



DEPARTMENT OF
SOFTWARE TECHNOLOGY

STINTSY

Machine Project

Major Details

Groupings:	At most 4 members in a group
Deadline:	December 1, 2023 (Friday) 6:00 PM
Demo Schedule:	December 4 to 7, 2022 (Week 14)
Percentage:	45%
Submission guidelines:	Submit the zip file to AnimoSpace
Filename format:	STINTSY-Project-<Section>-Group<#>.zip

Deliverables

Zip file containing:

- Jupyter Notebook file – ipynb file
- Other Python 3 files – py files
- Dataset files – csv files

Specifications

You are tasked to go through the process of selecting a dataset, describing the dataset, performing exploratory data analysis, data preprocessing and cleaning, model training, hyperparameter tuning, model selection, and extracting insights from the data.

The project is to be submitted as a Jupyter Notebook and, optionally, some Python 3 source files. The notebook should be a self-explanatory document containing a report of the entire process undertaken to come up with the generated insights from the raw dataset. It should contain markup cells explaining the processes undertaken in the project, as well as code cells showing all the code that was performed. Please make sure that the codes could be successfully run sequentially to replicate the processes done in the project. Offshoots (a new task different from the original task because the new task seems interesting) are also encouraged, but make sure that the original task has already been completed.

Outline for the Notebook

Section 1. Introduction to the problem/task and dataset

Each group should select one real-world dataset from the list of datasets provided for the project. Each dataset is accompanied with a description file, which also contains detailed description of each feature.

The target task (i.e., classification or regression) should be properly stated as well.

Section 2. Description of the dataset

In this section of the notebook, you must fulfill the following:

- State a brief description of the dataset.
- Provide a description of the collection process executed to build the dataset. Discuss the implications of the data collection method on the generated conclusions and insights. Note that you may need to look at relevant sources related to the dataset to acquire necessary information for this part of the project.
- Describe the structure of the dataset file.
 - What does each row and column represent?
 - How many instances are there in the dataset?
 - How many features are there in the dataset?
 - If the dataset is composed of different files that you will combine in the succeeding steps, describe the structure and the contents of each file.
- Discuss the features in each dataset file. What does each feature represent? All features, even those which are not used for the study, should be described to the reader. The purpose of each feature in the dataset should be clear to the reader of the notebook without having to go through an external link.

Section 3. List of requirements

List all the Python libraries and modules that you used.

Section 4. Data preprocessing and cleaning

Perform the necessary steps before using the data. In this section of the notebook, please take note of the following:

- If needed, perform preprocessing techniques to transform the data to the appropriate representation. This may include binning, log transformations, conversion to one-hot encoding, normalization, standardization, interpolation, truncation, and feature engineering, among others. There should be a correct and proper justification of the use of each preprocessing technique used in the project.
- Make sure that the data is clean, especially features that are used in the project. This may include checking for misrepresentations, checking the data type, dealing with missing data, dealing with duplicate data, and dealing with outliers, among others. There

should be a correct and proper justification of the application (or non-application) of each data cleaning method used in the project. Clean only the variables utilized in the study.

Section 5. Exploratory data analysis

Perform exploratory data analysis comprehensively to gain a good understanding of your dataset. In this section of the notebook, you must present relevant numerical summaries and visualizations. Make sure that each code is accompanied by a brief explanation. The whole process should be supported with verbose textual descriptions of your procedures and findings.

Section 6. Model training

Use machine learning models to accomplish your chosen task for the dataset. In this section of the notebook, please take note of the following:

- The project should train and evaluate at least 3 different kinds of machine learning models.
- Each model should be appropriate in accomplishing the chosen task for the dataset. There should be a clear and correct justification on the use of each machine learning model.
- Make sure that the values of the hyperparameters of each model are mentioned. At the minimum, the optimizer, the learning rate, and the learning rate schedule should be discussed per model.
- The report should show that the models are not overfitting nor underfitting.

Section 7. Hyperparameter tuning

Perform grid search or random search to tune the hyperparameters of each model. In this section of the notebook, please take note of the following:

- Make sure to elaborately explain the method of hyperparameter tuning.
- Explicitly mention the different hyperparameters and their range of values. Show the corresponding performance of each configuration.
- Report the performance of all models using appropriate evaluation metrics and visualizations.
- Properly interpret the result based on relevant evaluation metrics.

