

## Project 1: Midterm-Project

*Instructor: Yuan Yao**Due: 23:59 Saturday 18 Oct, 2025*

## 1 Project Requirement

This project as a warm-up aims to explore basic techniques in machine learning.

1. Pick up ONE (or more if you like) favourite dataset below to work. If you would like to work on a different problem outside the candidates we proposed, please email course instructor about your proposal.
2. Team work: we encourage you to form small team, up to **FIVE** persons per group, to work on the same problem. Each team just submit **ONE** report, *with a clear remark on each person's contribution*. The report can be in the format of either a *poster*, e.g.

`https://github.com/yuany-pku/2017\_math6380/blob/master/project1/DongLoXia\_poster.pptx`

or *technical report within 8 pages*, e.g. NIPS conference style (preferred format)

`https://nips.cc/Conferences/2019/PaperInformation/StyleFiles,`

with source codes such as Python (Jupyter) Notebooks with a detailed documentation.

3. For Kaggle contests, if please register your team with name in the format of `math5470_lastname`, so that we could easily find your results on Kaggle. For example, a team with Shawn Zhu and Kate Wong would be named by `math5470_Zhu_Wong`.
4. In the report, show your proposed scientific questions to explore and main results with a careful analysis supporting the results toward answering your problems. If possible, you should include your Kaggle contest score or rating in the report. Remember: scientific analysis and reasoning are more important than merely the performance tables. Separate source codes may be submitted through email as a GitHub link, or a zip file.
5. Submit your report by email or paper version no later than the deadline, to the following address (`datascience.hw@gmail.com`) with a title "MATH5470: Project 1"

## 2 Kaggle Contest: Home Credit Default Risk

Many people struggle to get loans due to insufficient or non-existent credit histories. And, unfortunately, this population is often taken advantage of by untrustworthy lenders.

Home Credit strives to broaden financial inclusion for the unbanked population by providing a positive and safe borrowing experience. In order to make sure this underserved population has a positive loan experience, Home Credit makes use of a variety of alternative data—including telco and transactional information—to predict their clients' repayment abilities.

While Home Credit is currently using various statistical and machine learning methods to make these predictions, they're challenging Kagglers to help them unlock the full potential of their data. Doing so will ensure that clients capable of repayment are not rejected and that loans are given with a principal, maturity, and repayment calendar that will empower their clients to be successful.

Visit the following website to join the competition.

<https://www.kaggle.com/c/home-credit-default-risk/>

**Requirements.** For Kaggle contests, if possible please register your team with name in the format of math5470.lastname, so that we could easily find your results on Kaggle. For example, a team with Shawn Zhu and Kate Wong would be named by math5470\_Zhu\_Wong.

## 3 Kaggle Contest: M5 Forecasting

Consider the following M5 forecasting challenge:

- Accuracy, Estimate the unit sales of Walmart retail goods. Can you estimate, as precisely as possible, the point forecasts of the unit sales of various products sold in the USA by Walmart?  
<https://www.kaggle.com/c/m5-forecasting-accuracy>

How much camping gear will one store sell each month in a year? To the uninitiated, calculating sales at this level may seem as difficult as predicting the weather. Both types of forecasting rely on science and historical data. While a wrong weather forecast may result in you carrying around an umbrella on a sunny day, inaccurate business forecasts could result in actual or opportunity losses. In this competition, in addition to traditional forecasting methods you're also challenged to use machine learning to improve forecast accuracy.

In this competition, you will use hierarchical sales data from Walmart, the world's largest company by revenue, to forecast daily sales for the next 28 days and to make uncertainty estimates for these forecasts. The data, covers stores in three US States (California, Texas, and Wisconsin) and includes item level, department, product categories, and store details. In addition, it has explanatory variables such as price, promotions, day of the week, and special events. Together, this robust dataset can be used to improve forecasting accuracy.