# 감성 분석을 통한 호감도 예측

2022.04.06

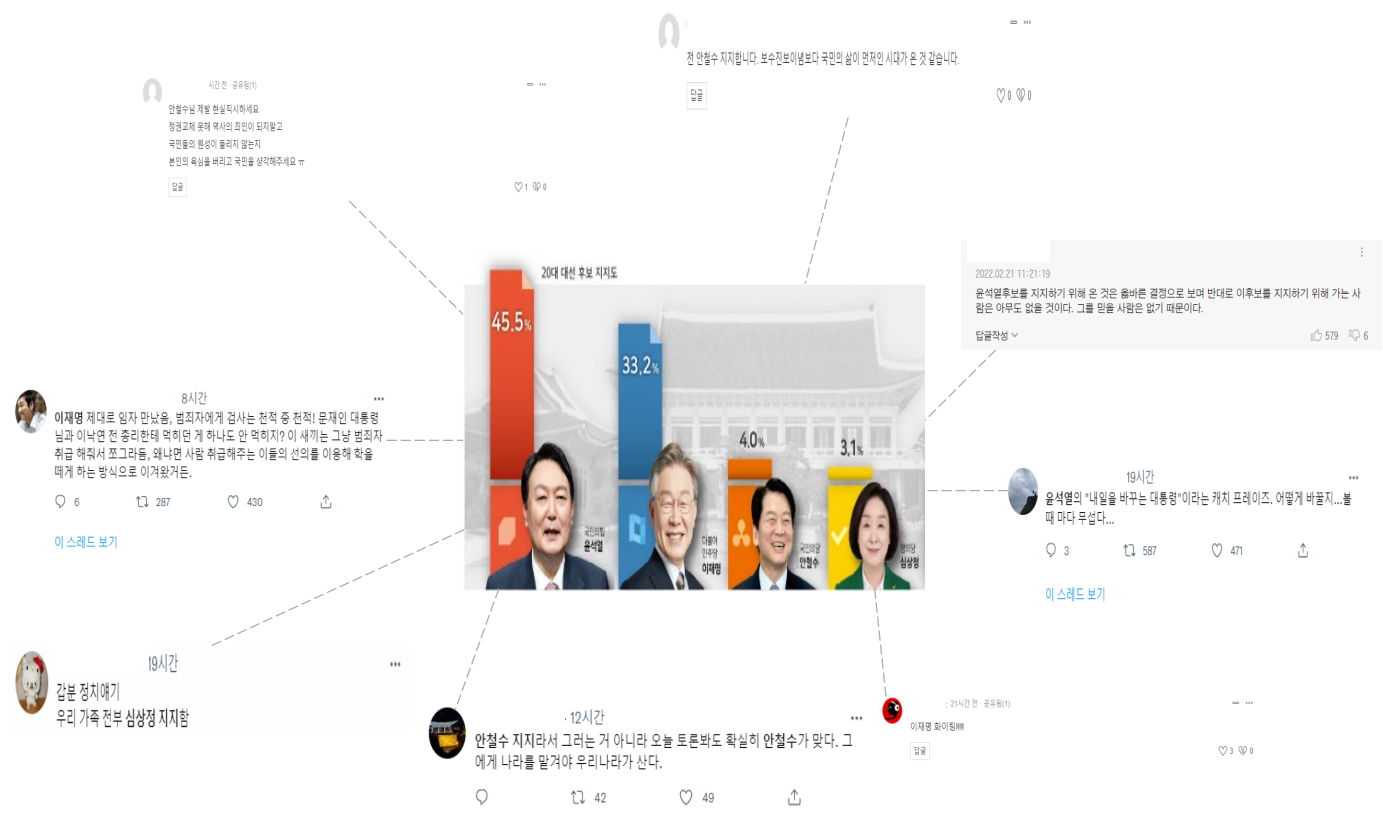김수영

# CONTENTS

**MLCL**
Kyungpook National University

01

Intro

MLCL

Kyungpook National University

# Project Name: 감성적인 투표

# Intro

**1K**

# Intro

# Goals

- Text–based unseen data 활용 연구

- 감성 분석을 통한 후보 호감도 예측

Kyungpook National University

# 02

# Data

Kyungpook National University

# Data Collection



**Date: 2/1 ~ 3/4**

**Data count: 2730K**

**Method: Web Crawling**

# Data Pre-processing

### Data classification

### Deduplication

2730K ⟹

| 이재명 | 336,049 |
|---|---|
| 윤석열 | 386,539 |
| 중복 | 104,414 |
| 총 합 | 827K |

⟹

| 이재명 | 11,675 |
|---|---|
| 윤석열 | 16,140 |
| 중복 | 4,248 |
| 총 합 | 32K |

Counter({('4tos****', '2022.03.09. 09:45', '국민들께 외칩니다...가르쳐 주십소∞ 거덜 은국안 개망 ...'): 2952, ('rrjj****', '2022.03.09. 09:36', '아무리 눈치가 없어도 파란코트 박근혜 파란마스크 홍준표 그래도 윤석열 찍는 모자란 대구시민은 없겠제?'): 2910, ('slsw****', '2022.03.09. 09:19', '전라도 20대 30대는 윤석열!!!!!!!'): 2910, ('hh18****', '2022.03.09. 09:43', '무조건윤석열'): 2562, ('jiny****', '2022.03.09. 09:17', '결국 코로나땜시 투표 분산효과만있고 투표율은 예전이나 지금이나 도찐개찐이네∼근데 무식한 대통령 나오면 안되는데 걱정이다 5년 주위에서 다 해쳐먹겠네∼'): 2562, ('kkhk****', '2022.03.09. 10:29', '윤석열 대통령님 같이 좌파 빨갱o 공산당 멸공합시다 !'): 1958, ('samy****', '2022.03.09. 10:10', '보수는 마지막에 강하다 !...대구 !....역시 대구 !....감사합니다 윤석열찍어주신. 대구시민들최고 !....'): 1958, ('ingc****', '2022.03.09. 08:16:29', '투표 독려해서 윤석열 당선시킵시다 '): 1957, ('duff****', '2022.03.09. 08:20:05', '열심히 투표해서 정권교체합시다 '): 1957, ('kj02****', '2022.03.08. 11:14',

Kyungpook National University

# Data

▷ Naver news comment data

| 이재명 | 11,675 |
|--------|--------|
