# CZ4041/CE4041: Machine Learning

## Course Project Description

Sinno Jialin PAN

School of Computer Science and Engineering, NTU, Singapore

Homepage: http://www3.ntu.edu.sg/home/sinnopan/

# Detailed Project Description

- This is a [group-based](group-based) course project
- Each group consists of [4-5](4-5) members
- Each group can choose either one of the <u>Kaggle competitions</u> or one of the <u>research topics</u> listed on the following two slides as the course project
- Teaching Assistant: Mr. Jianda CHEN ([JIANDA001@e.ntu.edu.sg](mailto:JIANDA001@e.ntu.edu.sg))

# Course Project Candidates

## Kaggle competitions:

- BigQuery-Geotab Intersection Congestion
  **url**: https://www.kaggle.com/c/bigquery-geotab-intersection-congestion/
  **dataset**: csv (54 MB)

- EEE-CIS Fraud Detection
  **url**: https://www.kaggle.com/c/ieee-fraud-detection
  **dataset**: csv (118 MB)

- Categorical Feature Encoding Challenge
  **url**: https://www.kaggle.com/c/cat-in-the-dat
  **dataset**: csv (21 MB)

- Severstal: Steel Defect Detection
  **url**: https://www.kaggle.com/c/severstal-steel-defect-detection
  **dataset**: image (2G)

- Northeastern SMILE Lab - Recognizing Faces in the Wild
  **url**: https://www.kaggle.com/c/recognizing-faces-in-the-wild
  **dataset**: image (381 MB)

- Instant Gratification
  **url**: https://www.kaggle.com/c/instant-gratification
  **dataset**: CSV (414 MB)

# Course Project Candidates (cont.)

- Research-based projects:
  - Semi-supervised Learning

  Recommended Datasets: http://sci2s.ugr.es/keel/semisupervised.php

  - Multi-label Classification

  Recommended Datasets: http://sci2s.ugr.es/keel/multilabel.php

  - Multi-instance Learning

  Recommended Datasets: http://sci2s.ugr.es/keel/category.php?cat=mul

  - Transfer Learning

  Recommended Datasets:
  https://www.kaggle.com/c/transfer-learning-on-stack-exchange-tags
  https://ai.bu.edu/visda-2018/

- Note: If you want to use other datasets to conduct the listed research topics, an approval is needed

# **Programming Languages**

- Programming Languages:
  - Any programming language can be used, e.g., Matlab, Python, C/C++, Java, R, etc
  - Any open-source ML toolbox can be used
- Note: for Kaggle competitions, directly using the source codes released by participants are not allowed (20% penalty will be made if found)

# Key Dates

- <u>Sent information on group members via email:</u>
  - by 21$^{st}$ Feb. 2020
- <u>Submit files, i.e., the project report, video, source codes, through NTULearn:</u>
  - by 11:59pm, 24$^{th}$ Apr. 2020

| FEBRUARY 2020 | | | | | | |
|---|---|---|---|---|---|---|
| SUN | MON | TUE | WED | THU | FRI | SAT |
| | | | | | | 1 |
| 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 | 13 | 14 | 15 |
| 16 | 17 | 18 | 19 | 20 | 21 | 22 |
| 23 | 24 | 25 | 26 | 27 | 28 | 29 |

www.theprintablecalendar.com

| APRIL 2020 | | | | | | |
|---|---|---|---|---|---|---|
| SUN | MON | TUE | WED | THU | FRI | SAT |
| | | | 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 | 16 | 17 | 18 |
| 19 | 20 | 21 | 22 | 23 | 24 | 25 |
| 26 | 27 | 28 | 29 | 30 | | |