Section 8. Model selection

Present a summary of all model configurations. Include each algorithm and the best set of values for its hyperparameters. Identify the best model configuration and discuss its advantage over other configurations.

Section 9. Insights and conclusions

Clearly state your insights and conclusions from training a model on the data. Why did some models produce better results? Summarize your conclusions to explain the performance of the models. Discuss recommendations to improve the performance of the model.

Section 10. References

Cite relevant references that you used in your project.

Final Project Presentation

Here are some guidelines regarding the final project presentation:

- Each group is given 30 minutes: 20 minutes to present, and 10 minutes for Q&A.
- Presentations will be done via Zoom.
- Groups are required to have a presentation ready. Open your Jupyter notebook as well.
- Open all the necessary files before your allotted presentation time slot. Do not wait until the presentation itself to load anything.
- All members should be present and should discuss a part in the final project presentation.
- Kindly read the rubrics to check different requirements and expectations on the project presentation.

Working With Groupmates

For this project, you are encouraged to work in groups of at most 4 members. Make sure that each member of the group has approximately the same amount of contribution for the project. Problems with groupmates must be discussed internally within the group, and if needed, with the lecturer.

Deliverables

Submit a zip file containing the source code files via AnimoSpace. All exploratory data analysis, machine learning, and core algorithms should be performed using Python 3 code and integrated into the Jupyter Notebook. Other code that you used for the project other than those in the Notebook should also be included in the submission of the project.

Academic Honesty Policy

Honesty policy applies. Please take note that you are NOT allowed to borrow and/or copy-and-paste – in full or in part – any existing related program code or solutions from the internet or other sources (such as printed materials like books, or source codes by other people that are not online). You should develop your own codes and solutions from scratch by yourselves.

The student handbook states that (Sec. 5.2.4.2):

“Faculty members have the right to demand the presentation of a student’s ID, to give a grade of 0.0, and to deny admission to class of any student caught cheating under Sec. 5.3.1.1 to Sec. 5.3.1.1.6. The student should immediately be informed of his/her grade and barred from further attending his/her classes.”

The student handbook also states that (Sec. 10.3):

A student caught cheating, as defined in Sec. 5.3.1.1., shall be penalized with a grade of 0.0 in the requirement or in the course, at the discretion of the faculty member, without prejudice to an administrative sanction. In cases of alleged cheating, the faculty member should report the incident to the Student Discipline Formation Office (SDFO).

RUBRIC FOR GRADING

Criteria		Ratings			Points
Description of the dataset and the task	COMPLETE 5 pts	INCOMPLETE 2 pts		NO MARKS 0 pt	5 pts
	<p>An overview or description of the data is provided, including how it was collected, and its implications on the types of conclusions that could be made from the data. A description of the variables, observations, and/or structure of the data is provided.</p> <p>The target task is well introduced and clearly defined.</p>	<p>An overview or description of the data is provided but lacks details. A description of variables, observations, and/or structure is present but is missing for some aspects of the dataset.</p> <p>The task is not clearly defined.</p>		<p>No overview or description of the data is provided.</p> <p>No description of variables, observations, and/or structure is provided.</p> <p>The task is not defined.</p>	
Exploratory data analysis	COMPLETE 10 pts	INCOMPLETE 7 pts	INCOMPLETE 4 pts	NO MARKS 0 pt	10 pts
	<p>The data is sufficiently explored to get a grasp of the distribution and the content of the data. Appropriate summaries and visualizations are presented. Insights into how the EDA can help the model training is mentioned.</p>	<p>Exploratory data analysis is not sufficiently performed. Summaries and visualizations are presented but have minor issues in terms of methods chosen.</p>	<p>Exploratory data analysis is rudimentary. Inappropriate methods of summarizing and visualizing data are frequently chosen.</p>	<p>No exploratory data analysis is attempted.</p>	
Data pre-processing and cleaning	COMPLETE 10 pts	INCOMPLETE 7 pts	INCOMPLETE 4 pts	NO MARKS 0 pt	10 pts
	<p>The necessary steps for pre-processing and cleaning are performed, including explanations for every step for each feature. If no preprocessing or cleaning is done, there should be a justification on why it is not needed.</p>	<p>Pre-processing and cleaning steps are performed but lacks explanation. Or, pre-processing and cleaning done are insufficient for less than half or half of the number of features.</p>	<p>Pre-processing steps do not match the ML model chosen. Or, pre-processing and cleaning done are insufficient for more than half of the number of features.</p>	<p>No pre-processing and cleaning are done, and no justification was provided as to why it was not done, or the justification is weak or incorrect.</p>	

