1                                                       Final Project

2                                                       Nate Koser

3                                                    Rutgers University

4                                              Final Project

## Introduction

6        This paper is an investigation of drunk driving patterns in South Korea for the year
7    2016. The data were obtained from kosis.kr, the Korean Statistical Information Service
8    website maintained by the National Statistical Office. The data include information on the
9    number of cases, deaths, and injuries resulting from drunk driving incidents organized by
10   time of day, month of the year, and level of intoxication.

## Hypotheses

12       The analysis presented here is an investigation of three main hypotheses. The first is
13   that the number of cases that occured will be a good predictor for time of day (day or night),
14   and that generally there are more incidents at night. The second is that there will be more
15   cases in months with major holidays. These months will include December/January - when
16   solar new year's day occurs; February - when lunar new year's day occurs; and September -
17   when the fall harvest holiday (similar to Thanksgiving in North America) occurs. The third
18   hypothesis to be tested involves high levels of intoxication. Within the `abv` variable, there is
19   a value ">.35%", indicating a blood alcohol content of over .35%. The hypothesis is that at
20   this extreme level of intoxication, the risk of death will be greater. Descriptive statistics,
21   plots, and model information for all three hypotheses will be provided in the following
22   sections.

## Data analysis

24   *Day and Night*

25       The following table shows a summary of some descriptive statistics regarding the
26   general occurance of drunk driving during the night and during the day. The numbers given
27   are log-adjusted.

| timeofday | CasesMean | CasesSd | InjuriesMean | InjuriesSd |
|---|---|---|---|---|
| day | 2.92 | 1.63 | 3.40 | 1.73 |
| night | 3.58 | 2.01 | 4.05 | 2.09 |

<sup>28</sup> We observe a higher mean incidence of cases during the night than during the day, as <sup>29</sup> well as a higher mean number of injuries during the night vs. the day. This is line with our <sup>30</sup> proposed hypothesis.

<sup>31</sup> The following is a plot of a binomial regression where a value of 1 indicates "day" and <sup>32</sup> a value of 0 indicates "night". The predictor is the number of cases that occured.
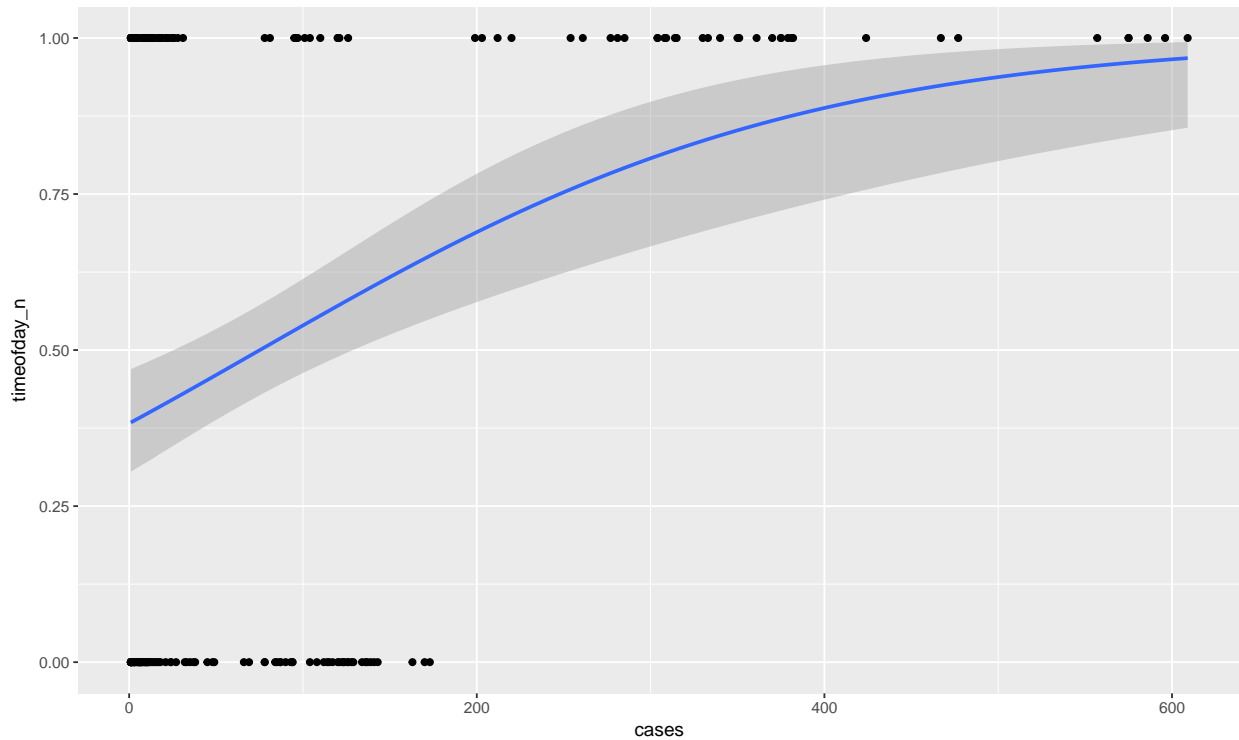