www.theprintablecalendar.com

# Submission (Kaggle)
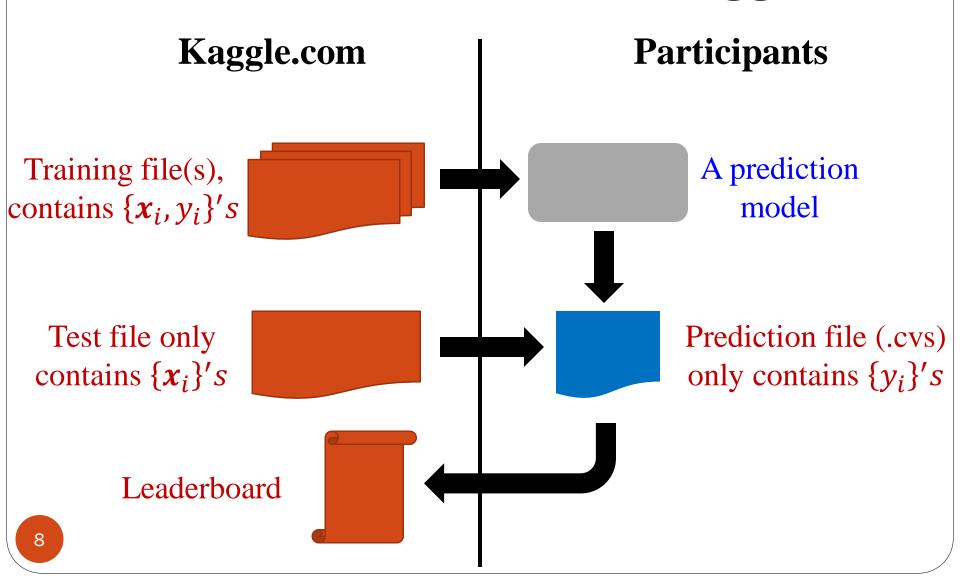
- <u>Submitted files:</u>
  1. A project report
  2. A presentation video
  3. The final .cvs file of your prediction results submitted to the specific completion in Kaggle you participate
  4. Your source codes (with a readme file)
- <u>Notes:</u>
  - The submitted .cvs is to double check whether the reported results are correct
  - The submitted source codes are to double check whether they are just copied from the ones released by some participants

# General Information of Kaggle

**Kaggle.com**  |  **Participants**

Training file(s), contains $\{\boldsymbol{x}_i, y_i\}'s$

A prediction model

Test file only contains $\{\boldsymbol{x}_i\}'s$

Prediction file (.cvs) only contains $\{y_i\}'s$

Leaderboard

# **Submission (Research)**

- <u>Submitted files:</u>
    1. A project report
    2. A presentation video
    3. Your source codes (with a readme file)
- <u>Notes:</u>
    - The submitted source codes are to double check whether the reported results are correct

# **Format and Content of Video**

- Presentation video:
  - To introduce your course project in a video of 10-15 minutes long
    - The video is a visual summary of your course report
  - You can use any tool to produce the video, e.g., simply using PowerPoint or other advanced tools
  - File size $\leq$ 8M
  - Some examples for reference:

  https://www.youtube.com/channel/UCSBrGGR7JOiSyzl60OGdKYQ

  https://www.youtube.com/channel/UC_sfvZvvPUbOQhDs_cqlx_A

# **Content of Project Report (Kaggle)**

- Specific roles of each group member

- An evaluation score and ranked position of your prediction results for the specific competition in Kaggle

- Problem statement (using your own words instead of copy-and-paste from Kaggle)

- Challenges of the problem

- Your proposed solution in detail (preprocessing, feature engineering/representation learning, methodologies, etc)

- Experiments to demonstrate why the solution you proposed is appropriate to solve the problem using experiments

