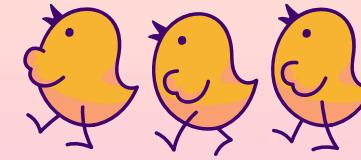


LangCon 2020



롯데의 딥러닝 모델은 어떻게 자기소개서를 읽고 있을까?

롯데정보통신

정보기술연구소 AI기술팀

김영준

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

기술 키워드 : 자기소개서 전처리, 딥러닝, TextRank 알고리즘, Bi-LSTM+CNN

2. 자기소개서의 어색한 문장을 찾을 수 있을까?

기술 키워드 : 딥러닝, 유사도, MaLSTM, Transformer

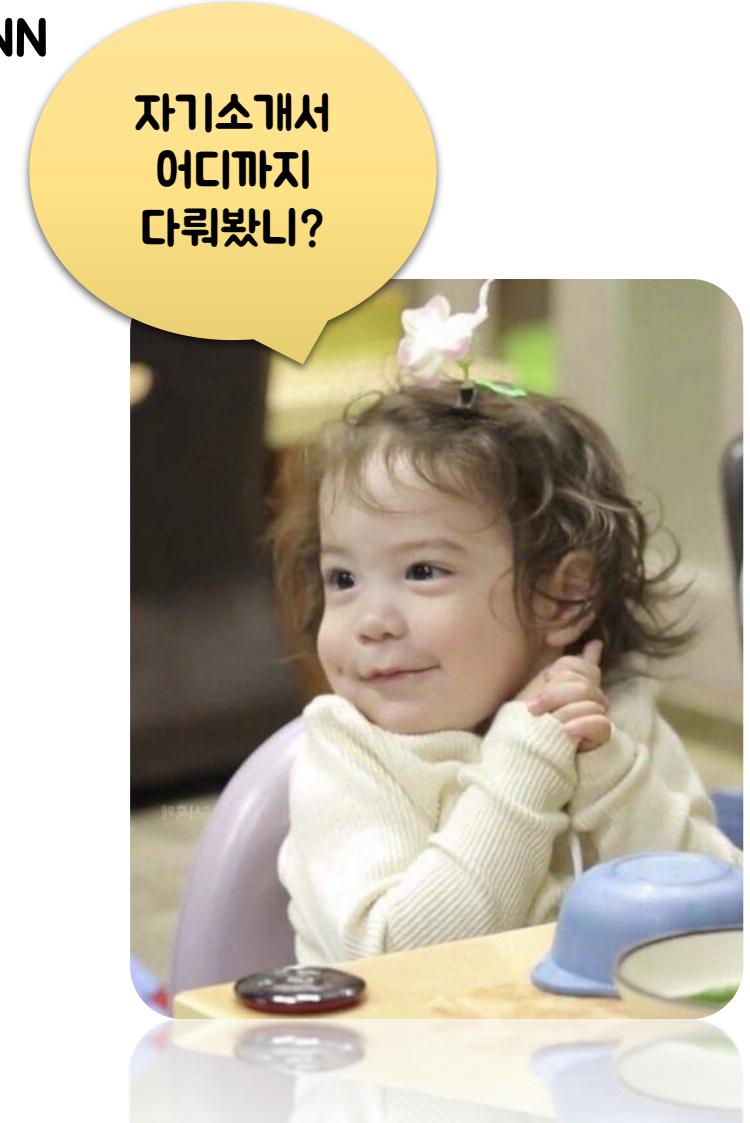
3. [Ctrl+C, Ctrl+V] 자기소개서 어떻게 찾을까?

기술 키워드 : 딥러닝, BERT

4. 오타 어떻게 찾을까?

기술 키워드 : Levenshtein distance 알고리즘

5. 롯데그룹 채용 안내



김영준

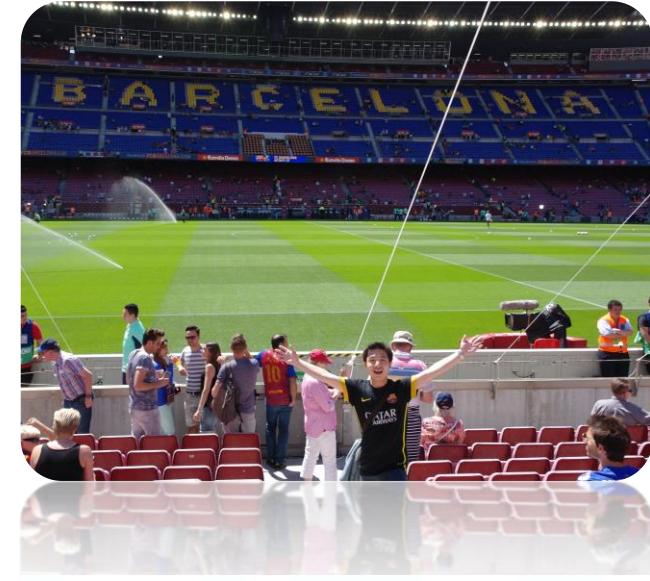
정보기술연구소 AI 기술팀
자연어 처리 담당
채용AI

성균관대학교
데이터사이언스융합학과 석사

- 최우수연구상
- 수석 졸업

딥러닝과 NLP 전문가를 꿈꾸며 이와 관련된 모든 일을 사랑합니다.
가방 속에 있는 IT 서적이 무거울 수록 마음이 편안해지는 특이한 취향이 있습니다.

data.juny@gmail.com



Predicting the Success of Job Application Using Personal Statements

ABSTRACT

A clear and well-documented \LaTeX document is presented as an article formatted for publication by ACM in a conference proceedings or journal publication. Based on the template, this article provides examples of the elements an author may use in the preparation of the documentation of their work.

CCS CONCEPTS

- Computer systems organization → Embedded systems; Redundancy; Robotics; Networks → Network reliability.

KEYWORDS

datasets, neural networks, gaze detection, text tagging.

ACM Reference Format:
2018. Predicting the Success of Job Application Using Personal Statements. In *Woodstock '18: ACM Symposium on Neural Data Detection*. June 03–05, 2018, Woodstock, NY, ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3220451.3220506>.

1 INTRODUCTION

Among the numerous resumes, research has been actively conducted to find suitable candidates. In 2013, Palagka and others introduced applications using machine learning and semantic matching technology to find the right candidate [1]. The application showed a methodology for extracting the necessary information from a linked-in profile and comparing it with job requirements.

Meanwhile, recruitment-related studies are focused on effectively comparing resumes and job requirements. In 2014, P. Monttuchi and others criticized the need for complex work to identify a suitable candidate [2]. In 2015, Monttuchi and others proposed a system to identify suitable candidates by taking into account how much the resume and job offer conform semantically [3]. In 2015, Carlos Diego et al. published a paper to find suitable candidates by comparing resumes and job requirements based on job requirements [5].

