

Part 1 Final Project Rubric

Deliverable: Jupyter Notebook or PDF Report

The Jupyter notebook should have a Project Topic and Data Cleaning/EDA section. If your work doesn't fit into one notebook (or you think it will be less readable by having one large notebook), make several notebooks or scripts in the GitHub repository and submit a report-style notebook or pdf instead.

Part of the project is to create a GitHub repository with your work. This repository needs to be specifically for this project. While Part 1 does not require you to submit the link, it would be a good idea to set up the repository and commit to it throughout the project.

Prompt	Points					
<u>Project Topic</u> Is there a clear explanation of what this project is about? Does it state clearly which type of problem? E.g. classification or regression	(0 pts) Not included in the project	(1 pts) States the type of problem (e.g. classification or regression). Missing explanation of what project is about.	(2 pts) Gives a clear explanation of what the project is about. Missing the type of problem (e.g. classification or regression).	(3 pts) Gives a clear explanation of what the project is about and clearly states the type of problem (e.g. classification or regression).		
<u>Project Topic</u> Is the goal of the project clearly stated? E.g. why it's important, what goal the author wants to achieve, or wants to learn.	(0 pts) Not included in the project	(1 pts) Needs improvement — attempts but doesn't get across the motivation or goal for the project	(2 pts) Very Good — clearly states the motivation or the goal for the project			
<u>Bonus (Optional Extra Credit)</u> <i>Is the project topic creative?</i> -- For the final project, it is a valid strategy to solve an existing online problem (e.g., go on the UCI Machine Learning Repository and work on the default task.) A bonus is available if you want to stretch and define your own project problem. For the bonus, you can use your own data or existing data (e.g., available online). If you are trying for the bonus, make sure to add a Bonus section in your Jupyter notebook/report and provide a brief explanation of why your topic meets the bonus criteria.	(0 pts) No extra credit attempted/did not define a creative project problem.	(2 pts) Includes a Bonus section in the Jupyter notebook/report with a brief description of how the learner defined their own creative project topic.				

<p>Data</p> <p>Is the data source properly cited and described? (including links, brief explanations)</p>	<p>(0 pts)</p> <p>Missing both of the following: Does not include a brief explanation of where the data is from/how it was gathered and does not include a citation (using the format of a style manual like APA) for either a public or unpublished dataset.</p>	<p>(2 pts)</p> <p>Includes one of the following: a brief explanation of where the data is from/how it was gathered or includes a citation (using the format of a style manual like APA) for either a public or unpublished dataset.</p>	<p>(4 pts)</p> <p>Includes both of the following: a brief explanation of where the data is from/how it was gathered and cites the dataset (using the format of a style manual like APA) for either a public or unpublished dataset.</p>			
<p>Data</p> <p>Is the data description explained properly? The data description should include the data size.</p> <ul style="list-style-type: none"> E.g. for tabulated data: number of samples/rows, number of features/columns, bytesize if a huge file, data type of each feature (or just a summary if too many features- e.g. 10 categorical, 20 numeric features), description of features (at least some key features if too many), whether the data is multi-table form or gathered from multiple data source. E.g. for images: you can include how many samples, number of channels (color or gray or more?) or modalities, image file format, whether images have the same dimension or not etc. E.g. sequential data: texts, sound file; please describe appropriate properties such as how many documents or words, how many sound files with typical length (are they the same or variable), etc. 	<p>(0 pts)</p> <p>Does not include any description of the data or the data size</p>	<p>(3 pts)</p> <p>Partially describes the data but does not refer to the data size or does not describe the data size appropriately for the type of data.</p>	<p>(6 pts)</p> <p>Describes the data including the data size appropriately for the type of data.</p>			

<p><u>Data Cleaning and EDA</u></p> <p><i>For Part 1, the learner should address at least one of the following questions.</i> For Part 2, we will have the same rubric (with increased points for this section) and expect the learner to show more effort on cleaning and EDA. For now, in Part 1, it is enough that they've tried something and started.:</p> <ol style="list-style-type: none"> Does it include clear explanations on how and why cleaning is performed? <ol style="list-style-type: none"> E.g. the author decided to drop a feature because it had too many NaN values and the data cannot be imputed. E.g. the author decided to impute certain values in a feature because the number of missing values were small and he/she was able to find similar samples OR, he/she used an average value or interpolated value, etc. E.g. the author removed some features because there are too many of them and they are not relevant to the problem, or he/she knows only a few certain features are important based on their domain knowledge judgment. E.g. the author removed a certain sample (row) or a value because it is an outlier. Does it have proper visualizations? (E.g. histogram, correlation matrix, etc.) Does it have adequate analysis? (E.g. analyzes visualizations, feature importance (if possible), etc.) Does it have conclusions or discussions? <ol style="list-style-type: none"> E.g. summary of data cleaning, findings, discussing foreseen difficulties and/or analysis strategy. 	<p>(0 pts)</p> <p>Doesn't adequately address at least one of the following as described in the rubric: includes clear explanations on how and why cleaning is performed or has proper visualizations or has adequate analysis or has conclusions/discussions.</p>	<p>(5 pts)</p> <p>Addresses at least one of the following as described in the rubric: includes clear explanations on how and why cleaning is performed or has proper visualizations or has adequate analysis or has conclusions/discussions.</p>				
--	---	--	--	--	--	--