

# DS-GA 1007: Programming for Data Science



Center for  
Data Science

## Final Project Guidelines

The final project comprises three steps, as detailed below. Throughout the project, groups will have to work collaboratively to manage expectations and meet goals. As discussed in class, the version control system [Git/GitHub](#) will be useful for the coding component and the typesetting platform [Overleaf](#) will be useful for the reporting component.

Groups should contain 1, 2, or 3 members. If you would like to be assigned at random to a group, then please contact the instructors. If you would like to determine your group, then please post to Forums under the *Project* thread to contact classmates.

**October 31<sup>st</sup>:** Project Proposal  
**November 28<sup>th</sup>:** Project Milestone  
**December 14<sup>th</sup>:** Project Report

### October 31<sup>st</sup>: Project Proposal

By October 31<sup>th</sup> groups must upload a one page pdf file on Gradescope containing:

- Title
- Summary of Plans
  - Description of Problem
  - General Approach
  - Suggested Experiments
- Group
  - Name and NetID of each member.
  - Member responsible for uploading submissions.

NOTE: Only the *member responsible for uploads* needs to upload the pdf file. In other words, each group should have only one pdf file uploaded on Gradescope.

### November 28<sup>th</sup>: Project Milestone

By November 28<sup>th</sup> groups must upload on Gradescope a two page pdf file containing:

- Title

- Group Members
  - Name and NetID of each member.
  - Member responsible for uploading submissions.
- Background
  - Description of Problem
  - Motivation for Problem
  - References
- Plans
  - Description of Methodology
  - Proposed Experiments
  - Some Relevant Datasets

NOTE: Only the *member responsible for uploads* needs to upload the pdf file. In other words, each group should have only one pdf file uploaded on Gradescope.

### **December 14<sup>th</sup>: Project Report**

Groups will not be responsible for a presentation. By December 14<sup>th</sup> groups must upload a notebook (.ipynb file format) and pdf on Gradescope.

The notebook should describe the problem, the methodology and experiments used to understand the problem, evaluation of results, and possible next steps. More specifically, the notebook should be structured as follows:

1. Title
2. Group Members
  - a. Name and NetID of each member.
  - b. Member responsible for uploading submissions
3. Abstract
4. Background
  - a. Description of Problem
  - b. Motivation for Problem
  - c. References
5. Results
  - a. Description of Methodology
  - b. Experiments Conducted
  - c. Description of Datasets
6. Discussion
  - a. Evaluation of Findings
  - b. Possible Next Steps

See the template below for more information.

The pdf should summarize the notebook. The summary should motivate the problem, explain some aspects of the approach and implementation, and describe the outcomes of the experiments. The pdf should be limited to four pages in [bulletin format](#). Groups can share their summaries with the class by electing to upload pdf's to <https://wp.nyu.edu/pdsf19/>

NOTE: Only the *member responsible for uploads* needs to upload the pdf file. In other words, each group should have only one pdf file uploaded on Gradescope.

### **Final Project Evaluation:**

The final project will be graded based on three main aspects:

1. adherence to guidelines
2. quality of the report
3. implementation of the code

While projects will not be assessed on the technicality of the problem, we will recognize efforts regarding

1. size and “cleanliness” of the datasets
2. sophistication of the algorithms
3. relevance of the problem to applications

A final report should:

1. clearly state the problem, pointing which are the hurdles and issues to solve it;
2. clearly present the methodology employed to solve the problem, pointing out:
  - a. the data sets used
  - b. the methods employed to (if necessary) handle missing data, transform data, combine data, etc.
  - c. the algorithms involved in the solution, as for example, SVM for classification, DBScan for clustering, etc.
  - d. present and discuss the results, highlighting the strengths and weaknesses of the proposed methodology
  - e. make some conclusion, emphasizing whether the chosen approach was success and, if not, why.

# **NOTEBOOK TEMPLATE**

**Title:** Project Title Here

**Authors:**

Name1, NetID1  
Name2, NetID2  
Name3, NetID3  
Name4, NetID4

---

**Abstract:** This project focus on ... we approached it using ... and the results shows that our approach is a good alternative.

## **1. Introduction and Motivation**

Presentation of the problem, its importance, and which are the difficulties involved on it.

## **2. Methodology**

- The data sets involved and how they were “cleaned”.
- Mathematical and computational methods employed to solve the problem.
- Particular design decisions that you deem important when handling the problem.

## **3. Results**

Description of the results, pointing how well the problem has been solved. Figures showing the results are, in general, better than tables.

## **4. Discussion**

Which are the strengths and weaknesses of the proposed methodology. Are the results good enough? Can they be improved? Which are the limitations?

## **5. Conclusion**

Summary of the problem, findings and limitations. Future work directions.

## **References:**

[Joia et al. 2011] Joia, P., Coimbra, D., Cuminato, J. A., Paulovich, F. V., and Nonato, L. G Local affine multidimensional projection. *IEEE Transactions on Visualization and Computer Graphics*, 17(12):2563-2571, 2011.