

Mini-Project 1. Machine Learning Basics

*Instructor: Yuan Yao**Due: 23:59 Sunday 4 Oct, 2018*

1 Mini-Project Requirement and Datasets

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 **THREE** 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. 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.
4. Submit your report by email or paper version no later than the deadline, to the following address (datascience.hw@gmail.com) with a title "MATH4995: Project 1"

2 Kaggle Contest: Predict Survival on the Titanic

The sinking of the Titanic is one of the most infamous shipwrecks in history.

On April 15, 1912, during her maiden voyage, the widely considered “unsinkable” RMS Titanic sank after colliding with an iceberg. Unfortunately, there weren’t enough lifeboats for everyone onboard, resulting in the death of 1502 out of 2224 passengers and crew.

While there was some element of luck involved in surviving, it seems some groups of people were more likely to survive than others.

In this challenge, we ask you to build a predictive model that answers the question: “what sorts of people were more likely to survive?” using passenger data (*i.e.* name, age, gender, socio-economic class, etc). Visit the following website to join the Kaggle contest:

<https://www.kaggle.com/c/titanic>

3 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/>