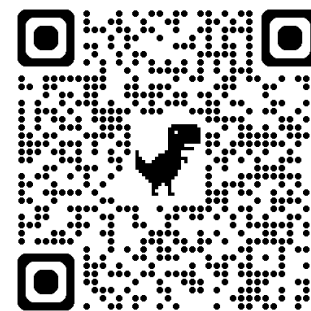# 自然語言處理與應用
# Natural Language Processing and Applications

## Project Introduction

Instructor: 林英嘉 (Ying-Jia Lin)
2025/04/07

Course GitHub

Slido # NLP_0407

# Outline

- Recap: Language Generation

- Decoding Strategies

  - Greedy Decoding

  - Beam Search

  - Top-k / Top-p Sampling

- Evaluations

NLP

# What is Kaggle?

- Kaggle is a platform that provides:

  - Real-word datasets for machine learning

  - Competitions with prizes (sometimes with money)

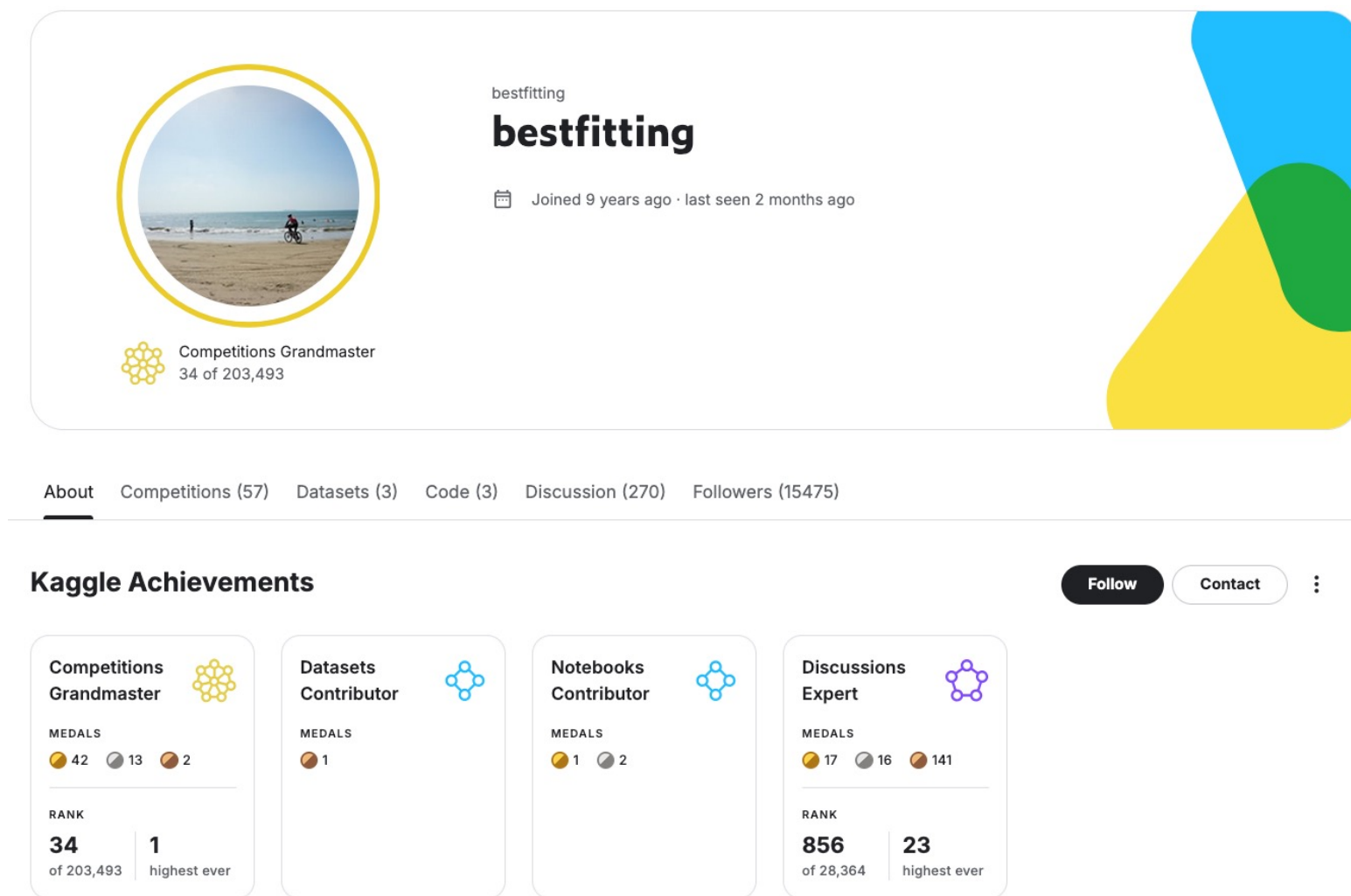  - Discussion forum with <span style="color:red">a lots of code examples</span>