Some studies have tried to rank resumes and job requirements. In 2017, Zain and others proposed how to use the Integrated Knowledge Graph as a hybrid approach that automatically ranks the conceptual-based classification between resumes

However, existing studies, which are considered based on resumes, are suitable for hiring experienced players, but are inappropriate for hiring new employees. This is because most of the new employees don't have any experience, which is a key piece of information on their resumes, or have worked briefly at previous even if they do. Second, self-introduction is more important than resumes for new employees. This is because the potential for growth is more important than the ability to perform tasks right now because there is a small gap in the ability of applicants to perform their jobs. Therefore, hiring new employees who lack experience and have important future growth potential means more to self-introduction than resumes.

The paper makes the following contributions to the field of recruitment research.

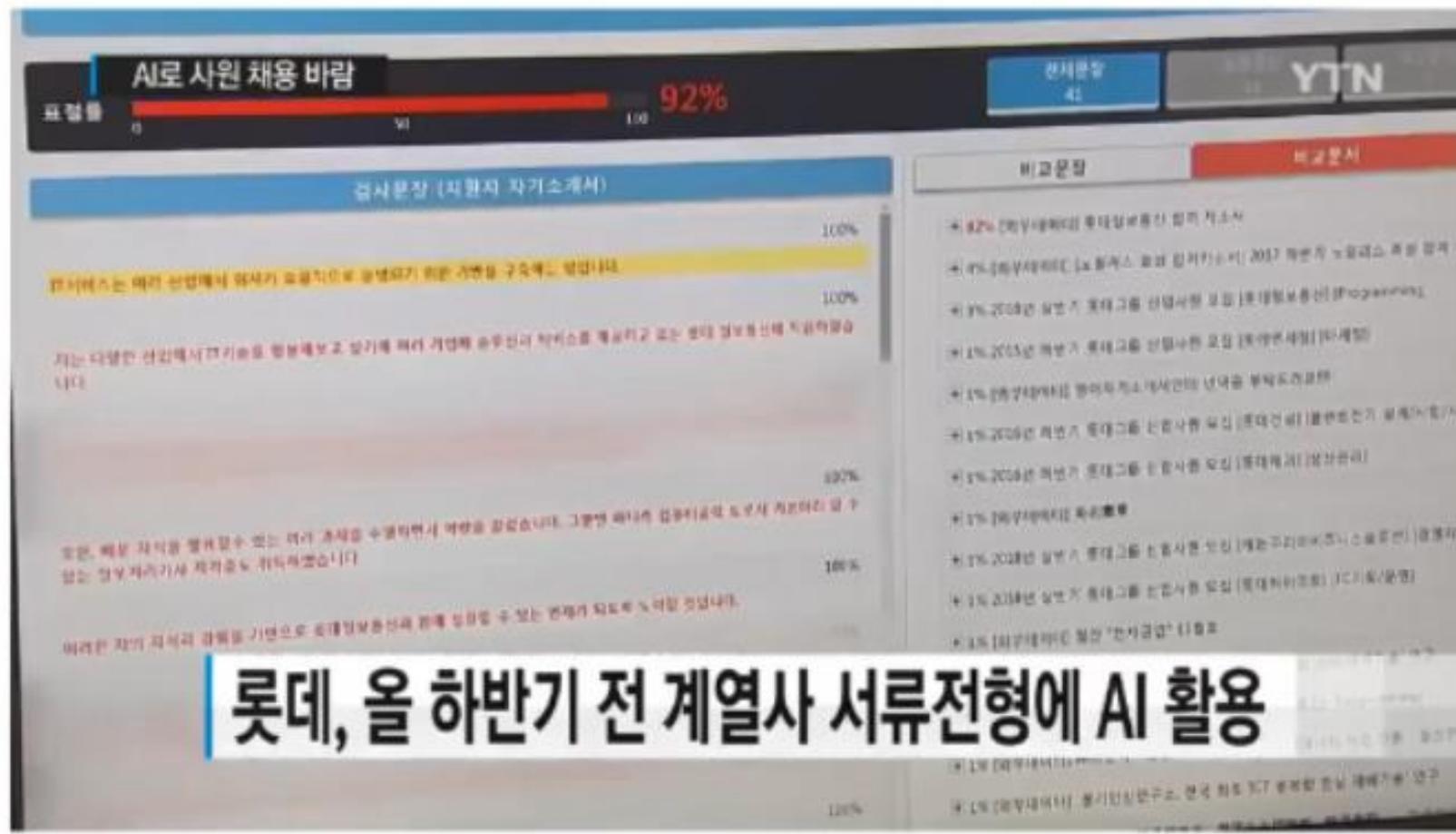
Premises: We can't expect all of our work for this project or resumes in general will be printed on the copies we are made or distributed for prints or commercial advertising and that expect these to be sold and the others to be sold. We can't expect all of our work for this project or resumes in general will be printed on the copies we are made or distributed for prints or commercial advertising and that expect these to be sold and the others to be sold. Abstaining with words is printed. To copy otherwise, or republish, in part or in whole, without the prior written permission and/or a fee, requires permission from permissions.org.

Keywords: resume, job application, self-introduction, deep learning, neural network, cosine similarity.

© 2018 Association for Computing Machinery.
978-1-4503-5626-0/18/06...\$15.00
<https://doi.org/10.1145/3220451.3220506>

This paper is licensed under a Creative Commons Attribution Non-Commercial-ShareAlike 4.0 International License.

CC-BY-NC-SA

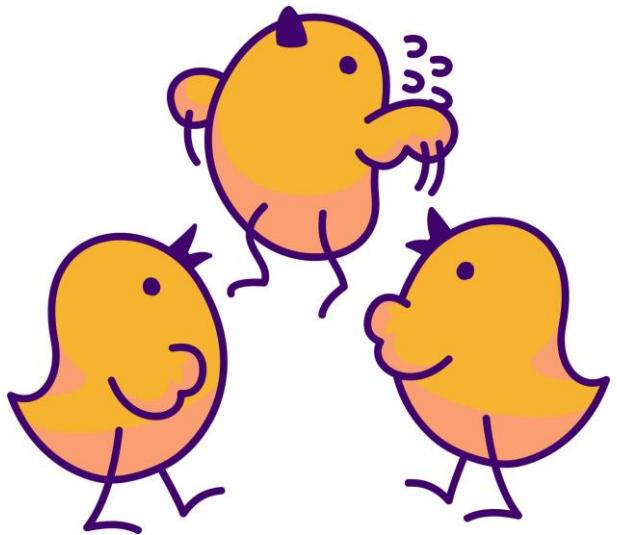


롯데그룹의
[채용AI]는
YTN, 전자신문 등
각종매체에서 소개된
AI 시스템입니다.



- ✓ 자기소개서는 여전히 사람이 읽고 있습니다!
- ✓ 저는 [채용AI] 기술 담당자이지 채용 담당자가 아닙니다!
Q. '어떻게 하면 딥러닝이 합격으로 분류하나요?'
A. '저도 궁금해요!'
- ✓ 본 세션은 딥러닝을 이미 알고 계신 분이 실무에 어떻게 활용되는지 궁금하신 분이 들으면 좋아요!