| 윤석열 | 16,140 |
| 중복 | 4,248 |
| 총 합 | 32,063 |

|  | 윤석열 | 이재명 | 총합 |
|--------|--------|--------|--------|
| Train | 9,000 | 9,000 | 18,000 |
| Test | 1,000 | 1,000 | 2,000 |

▷ Naver movie review data

| Train | 150K |
|-------|------|
| Test | 50K |

|  | id | document | label |
|---|------|----------|-------|
| 0 | 6270596 | 굳 ㅋ | 1 |
| 1 | 9274899 | GDNTOPCLASSINTHECLUB | 0 |
| 2 | 8544678 | 뭐야 이 평점들은.... 나쁘진 않지만 10점 짜리는 더더욱 아니잖아 | 0 |
| 3 | 6825595 | 지루하지는 않은데 완전 막장임... 돈주고 보기에는.... | 0 |
| 4 | 6723715 | 3D만 아니었어도 별 다섯 개 줬을텐데.. 왜 3D로 나와서 제 심기를 불편하게 하죠?? | 0 |

Kyungpook National University

03

# Experiment

# Method1. Pre-trained model (KoBERT)

**KoBERT**

mask 1　　mask 2

↑　　↑

**Pre-training**

↑

**masked sentence**
수영이는 <mask>을 먹었다.
고기는 정말 <mask>.

| 데이터 | 문장 | 단어 |
|---|---|---|
| 한국어 위키 | 5M | 54M |

➡

긍정: 1  부정: 0

↑

**Fine-tuning**

↑

sentence

# Method2. Pre-trained model + Self-training

# Method2. Pre-trained model + Classic Self-training

**Algorithm 1** Classic Self-training

1: Train a base model $f_{\boldsymbol{\theta}}$ on $L = \{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{l}$
2: **repeat**
3:      Apply $f_{\boldsymbol{\theta}}$ to the unlabeled instances $U$
4:      Select a subset $S \subset \{(\boldsymbol{x}, f_{\boldsymbol{\theta}}(\boldsymbol{x})) | \boldsymbol{x} \in U\}$
5:      Train a new model $f_{\boldsymbol{\theta}}$ on $S \cup L$
6: **until** convergence or maximum iterations are reached

From [1]

[1] He, Junxian, et al. "Revisiting self-training for neural sequence generation." arXiv preprint arXiv:1909.13788 (2019).