- Conclusion: what you have learned from the project

# **Content of Project Report (Research)**

- Specific roles of each group member
- A review on the specific research topic
- Your new proposed method if applicable
- Comparison experiments on state-of-the-art methods (and your proposed method if applicable)
- Analysis on pros and cons of the compared methods
- Conclusion: you own insights on the research project
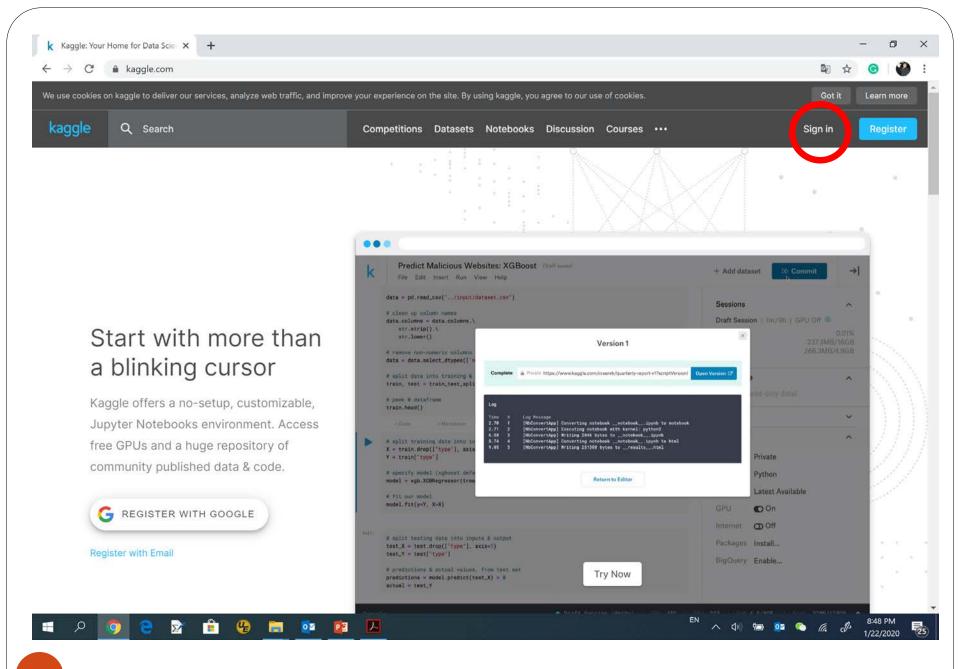
# **Format and Assessment on Project Report**
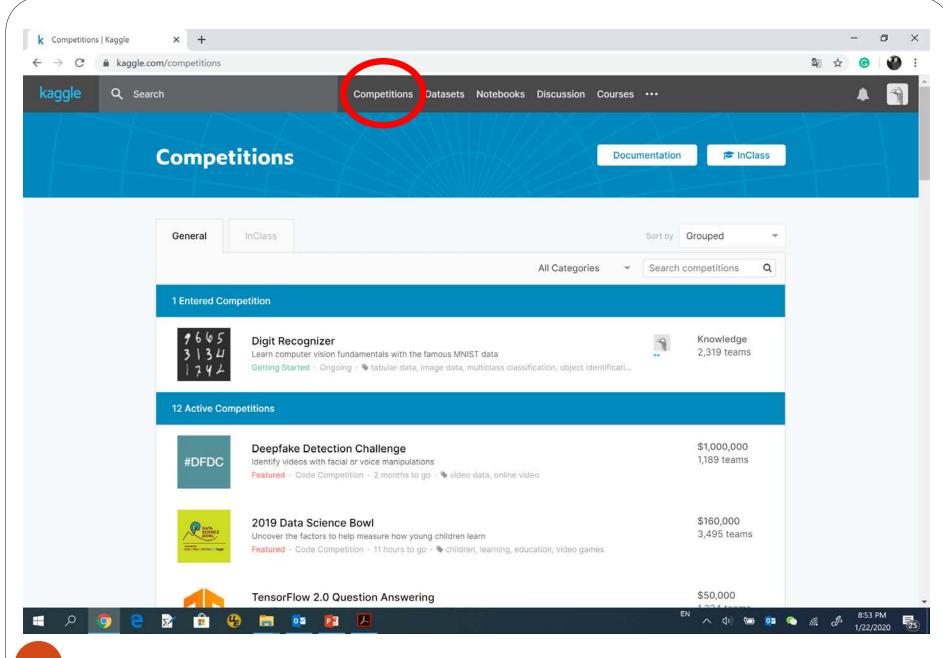
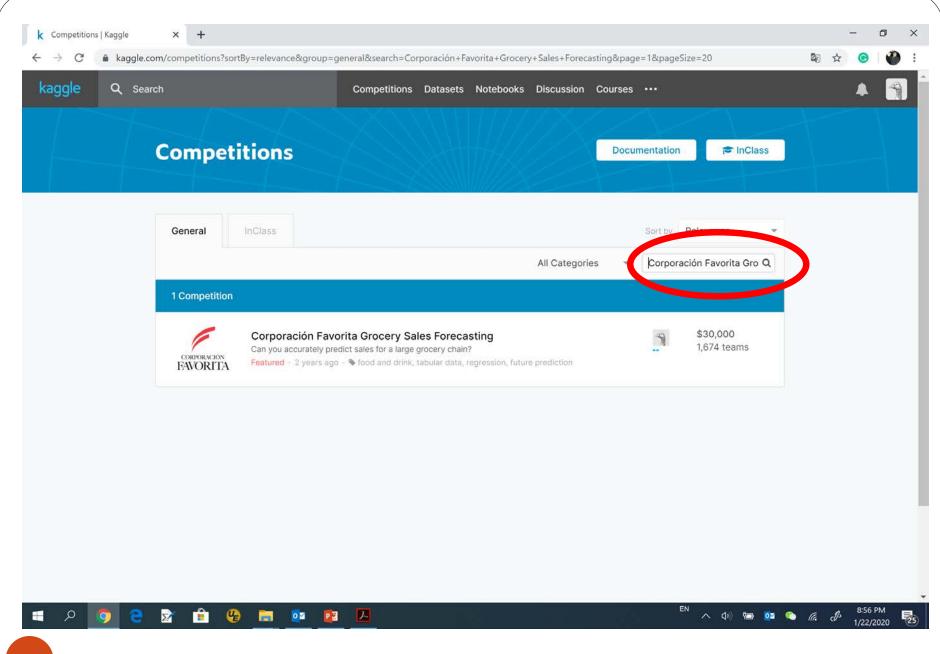- Report format:
  - 12 point font, single space, 20-25 pages