NLP

3

# Kaggle submission types

- Traditional competitions

  - Upload s**ubmission file** (e.g., *.csv)

- Code / Notebook competitions    本課程 projects 採用此方式

  - Upload **code** (e.g., *.ipynb)

NLP

# If you play Kaggle a lot …

Figure source: https://www.kaggle.com/bestfitting

# Outline of tasks

| Platform | Competition Name | Data Type | Task Type |
|---|---|---|---|
| Kaggle | LLM Classification Finetuning | text | Text classification |
| | NBME - Score Clinical Patient Notes | text | Token classification (like Named Entity Recognition) |
| | LLM - Detect AI Generated Text | text | Text classification |

NLP

# LLM Classification Finetuning (intro)

**Chatbot Arena:** https://lmarena.ai/?leaderboard
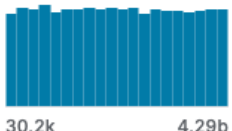
# LLM Classification Finetuning (data example)

**train.csv** (184.18 MB)    約55K rows

NLP

# LLM Classification Finetuning (data example)



**test.csv** (10.66 kB)

Detail    Compact    Column                                4 of 4 columns ⌄

| 🔑 id ≡ | 🔤 prompt ≡ | 🔤 response_a ≡ | 🔤 response_b ≡ |
|---|---|---|---|
| **3** total values | **3** unique values | **3** unique values | **3** unique values |
| 136060 | ["I have three oranges today, I ate an orange yesterday. How many oranges do I have?"] | ["You have two oranges today."] | ["You still have three oranges. Eating an orange yesterday does not affect the number of oranges you... |
| 211333 | ["You are a mediator in a heated political debate between two opposing parties. Mr Reddy is very hun... | ["Thank you for sharing the details of the situation. As a mediator, I understand the importance of ... | ["Mr Reddy and Ms Blue both have valid points in their arguments. On one hand, Mr Reddy is correct t... |
| 1233961 | ["How to initialize the classification head when I do transfer learning. For example, I have a pre-t... | ["When you want to initialize the classification head for transfer learning, you can follow these st... | ["To initialize the classification head when performing transfer learning, follow these steps: \n\n1.... |

No more data to show

此競賽為 Code competition
因此 test set 實際上不公開

**NLP**

# LLM Classification Finetuning (submission)

- 上傳到 Kaggle Leaderboard 需要實作 .ipynb
- 你的 code 需要能 predict test.csv 並產生檔案名稱為 submission.csv 的檔案，Kaggle 才能幫你執行程式碼並打分數，可參考 LMSYS: KerasNLP Starter

NLP

# LLM Classification Finetuning (evaluation)

- Binary cross-entropy，對於任一個 class 而言 (win_a, win_b, tie):

$$L_{\log}(y, p) = -\big(y\log(p) + (1 - y)\log(1 - p)\big)$$

$p$: 預測該類別為1的機率
$y$: 正確答案的類別

- Leaderboard 比的是 test set的 Loss 值 (越小越好)

# NBME - Score Clinical Patient Notes

- 找出 patient note 中的重要特徵，例如：

  - 輸入一篇 patient note，輸出為 "diminished appetite"

- 競賽的實際範例：

```
>>> import pandas as pd
>>> df = pd.read_csv("patient_notes.csv")
>>> df[df["pn_num"]==16].pn_history.values[0][696:724]
'dad with recent heart attcak'
```

