# DSC 190 – Intro to DM: Course Project Choices

**Jingbo Shang** (jshang@ucsd.edu)

CSE & HDSI, UCSD

Jan 20, 2022

# Project Choice 1 – Overview

❑ Team: 1 to 4 people
  ❑ We will have a really high expectation to teams of 4 people → Nearly a top conference submission quality
❑ Report
  ❑ At least 4 pages
    ❑ Double-column, 11 pt
    ❑ Roughly 2.5-3 thousand words + figures, tables, and equations
❑ Code
  ❑ A GitHub repo
  ❑ Working demo (bonus points up to 5%)
    ❑ E.g., simple UI + your model: TAs or other students can give their own inputs and check the prediction

# Report Template

- ❑ ACM Proceedings Templates
  - ❑ https://www.acm.org/publications/proceedings-template
  - ❑ Overleaf: https://www.overleaf.com/latex/templates/acm-conference-proceedings-master-template/pnrfvrrdbfwt

# Five Components

- ❑ Dataset (5%)
- ❑ Predictive Task (5%)
- ❑ Model (5%)
- ❑ Literature (5%)
- ❑ Results (5%)

# Dataset

❑ Identify a dataset
❑ Perform an exploratory data analysis
  ❑ Basic Statistics
  ❑ Properties
  ❑ Interesting findings
❑ All these should motivate your model design/choice

❑ The dataset should be large enough (e.g., 50,000 instances in total)

# Dataset – Example

- ❑ Heroes of the Storms (HOTS)
    - ❑ A 5v5 online video game

- ❑ It has massive log data available!
    - ❑ https://www.hotslogs.com/info/API
    - ❑ API & 30-day logs download

- ❑ EDA
    - ❑ Any frequent combinations of heroes?
    - ❑ Which hero is the "weakest"?
    - ❑ …
    - ❑ Some analysis: https://github.com/shangjingbo1226/HOTS-Analysis/blob/master/combo-analysis.ipynb

# Predictive Task

❑ Identify a predictive task based on your dataset

❑ How will you evaluate different models in this task?

❑ What are the baseline models you want to compare with?

   ❑ Why do you think they are appropriate?

   ❑ Why do you think your model can outperform them? Or say, what are their drawbacks?

# Predictive Task – Example

❑ Heroes of the Storms (HOTS)
  ❑ A 5v5 online video game

❑ It has massive log data available!
  ❑ https://www.hotslogs.com/info/API
  ❑ API & 30-day logs download

❑ Given the hero picks of two teams (of similar levels), which team will win?
  ❑ A classification problem!

# Predictive Task – Example

- Given the hero picks of two teams (of similar levels), which team will win?
  - A classification problem!

- Evaluation
  - Accuracy, F1

- Baselines
  - Logistic regression: Assume it's a linear combination of selected heroes
  - Naïve Bayes: Assume the selected heroes are conditionally mutual independent to each other

# Model

- ❑ What is the model that you propose to attack this task?
  - ❑ It's fine to use models that were described in class here
  - ❑ Explain and justify your choice/proposal What are the features you designed for your model?
  - ❑ Any unsuccessful tries?
- ❑ How will you optimize your model?
  - ❑ It's fine here to call any 3$^{rd}$-party libs
- ❑ Did you encounter any troubles?
  - ❑ Scalability? Overfitting?

# Model – Example

❑ How can we represent the input?
   ❑ One hot encoding?
   ❑ Heroes have no order?
   ❑ Meta-data (e.g., hero types, user ratings, …)?
   ❑ Battlegrounds (Maps)?
   ❑ …

❑ Application?
   ❑ Real-time demo
   ❑ Retrieve all users' info (rating, history, …)
   ❑ Suggest some picks
   ❑ How to update when there is a new hero or a rework?

# Literature

❑ Has your dataset/task been studied by others before?

❑ How the dataset was used?

❑ Are you working on a brand-new task?

❑ How are other people attacking the same/similar tasks?

❑ What is state-of-the-art method in this task or related tasks?

❑ Are your conclusions similar or different from existing work?

❑ What's the major novelty of your work?

❑ …

# Literature – Example

❑ There are some apps doing this task already
    ❑ Looking into that and check their performance, if possible

❑ There might be similar tasks (e.g., DoTA AI?)
    ❑ What are their solutions?

❑ …

# Results

❑ Does your proposed method outperform the baselines?
   ❑ Why your model can outperform?
   ❑ Or why your model fails?
❑ Whether the gap is significant?
❑ Are all features you designed effective?
❑ How shall one set the hyper-parameters of your model?
❑ What are the major takeaways (i.e., conclusions)?
❑ …

# Results – Example

- ❑ Performance comparison different methods
  - ❑ Baselines + Your proposed model
- ❑ Ablation study
  - ❑ What if some of the features/designs of your proposed model degenerates?
- ❑ Case Study
  - ❑ Some interesting cases when your model performs very well/poor
- ❑ Parameter Sensitivity
  - ❑ How do you decide hyperparameters?
  - ❑ Is the result sensitive to these hyperparameters?
- ❑ …

# Compare with Project Choice 2

- ❑ Choice 2 is Individual
- ❑ Implement ~4 models learned from this course from scratch.
  - ❑ Skeleton codes will be provided. Your work is more like "filling in blanks"
- ❑ Write a report (about 3~4 pages) describing your interesting findings and takeaways.
  - ❑ 11 pt, single column should be enough
  - ❑ Slightly less writing

# Project Choice 2 - Model Choices

❑ Model Choices:
  ❑ Linear Regression (1 Point)
  ❑ Logistic Regression (1 Point)
  ❑ Naïve Bayes (1 Point)
  ❑ K-means (1 Point)
  ❑ Gaussian Mixture (2 Points)
  ❑ Decision Tree + Random Forest (4 Points)
  ❑ Plain Matrix Factorization (4 Points)

❑ Choose any combination s.t. total points >= 6
  ❑ More points? ➔ Up to 5% bonus!
❑ 4 points ➔ 1 page; Otherwise, half page each

# Make your choice!  Q&A

❑ We have released all models

❑ Make your choice by the end of **Jan 31**

   ❑ We will create a multiple-choice question in Canvas or Gradescope