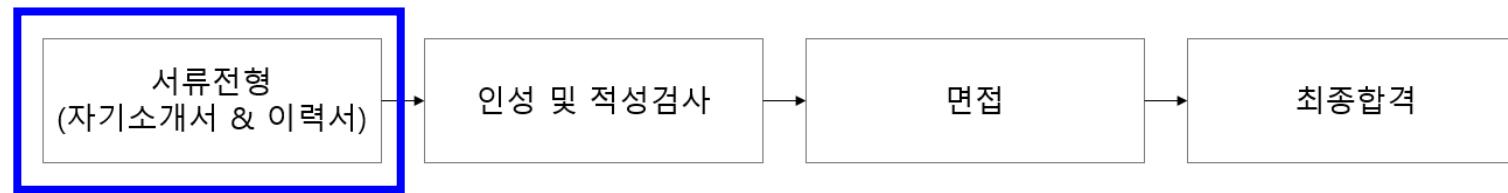


1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

기술 키워드 : 자기소개서 전처리, 딥러닝, TextRank, Bi-LSTM+CNN

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference



[신입사원 채용 프로세스(Process)]

- 2019년 하반기 롯데 신입사원 지원자 약 49,000명.
- 신입사원 채용은 이력서보다 자기소개서를 중요하게 평가.
 - 서류전형이 마감되면 인적성 검사일까지 10일안에 서류전형을 끝내야 한다.
 - 경력이 없는 신입사원의 업무 수행능력은 비슷하기 때문에 지원자의 잠재력을 자기소개서로 평가한다.
- 매 채용시 지원자는 동일한 질문을 제한된 글자 수 내에서 답변을 하고 있다. 그렇다면 **딥러닝**으로 우수한 자기소개서를 **이항분류(Binary Classification)** 할 수 있지 않을까?

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference

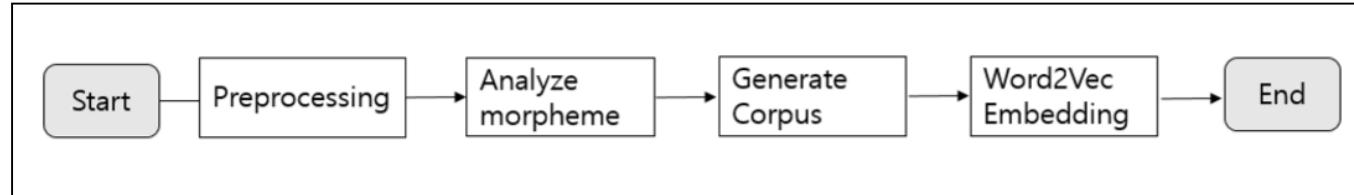
자기소개서 질문 및 답변 예시

Index	범주	질문 및 답변	제한 글자수 (공백포함)
1	지원동기	지원동기를 구체적으로 기술해주세요. [핵심은 사람이다] 학교 선배들, 지인 중 유통업에 종사하는 사람들이 회사에 대한 만족도가 높고...	500
2	성장과정	성장과정을 구체적으로 기술해주세요. [만화책을 위해 시작한 아르바이트] 중학교 1학년 때는 순정만화를 참 좋아했습...	800
3	사회활동	학업 이외에 관심과 열정을 가지고 했던 다양한 경험 중 가장 기억에 남는 것을 구체적으로 기술해주세요. [춤과 노래하기를 잊게 만든 브랜딩 마케팅] 브랜딩 마케팅은 학문이지만 저에게는...	800
4	직무경험	희망직무 준비과정과 희망직무에 대한 본인의 강점과 약점을 기술해주세요. [풍부한 프로젝트 경험] 업무 수행에 필요한 Java, JSP, Oracle Database 를 꾸준히...	800
5	입사 후 포부	입사 후 10년 동안의 회사생활 시나리오와 그것을 추구하는 이유를 기술해주세요. [고객의 니즈를 충족시키는 MD] 옴니 채널을 통해 유통시장에 변화를 불어넣고...	500

- 딥러닝의 시작은 **데이터 수집!**
- 지원자에게 연구용으로 데이터가 활용됨을 알리고 동의를 받았어요!
- 덕분에 개인정보가 비식별화된 데이터를 모을 수 있었답니다.

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference



[자기소개서 전처리 및 임베딩(Embedding) 순서도]

- 자기소개서의 한글, 영문, 숫자를 제외한 특수문자, 따옴표 등은 정규표현식을 사용하여 모두 제거한다.
- Okt 형태소 분석기로 자기소개서를 형태소 단위로 토큰나이징(Tokenizing)한다.
- 토큰나이징(Tokenizing)된 자기소개서 → 자기소개서 말뭉치(Corpus)
- Word2Vec Embedding은 CBOW 보다 성능이 좋은 중심단어로 주변 단어를 예측하는 Skip-gram
⇒ 100,365개 형태소(*128 차원으로 Embedding / 10 번 이상 등장한 단어로 제한 / Window 크기 : 10)

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference

[고양이 + 애교 = 강아지]

[손흥민 + 아이버슨(농구선수) = 박지성]

(*자기소개서에서 존경하는 인물로 자주 등장하는 인물)

[남자 - 여자 = 다른 사람]

```
w2v_model.wv.most_similar(positive=['고양이', '애교'], topn=1)
```

```
[('강아지', 0.7390154600143433)]
```

```
w2v_model.wv.most_similar(positive='남자', negative='여자', topn=1)
```

```
[('다른 사람', 0.5418267250061035)]
```

```
w2v_model.wv.most_similar(positive=['손흥민', '아이버슨'], topn=1)
```

```
[('박지성', 0.7613238096237183)]
```

```
w2v_model.wv.most_similar("부모님")
```

```
[('어머니', 0.8601336479187012),  
 ('어머님', 0.7858796715736389),  
 ('아버지', 0.7813066244125366),  
 ('아버님', 0.7686976194381714),  
 ('외할머니', 0.7564709782600403),  
 ('고모', 0.7508319616317749),  
 ('집안', 0.7503125667572021),  
 ('키워 주신', 0.7484757900238037),  
 ('키우셨습니다', 0.7364680171012878),  
 ('동생', 0.734329104423523)]
```

```
w2v_model.wv.most_similar("프로그래밍")
```

```
[('java', 0.8812177777290344),  
 ('프로그래밍언어', 0.878530740737915),  
 ('코딩', 0.8749273419380188),  
 ('자료구조', 0.8630567789077759),  
 ('자바', 0.8379415273666382),  
 ('액체지향', 0.828315258026123),  
 ('python', 0.8225173950195312),  
 ('소프트웨어', 0.8223742842674255),  
 ('임베디드', 0.8213320374488831),  
 ('운영체제', 0.8109732866287231)]
```

```
w2v_model.wv.most_similar("손흥민")
```