# NBME – patient note example

```
>>> df[df["pn_num"]==16].pn_history.values[0]
```

HPI: 17yo M presents with palpitations. Patient reports 3-4 months of intermittent episodes of "heart beating/pounding out of my chest." 2 days ago during a soccer game had an episode, but this time had chest pressure and felt as if he were going to pass out (did not lose conciousness). Of note patient endorses abusing adderall, primarily to study (1-3 times per week). Before recent soccer game, took adderrall night before and morning of game. Denies shortness of breath, diaphoresis, fevers, chills, headache, fatigue, changes in sleep, changes in vision/hearing, abdominal paun, changes in bowel or urinary habits. \r\nPMHx: none\r\nRx: uses friends adderrall\r\nFHx: mom with "thyroid disease," dad with recent heart attcak\r\nAll: none\r\nImmunizations: up to date\r\nSHx: Freshmen in college. Endorses 3-4 drinks 3 nights / week (on weekends), denies tabacco, endorses trying marijuana. Sexually active with girlfriend x 1 year, uses condoms

NLP

## NBME - Score Clinical Patient Notes (data example - train.csv)

但有的feature 可能沒有標註

| id | case_num | pn_num | feature_num | annotation | location |
|---|---|---|---|---|---|
| Unique identifier for each patient note / feature pair. | The case to which this patient note belongs. | The patient note annotated in this row. | The feature annotated in this row. | The text(s) within a patient note indicating a feature. A feature may be indicated multiple times within a single note. | Character spans indicating the location of each annotation within the note. |
| 14300 unique values | 0 — 9 | 16 — 95.3k | 0 — 916 | [] 31%<br>['F'] 2%<br>Other (9597) 67% | [] 31%<br>['0 5'] 1%<br>Other (9719) 68% |
| 00016_000 | 0 | 00016 | 000 | ['dad with recent heart attcak'] | ['696 724'] |
| 00016_001 | 0 | 00016 | 001 | ['mom with "thyroid disease'] | ['668 693'] |
| 00016_002 | 0 | 00016 | 002 | ['chest pressure'] | ['203 217'] |
| 00016_003 | 0 | 00016 | 003 | ['intermittent episodes', 'episode'] | ['70 91', '176 183'] |
| 00016_004 | 0 | 00016 | 004 | ['felt as if he were going to pass out'] | ['222 258'] |
| 00016_005 | 0 | 00016 | 005 | [] | [] |
| 00016_006 | 0 | 00016 | 006 | ['adderall', 'adderrall', 'adderrall'] | ['321 329', '404 413', '652 661'] |
| 00016_007 | 0 | 00016 | 007 | [] | [] |
| 00016_008 | 0 | 00016 | 008 | [] | [] |
| 00016_009 | 0 | 00016 | 009 | ['palpitations', 'heart beating/ pounding'] | ['26 38', '96 118'] |
| 00016_010 | 0 | 00016 | 010 | ['3-4 months of'] | ['56 69'] |
| 00016_011 | 0 | 00016 | 011 | ['17yo'] | ['5 9'] |
| 00016_012 | 0 | 00016 | 012 | ['M'] | ['10 11'] |

每筆資料都有 feature_num代號

# NBME - Score Clinical Patient Notes (data example - features.csv)

| # feature_num ≡ | # case_num ≡ | ⩍ feature_text ≡ |
|---|---|---|
| A unique identifier for each feature. | The case to which this patient note belongs. | A description of the feature. |
|  0                              916 |  0                              9 | Female 5%<br>Male 2%<br>Other (133) 93% |
| 000 | 0 | Family-history-of-MI-OR-Family-history-of-myocardial-infarction |

- 競賽的實際範例：

```
>>> import pandas as pd
>>> df = pd.read_csv("patient_notes.csv")
>>> df[df["pn_num"]==16].pn_history.values[0][696:724]
'dad with recent heart attcak'
```

NLP

# NBME - Score Clinical Patient Notes (submission and evaluation)

## sample_submission.csv (93 B)

Detail  Compact  Column

**About this file**

A sample submission file in the correct format.

| ⚠ id | ⚠ location |
|---|---|
| Unique identifier for this instance, a feature within a patient note. | Character spans indicating the location(s) of the feature within the note. |
| **5** unique values | [null] 40% / 0 100 20% / Other (2) 40% |
| 00016_000 | 0 100 |
| 00016_001 | |
| 00016_002 | 200 250;300 400 |
| 00016_003 | |
| 00016_004 | 75 110 |

No more data to show

- 只需要預測 location 即可
- 上傳到 Kaggle Leaderboard 需要實作 .ipynb
- 你的 code 需要能 predict test.csv 並產生檔案名稱為 submission.csv 的檔案，Kaggle 才能幫你執行程式碼並打分數，可參考 NBME / Deberta-base baseline [inference]