Model training	COMPLETE 10 pts Appropriate models are used to accomplish the machine learning task. Justification of choosing the models is discussed. All models are trained correctly. The report shows that models are not overfitting nor underfitting.	INCOMPLETE 7 pts A lot of various models are used without proper justification of why they are chosen. Or, models are not trained correctly. The report shows no evidence proving that the model is not overfitting nor underfitting.	INCOMPLETE 4 pts Only one model is generated. Or all models that are chosen are not appropriate for the task.	NO MARKS 0 pt No model training is done.	10 pts
Model selection and hyperparameter tuning	COMPLETE 15 pts Appropriate data-driven error analysis is made, and changes to the model selection and hyperparameters are performed to improve model performance. The study exhausts improvements that can be done to the model.	INCOMPLETE 10 pts Model selection and hyperparameter tuning is done exhaustively but without proper justification or analysis. Or improvements to the models are not exhausted.	INCOMPLETE 5 pts Model selection and hyperparameter tuning is done, but no efforts to further improve the model are done.	NO MARKS 0 pt No model selection and hyperparameter tuning are done.	15 pts
Results	COMPLETE 10 pts Multiple appropriate evaluation metrics and visualizations are used to report the performance of all models. Results are correctly interpreted.	INCOMPLETE 5 pts Incorrect evaluation metric or visualization is used to report the performance of at least one models. Or, results are incorrectly interpreted.		NO MARKS 0 pt No evaluation metric nor visualization is used to report the performance of the models.	10 pts
Insights and conclusions	COMPLETE 5 pts The study is concluded by effectively summarizing the efforts of the authors. Recommendations on how the model could be further improved are provided.	INCOMPLETE 2 pts The study is concluded but misses key insights performed in the study. Or recommendations on how the model could be further improved are provided without clear justification.		NO MARKS 0 pt No insights or conclusions are presented.	5 pts

Notebook	COMPLETE 5 pts The report discusses all steps in the machine learning process.	INCOMPLETE 2 pts The report discusses some steps in the machine learning process.	NO MARKS 0 pt No steps are discussed in the report.	5 pts	
Presentation manner	COMPLETE 5 pts The presenter seldomly looks at notes. The presenter helps the audience visualize through gestures and movements. The presenter displays a relaxed, self-confident nature about self, with no mistakes.	INCOMPLETE 2 pts The presenter looks at his notes most of the time. The presenter uses very little movement or descriptive gestures. The presenter displays mild tension; has trouble recovering from mistakes.	NO MARKS 0 pt The presenter reads the entire report from his notes. The presenter does not use movement or descriptive gestures. The presenter displays tension and nervousness; has trouble recovering from mistakes.	5 pts	
Presentation organization	COMPLETE 5 pts Information is presented in a logical and interesting sequence which the audience can follow.	INCOMPLETE 2 pts Audience has difficulty following the presentation because the presenter jumps around different topics.	NO MARKS 0 pt Audience cannot understand the presentation because there is no logical sequence of information.	5 pts	
Presentation Q&A	COMPLETE 20 pts The group convincingly answered all questions about both the code and the machine learning process.	INCOMPLETE 14 pts The group convincingly answered more than half or half of the number of questions about both the code and the machine learning process.	INCOMPLETE 6 pts The group convincingly answered less than half of the number of questions about both the code and the machine learning process.	NO MARKS 0 pt The group failed to answer any question about the code and the machine learning process.	20 pts
Total points:					100

Note: Each member of the group is expected to have a good understanding of the group's submission, even the parts that were not directly delegated to them. Failure to answer the questions during the demo, in a such a way that suggests that one or more group members did not sufficiently understand the work that was delivered, will result in a grade of 0 for those members for the entire project.