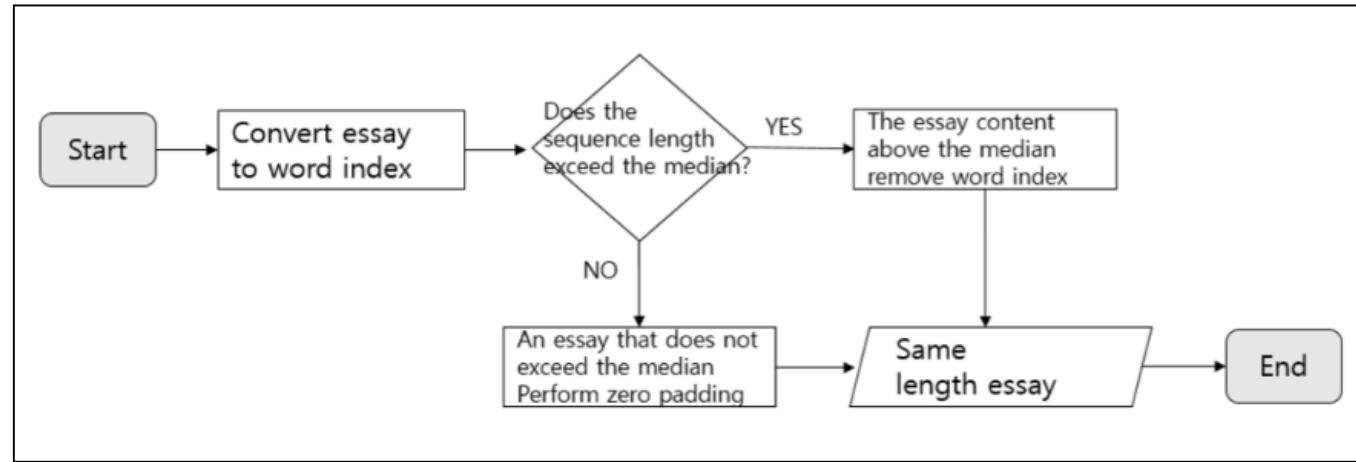
```
[('축구선수', 0.6439859867095947),  
 ('안정환', 0.629976212978363),  
 ('맨유', 0.6285330057144165),  
 ('박지성', 0.6210668683052063),  
 ('선수', 0.6134283542633057),  
 ('아이버슨', 0.6086547374725342),  
 ('토트넘', 0.591521143913269),  
 ('유나이티드', 0.5879682898521423),  
 ('프리미어리그', 0.580830454826355),  
 ('야구선수', 0.5790668725967407)]
```

임베딩(Embedding)01
잘 됐나 한 번
볼까요?



간합 보까요?

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?



[자기소개서 단어 인덱스 변경 순서도]

- 딥러닝 모델에 입력하기 위해 형태소를 사전의 **인덱스**로 변경한다.
- 각각의 질문에 대한 답변 별로 시퀀스(Sequence) 길이의 **중앙값**을 구한다.
(*평균이 아닌 중앙값을 기준으로 삼은 이유는 중앙값이 평균보다 이상치(Outlier)에 덜 민감하기 때문이다.)
- 답변 별로 동일한 길이의 시퀀스를 만들기 위해 중앙값보다 **짧은 답변**은 자기소개서 앞쪽에 **제로 패딩** 중앙값보다 **긴 답변**은 중앙값 이후 단어 인덱스를 **삭제**한다.

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference



그러나 베이스 모델의
정확도는 기대에 미치지
못했습니다.
무엇이 문제였을까요?

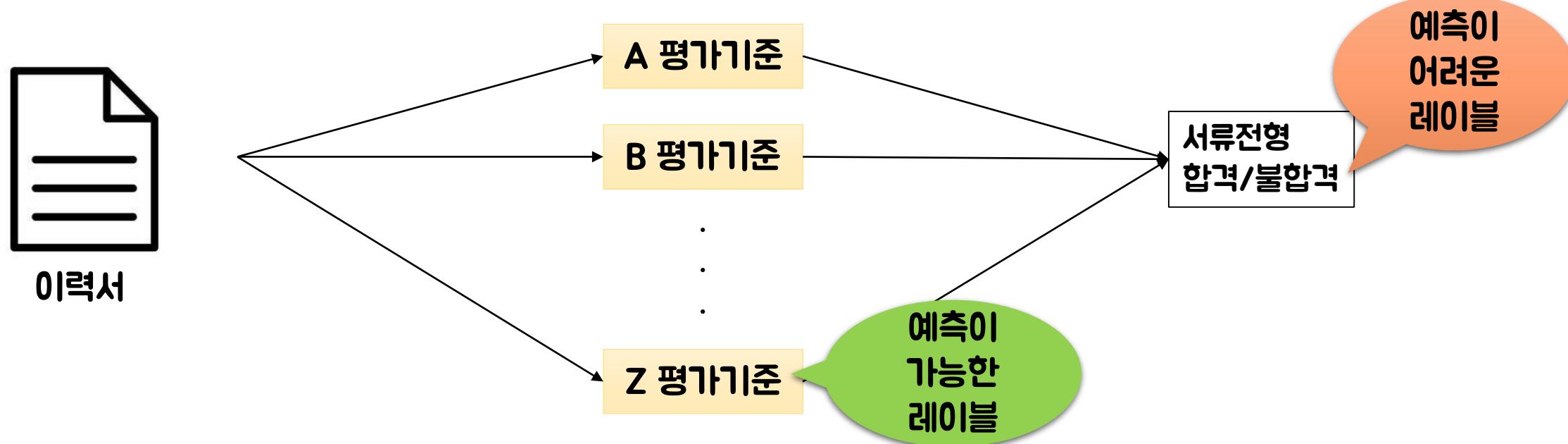
문제는 레이블링이었습니다.
서류전형 합격이라고 모두
우수한 자기소개가 아니고,
서류전형 탈락이라고 모두
우수하지 못한 자기소개서가
아니었던 것입니다.

