

**STA302 Fall 2024 Final Project Part 1**  
**Research Proposal and Data Introduction**  
**Due: October 4, 2024 by 8:00PM ET**  
**Latest Acceptance: October 11, 2024 by 8:00PM ET**

---

*Please note that if you intend to use the NQA extension time, **you should not submit any documents prior to the posted deadline** as Quercus will not allow any changes or additions to the submission after the initial deadline. Instead, make sure you have all your documents prepared and ready to submit all at once.*

---

Goal of the Assessment:	Learning Outcomes being Assessed:
<ul style="list-style-type: none"><li>• To have the opportunity to work on a topic of interest to them and to be creative about this topic.</li><li>• To experience the process of conducting a small literature review and incorporating knowledge gained into analysis.</li><li>• To think about whether a research question and/or a dataset is appropriate for use with linear regression.</li><li>• To create a draft of the components to be included in an introduction section of a report, as well as summary figures and/or tables for results section.</li></ul>	<ul style="list-style-type: none"><li>• Apply multiple linear models on various datasets using R statistical software.</li><li>• Differentiate the relationships modelled using qualitative predictors, interactions between predictors, and continuous predictors.</li><li>• Create appropriate residuals plots to evaluate model assumptions for a given data set using software.</li><li>• Recognize distinct patterns in appropriate residual plots and correctly conclude which assumption is violated.</li><li>• Report the results of a residual plot analysis and recommend a course of action.</li></ul>

**Instruction Summary:**

1. Locate open-source data in an area of interest to the group that meets the data requirements listed below. Some examples could be (but are certainly not limited to) sports, medicine, public health, economics, video games, literature, etc. Students/groups will also need to argue for why their dataset is suitable to be used with a linear regression model.
2. Define an explicit research question using the information in that dataset. Note that students/groups will need to argue for why linear regression is appropriate to answer this question with this dataset.
3. Locate three **peer-reviewed** academic papers related to the specific research question or topic of interest. Students/groups will need to describe how each article relates back to their proposed research question.
4. Select at least 5 variables from the dataset to be predictors in a preliminary multiple linear regression model, with at least one of these five being categorical in nature. These predictors must have been mentioned and summarized in the three academic papers above. The model will then be fit and a complete residual analysis to assess model assumptions will be done.
5. Provide a table that numerically summarizes each variable used in their preliminary model, with an informative caption that highlights any interesting features of the variables (e.g., skews, possible outliers or non-sensical observations, high spread, missing values).

### Dataset Requirements:

- Dataset must be open-source and the website where it was found/downloaded from must be provided.
- MUST contain at least 1000 observations (i.e., rows).
- MUST contain 1 response variable suitable for linear regression and at least 9 predictor variables, one of which must be categorical. Categorical variables with multiple levels count as 1 variable here.
  - Since at least one predictor will need to be categorical, you may convert one of your numerical variables to categorical if no such variable is available in your downloaded dataset. However, you will need to justify your choice of variable and categorization in the proposal.
- Should **NOT** be from an educational resource, such as a textbook dataset. If you're not sure, please ask the instructor or one of the TAs.
- Should **NOT** be one of the following datasets: Boston Housing dataset or Red Wine Quality dataset.
- If the dataset was found in a data repository (e.g., Kaggle, UCI Repository, etc.), you MUST ensure that your research question is novel and different from the original usage of the data.

### Proposal Format:

Your group will create a written proposal that should introduce your research question and data, summarize existing knowledge in that area, fit a preliminary model based on the existing knowledge, and conduct a residual analysis of the model. The proposal **must include** the following sections and **must not exceed** the word count in each case:

