

Final exam, STATS 401 W18

Name:

UMID:

Instructions

- You have a time allowance of 120 minutes. The exam is closed book and closed notes. Any electronic devices (including calculators) in your possession must be turned off and remain in a bag on the floor.
- If you need extra paper, please number the pages and put your name and UMID on each page.
- Responses will be assessed on quality of explanation as well as whether they lead to a correct answer.
- You may use the following formulas. Proper use of these formulas may involve making appropriate definitions of the necessary quantities.

(1) $\mathbf{b} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{y}$

(2) $\text{Var}(\hat{\beta}) = \sigma^2 (\mathbb{X}^\top \mathbb{X})^{-1}$

(3) $\text{Var}(\mathbb{A}\mathbf{Y}) = \mathbb{A} \text{Var}(\mathbf{Y}) \mathbb{A}^\top$

(4) $\text{Var}(X) = \text{E}[(X - \text{E}[X])^2] = \text{E}[X^2] - (\text{E}[X])^2$

(5) $\text{Cov}(X, Y) = \text{E}[(X - \text{E}[X])(Y - \text{E}[Y])] = \text{E}[XY] - \text{E}[X] \text{E}[Y]$

(6) The binomial (n, p) distribution has mean np and variance $np(1 - p)$.

(7) $f = \frac{(\text{RSS}_0 - \text{RSS}_a)/(q - p)}{\text{RSS}_a/(n - q)}$.

Problem	Points	Your Score
1	8	
2	4	
3	6	
4	10	
5	8	
6	8	
Total	44	

All the questions in this exam refer to the field goal kicking data provided in the R dataframe `goals`. These data record the results of field goal attempts for the kickers who played in all the 2002–2006 National Football League (NFL) seasons. The primary question of interest is whether a kicker who exceeds expectations in one season is likely to do better, or worse, than expected in the following season.

Name. The name of the field goal kicker.

Yeart. The year t corresponding to the row in the dataset.

Teamt. An abbreviation of the name of the team for the kicker in year t .

FGAt. Field goal attempts in year t .

FGt. Percentage of field goal attempts that were successful in year t .

Team.t.1. An abbreviation of the name of the team for the kicker in year $t - 1$.

FGAtM1. Field goal attempts in year $t - 1$.

FGtM1. Percentage of field goal attempts that were successful in year $t - 1$.

Throughout the exam, you may write y_i for the field goal percentage recorded on the i th row of the data file, for $i = 1, \dots, n$ with $n = 4k$ corresponding to four data points on each of $k = 19$ kickers. You may also write $y_{i,j}$ for the j th measurement on kicker i , for $i = 1, \dots, k$ and $j = 1, \dots, 4$. You may use this notation without explanation. Other additional notation you use should be defined as appropriate.

```
head(goals)
```

```
##           Name Yeart Teamt FGAt  FGt Team.t.1. FGAtM1 FGtM1
## 1 Adam Vinatieri 2003    NE   34 73.5         NE    30 90.0
## 2 Adam Vinatieri 2004    NE   33 93.9         NE    34 73.5
## 3 Adam Vinatieri 2005    NE   25 80.0         NE    33 93.9
## 4 Adam Vinatieri 2006    IND  19 89.4         NE    25 80.0
## 5   David Akers  2003    PHI  29 82.7         PHI    34 88.2
## 6   David Akers  2004    PHI  32 84.3         PHI    29 82.7
```

1. Factors and their coding in R.

We will start the analysis by fitting a basic model, seen earlier in class and homework, specified in R code as

```
lm1 <- lm(FGt~Name+FGtM1, data=goals)
```

- (a) [5 points]. Write down the sample model fitted by `lm1` in subscript form.

(b) [3 points]. Write down the first 6 rows of the design matrix for `lm1`. You may use dots (`...`) to abbreviate entries following a repeated pattern, but if you do this it must be clear what they represent.

```
coef(summary(lm1))["FGtM1",]
```

```
##      Estimate   Std. Error    t value   Pr(>|t|)
## -5.037008e-01  1.127613e-01  -4.466963e+00  3.899977e-05
```

2. Model interpretation. [4 points].

A direct interpretation of the estimated coefficient for the previous year field goal percentage from `lm1` (shown above) is that field goal kickers who kick well one season tend to kick relatively poorly the next season. Explain why general principles for the interpretation of observational studies should make us cautious about jumping to that conclusion.