★딥러닝 정확도 15% 이상 향상



1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference

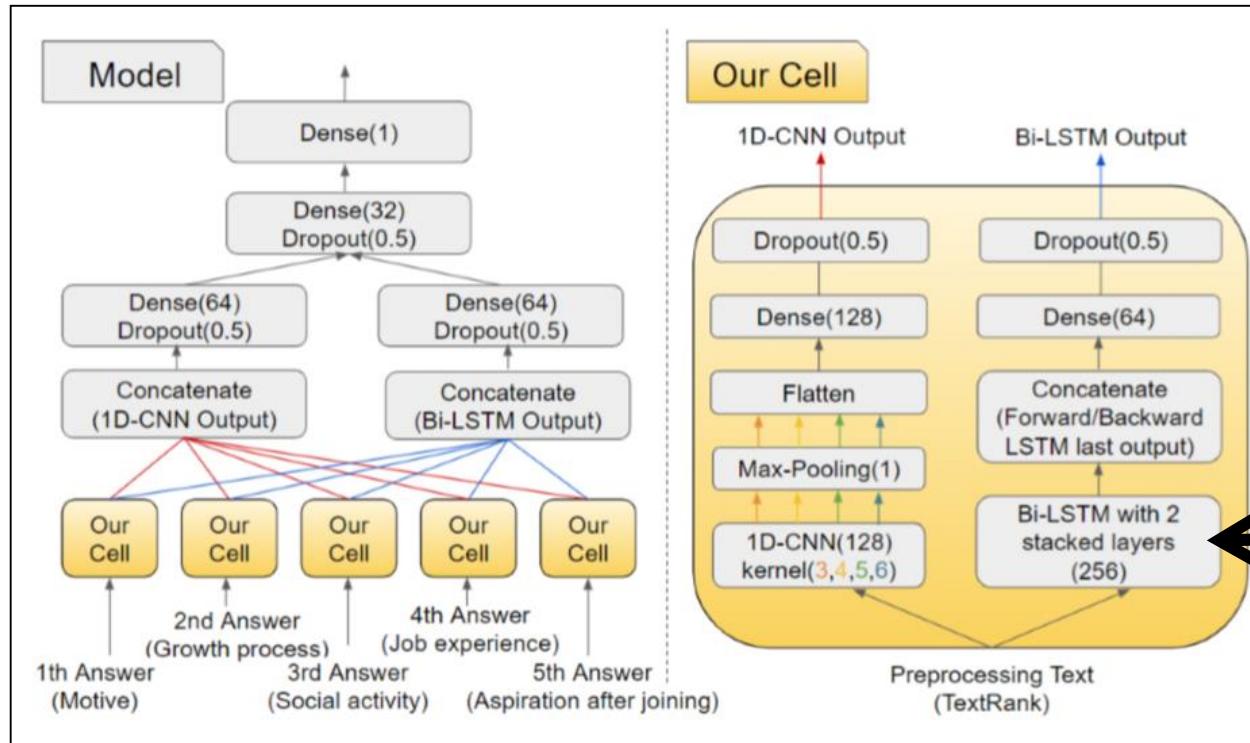


- 자기소개서를 읽고 우수함을 레이블링(Labeling) 하기에는 평가자의 주관적 요소가 들어갈 수 있을뿐더러 많은 비용을 필요로 함.
- 그룹사의 필요인재부합도(PFA)라고 하는 평가기준을 레이블링(Labeling) 하였을 때 높은 정확도를 보임.
(레이블링에 사용된 평가기준은 스펙을 평가하는 알고리즘.)
- 해당 평가기준에서 높은 점수를 받은 지원자의 자기소개서와 낮은 점수를 받은 지원자를 추출.
(높지도 낮지도 않은 애매한 점수를 받은 자기소개서는 학습을 어렵게하는 노이즈로 판단하여 제거함.)
- 딥러닝 모델로 스펙을 보지 않고 자기소개서만으로 우수한 자기소개서 여부를 가려낼 수 있게 됨.
(*딥러닝이 예측한 자기소개서를 실제로 채용 담당자분들이 읽고 겸증한 결과 사용 가능한 평가기준이라고 판단. → 44개 계열사 적용)

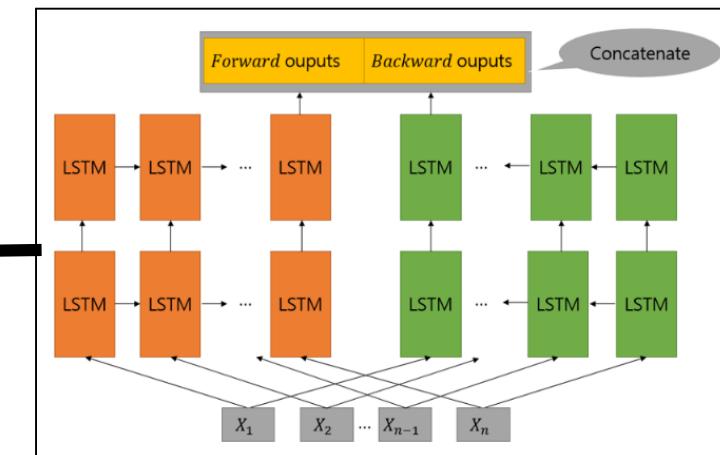
1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference

- 딥러닝 모델이 5개의 답변을 수용할 수 있도록 5개의 셀(Cell)을 두었다.
- 1D-CNN은 3, 4, 5, 6의 커널(Kernel)이 존재해서 3개, 4개, 5개, 6개의 연속된 시퀀스(Sequence)의 특징을 추출한다.
- Gradient vanishing을 방지하기 위해 미분값이 상대적으로 큰 활성화 함수(Activation Functions) 하이퍼볼릭 탄젠트 함수(tanh)를 사용한다.(Xavier Initialization)
- 마지막 노드는 0~1 사이의 값을 출력하기 위해 시그모이드(sigmoid) 함수를 사용한다.
- 손실함수는 크로스 엔트로피(Cross-entropy)를 사용한다.

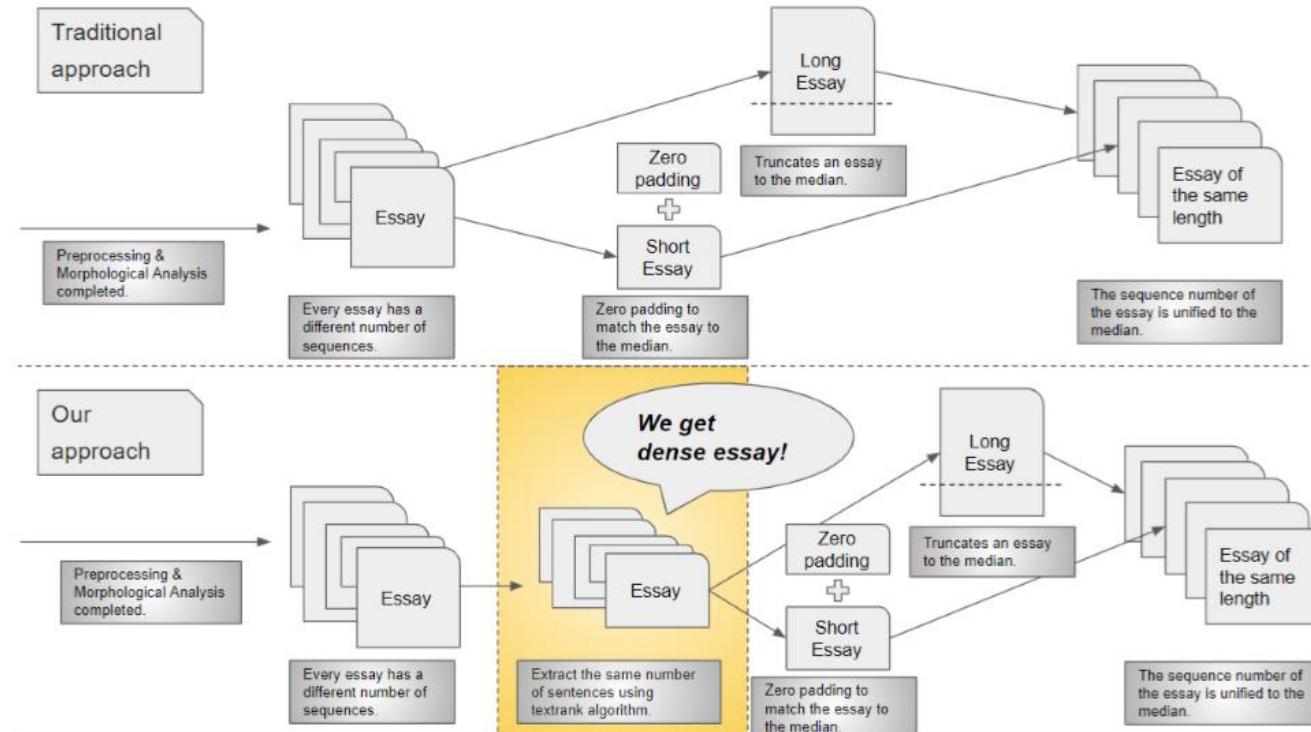


[자기소개서 분류 딥러닝 모델 예시]