MLCL
Kyungpook National University

# Data

| | Pre-trained model (KoBERT) | Self-training |
|---|---|---|
| Labeled data (movie review) | 200K | 200K |
| Unlabeled data (news comment) | 18K | 18K |
| Test (news comment) | 2K | 2K |

MLCL Kyungpook National University

# Parameter

| | Pre-trained model (KoBERT) | Self-training |
|---|---|---|
| Batch size | 64 | 64 |
| Epoch | 5 | 5 |
| Learning rate | 1e-4 | 1e-4 |
| Dropout | 0.5 | 0.5 |
| Iteration | - | 4 |
| Optimizer | AdamW | AdamW |
| Criterion | CrossEntropyLoss | CrossEntropyLoss |

MLCL
Kyungpook National University

# Experiment Result

## Method1. Pre-trained model (KoBERT)

**68.16 %**

## Method2. Classic Self-training

# 04

# Analysis

MLCL

Kyungpook National University

# Analysis

1. Self-training with Noisy Student

2. Data ratio

3. Labeled data change

4. Test data change and F1-score measure

MLCL
Kyungpook National University

# 1. Self–training with Noisy Student



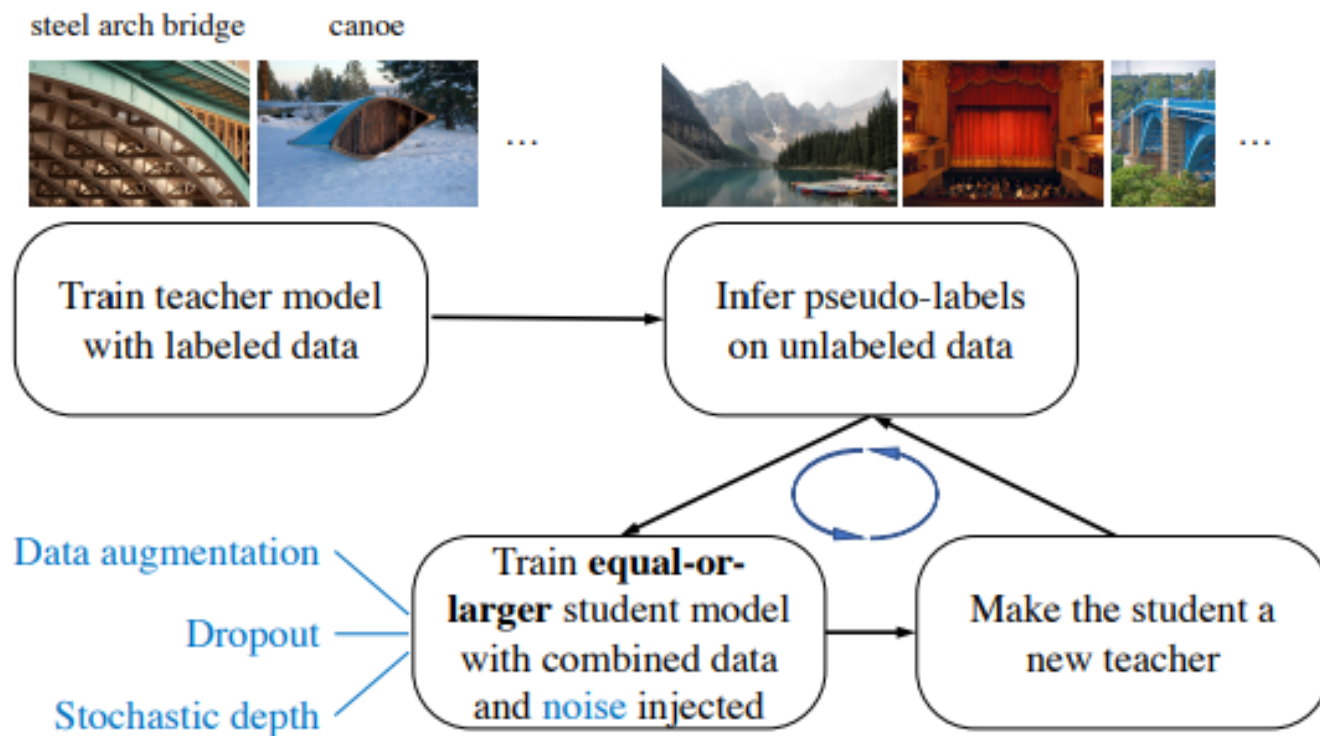| Model / Unlabeled Set Size | 1.3M | 130M |
|---|---|---|
| EfficientNet-B5 | 83.3% | 84.0% |
| Noisy Student Training (B5) | **83.9%** | **85.1%** |
| student w/o Aug | 83.6% | 84.6% |
| student w/o Aug, SD, Dropout | 83.2% | 84.3% |
| teacher w. Aug, SD, Dropout | 83.7% | 84.4% |

From [1]

[1] Xie, Qizhe, et al. "Self-training with noisy student improves imagenet classification." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.