# NBME - Score Clinical Patient Notes (evaluation)

- Metric: Micro F1-score

**Example**

Suppose we have an instance:

```
| ground-truth | prediction     |
|--------------|----------------|
| 0 3; 3 5     | 2 5; 7 9; 2 3  |
```

These spans give the sets of indices:

```
| ground-truth | prediction  |
|--------------|-------------|
| 0 1 2 3 4    | 2 3 4 7 8   |
```

We therefore compute:

- TP = size of {2, 3, 4} = 3
- FN = size of {0, 1} = 2
- FP = size of {7, 8} = 2

# Confusion Matrix

|  | **Actually positive** | **Actually negative** |
|---|---|---|
| **Predicted positive** | True positive (TP) | False positive (FP) |
| **Predicted negative** | False negative (FN) | True negative (TN) |

- Precision = TP / (TP + FP)
  - 模型預測的TP比例
- Recall = TP / (TP + FN)
  - True Positive Rate (TPR)
- F1-score (廣義) = 2(Precision*Recall) / (Precision+Recall)

# Macro vs. Micro

- 假設總共有100筆資料：

  - 每一筆都算出一個 F1-score，最後取平均 => Macro F1-score

  - 加總全部100筆的 TP, FN 以及 FP 之後，再算出 F1-score => Micro F1-score

# LLM - Detect AI Generated Text

- The competition dataset comprises about 10,000 essays

  - Some written by students and some generated by a variety of large language models (LLMs).

- The goal of the competition is to determine whether or not essay was generated by an LLM.

# LLM - Detect AI Generated Text (data)



可能有幫助的外部資料集：https://www.kaggle.com/datasets/alejopaullier/daigt-external-dataset

NLP

# LLM - Detect AI Generated Text (submission and evaluation)

## Evaluation

Submissions are evaluated on **area under the ROC** curve between the predicted probability and the observed target.

## Submission File

For each `id` in the test set, you must predict a probability that that essay was `generated`. The file should contain a header and have the following format:

```
id,generated
0000aaaa,0.1
1111bbbb,0.9
2222cccc,0.4
...
```
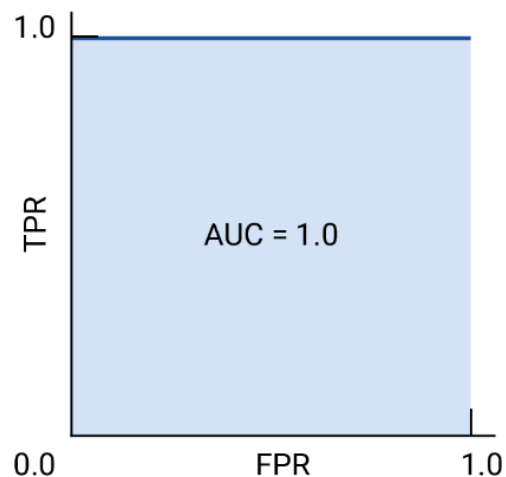
- 上傳到 Kaggle Leaderboard 需要實作 .ipynb
- 你的 code 需要能 predict test_essays.csv 並產生檔案名稱為 submission.csv 的 檔案，Kaggle 才能幫你執行程式碼並打分數，可參考 0.960| Phrases are keys

# ROC Curve

ROC: Receiver operating characteristic curve (接收者操作特徵曲線)

AUC: ROC 的底面下面積 (越大越好)



圖 1：ROC 和 AUC，這是一個完美假設的模型。

- TP / (TP + FN): True Positive Rate (TPR)
  - 又稱作 Recall
- FP / (FP+TN): False Positive Rate (FPR)

NLP

23

# Project checkpoints (暫定)

- Week 9: 確定各組的題目

- Week 11: 進度報告 PPT (5 pages)

- Week 13: 進度報告 PPT (5+5 pages), Presentations (selected teams)

- Week 15 – Week 16: Final presentations for all teams (maybe poster)

- Week 16 結束前: 繳交書面報告以及程式碼

# 期末 Project 規定 (暫定)

- 需要上傳 Kaggle Leaderboard
  - 請留意，每個題目都是 Code competition，code 必須沒有 bug 才能上傳
- 每次報告都需要列出每位組員的貢獻內容，以及組員間的工作比重 (%)

# Thank you!

Instructor: 林英嘉

✉ yjlin@cgu.edu.tw