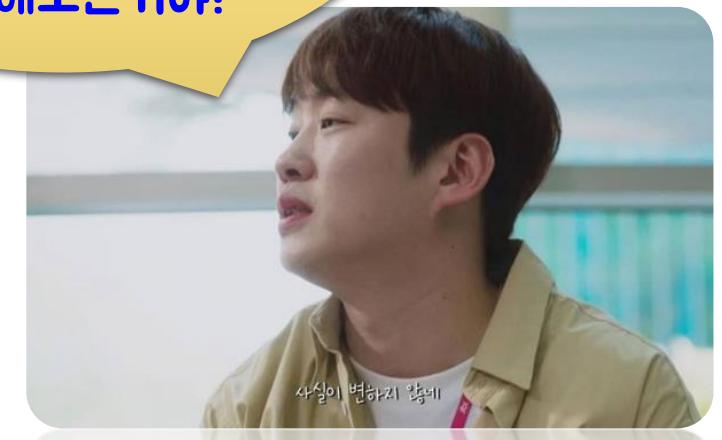


1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference



(아련하게) 어떻게
하면 정확도를 더 올릴
수 있을까...?
그래! 자기소개서를
요약해보는거야!



- **TextRank**를 활용하여 중요도가 낮은 문장을 제거한다.
*이는 딥러닝 모델의 정확도를 높일 수 있다. 이는 기계학습의 **특징 선택(Feature selection)**의 효과와 유사하다. 자연어는 한 문장 안에서 특정 형태소를 제거하면 의미가 훼손되기 때문에 전통적인 **특징 선택(Feature selection)** 기법을 적용하기 어렵다. 따라서 의미를 최대한 보존하면서 불필요한 노이즈(Noise)를 제거하기 위하여 중요도가 낮은 문장을 제거하는 전처리 과정을 추가했다.

1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference

[TextRank 활용한 자기소개서 요약 예시]

순서	TextRank 점수 순위	TextRank 점수	문장
1	2	1.45233	1년 동안의 교내 영어 학습 동아리 회장으로서...
2	3	1.25027	이른 오전 매일 7시 30분에 시작하는 동아리...
3	5	1.09548	동아리의 분위기는 구성원들의 활동에 의해 결정된다고 생각합니다. 그래서 구성원들의 자발적인 참여를 이끌어내기 위해 80% 이상의 출석률을 목표로 설정했습니다. 이 목표를 이루기 위해 릴레이 모닝콜을 제안했습니다. 학교와 집과의 거리가 먼 사람부터 가까운 사람에게 한명씩 모닝콜을 하는 시스템입니다. 그 결과 출석률을 높이는데 큰 도움이 되었습니다. 또한 사소한 대화로 시작하여 작은 소모임을 결성하여 친밀감을 형성을 이루려고 노력했습니다. 그래서 신뢰를 바탕으로 자발적 참여를 유도하여 화기애애한 분위기로 바꿀 수 있었습니다. 이러한 노력 끝에 수동적 참여가 아닌 스스로 흥미를 가진 능동적 참여를 이끌어 냈습니다. 그래서 한 학기 동안 낙오자가 없고 다른 어떤 동아리보다도 결속력 있는 영어학습동아리를 만들 수 있었습니다.
4	1	1.49881	그래서 구성원들의 자발적인 참여를...
5	8	0.93056	이 목표를 이루기 위해 릴레이 모닝콜...
6	9	0.91456	학교와 집과의 거리가 먼 사람부터 가까운...
7	13	0.63766	그 결과 출석률을 높이는데 큰 도움이 되었습니다.
8	11	0.84565	또한 사소한 대화로 시작하여 작은 소모임을 결성...
9	10	0.86788	그래서 신뢰를 바탕으로 자발적인 참여를 유도하여...
10	14	0.62134	이러한 노력 끝에 수동적 참여가 아닌 스스로 흥미...
11	12	0.67191	그래서 한 학기 동안 낙오자가 없고 다른 어떤...
12	6	1.01100	4년 이 지난 지금도 구성원들과 그 당시 추억...
13	7	0.98664	이처럼 구성원들과의 무한한 신뢰 온정과 수용...
14	4	1.21590	앞으로 롯데 정보통신에서 구성원들을 아우르는...

1년 동안의 교내 영어학습동아리 회장으로서 구성원들과 소통하고 협력하여 능동적 참여를 이끌어 낸 경험이 있습니다. 이른 오전 매일 7시 30분에 시작하는 동아리이기 때문에 구성원들의 늦잠으로 인해 출석률이 낮아지는 것에 대한 고민을 했습니다. 동아리의 분위기는 구성원들의 참여율에 의해 결정된다고 생각했습니다. 그래서 구성원들의 자발적인 참여를 이끌어내기 위해 80% 이상의 출석률을 목표로 설정했습니다. 이 목표를 이루기 위해 릴레이 모닝콜을 제안했습니다. 학교와 집과의 거리가 먼 사람부터 가까운 사람에게 한명씩 모닝콜을 하는 시스템입니다. 그 결과 출석률을 높이는데 큰 도움이 되었습니다. 또한 사소한 대화로 시작하여 작은 소모임을 결성하여 친밀감을 형성을 이루려고 노력했습니다. 그래서 신뢰를 바탕으로 자발적 참여를 유도하여 화기애애한 분위기로 바꿀 수 있었습니다. 이러한 노력 끝에 수동적 참여가 아닌 스스로 흥미를 가진 능동적 참여를 이끌어 냈습니다. 그래서 한 학기 동안 낙오자가 없고 다른 어떤 동아리보다도 결속력 있는 영어학습동아리를 만들 수 있었습니다.