3. Model diagnostics.

One possible explanation behind some, or all, of the negative association between kicking percentages in subsequent years could be that coaches who have lower expectation of the abilities of the kicker tend to refrain from hard field goal attempts the following season, pushing up the next season's success rate average. Correspondingly, a coach emboldened by successful kicking may follow this up with choosing to kick in challenging situations. To investigate this, we can consider a linear model where the number of field goal attempts in year t is explained by the field goal success rate in year $t - 1$.

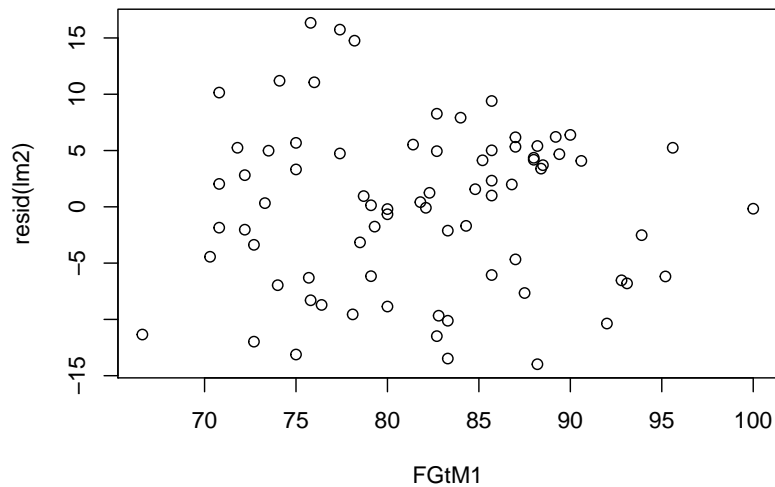
```
lm2 <- lm(FGAt~Name+FGtM1, data=goals)
anova(lm2)
```

```
## Analysis of Variance Table
##
## Response: FGAt
##           Df Sum Sq Mean Sq F value Pr(>F)
## Name      18  623.0  34.613   0.5027 0.9459
## FGtM1      1    1.8   1.823   0.0265 0.8713
## Residuals 56 3855.7  68.851
```

- (a) [4 points]. Interpret the results of this fitted linear model in the context of question of primary interest in the data analysis. You are not asked to give all the details for a hypothesis test or confidence interval. That will come in later questions; here, it is enough to describe briefly the statistical reasoning behind your interpretation.

We should always investigate the data graphically in addition to fitting a model.

```
plot(resid(lm2)~FGtM1, data=goals)
```



(b) [2 points]. Comment on your interpretation of the above residual plot, and how it relates to your answer to (a).

One other possibility proposed in class to explain the unexpected results of our first model is that kickers must do well in the earlier years included in the dataset, since they necessarily maintained their position on the team throughout the 2002–2006 interval. The following model investigated the evidence for the magnitude of this effect.

```
lm3 <- lm(FGt~Name+FGtM1+factor(Yeart), data=goals)
anova(lm3)
```

```
## Analysis of Variance Table
##
## Response: FGt
##           Df Sum Sq Mean Sq F value    Pr(>F)
## Name      18 1569.68   87.20  2.1577  0.01573 *
## FGtM1      1  769.99  769.99 19.0520 5.923e-05 ***
## factor(Yeart) 3   18.97    6.32  0.1564  0.92508
## Residuals  53 2141.99   40.41
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4. An investigation using an F-test.

(a) [5 points]. Write out in full, using subscript form, the alternative hypothesis, H_a , for using `lm3` to test whether the field goal average changes over time.

(b) [5 points]. Carry out an F test of the hypothesis H_a against a suitably constructed null hypothesis, H_0 , giving explanation of how this test is constructed. What do you conclude?

5. A confidence interval.

- (a) *[5 points]*. Using the model in Question 1 and the R output on `lm1`, explain how R obtains the estimated coefficient of goal kicking percentage in year $t - 1$ as a predictor of goal kicking percentage in year t . Also, using the probability model implicitly assumed in the analysis of Question 1, explain how to construct a 95% confidence interval for the true coefficient.

- (b) *[3 points]*. A confidence interval is only as trustworthy as the model that it is derived from. Explain to what extent you feel the confidence interval is justified based on the analysis available in this exam. Propose any supplementary analysis you would do to strengthen this inference.

6. Collinearity.

Suppose someone suggests that the rest of the team may also be an important component of field goal success. This leads you to try adding to the model a factor for the team in year t with the following consequence.