*Figure 1*

<sup>33</sup> Visual inspection reveals a general upward trend towards nighttime as the number of <sup>34</sup> cases increases. In fact, the max number of cases occuring in any one daytime period (~180) <sup>35</sup> appears to be much lower than the max number of cases occuring on the busiest night (>600). <sup>36</sup> The shape of the curve reflects this and it is once again amenable to the stated hypothesis.

| timeofday | CasesMean | CasesSd | InjuriesMean | InjuriesSd |
|---|---|---|---|---|
| day | 2.92 | 1.63 | 3.40 | 1.73 |
| night | 3.58 | 2.01 | 4.05 | 2.09 |

[28] We observe a higher mean incidence of cases during the night than during the day, as [29] well as a higher mean number of injuries during the night vs. the day. This is line with our [30] proposed hypothesis.

[31] The following is a plot of a binomial regression where a value of 1 indicates "day" and [32] a value of 0 indicates "night". The predictor is the number of cases that occured.



*Figure 1*

[33] Visual inspection reveals a general upward trend towards nighttime as the number of [34] cases increases. In fact, the max number of cases occuring in any one daytime period (~180) [35] appears to be much lower than the max number of cases occuring on the busiest night (>600). [36] The shape of the curve reflects this and it is once again amenable to the stated hypothesis.

37    A generalized linear model was fit to the data using a binomial regression with a logit

38  linking function. The following table is a summary of the results of this model fit.

| term | estimate | std. error | z value | p value |
|------|----------|-----------|---------|---------|
| (Intercept) | -0.48 | 0.18 | -2.66 | < 0.01 * |
| cases | 0.01 | 0.001 | 4.25 | < 0.001 ** |

39    We see that there is a significant effect of cases as a predictor for night or day in the

40  data ($p$-value <0.001). Taking the inverse logit of the estimates for the night-day intercept

41  and cases reveals a 62% chance of the incident having occured at night when the number of

42  cases equals 100, and an increase in this chance by ~20% when the number of cases is 200.

43  Increasing the number of cases to 300 reflects a further ~10% increase in the probability of

44  the event occuring at night, up to 92.55%. These results are reflective of the trend observed

45  in the plot above.

46    *Months of the Year*

47    The following table shows a summary of some descriptive statistics regarding the

48  general occurance of drunk driving during different months of the year. The numbers given

49  are log-adjusted.