*TextRank
점수가 가장 낮은
문장이 요약된다.

$$\text{similarity}(S_1, S_2) = \frac{|\{m_k | m_k \in S_1 \text{ & } m_k \in S_2\}|}{\log(|S_1|) + \log(|S_2|)}$$

$$\text{TR}(V_i) = (1 - d) + d * \sum_{V_j \in \text{In}(V_i)} \frac{w_{ji}}{\sum_{V_k \in \text{Out}(V_j)} w_{jk}} \text{TR}(V_j)$$

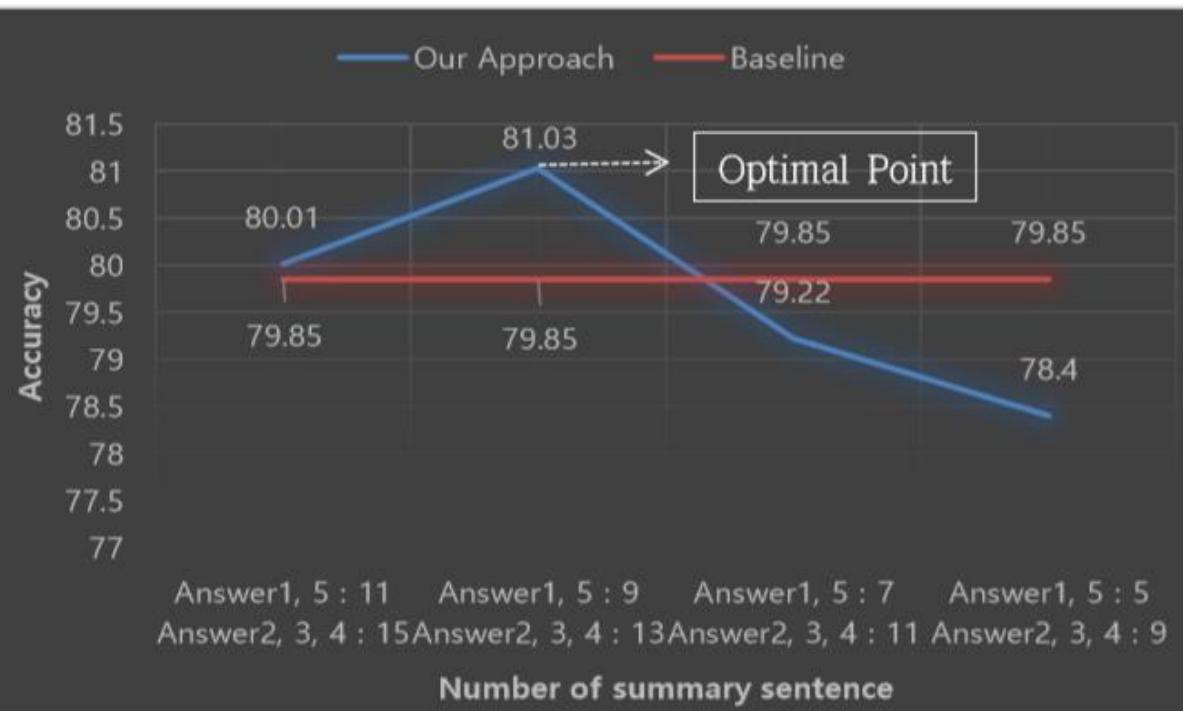
*참고 논문

- [1] Mihalcea, R., & Tarau, P. (2004). Textrank: Bringing order into text. In Proceedings of the 2004 conference on empirical methods in natural language processing
- [2] Federico Barrios, Federico López, Luis Argerich, and Rosa Wachenchauzer. (2016). Variations of the similarity function of TextRank for automated summarization. arXiv:1602.03606[cs.CL]

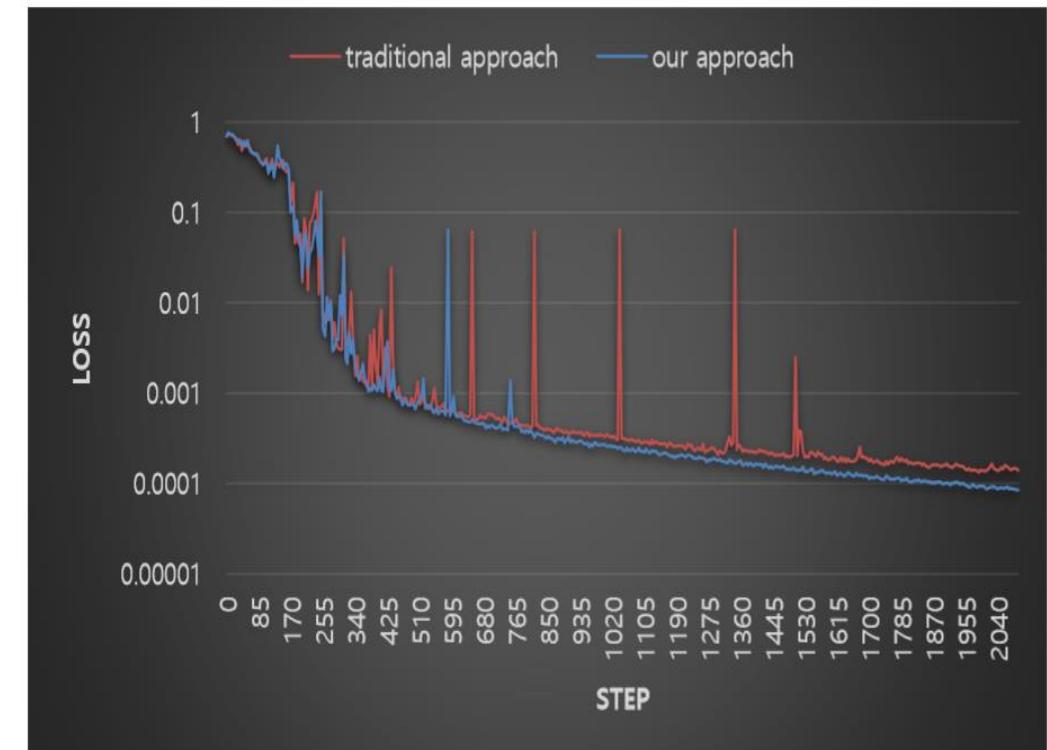
1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

- **(정확도 향상)** 제한글자가 500자인 자기소개서는 문장 수는 **9** 문장, 제한글자가 800자인 답변은 평균 문장 수는 **13** 문장일 때 성능이 가장 좋았습니다. 이 이상 요약할 때 정확도가 기존 방법보다 떨어지는 것으로 미루어 합격 및 불합격을 결정짓는 중요한 문장마저 지워진다고 할 수 있다.
- **(학습 속도 향상)** 요약된 자기소개서를 학습하는 딥러닝 모델은 1,865Step 부터 Loss 가 0.0001 이하로 떨어지지만, 자기소개서 원문을 학습하는 딥러닝 모델은 Loss 가 에폭(Epoch)을 모두 수행해도 0.0001 이하로 떨어지지 않았다.

