

**STA302 Summer 2025 Final Project Part 1**  
**Research Proposal and Data Introduction**  
**Due: May 21, 2025, by 11:00PM ET**

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 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, ecology, finance, 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. 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**. The model will then be fit and a complete **residual analysis** to assess model assumptions will be done.
4. **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, 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**, 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, 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.
- **Preliminary results (300 words):** **fit a preliminary model using at least 5 predictors**, 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.
- ✓ 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 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="#">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="#">Kabble</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 Resources portion on Quercus.

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 <b>response variable</b> of interest is <b>clearly identifiable</b>, and the <b>predictors</b> hypothesized to be related to the response are <b>explicitly stated</b> (or at minimum groups of common predictor characteristics are listed).</li> <li>It is phrased using <b>clear language</b> and <b>familiar terminology</b> and makes a <b>clear hypothesis</b> about the population relationship.</li> <li>It is directly <b>connected to the stated importance/relevance</b> of the project topic.</li> </ul>	All three criteria are met.	Only two criteria are met.	Only one or fewer criteria are met.
<p>Suitability of linear regression:</p> <ul style="list-style-type: none"> <li>Uses <b>appropriate terminology</b> from the course materials.</li> <li>Provides a <b>reasonable justification</b> for <b>why and how</b> estimating a <b>linear</b> trend will answer the research question proposed.</li> <li>Provides a reasonable justification for <b>whether the focus of the model will be on interpretability (description) or precision/accuracy (prediction)</b>.</li> </ul>	All three criteria are met.	Only two criteria are met.	Only 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 <b>sourced/downloaded</b> from is explicitly mentioned with a corresponding citation in the bibliography.</li> <li>The <b>original usage or purpose of the dataset is described</b>, and it is explicit <b>how that usage differs from the current research proposal</b>.</li> <li><b>How the data were originally collected</b> by the curator of the dataset is described and a <b>corresponding reference is cited</b> from the bibliography.</li> </ul>	All three criteria are met.	Only two criteria are met.	Only one or fewer criteria are met.
<p>Response variable summary:</p> <ul style="list-style-type: none"> <li>An appropriate and suitably presentable <b>numerical or graphical summary</b> is used to statistically describe the response variable.</li> <li>A <b>written description</b> of the response variable highlights <b>important features</b> of the <b>response distribution</b>, in the context of <b>what is being measured/the research question</b>.</li> <li>A justification for <b>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</b>.</li> </ul>	All three criteria are met.	Only two criteria are met.	Only one or fewer criteria are met.

<p>Predictor variable summaries:</p> <ul style="list-style-type: none"> <li>An appropriate and suitably presentable <b>numerical or graphical summary</b> is used to statistically describe the chosen predictor variables.</li> <li>Important/interesting variable <b>characteristics (e.g. skews, abnormal values)</b> or lack thereof are, <b>in the context of what is being measured/the research question</b>.</li> <li>A justification for <b>why the chosen predictor variables are relevant to answering the research question</b>, making explicit reference to any modifications to variables that have been made.</li> </ul>	All three criteria are met.	Only two criteria are met.	Only 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><b>All plots needed for a complete residual analysis have been presented, are correct, and are easily readable</b> with appropriate axes and labels.</li> <li>Each <b>assumption and condition are assessed</b> and a <b>conclusion for each is provided</b>.</li> <li><b>Correct details are provided</b>, with <b>reference to the appropriate plot</b>, to describe how such a conclusion was made for each <b>assumption and condition</b>.</li> </ul>	All three criteria are met.	Only two criteria are met.	Only one or fewer criteria are met.
<p>Preliminary model discussion:</p> <ul style="list-style-type: none"> <li><b>Model estimates</b> from preliminary model are presented in an <b>easily readable, understandable, and professional</b> way.</li> <li>A discussion on <b>what these estimates tell the reader</b> about a possible <b>answer to the research question</b> is provided in context, highlighting the <b>effect of at least one numerical and one categorical predictor explicitly</b>.</li> </ul>	All two criteria are met.	Only one criterion is met.	None of the criteria are met.
<b>Overall Proposal Formatting</b>			
<ul style="list-style-type: none"> <li>The <b>bibliography</b> and in-text citations are formatted correctly using a consistent style.</li> <li><b>Word counts</b> for each section are met or are <b>no more than 15 words in excess</b>.</li> <li><b>Headers and paragraphs</b> are used effectively to increase readability and separate ideas for increased comprehension.</li> <li><b>No R code or R output (other than plots)</b> are displayed in the written proposal.</li> </ul>	All four criteria are met.	Only three criteria are met.	Only two or fewer criteria are met.
			<b>Total Points: /16</b>