- **Contributions:** each group member's name is listed and a description of their contribution to the proposal is outlined (this does not count towards the word limit).
- **Introduction (350 words):** introduce the relevance/importance of the topic, state the research question of interest, summarize the results of three peer-reviewed research papers with a focus on their connection to the research question, and describe why linear regression is a suitable statistical tool to answer the research question.
  - i.e., why should someone be interested in your project, what are you trying to answer what is already known about this question, and why should you use linear regression.
- **Data description (300 words):** state where the data was found, explain how the data was originally collected (not how you found the data but how the original curator of the data collected it), describe the response variable (both statistically and with a written description of what it measures and why it meets the requirements for use in a linear model), summarize numerically or graphically (in a single figure/table) each predictor in your dataset that will be used in the preliminary model, and interpret the descriptive statistics in the context of what the predictors measure and how it relates to the research question.
  - NOTE: if you had to convert a numerical predictor to a categorical predictor to meet the data requirements, you must justify your choice and the chosen categories in this section.
- **Ethics discussion (100-200 words, only for L0101/L0201/L2001/L2002 students):** Would you consider your dataset to be trustworthy, given the criteria discussed in the ethics module? Justify briefly using material and terminology discussed in the first ethics module.
  - Bonus exercise: We encourage you to think about whether your dataset was collected ethically, and whether you are making ethically appropriate use of it, given the issues raised in the ethics module (you do not need to write anything about this question).
- **Preliminary results (300 words):** fit a preliminary model using 5 predictors noted in the literature, conduct a full analysis of the linear regression assumptions noting any violations and what led to your conclusions. Discuss whether your preliminary model results are similar or different to results in the literature and why.

- NOTE: Place residual plots into the document in a grid (i.e., 2-3 plots placed horizontally in a single figure) so that multiple plots will display in a single figure for improved readability (see Resources below).
- **Bibliography:** an appropriately formatted list of resources and literature cited in the proposal (not included in work count). APA format is acceptable.

### What to Submit:

Only ONE member of the group should submit ALL required submission components. A complete submission to Quercus will include:

- ✓ Your group's completed Group Teamwork Agreement, saved as a PDF (see [Quercus Final Project](#) page).
- ✓ The completed proposal, saved as a PDF.
- ✓ The Rmd file containing the code used to subset and clean the data, fit the model, produce a summary table, and conduct the residual analysis for checking assumptions.
- ✓ The original and cleaned (where appropriate) datasets as CSV files, uploaded to a cloud-based storage service (e.g., OneDrive), with the **shareable link** included as a **submission comment** on Quercus.

Failure to meet these submission requirements, including incorrect format of components, missing components, and cloud links that do not allow shared access will result in a **one-mark deduction** on the grade of the proposal.

### Resources:



Should your group have difficulty locating a suitable dataset that meets the group's interest and the dataset requirements, your group can consider using one of the datasets in the table below. You may also consider consulting the library resources for help performing your literature search and citing the results. Should your group use R Markdown to produce the proposal, the R Markdown resources will help you format your document and make it more presentable.

Dataset Resources	Library Resources	R Markdown Resources
<ul style="list-style-type: none"> <li>• <a href="#">Ames Housing dataset</a></li> <li>• <a href="#">NHANES survey dataset</a></li> <li>• <a href="#">AirBnB dataset</a> (needs you to create a free account)</li> <li>• <a href="#">Million Song dataset</a></li> <li>• <a href="#">NBA player dataset</a></li> </ul>	<ul style="list-style-type: none"> <li>• <a href="#">How to search for academic articles</a></li> <li>• <a href="#">Using search operators to find articles</a></li> <li>• <a href="#">Limiting search to peer-reviewed articles</a></li> <li>• <a href="#">Why and how to cite your references</a></li> <li>• <a href="#">Help getting the correct citation format</a></li> <li>• <a href="#">Exporting a citation</a></li> </ul>	<ul style="list-style-type: none"> <li>• <a href="#">Settings for displaying or not displaying R code in knitted document</a></li> <li>• <a href="#">Adding captions and other plotting features</a></li> <li>• Including multiple plots in a grid using <a href="#">patchwork</a> or <a href="#">base R plot</a> commands</li> <li>• Creating tables in RMarkdown using <a href="#">Kable</a> or <a href="#">manually</a></li> <li>• <a href="#">Exporting plots in RStudio</a></li> </ul>



You may also wish to consider the writing resources posted on the [General Resources](#) Quercus page. Alternatively, keep an eye on the course announcements for dedicated writing office hours with our English Language Learning TA, Dory.



For some advice in formulating a research question and searching the academic literature, see our [Tip Sheet for Creating a Research Question](#), designed by Dory Abelman.