**Kaggle competitions**

- Leaderboard performance
- Convincingness
- Solution novelty
- Writing

**Research-based projects**

- Literature review
- Comparison analysis
- Methodology novelty
- Writing

This assessment is to evaluate whether the organization of report is clear and easy to follow, whether the report contains a lot of typos

# **Assessments – Kaggle**

- **<u>Leaderboard Performance:</u>** though all the listed Kaggle competitions are completed, you can still submit your results to Kaggle to obtain an evaluation score and a ranking position

- The performance assessment is based on the relatively ranking of your results on the specific competition (i.e., top 20%, top 40%, top 60%, top 80%, and top 100%)
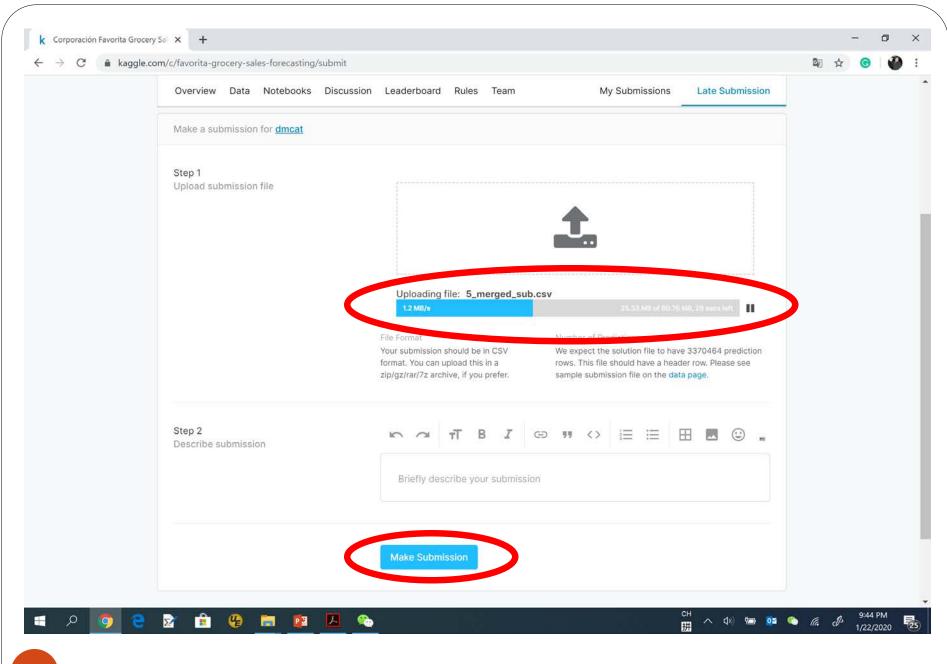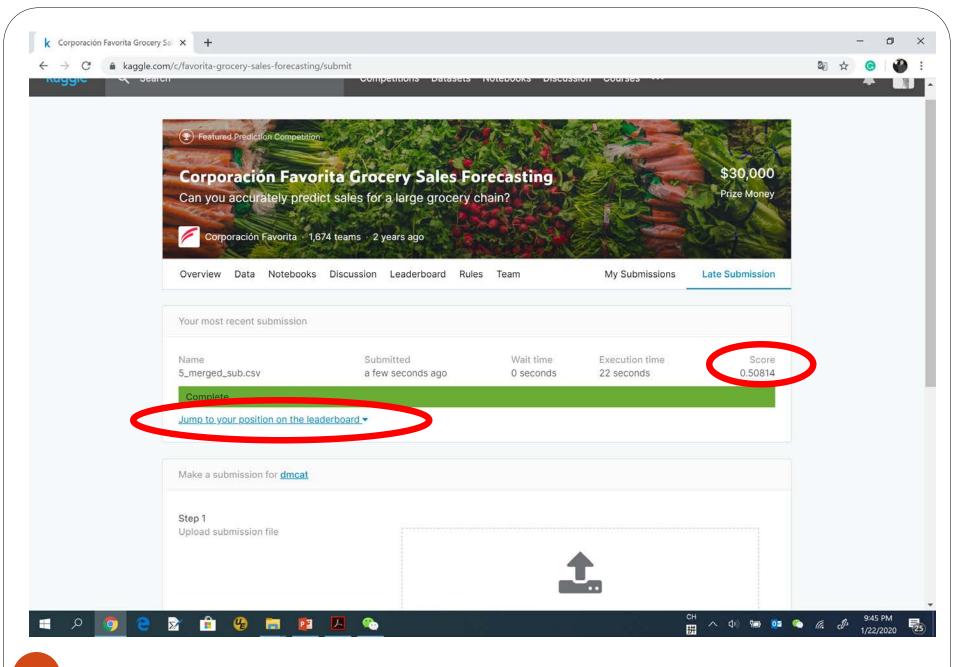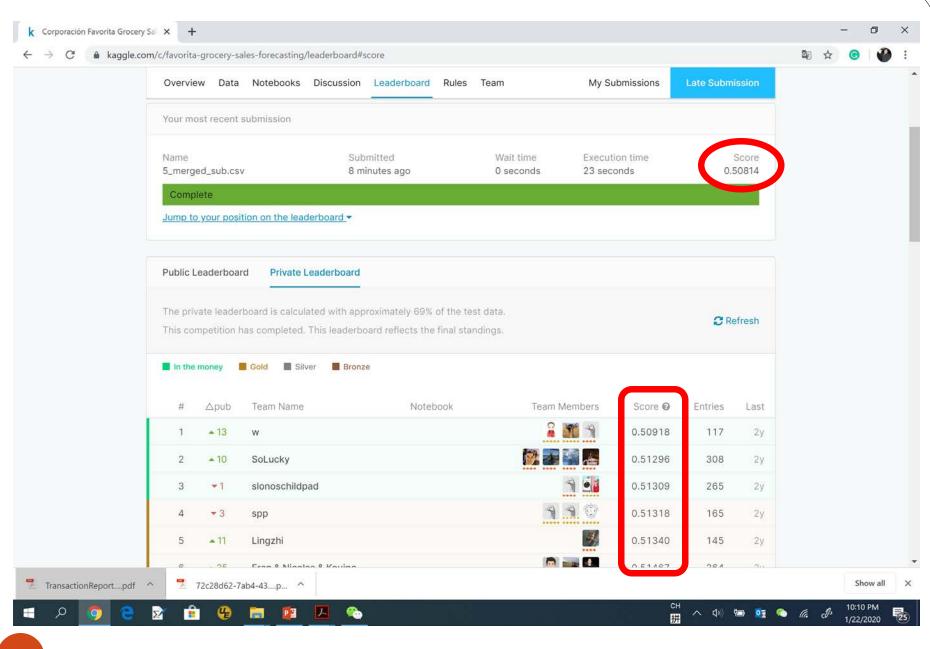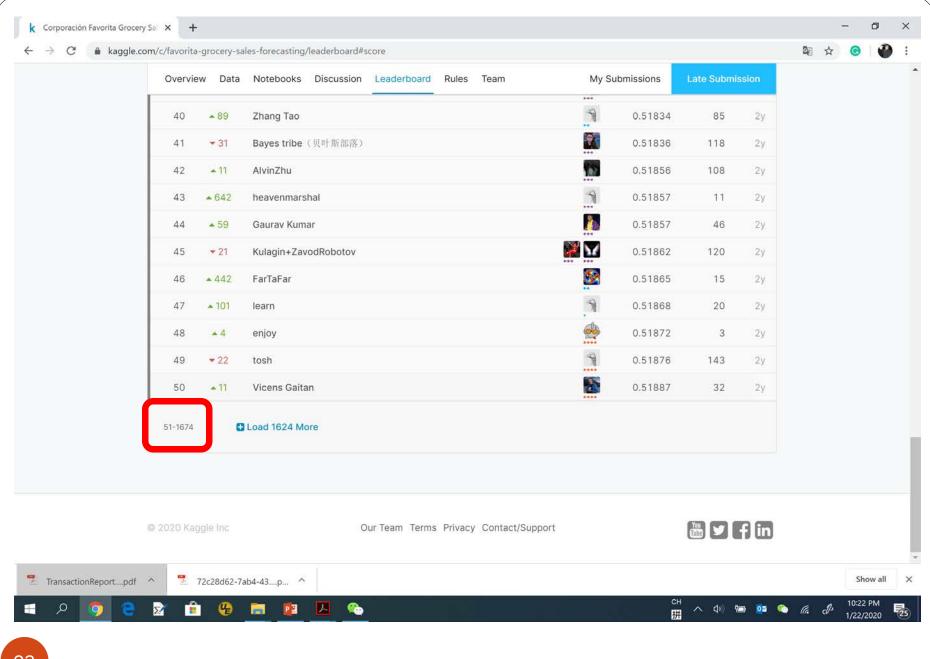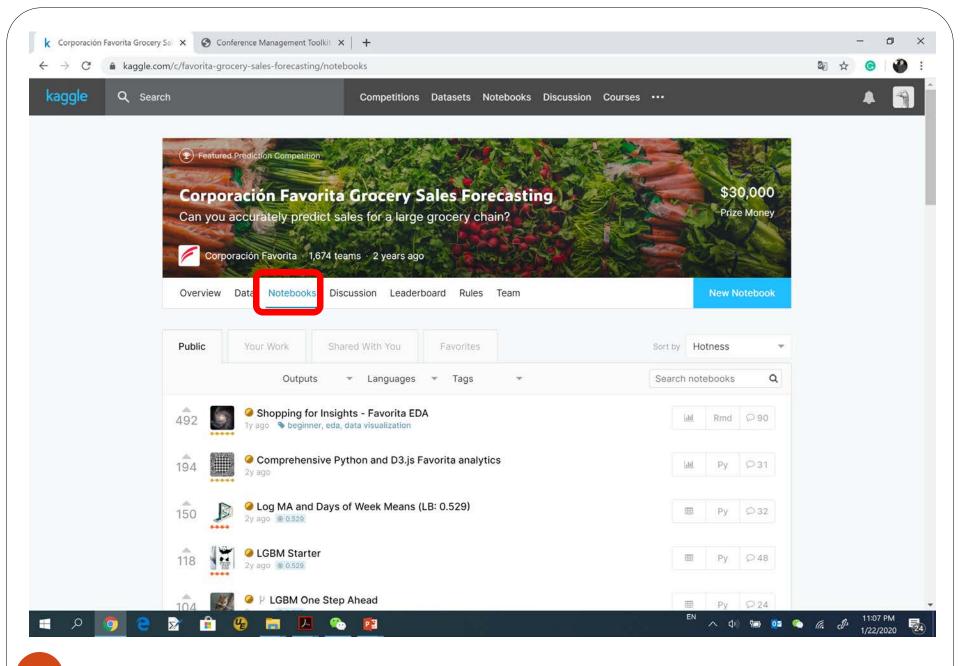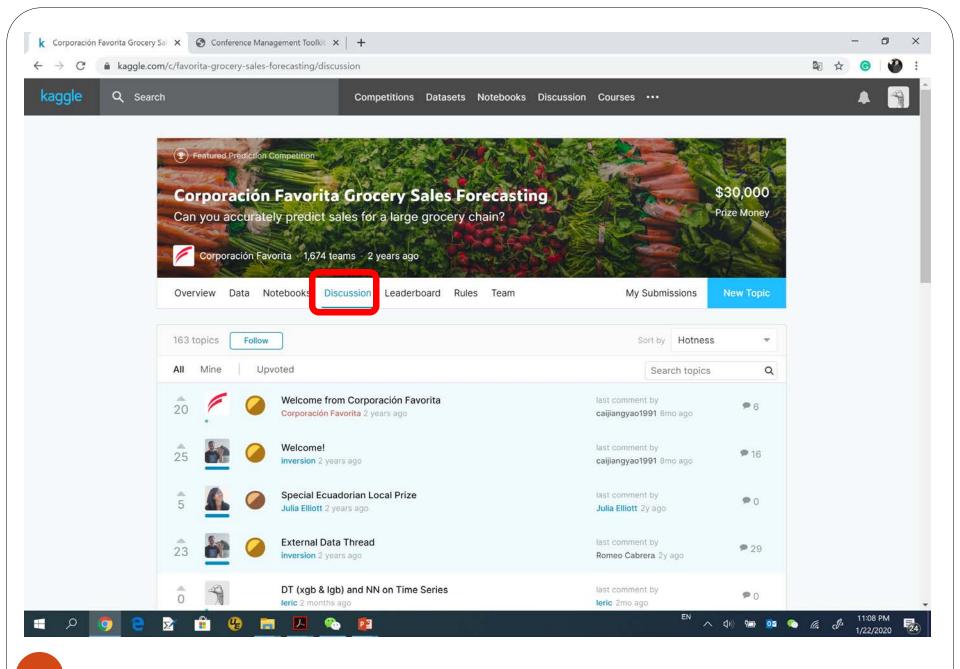
# Assessments – Kaggle (cont.)