Kyungpook National University

# 1. Self–training with Noisy Student

steel arch bridge          canoe



| Model / Unlabeled Set Size | 1.3M | 130M |
|---|---|---|
| EfficientNet-B5 | 83.3% | 84.0% |
| Noisy Student Training (B5) | **83.9%** | **85.1%** |
| student w/o Aug | 83.6% | 84.6% |
| student w/o Aug, SD, Dropout | 83.2% | 84.3% |
| teacher w. Aug, SD, Dropout | 83.7% | 84.4% |

From [1]

Data augmentation

Dropout

Stochastic depth

KoBERT
dropout = 0.2

number of data

(1) 4500
(2) 9000
(3) 13500
(4) 18000

[1] Xie, Qizhe, et al. "Self-training with noisy student improves imagenet classification." *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2020.

Kyungpook National University

# 1. Self–training with Noisy Student

Pre-trained teacher model
(KoBERT)

→

**Fine-Tuning**
teacher model
with labeled data

→

**Infer pseudo labels**
on unlabeled data

**Teacher = Student**

Make **the new student model with noisy**

←

Train the student model with **combined data**

[1] Zhang, Yu, et al. "Pushing the limits of semi-supervised learning for automatic speech recognition." *arXiv preprint arXiv:2010.10504* (2020).