| month | CasesMean | CasesSd | InjuriesMean | InjuriesSd |
|-------|-----------|---------|--------------|------------|
| April | 3.43 | 1.93 | 4.00 | 1.83 |
| August | 3.26 | 1.86 | 3.68 | 2.02 |
| December | 3.30 | 2.05 | 3.68 | 2.25 |
| February | 3.16 | 1.99 | 3.78 | 1.90 |
| January | 3.22 | 1.93 | 3.61 | 2.13 |
| July | 3.06 | 2.01 | 3.55 | 2.02 |
| June | 3.10 | 1.80 | 3.55 | 1.95 |
| March | 3.31 | 1.94 | 3.82 | 1.96 |

| month | CasesMean | CasesSd | InjuriesMean | InjuriesSd |
|---|---|---|---|---|
| May | 3.38 | 1.80 | 3.93 | 1.79 |
| November | 3.36 | 2.01 | 3.74 | 2.20 |
| October | 3.34 | 1.87 | 3.81 | 2.00 |
| September | 3.25 | 1.75 | 3.74 | 1.88 |

50    The months described in the hypothesis above do not appear to have any special status

51  in terms of the number of cases or injuries that occur. In particular, Janurary and

52  February's cases and inuries are among the lowest, while April sticks out as being the most

53  active. Also of note is the standard deviation values, which are more than 50% of the

54  corresponding mean for most values in the table. This indicates data that is more spread

55  out, and less concentrated around any one point (the mean).

56    *Figure 2* is a boxplot of the log cases of drunk driving with the month as the predictor.
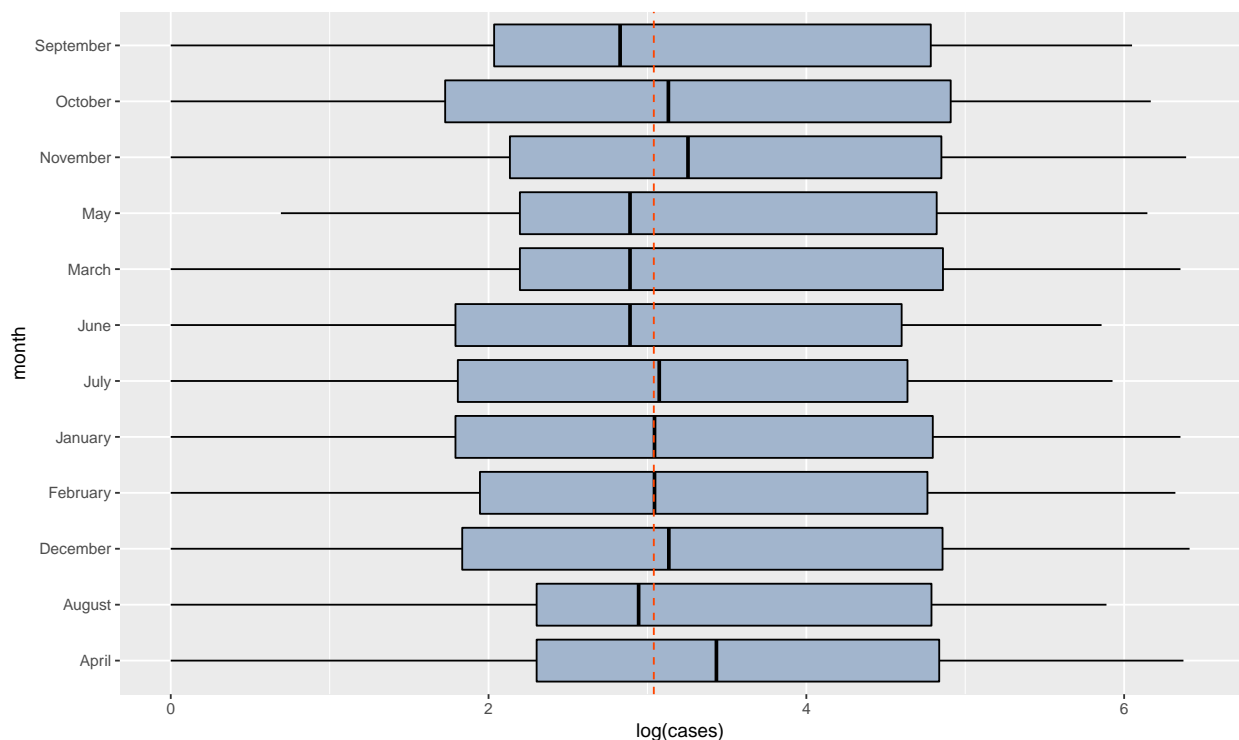


*Figure 2*

57   The median of the log cases for all months is 3.04, and the vertical red line in the plot

58   marks this value. Visual inspection reveals that the median of the log case value for January

59   and February is right on the the red line. The value for December is slightly above it, and

60   the value for September appears to be lower. This is not in line with the proposed

61   hypothesis. Unexpectedly, April shows the biggest departure from the median value - it is

62   the highest above the red line.