```
lm4 <- lm(FGt~Name+Teamt+FGtM1, data=goals)
summary(lm4)
```

```
##
## Call:
## lm(formula = FGt ~ Name + Teamt + FGtM1, data = goals)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11.0807  -3.2025  -0.4982   4.0692  13.2308
##
## Coefficients: (17 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    126.7703    10.6630  11.889 < 2e-16 ***
## NameDavid Akers    -3.6917     4.7822  -0.772  0.4436
## NameJason Elam     -2.0890     4.8118  -0.434  0.6660
## NameJason Hanson     3.1180     4.7613   0.655  0.5154
## NameJay Feely      -5.2243     5.7213  -0.913  0.3654
## NameJeff Reed      -7.3385     4.7801  -1.535  0.1308
## NameJeff Wilkins    3.2869     4.7674   0.689  0.4936
## NameJohn Carney    -5.0437     4.8041  -1.050  0.2986
## NameJohn Hall      -7.5838     4.8506  -1.563  0.1240
## NameKris Brown    -12.4942     4.9275  -2.536  0.0143 *
## NameMatt Stover     9.7595     4.7649   2.048  0.0456 *
## NameMike Vanderjagt  3.6936     7.2192   0.512  0.6111
## NameNeil Rackers   -5.6610     4.7785  -1.185  0.2415
## NameOlindo Mare   -12.1338     4.8506  -2.501  0.0156 *
## NamePhil Dawson    4.5452     4.7621   0.954  0.3443
## NameRian Lindell   -3.9423     4.8153  -0.819  0.4167
## NameRyan Longwell  -5.2597     7.3294  -0.718  0.4762
## NameSebastian Janikowski -3.0388     4.7995  -0.633  0.5294
## NameShayne Graham   3.1111     4.7677   0.653  0.5169
## TeamtATL           -8.4916     6.2682  -1.355  0.1814
## TeamtBAL              NA          NA      NA      NA
## TeamtBUF              NA          NA      NA      NA
## TeamtCIN              NA          NA      NA      NA
## TeamtCLE              NA          NA      NA      NA
## TeamtDAL           -2.9588    10.1814  -0.291  0.7725
## TeamtDEN              NA          NA      NA      NA
## TeamtDET              NA          NA      NA      NA
## TeamtGB              5.3209     7.3222   0.727  0.4707
## TeamtHOU              NA          NA      NA      NA
## TeamtIND             3.9384     7.2302   0.545  0.5883
## TeamtMIA              NA          NA      NA      NA
## TeamtMIN              NA          NA      NA      NA
## TeamtNE              NA          NA      NA      NA
## TeamtNO              NA          NA      NA      NA
## TeamtNYG              NA          NA      NA      NA
## TeamtOAK              NA          NA      NA      NA
```



```

## TeamtPHI          NA          NA          NA          NA
## TeamtPIT          NA          NA          NA          NA
## TeamtSTL          NA          NA          NA          NA
## TeamtWAS          NA          NA          NA          NA
## FGtM1             -0.5164      0.1170     -4.414  5.15e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.234 on 52 degrees of freedom
## Multiple R-squared:  0.551, Adjusted R-squared:  0.3524
## F-statistic: 2.774 on 23 and 52 DF, p-value: 0.00117

```

(a) [4 points]. Explain why all but four of the coefficients for the team factors take value NA.

The following results show that if we put the kicker into the model first, then the team appears insignificant from an F test. However, if we put team first then it is significant and kicker becomes insignificant.

```
anova(lm(FGt~Name+Teamt+FGtM1, data=goals))
```

```

## Analysis of Variance Table
##
## Response: FGt
##           Df Sum Sq Mean Sq F value    Pr(>F)
## Name      18 1569.68   87.20  2.2440  0.0121 *
## Teamt     4  153.02   38.25  0.9844  0.4242
## FGtM1     1  757.14  757.14 19.4831 5.147e-05 ***
## Residuals 52 2020.79   38.86
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```
anova(lm(FGt~Teamt+Name+FGtM1, data=goals))
```

```
## Analysis of Variance Table
##
## Response: FGt
##           Df Sum Sq Mean Sq F value    Pr(>F)
## Teamt     21 1721.49   81.98  2.1094  0.01508 *
## Name       1    1.20    1.20  0.0310  0.86100
## FGtM1      1  757.14  757.14 19.4831 5.147e-05 ***
## Residuals 52 2020.79   38.86
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- (b) [4 points]. Explain why the significance of the effect of the team and the kicker depends on the order in which the variables occur in the model. Can the data distinguish whether the goal kicking percentage is best explained by team or by kicker or by both?

Acknowledgments: The `goals` data were presented by *A Modern Approach to Regression with R* by S. J. Sheather, and originally come from <http://www.rortimes.com/nfl/stats>.

License: This material is provided under an [MIT license] (<https://ionides.github.io/401w18/LICENSE>)
