

Introduction to Machine Learning, Spring 2023

Course Project

(Due Sunday, Jun. 11 at 11:59pm (CST))

Instructor: Lu Sun

April 13, 2023

One of main goals of CS182 is to prepare you to utilize machine learning techniques to solve real-world problems, and the course project provides a good opportunity for you to start in this direction.

1 Guidelines

As a part of evaluation of CS182, you are required to complete a course project based on this material. Your project must be closely related to what you have learnt in this course, and you cannot use the ideas or results developed in previous semesters or other courses. Also, you should not submit a project that is largely collaborated with people outside this course.

In short, a typical project consists of picking an interesting dataset or application, applying one or more well-known machine learning algorithms as baselines, and extending these baselines in creative and innovative ways. The general guidelines are listed as follows.

- Projects should be completed in groups, each of which is composed of **1-2** students (we recommend groups of 2 students). Only one group member is supposed to submit the project, tag the rest of group members, and make sure the member-specific contributions.
- Each project consists of three major parts: presentation with slides, one final writeup and the source code.
 - **Presentation with slides:** Each group should prepare and give an in-class presentation (**5 min**) with slides. The presentation should focus on your high-level idea, and leave any technical details to the writeup. (30 points)
 - **Final writeup:** You are expected to submit a final writeup summarizing your findings, ideas and contributions. Final writeup must be in the form of a NeurIPS paper. and **no longer than 5 pages** in total: 4 pages for the body of the writeup, and an additional page for references. Only the hard copy in PDF, rather than the source latex code, should be submitted. (60 points)
 - **Source code:** For the sake of convenience, it is highly recommended to use Python to implement your ideas and algorithms in the project. However, any other programming languages are allowed. It is your responsibility to make sure the source code is executable and contains no severe bugs. Please submit the code in a **zip** file, which should not exceed 5MB. (10 points)

- All the projects must be submitted before **11:59pm (CST), June 11th, 2023**. There will be **no late days** for the final writeup. Note that this course project is **30%** of the final grade.

2 Project Types

Basically speaking, there are three types of projects:

1. **Application project.** This is the easiest and most common project. Select one application that interests you, and explore how best to apply existing learning algorithms to solve it. If you want to choose this project type, please make sure the application is somewhat new, and compare a sufficient number (≥ 5) of cutting-edge algorithms in your projects.
2. **Algorithmic project.** It aims to solve a problem or a family of problems by developing a new learning algorithm, or a novel variant of an existing algorithm. The proposed algorithm is expected to be comparable with state-of-the-art algorithms on some datasets or specific problem settings, and should not be the exactly same with any existing algorithms.
3. **Theoretical project.** This is purely theoretical, and is usually the most difficult project. Prove some interesting or non-trivial properties of a new or an existing algorithm. If you decide to challenge this type of project, please let me know before you start your project.

Of course, it will be also fine if you would like to complete your project by combining the three elements of applications, algorithms and theoretical analysis.

3 Project Topics

The first task for you is to pick one interesting project topic. If you are looking for topics, or would like to discuss them, please feel free to talk with TAs and/or me. We are happy to give you some advice and suggest some project ideas.

There are many avenues that you may pursue for this project, and we encourage you to be brave and creative even if you don't think you'll necessarily get "good" results. Here are some preliminary ideas¹:

- Extend classical supervised learning algorithms in the setting of semi-supervised or active or reinforcement learning.
- Develop algorithms that stores both global/linear and local/non-linear information during learning.
- Propose methods that overcome the limitations of existing methods in terms of classification accuracy or computational complexity or theoretical analysis.
- Extend existing binary classification/regression models to handle multi-class or multi-task or multi-label or multi-view/modal/source or multi-instance problems.
- Propose a variant of top methods to address real-world problems in modern applications, such as missing feature values, imbalanced labels, large-scale sample space, high-dimensional feature space, and so on.

¹Note that this list is by no means comprehensive, and you can pick any topic that is related to our course and interests you

- Solve clustering or dimensionality reduction problems by revising current unsupervised methods to handle semi-supervised or supervised applications.
- Design a new deep learning framework to solve a specific, practical problem.

In addition to the ideas listed above, you might also refer to some recent machine learning research papers. Two top-tier conferences in machine learning are ICML, NeurIPS and ICLR.

4 After CS182

An excellent CS182 project will be publishable or nearly-publishable piece of work. After completing CS182, if you would like to continue working on your project along this direction as your graduation project, or submit your work to a machine learning (or other appropriate non-machine learning) conference or journal, please feel free to talk with me. I am happy to give you some guidance and support your further work.

A Some Suggestions for Technical Writeups

The final writeup is an important part of the project, as well as the course learning exercise. Please spend enough time on completing your writeup so that it is well motivated, precisely described, and clearly presented. Typically, a technical report/paper comprises of four sections: introduction, methodology, experiment and conclusion. When writing these sections, it is better for you to keep the following ideas in mind.

- The **introduction** section should set up the problem (e.g., why is it interesting? important?) and provide some context for what work has been done in the past (e.g., what is known? what is the open problem? what are the deficiencies of existing methods?).
- The **methodology** and **experiment** sections should describe clearly what you did (e.g., enough details for someone to re-implement your algorithms, if needed.), and provide some summary tables or figures that illustrate your experimental results, along with some descriptive interpretations.
- The **conclusion** section should concisely summarize your answers for the following questions:
 - What is your problem?
 - What have you done?
 - Why did it work well? (or why didn't?)
 - What might you do if you want to continue working on this project?

B Some Data Repositories

- Kaggle Datasets
- UCL ML Repository
- UCL KDD Repository
- CIFAR Dataset

- Caltech 101 Dataset
- NUS-Wide Dataset
- MirFlickr Dataset
- Amazon Dataset
- MovieLens Dataset
- Multi-Label Repository