- **<u>Solution Novelty:</u>** as on Kaggle.com, most participants or winners may discuss their solutions on the forums of the specific competitions.
  - If you propose a new and effective solution, you can get bonus. You are encouraged to propose your own solutions based on your own understandings on the competitions

# Assessments – Kaggle (cont.)

- **<u>Convincingness:</u>** the goal of the project report is to convince readers that your proposed solution is proper to solve the specific machine learning task. In your report, you need to conduct experiments to verify your proposed ideas

# Assessments – Kaggle (cont.)

- Weight priority:

    Convincingness = Writing > Leaderboard
    Performance = Solution Novelty

# Assessments – Research

- **<u>Literature Review:</u>** as this is a research project, figuring out what have been done in the literature is important. You should provide a comprehensive review on the specific research topic studied in your project

# Assessments – Research (cont.)

- **<u>Comparison Analysis:</u>** you need to implement various state-of-the-art methods for the research topic studied in your research project, and analyze their cons and pros with your own insights

# Assessments – Research (cont.)

- **<u>Methodology Novelty:</u>**  if you propose a new and effective method for the specific research topic, even though it might be incremental, you can get bonus. You are encouraged to propose your own methods based on your understandings on the research topic

# Assessments – Research (cont.)

- Weight priority:

  Literature Review = Writing = Comparison Analysis > Methodology Novelty