# CIS 530 Advanced Data Mining

## Course Project Instructions

Shuwen Wang, Computer & Info Science Department

# Project Choice 1 - Overview (100 points)

- ❏ Team: 1 to 4 people
  - ❏ We will have a really high expectation to teams of 4 people → Nearly a top conference submission quality
- ❏ Presentation: In-person on April 23, April 28, April 30 (if needed)
  - ❏ No presentation, No credit
  - ❏ Each team will have 10 minutes to present
- ❏ Report Due **11:59pm May 5** No report, No credit
  - ❏ At least 4 pages
    - ❏ Double-column, 11 pt
    - ❏ Roughly 2.5-3 thousand words + figures, tables, and equations
- ❏ Code Due **11:59pm May 5** No code, No credit
  - ❏ A GitHub repo
  - ❏ Or A PDF

# 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
- ❑ Predictive Task
- ❑ Model
- ❑ Literature
- ❑ Results

# 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 (at least 10,000 instances in total)

# Dataset – Example

❑ Grammy Winners and Nominees from 1965 to 2024
  ❑ More than 25,000 records

❑ This dataset compiles historical Grammy Award winners and nominees, including details such as the award category, artist, song or album, and year of recognition. The data was sourced from grammy.com and can be used for music trend analysis, industry research, and historical insights.

❑ EDA
  ❑ Most awarded categories?
  ❑ Top Grammy Award winners?
  ❑ …

# 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

❑ Grammy Winners and Nominees from 1965 to 2024
  ❑ More than 25,000 records

❑ Whether a nomination is a winner or not
  ❑ A classification problem!

# Predictive Task – Example

❑ Whether a nomination is a winner or not

❑ Evaluation
  ❑ Accuracy, F1, AUC

❑ Baseline
  ❑ Logistic regression: Assume it's a linear combination of selected features

# 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?

# 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?

❑ …

# 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
- ❑ 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?
- ❑ …