63   As December is the month among our hypothesized months with the highest mean and

64   median cases, it makes sense to isolate it as a "best case scenario" and fit models considering

65   it. To this end, I fit a binomial generalized linear model to the data, where a response of 1

66   indicates that the event occured in December, and a response of 0 indicates that the even

67   occured in some other month. The results are summarized in the table below.

| term | estimate | std. error | z value | p value |
|------|----------|-----------|---------|---------|
| (Intercept) | -2.56 | 0.32 | -7.95 | $< 0.001$** |
| cases | 0.001 | 0.002 | 0.5 | 0.6 |

68   No significant effect of cases was found on an event occuring in December or not. An

69   investigation of inverse logit values reveals an 8% chance of having occured in December

70   when the number of cases is equal to 100. An increase to 200 cases comes with a positive

71   change of <1% probability of being in December, and a further increase to 300 cases has the

72   same result, raising the probability to only 9.4%. None of the methods employed here

73   confirm our hypothesis that drunk driving incidents are more frequent in months with a

74   major holiday.

75   *High Blood Alcohol Content*

76   The following table shows a summary of some descriptive statistics regarding drunk

77   driving at a measured blood alcohol content of over and under .35% and the mean deaths

78   per case for each category. The numbers given are log-adjusted.

| BAC | Deaths as percent of cases |
|---|---|
| over .35 | 0.03 |
| under .35 | 0.02 |

79    The numbers indicate that there is a slightly higher chance of death at the advanced

80    BAC of >.35% than there is for lower BAC values.

81    The following is the plot of a binomial regression where a response of 1 indicates that

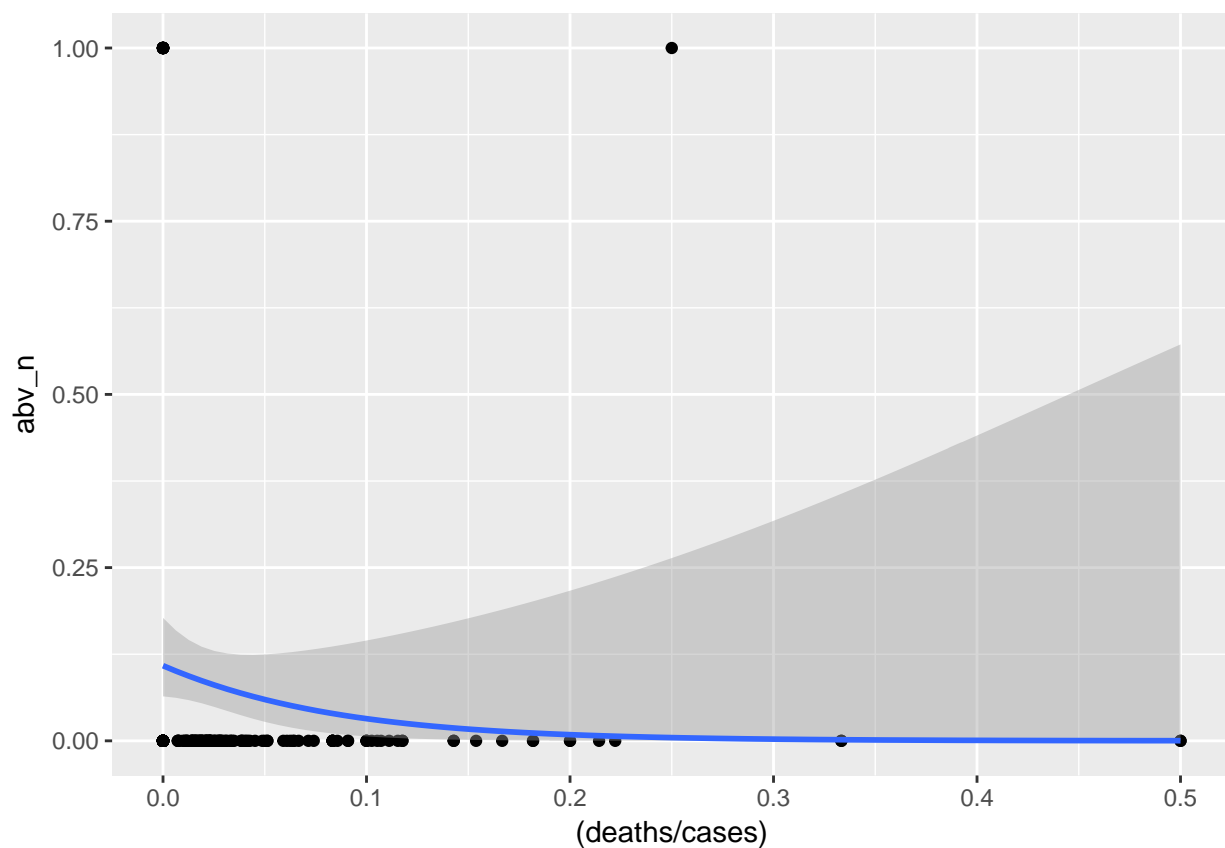82    the event occured at a BAC of >.35%, and a response of 0 indicates otherwise.



