

Data Analysis Project 2

36-401, Modern Regression, Fall 2015, Section B

Due at 3:00 pm on Tuesday, 24 November 2015

This exam is week-long take-home data analysis exam. You are allowed to use your textbook as well as other reference books you feel you might need. You should use the statistical software R to perform your analysis. **You are under no circumstances allowed to consult with any person other than your professor and your teaching assistants.** You are expected to comply with the CMU policy on academic integrity. Unauthorized help will result in failing the exam, and possibly more severe disciplinary action.

Please submit two files to Blackboard: one is the PDF of your report. If you use R Markdown (or knitr) to prepare your report, also submit your .Rmd (or .Rnw) file. If you did not use R Markdown, submit a separate plain-text file with all of your R code, clearly commented so that it is clear which parts of your code go with which parts of your report. **DO NOT SUBMIT WORD FILES; THEY WILL NOT BE GRADED.**

While there are wrong answers, there are many possible right answers. Any data analysis decisions or conclusions that you make should be justified and explained. Your job is to correctly analyze the data, not force the analysis to match a pre-conceived idea. Discuss whether or not your results match your hypothesis.

1 Formatting Instructions

Your answer should be written in the format of a scientific report. **You have a seven page length limit.** Nothing over the limit will be read. Do not try to game this: fonts should be no smaller than 10 points, margins should be reasonable, graphs and tables should be embedded in the report and count against the length.

2 Data and Research Problem

Bike sharing systems are variants of traditional bicycle rentals, where the process of renting and returning is heavily automated; typically, bikes can be rented at one location and returned at another without ever having to deal with a human being. There are currently several hundred bike sharing programs in many different cities, and a great deal of interest in their potential.

You are approached by a bicycle rental company that would like to predict the daily level of bicycle rentals from environmental and seasonal variables. Since you do not know any other kind of model (yet), you will try to give them such predictions using a multiple linear regression model. The data set, <http://www.stat.cmu.edu/~cshalizi/mreg/15/dap/2/bikes.csv>, is derived from a two-year usage log of a Washington, D.C. bike-sharing system called Capital Bike Sharing (CBS) (<http://www.capitalbikeshare.com/>).

The variables are as follows:

date The full date, in year-month-day format

season Season of the year, 1 to 4

yearr Year, 0=2011, 1=2012

month Month (1 to 12)

holiday Whether the day is holiday or not

weekday Day of the week (coded by 0–6)

workingday 1 for working days, 0 for weekends and holidays

weather Weather, coded as follows:

1. Clear to partly cloudy
2. Mist but no heavier precipitation
3. Light rain or snow, possibly with thunder
4. Heavy rain or snow

temp Normalized temperature in Celsius. The values are derived via $(t - t_{min}) / (t_{max} - t_{min})$, $t_{min} = -8$, $t_{max} = +39$

atemp Normalized feeling temperature in Celsius. The values are derived via $(t - t_{min}) / (t_{max} - t_{min})$, $t_{min} = -16$, $t_{max} = +50$

hum Normalized humidity (= actual humidity divided by 100)

windspeed Normalized wind speed (= actual wind speed in miles per hour divided by 67)

registered Number of bike rentals that day by registered users.

count Count of total bike rentals that day, including both casual and registered users

The response variable of interest is **count**, the total number of rentals each day.

Note: You may find your results easier to interpret if you transform some variables. Also, using all of the predictor variables may result in collinearity problems.

Specific Analytical Questions

The following points are of special interest to the client, and you should be sure to address them in your report.

- Does having more registered users renting bikes on a given day predict higher total bike rentals?
- What is the relationship between temperature and the number of bikes rented?
- Is the relationship between temperature and the number of bikes rented the same in the two years?
- What is the relationship between humidity and the number of bikes rented?
- Is the relationship between humidity and rentals different under different weather conditions?