Criteria of Assessment	Excellent (2 points)	Satisfactory (1 point)	Needs Revision (0 points)
<b>Introduction Section</b>			
<p>Proposed research question:</p> <ul style="list-style-type: none"> <li>The response variable of interest is clearly identifiable, and the predictors hypothesized to be related to the response are explicitly stated (or at minimum groups of common predictor characteristics are listed).</li> <li>It is phrased using clear language and familiar terminology and makes a clear hypothesis about the population relationship.</li> <li>It is directly connected to the stated importance/relevance of the project topic.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<p>Literature summary:</p> <ul style="list-style-type: none"> <li>Three legitimate peer-reviewed articles are summarized.</li> <li>The main result of each article is summarized concisely and in the context of the original study population.</li> <li>A strong and explicit connection is made between each article's results and the proposed research question.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<p>Suitability of linear regression:</p> <ul style="list-style-type: none"> <li>Uses appropriate terminology from the course materials.</li> <li>Provides a reasonable justification for why and how estimating a linear trend will answer the research question proposed.</li> <li>Provides a reasonable justification for whether the focus of the model will be on interpretability (description) or precision/accuracy (prediction).</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<b>Data Description Section</b>			
<p>Description of data source:</p> <ul style="list-style-type: none"> <li>Where the data was sourced/downloaded from is explicitly mentioned with a corresponding citation in the bibliography.</li> <li>The original usage or purpose of the dataset is described, and it is explicit how that usage differs from the current research proposal.</li> <li>How the data were originally collected by the curator of the dataset is described and a corresponding reference is cited from the bibliography.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<p>Response variable summary:</p> <ul style="list-style-type: none"> <li>An appropriate and suitably presentable numerical or graphical summary is used to statistically describe the response variable.</li> <li>A written description of the response variable highlights important features of the response distribution, in the context of what is being measured/the research question.</li> <li>A justification for why the chosen response variable is suitable to be used in a linear regression model is provided and is correct, based on the statistical summary presented.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.

<p>Predictor variable summaries:</p> <ul style="list-style-type: none"> <li>An appropriate and suitably presentable numerical or graphical summary is used to statistically describe the chosen predictor variables.</li> <li>Important/interesting variable characteristics (e.g. skews, abnormal values) or lack thereof are, in the context of what is being measured/the research question.</li> <li>A justification for why the chosen predictor variables are relevant to answering the research question, making explicit reference to the summarized literature and to any modifications to variables that have been made.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<b>Ethics Discussion Section (L0101/L0201/L2001/L2002 only)</b>			
<ul style="list-style-type: none"> <li>Answer correctly references some of the criteria discussed in the first ethics module.</li> <li>Response makes a reasonable and clear attempt to argue for its conclusion.</li> <li>Meets minimum and maximum word count.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<b>Preliminary Model Results Section</b>			
<p>Residual analysis of preliminary model:</p> <ul style="list-style-type: none"> <li>All plots needed for a complete residual analysis have been presented, are correct, and are easily readable with appropriate axes and labels.</li> <li>Each assumption and condition are assessed and a conclusion for each is provided.</li> <li>Correct details are provided, with reference to the appropriate plot, to describe how such a conclusion was made for each assumption and condition.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<p>Preliminary model discussion:</p> <ul style="list-style-type: none"> <li>Model estimates from preliminary model are presented in an easily readable, understandable, and professional way.</li> <li>A discussion on what these estimates tell the reader about a possible answer to the research question is provided in context, highlighting the effect of at least one numerical and one categorical predictor explicitly.</li> <li>A comparison is made between the preliminary model results and those summarised from the literature, and it is discussed why these may be similar or different.</li> </ul>	All three criteria are met.	Only two criteria are met.	One or fewer criteria are met.
<b>Overall Proposal Formatting</b>			
<ul style="list-style-type: none"> <li>The bibliography and in-text citations are formatted correctly using a consistent style.</li> <li>Word counts for each section are met or are no more than 15 words in excess.</li> <li>Headers and paragraphs are used effectively to increase readability and separate ideas for increased comprehension.</li> <li>No R code or R output (other than plots) are displayed in the written proposal.</li> </ul>	All four criteria are met.	Only three criteria are met.	Two or fewer criteria are met.
<b>Total Points:</b>			<b>/20</b>