

Final Project

DS 4021 – Instructor: Javier Rasero; **TAs:** Jeffrey Zhang, Ethan Meidinger

Due: See Canvas

Submission format: GitHub repository (submitted by link in Canvas)

Group Assignment

General Description: Submit to Canvas the link to your group's GitHub repository for this project. This link should be **the same** that you created for the prerequisite assignment, and will contain all materials related to your final project, including Jupyter notebooks, documentation, figures, and a final report summarizing the results.

We will use GitHub's commit history to verify submission times. Any commits after the deadline will be considered late, and the final project **will receive no credit**.

Why am I doing this? The goal of the final project is to put into practice everything you have learned throughout this course. Using the same dataset and research question that you selected previously, you will perform exploratory visualization, descriptive statistics, and the training and optimization of four predictive models: a linear model with penalization, a Support Vector Machine, a neural network implemented in PyTorch, and an ensemble method. Proper validation through cross-validation, model comparison, and final evaluation on the held-out test set are required.

What will we do? Your group will complete a full machine learning analysis on your selected dataset and write a concise report presenting the main results and findings. Work should be divided among team members.

How will this be graded? This final project is worth 100 points.

- The notebooks portion is worth 80 points.
- The final report is worth 20 points.

See the table below for a detailed description of where the points come from.

Important: Although only the notebooks and the final report count toward the grade, your repository must include all required files and folders **outlined below** in order to receive full credit.

Formatting	<ul style="list-style-type: none">• One Github Repository (submitted via link on canvas)• The top level page of the repository should contain:<ul style="list-style-type: none">○ A README.md file (which auto displays)○ A NOTEBOOKS folder○ A Data Folder○ A OUTPUT folder○ A final_report.pdf file
README.md	<ul style="list-style-type: none">• <u>Goal:</u> The README helps others quickly understand your project. Keep it short, clear, and written for someone who has never seen your files before.• Use markdown headers to divide content.• Make an H2 (##) section explaining the contents of the repository <p>It must include two simple sections:</p>

	<ul style="list-style-type: none"> Section 1: Software and platform section List the main tools needed to run your notebooks. For example: <ul style="list-style-type: none"> Python version Jupyter notebooks Packages (pandas, numpy, scikit-learn, matplotlib, seaborn, PyTorch) Operating system (Mac, Windows, or Linux) Section 2: A Map of your documentation. In this section, you should provide an outline or tree illustrating the hierarchy of folders and subfolders contained in your Project Folder, and listing the files stored in each folder or subfolder.
NOTEBOOKS folder	<ul style="list-style-type: none"> <u>Goal:</u> This folder contains all the jupyter notebooks for your project. <p>These are the required notebooks you will need to upload:</p> <ol style="list-style-type: none"> Descriptive Analysis Notebook (10 points): <ul style="list-style-type: none"> Exploratory visualizations Summary statistics Models optimization and training Notebook (s) (60 points) Must include complete training and optimization of: <ul style="list-style-type: none"> A Penalized (Ridge, Lasso or ElasticNet) linear model (Linear Regression or Logistic Regression). Support Vector Machine Ensemble model (e.g. Random Forest or Gradient Boosting) Neural network implemented in PyTorch <p>You may use one combined notebook or separate notebooks for each model.</p> <p><u>REMEMBER:</u> For the optimization and training stage, you must not use the test set you put aside in the prerequisite Final Project assignment.</p> <p>Using the test set for tuning or training will result in major penalties.</p> <ol style="list-style-type: none"> Final Test-Set Evaluation Notebook (10 points) <ul style="list-style-type: none"> Select the best model (and by best, this also includes its optimal hyperparameters) from the optimization phase and fit it to the entire training set. Evaluate it on the test set you had initially put aside. Report final results. <p><u>REMEMBER:</u> Here is where you must use the test set you put aside in the prerequisite Final Project assignment.</p> <p>Not evaluating on the correct test set will result in zero points for this part.</p> <p><u>Notebook Requirements</u></p> <ul style="list-style-type: none"> Use cross-validation to tune your models. Try different hyperparameter combinations.

	<ul style="list-style-type: none"> • Embed preprocessing inside the cross-validation to avoid any data leakage. You may accomplish this using pipelines. Failing to do this will incur in penalties. • Assess and report the final model's performance using more than one metric. e.g. Accuracy and confusion matrix for <i>classification</i>, Coefficient of determination and Mean Square Error for <i>regression</i>. • Cross-validation and hyperparameter tuning must be clearly documented. • Notebooks must run top-to-bottom without errors. • Code must be accompanied by explanatory text. • Figures must be readable and well labeled. • Random seeds must be set whenever possible for reproducibility. Scikit-learn objects typically include random_state as an argument for this. For Pytorch, you may use <code>torch.manual_seed</code>.
Data folder	<ul style="list-style-type: none"> • This is the same folder where you originally uploaded the test set. • For the final project, this folder should also contain the training set. • If any of these pieces of data are too large for GitHub, upload it to a platform such as Google Drive and include a link with appropriate download permissions.
OUTPUT folder	<ul style="list-style-type: none"> • <u>Goal:</u> Store all outputs generated during the project. For example: Figures, Tables, any materials used for your final presentation, etc.
final_report.pdf	<ul style="list-style-type: none"> • <u>Goal:</u> to summarize everything into a final report file. • A PDF file 1 or 2 pages long. • 20 points towards the final grade. <p>Structure</p> <ul style="list-style-type: none"> • One introductory paragraph with <ul style="list-style-type: none"> ◦ Research question. You may reinstate or refine the same research question you used for the prerequisite document. Remember that your research question needs to be testable. For example, something like: “<i>which factors contribute most to...</i>”, “<i>which features are the most predictive...</i>” might now be directly addressable for some models (e.g. Neural Networks or SVM with non-linear Kernel). A good and testable research question could just be: “<i>Can we successfully predict [SOMETHING]?</i>” “<i>What model best predicts [SOMETHING]?</i>” ◦ Motivation • One or two methods paragraphs <ul style="list-style-type: none"> ◦ Dataset description. ◦ Summary of any data cleaning/curation you needed to perform. ◦ Outline each model used, including the preprocessing steps each needed to include. ◦ Cross-validation strategy. ◦ Metrics used to assess model performance. • One–two results paragraphs <ul style="list-style-type: none"> ◦ Key findings and model results. ◦ Include one table showing the cross-validated performance of each best model ◦ Report the performance for the selected best model on the test set, i.e. the piece of data you prepare for this at the very beginning of the project.

- | | |
|--|--|
| | <ul style="list-style-type: none">● One–two discussion paragraphs<ul style="list-style-type: none">○ Interpretation of results○ Why some models may perform better than others○ Limitations of the project |
|--|--|