Kyungpook National University

# 1. Self–training with Noisy Student - Result

# 2. Data ratio

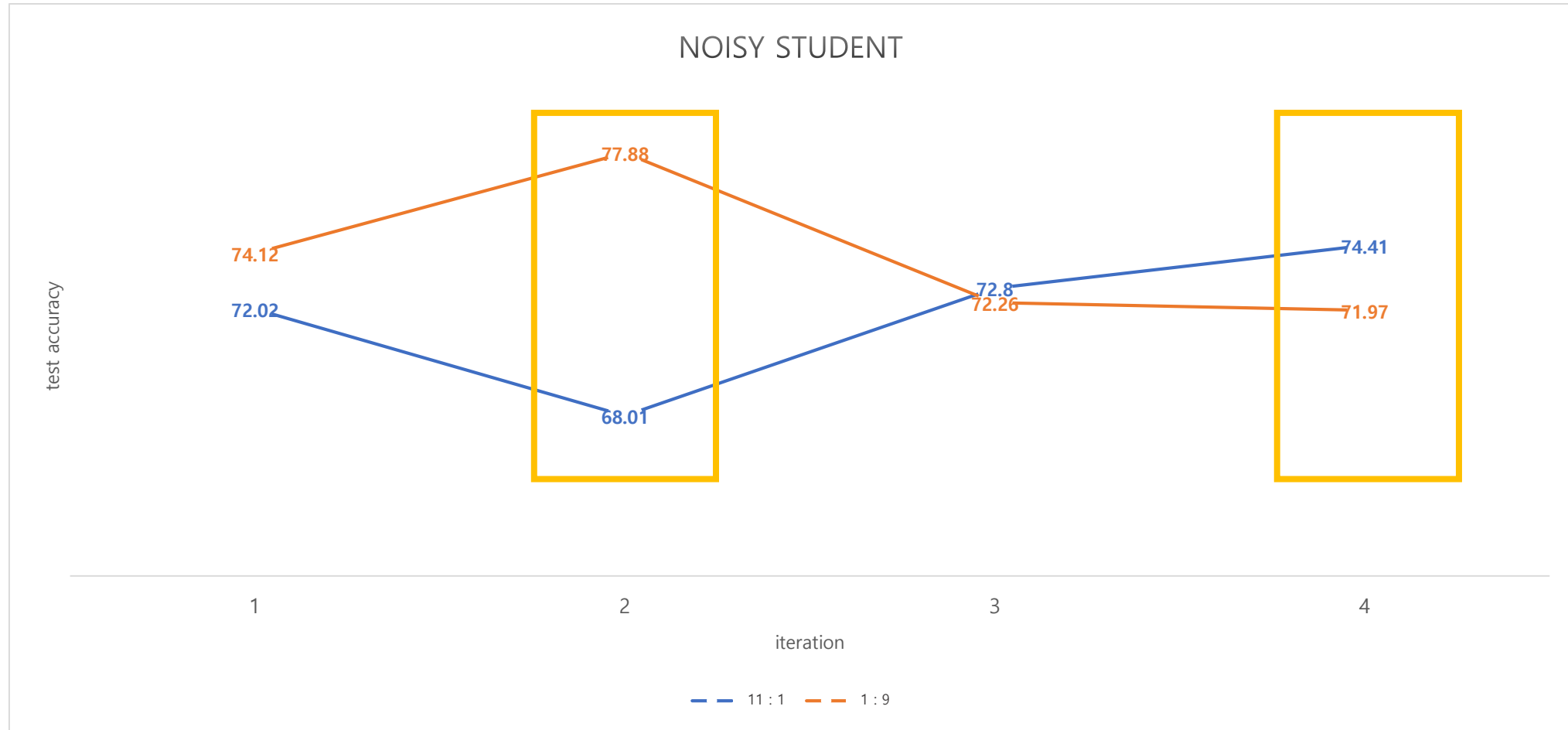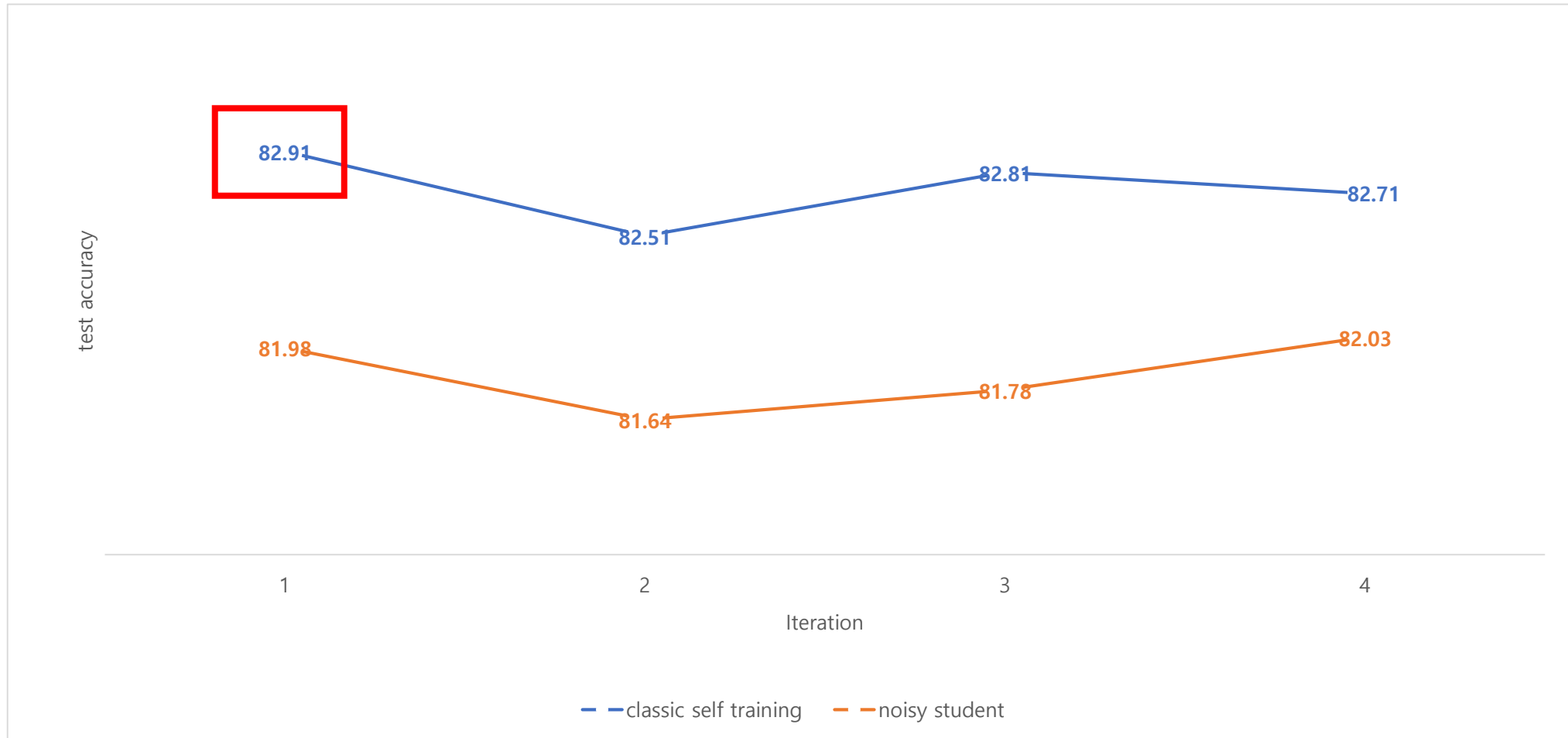| Labeled data: 200K | | Labeled data: 2K |
| Unlabeled data: 18K | $\Rightarrow$ | Unlabeled data: 18K |
| 11 : 1 | | **1 : 9** |