Suggested Outline

1. *Introduction* Write four to five sentences introducing the research problem and describing the specific research hypothesis. Cite any information sources in parentheses or foot- or end- notes.
2. *EDA* How many observations do you have?
 - Examine the (predictor and response) variables univariately and multivariately.
 - Provide graphical displays or numerical summaries for all variables.
 - You need EDA for all pairs of numerical variables and at least for the categorical variables and the response variable. Describe your results.
 - Which variables seem associated with the bicycle usage?
3. *Initial modeling* Start by building a multivariate linear regression to the data predicting usage from the predictor variables. Address the specific questions of the client when building the model. Be sure to justify the choices you made in building this initial model.
4. *Diagnostics/model selection*
 - Are the basic assumptions met for your multivariate linear regression model? Why or why not?
 - What transformations do you choose (if any)? Why?
 - Are there any outliers in your sample overly influencing your model? Identify any outlier candidates and decide whether or not to remove them. Give details.

- Do you exclude any variables? Why? All exclusions/inclusions must be justified.
5. *Final model inference/results* Create a table that summarizes your final model (coefficients, standard errors, confidence intervals, p-values). Provide interpretations of all your coefficients in the context of the problem. Be sure to address the specific questions of the client (above).
 6. *Discussions/results*: What are your conclusions? Identify a few key findings, and discuss, with reference to the supporting evidence. Can you come up with explanations for the patterns you have found? Suggestions or recommendations for the client? How could your analysis be improved? (6–8 sentences)

Rubric

Words (10) The text is laid out cleanly, with clear divisions and transitions between sections and sub-sections. The writing itself is well-organized, free of grammatical and other mechanical errors, divided into complete sentences logically grouped into paragraphs and sections, and easy to follow from the presumed level of knowledge.

Numbers (5) All numerical results or summaries are reported to suitable precision, and with appropriate measures of uncertainty attached when applicable.

Pictures (5) Figures and tables are easy to read, with informative captions, axis labels and legends, and are placed near the relevant pieces of text.

Code (10) The code is formatted and organized so that it is easy for others to read and understand. It is indented, commented, and uses meaningful names. It only includes computations which are actually needed to answer the analytical questions, and avoids redundancy. Code borrowed from the notes, from books, or from resources found online is explicitly acknowledged and sourced in the comments. Functions or procedures not directly taken from the notes have accompanying tests which check whether the code does what it is supposed to. The text of the report is free of intrusive blocks of code. If you use R Markdown, all calculations are actually done in the file as it knits, and only relevant results are shown¹. If you do not use R Markdown, the code in your appendix must generate exactly the results you show in your report, and must have comments making it clear which parts of your code go with which results.

Exploratory Data Analysis (15) Variables are examined individually and bivariate. Features/observations are discussed with appropriate figure or tables. The relevance of the EDA to the modeling is clearly explained.

¹See the model report for DAP 1 for examples of how to do this.

Model formulation and checking (30) The initial model's formulation is clearly related to the substantive questions of interest. The model's assumptions are checked by means of appropriate diagnostic plots or formal tests; if the model is re-formulated, the changes are both well-motivated by the diagnostics, and still allow the model to answer the original substantive question. Limitations from un-fixable problems are clearly noted.

Estimation, Inference and Uncertainty (15) The actual estimation of model parameters or predictions is technically correct. All calculations based on estimates are clearly explained, and also technically correct. All estimates or derived quantities are accompanied with appropriate measures of uncertainty.

Conclusions (10) The substantive questions are answered as precisely as the data and the model allow. The chain of reasoning from estimation results about models, or derived quantities, to substantive conclusions is both clear and convincing. Contingent answers ("if X , then Y , but if Z , then W ") are likewise described as warranted by the model and data. If uncertainties in the data and model mean the answers to some questions must be imprecise, this too is reflected in the conclusions.

Extra credit (10) Up to ten points may be awarded for reports which are unusually well-written, where the code is unusually elegant, where the analytical methods are unusually insightful, or where the analysis goes beyond the required set of analytical questions.

Writing Advice

Your language should be very clear and precise. Do not make claims for which you have no evidence. Do not say "will" or "would" when you really mean "may" or "might". Do not use language that implies causation; you are studying associations between variables only. When revising, look for "throat-clearing" phrases ("This is because", "this is due to", "the reason for this is", "this means that", "I believe that", "I think the reason is that") and either cut them or replace them with more concrete and informative phrases. Make sure pronouns have clear referents ("these confidence intervals show" vs. "this shows").

You may assume that the reader has a general familiarity with the contents of 36-225 and 36-226, and with the models and methods we have covered so far in the course, but will need to be reminded of any details. The reader should not be assumed to have any prior familiarity with the data set.