

Team Project Instructions

Introduction

Team Project Instructions (document) Throughout this project, you will act out a scenario where you and your team are analysts retained by a company (large or small) or funding source (e.g., a VC firm or incubator) who have asked you to use data science techniques to address a chosen business problem. You and your team will design the data science task, analyze the data, and describe your results.

Your own data and results need not be on par with actual industry results—the goal is for you to get a realistic hands-on experience given the constraints and techniques you've learned. Don't worry too much about coming up with a novel idea; it's more important to develop the idea well (within the scope of what we've discussed in class).

You must choose a classification problem (we will talk about this in Module 2) and use the “data science process” to structure your research and write-up. Keep in mind that it may be ineffective to simply proceed linearly through the steps, and this may need to be reflected in your analysis.

You should interact with the instructor from the preparation of your initial ideas through your write-up, just as a consulting group would interact with a firm or funding source in preparing a report. Feel free to ask the instructor to help to fill in any gaps between the material available and what you would like to find out.

Schedule of Deliverables

Use this schedule of deliverables as a birds-eye view of what items will be due in each module.

Module	Deliverable
Module 1	Nothing due
Module 2	Dataset and Business Question
Module 3	Proposal
Module 4	Nothing due
Module 5	Nothing due
Module 6	One-page outline
Module 7	Nothing Due
Module 8	Report Presentation Supporting Data and Artifacts

Deliverables

Each deliverable will build upon the last, culminating in the final report and presentation. Please take time to review the instructions for each deliverable.

M2 Dataset and Business Question

Teams will select their data set and main business question via a Google document sent out by the instructor. You will not need to submit anything via Canvas, however there will still be work to do.

Possible datasets:

- [UCI Machine Learning Repository](#)
- [Kaggle](#)
- [KD nuggets](#)
- [R Bloggers](#)
- [FiveThirtyEight](#)

M3 Proposal

Teams will complete a proposal that addresses any feedback from the Dataset and Business Question deliverable and that also addresses the following:

- What are your business problem and business question?
- What is the data instance? (That is, what does each data point represent? A customer, a country, a product, etc.)
- What might be the target variable (the variable of interest to be predicted)?
- What is your proposed method to test or examine the problem?
- How will the results be used and how will they provide business value?

M6 Outline

Teams will submit an outline of their work so far. Please make sure to include feedback from the proposal submission and the following:

- Business problem and question
- Data understanding and prep process
- Identify target variable
- Work plan for the next two weeks

M8 Report

The write-up should be a maximum of 10 pages, double-spaced, plus any appendices you would like to include (where appropriate). Use external sources where appropriate with clear citations and a bibliography.

Your report should follow the rubric and contain all the steps in the data science process:

- Business understanding and business question
- Data understanding and data preparation, including identifying target variable, data visualization, some descriptive summary statistics, attribute understanding, etc.
- Modeling and results
- Model tuning and evaluation

- Discussion and limitations, including deployment related issues, the potential hazards, and bias your finding might result in if deployed
- All group members should contribute to the analysis and report.

M8 Presentation

Each team will record a 10–15-minute (no longer than 15-minute) presentation of their research to their employer.

All group members must contribute (have face-time) during the presentation.

M8 Supporting Data and Artifacts

This includes the PowerPoint deck used in the presentation, as well as data sources

Rubric

Assessment Criteria	Very Good	Good	Not Good Enough
Business Understanding (15 pts)	<ul style="list-style-type: none"> Thorough and clear discussion of the business question and objective. The target variable was appropriate for the business problem and was clearly defined to be readily used for analysis. 	<ul style="list-style-type: none"> Demonstrate good understanding of business problem and reasonably clear objective of the project. The target variable was appropriate for the business problem. 	<ul style="list-style-type: none"> The business question and objective are not clear. No or limited discussion of the background of the business problem. The target variable is not clearly defined and is irrelevant to the business problem.
Data Understanding & Visualization (20 pts)	<ul style="list-style-type: none"> Clear and effective description of data. Appropriate selection of data for the business problem at hand. The data collection and preparation procedures (and the sources) were clearly described. Creative and very effective data visualization that directly guides analysis. 	<ul style="list-style-type: none"> Clear and effective description of data and data collection procedure. Appropriate selection of data for the business problem at hand. The data collection and preparation procedures (and the sources) were clearly described. Effectively utilized data visualizations to describe the data and guide some analysis. 	<ul style="list-style-type: none"> The data are not clearly described. The data are not appropriate to address the business question being analyzed. The data collection and preparation procedures (and the sources) are not clearly described. No or very limited attempt to visualize data.
Modeling (15 pts)	<ul style="list-style-type: none"> Insightful and thorough analysis of all possible issues in the modeling stage. The team provided sufficient discussion on the choice of model (e.g., alternative model specifications, the pros and cons of their model). The application of algorithm and the interpretations of the results were accurate and clearly explained. 	<ul style="list-style-type: none"> The choice of model is well discussed. Reasonably complete analysis of the issues. The model is applied appropriately and the interpretation of the result is accurate. 	<ul style="list-style-type: none"> The team provided insufficient discussion on the choice of model (e.g., alternative model specifications, the pros and cons of their model). The model is applied inappropriately and the interpretation of results is flawed.

		<ul style="list-style-type: none"> The team provided sufficient discussion on the choice of model (e.g., alternative model specifications, the pros and cons of their model). The application of algorithm and the interpretations of the results were accurate and clearly explained. 	
Evaluation (15 pts)	<ul style="list-style-type: none"> The team clearly demonstrated the generalization performance of their model. (I.e., how was the model evaluated?) The analysis of expected benefit follows logical development and is supported by data. 	<ul style="list-style-type: none"> The team clearly demonstrated the generalization performance of their model. (I.e., how was the model evaluated?) The analysis of expected benefit follows logical development and is supported by data. 	<ul style="list-style-type: none"> No or very limited attempt to evaluate model performance. Poor flow of reasoning or logic.
Deployment (15 pts)	<ul style="list-style-type: none"> Clear and thorough demonstration of the use scenario of the result. Well-reasoned and thoughtful guidelines and recommendations for deployment. Recognize obstacles, challenges, and risks. Thoughtful discussion of the potential issues associated with the proposed plan. Comprehensive discussion on potential mitigation strategies. 	<ul style="list-style-type: none"> Reasonably complete demonstration of the use scenario of the result. Demonstrated recognition of potential issues associated with the proposed plan and provided potential mitigation strategies. 	<ul style="list-style-type: none"> Superficial, obvious, or inappropriate demonstration of the use scenario of the result. No or very limited awareness of potential issues. No or very limited offering of strategies to address issues.
Presentation/Report (20 points)	<ul style="list-style-type: none"> Exceptionally well organized and easy-to-follow structure. 	<ul style="list-style-type: none"> Reasonably well organized and easy-to-follow structure. The presenter made good use of time. 	<ul style="list-style-type: none"> Overall lack of organization and structure of the content. The presentation or report was unstructured and difficult to follow.

			<ul style="list-style-type: none">• The presentation was over time limit.
--	--	--	---