*Figure 3*

83    The plot indicates that there is a slight negative trend - as the number of deaths per

84    cases increases, the chance that the incident in question occured at a BAC value of >.35%

85    decreases. This is not what the descriptive statistics in the table above indicated. This may

86  be due to the scarcity of cases occuring at such a high BAC - the larger number of cases

87  where a larger propotion of people died at a lower BAC could offset the fewer number of

88  cases at a higher BAC.

89      The following table is the result of fitting a generalized linear model with a binomial,

90  logit linking function to the data.

| term | estimate | std. error | z value | p value |
|---|---|---|---|---|
| (Intercept) | -2.11 | 0.29 | -7.21 | < 0.001** |
| deathchance | -13.04 | 9.37 | -1.39 | 0.16 |

91      There is what could be considered a non-significant trend (p = .16) for chance of death

92  when being invovled in an incident at a BAC of over .35%. This trend is negative, as can be

93  observed in the plot. Taking the inverse logit indicates a 10.8% chance of the incident

94  occuring at a BAC of over .35% when the number of deaths per case is equal to zero. When

95  the deaths per case increases to .5, the chance of having occured at BAC >.35% is a mere

96  .02%. This does not confirm the hypothesis that more deaths occur at a higher blood alcohol

97  content.

## Discussion

99      In the preceding analysis, three hypotheses were proposed and analyzed. The first

100  hypothesis was that there are more cases of drunk driving during the night than during the

101  day. Statistical analysis through generalized linear models revealed that this is the case -

102  there was a significant effect of number of cases on time of occurance ($p$-value $< 0.001$).

103      The second hypothesis was that more drunk driving incidents occur in months that

104  have major holidays. These included September, December, January, and February.

105  Statistical anaylsis showed that there is no such effect. For the month of December, which

106  was determined ot be the most hopeful in terms of the hypothesis, there was only a slight

107  positive trend which was not found to be statistically significant ($p$-value $< 0.6$).

108        The third hypothesis was that at the high BAC level of $>.35\%$, there would be a

109   higher number of deaths per incident of drunk driving. While examination of the mean

110   initially suggested that this may be the case, there was no statistically significant effect

111   found in support of the hypothesis. In fact, a general negative trend was found ($p$-value .16).

112   This non-intuitive result may be due to the dearth of cases at this high BAC - visual

113   inspection of the data set reveals that they are only a tiny portion of the overall number of

114   cases. It is possible that incorporation of more data would reverse this trend.

115   _____