▷ Infer pseudo labels on **accumulated** unlabeled data

(1) 500 : 4500

(2) 1000 : 4500 + 4500

(3) 1500 : 9000 + 4500

(4) 2000 : 13500 + 4500

# 2. Data ratio - Result



NOISY STUDENT

test accuracy

74.12
72.02
77.88
68.01
72.8
72.26
74.41
71.97

iteration

1        2        3        4

— — 11 : 1    — — 1 : 9

# 3. Labeled Data Change (1 : 9)

# 4. Test data change and F1-score measure

▷ Test data change

| 0 | 1 |
|---|---|
| 71.6% | 28.4% |

➡️

| 0 | 1 |
|---|---|
| 58.05% | 41.95% |

▷ F1-score



Recall        F1 score        Precision

Real Answer                    Prediction Confidence

Kyungpook National University

# 4. Test data change and F1-score measure



ACCURACY

- 81.54
- 80.41
- 80.37
- 79.83
- 79.58
- 79.34
- 79.19
- 78.66

iteration

— classic training   — noisy student

Pre-trained model: 78.85

F1-SCORE

- 80.61
- 79.13
- 78.81
- 78.49
- 78.2
- 78.19
- 78.01
- 77.33

iteration

— classic training   — noisy student

Pre-trained model: 77.20

MLCL Kyungpook National University

# 4. Test data change and F1-score measure

|  | Positive | Negative |
|---|---|---|
| **True** | 2.8 | 1.4 |
| **False** | 0.37 | 5.4 |

$$\text{recall} = \frac{2.8}{2.8 + 5.4} = 0.34$$

$$\text{precision} = \frac{2.8}{2.8 + 0.37} = 0.88$$

06

# Conclusion

MLCL

**Kyungpook National University**

# Conclusion

▷ Labeled data : dissimilar 200K < **similar** 2K

▷ Adding **noise** in Self-training can prevent learning incorrect pseudo labels

▷ The use of self-training is helpful in the **fine-tuning** stage.

MLCL
Kyungpook National University
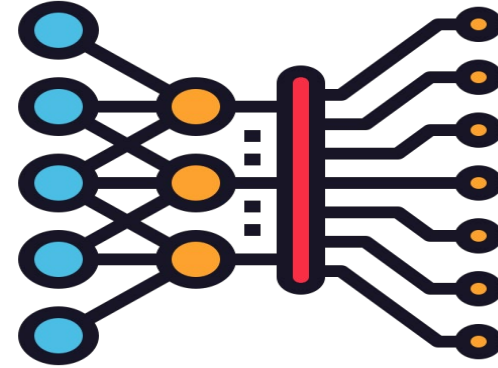
후보 호감도 예측

긍정: 23.6%   부정: 76.4%   >   긍정: 20.8% 부정: 79.2%

# Self-training Effect

- Time – 27 hours (18K x 90m)

- Cost – 2,700,000 (18K x150)

- Artificial intelligence is used to avoid borrowing human hands, but data labeling consumes a lot of manpower

- Self-training can overcome the limits of deep learning

- Unseen data input the field can be utilized in real time

07

# Future works

Kyungpook National University

# Future works

▷ **Self-training performance**

- Select confidence score

- KcBERT

- Experiment with large amounts of data

- Zoph, Barret, et al. "Rethinking pre-training and self-training." *Advances in neural information processing systems* 33 (2020): 3833-3845.

▷ **Multi-view Algorithm: Co-training**
- Several models work together to learn

Q & A

MLCL
Kyungpook National University

Thank you

Kyungpook National University