.1494	0.3304
Weight	23.1	1	0.9606	0.3349
Location:Weight	4892.8	2	101.6878	4.325e-14 ***
Residuals	721.7	30		

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

```
> summary(FRModel)
```

Call:

```
lm(formula = Noise ~ Location * Weight, data = RangeData)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	53.65923	2.18099	24.603	< 2e-16 ***
LocationIndoor	-4.25527	3.04922	-1.396	0.173
LocationOutdoor	-3.74754	3.07503	-1.219	0.232
Weight	-0.01849	0.01887	-0.980	0.335
LocationIndoor:Weight	0.19156	0.02790	6.867	1.27e-07 ***
LocationOutdoor:Weight	0.39098	0.02742	14.259	6.69e-15 ***

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

Residual standard error: 4.905 on 30 degrees of freedom

Multiple R-squared: 0.9546, Adjusted R-squared: 0.947

F-statistic: 126.1 on 5 and 30 DF, p-value: < 2.2e-16

```
> emtrends(FRModel, pairwise ~ Location, var = "Weight")
```

\$emtrends

Location	Weight.trend	SE	df	lower.CL	upper.CL
Control	-0.0184936	0.01886877	30	-0.05702877	0.02004157
Indoor	0.1730671	0.02054886	30	0.13110076	0.21503351
Outdoor	0.3724846	0.01989608	30	0.33185137	0.41311778

Confidence level used: 0.95

\$contrasts

contrast	estimate	SE	df	t.ratio	p.value
Control - Indoor	-0.1915607	0.02789778	30	-6.867	<.0001
Control - Outdoor	-0.3909782	0.02742051	30	-14.259	<.0001
Indoor - Outdoor	-0.1994174	0.02860261	30	-6.972	<.0001

P value adjustment: tukey method for comparing a family of 3 estimates.