[핵심문장 수별 딥러닝 모델 정확도 변화]



[학습 스텝별 손실변화 추이]



1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

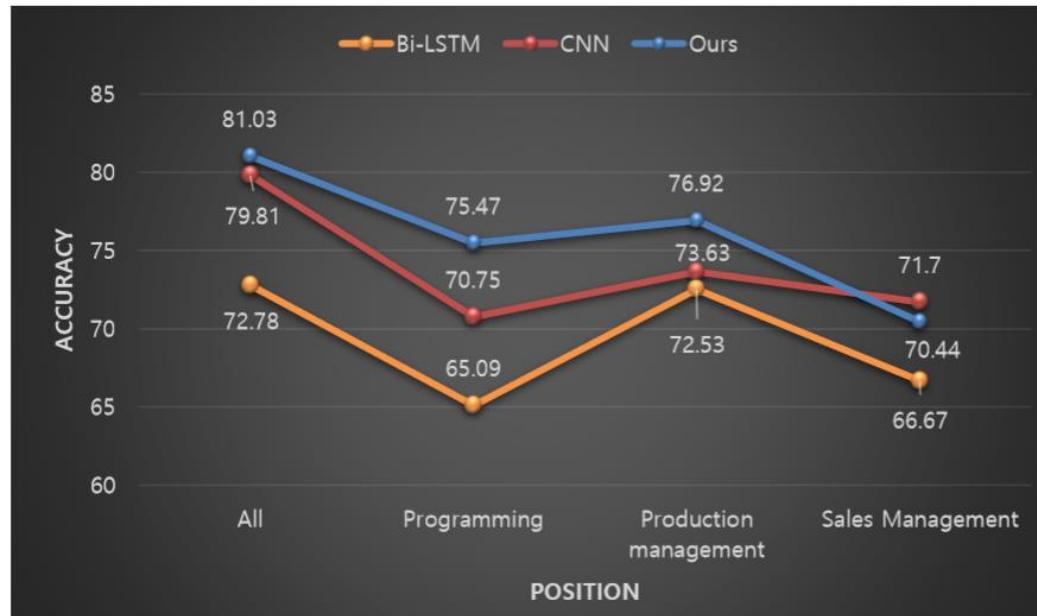
- Train data : 2016년 ~ 2018년 상반기 자기소개서 / Test data : 2018년 하반기 자기소개서
- Ours의 F1-Score는 각각 0.8, 0.82로 합격 및 불합격 중 한쪽으로 치우치지 않고 균형 잡힌 정확도를 보였다.
- 직무를 구분했을 때 오히려 정확도가 내려갔다. 원인은 첫째, 직무를 구분하면서 학습 데이터 수가 줄었다. 둘째, 직무를 구분하면서 학습 데이터(Train Data)의 합격 및 불합격의 비율이 1:1이 아니게 되었다. 셋째, 경력이 아닌 신입사원이기 때문에 직무 특성이 자기소개서에 크게 드러나지 않는다.
- 자기소개서는 데이터가 부족한 도메인(Domain)인 만큼 딥러닝의 특성상 직무별로 데이터가 충분하지 않다면 직무를 구분하지 않는 것이 오히려 좋은 성능을 낼 수 있다.

Table 4: Comparison of Performances

Position	Model	Precision		Recall		F1-score		Accuracy
		Fail(0)	Pass(1)	Fail(0)	Pass(1)	Fail(0)	Pass(1)	
All	Bi-LSTM	.69	.77	.76	.70	.72	.73	72.78%
	1D-CNN	.72	.88	.87	.75	.79	.81	79.81%
	Ours	.74	.88	.87	.76	.80	.82	81.03%
Programming	Bi-LSTM	.69	.77	.76	.70	.72	.73	72.78%
	1D-CNN	.72	.88	.87	.75	.79	.81	79.81%
	Ours	.74	.88	.87	.76	.80	.82	81.03%
Production Management	Bi-LSTM	.69	.77	.76	.70	.72	.73	72.78%
	1D-CNN	.72	.88	.87	.75	.79	.81	79.81%
	Ours	.74	.88	.87	.76	.80	.82	81.03%
Others	Bi-LSTM	.69	.77	.76	.70	.72	.73	72.78%
	1D-CNN	.72	.88	.87	.75	.79	.81	79.81%
	Ours	.74	.88	.87	.76	.80	.82	81.03%

Table 2: Dataset Distribution for Training and Test

Job Description	Training dataset		Test dataset	
	Failed(0)	Passed(1)	Failed(0)	Passed(1)
Sales Management	852(41.6%)	1,195(58.4%)	93(58.5%)	66(41.5%)
Production Management	583(62.3%)	353(37.7%)	66(72.5%)	25(27.5%)
Programming	237(45.1%)	288(54.9%)	24(22.6%)	82(77.4%)
Others	3,670(51.1%)	3,506(48.9%)	1,138(52%)	1,052(48%)
Total	5,342(50%)	5,342(50%)	1,321(51.9%)	1,225(48.1%)

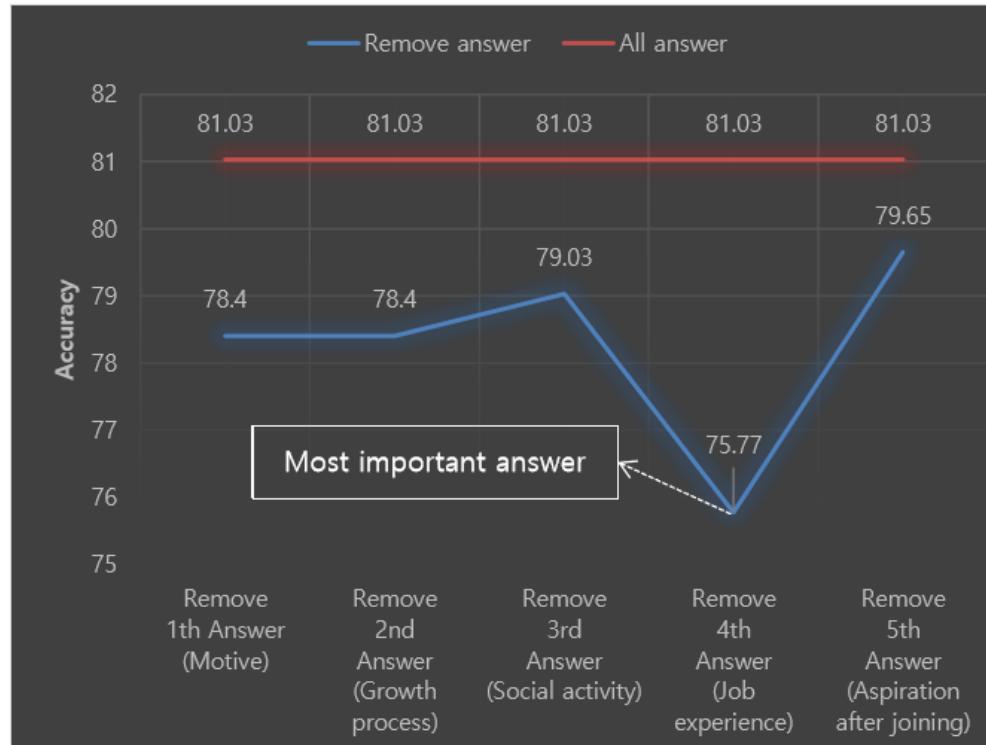
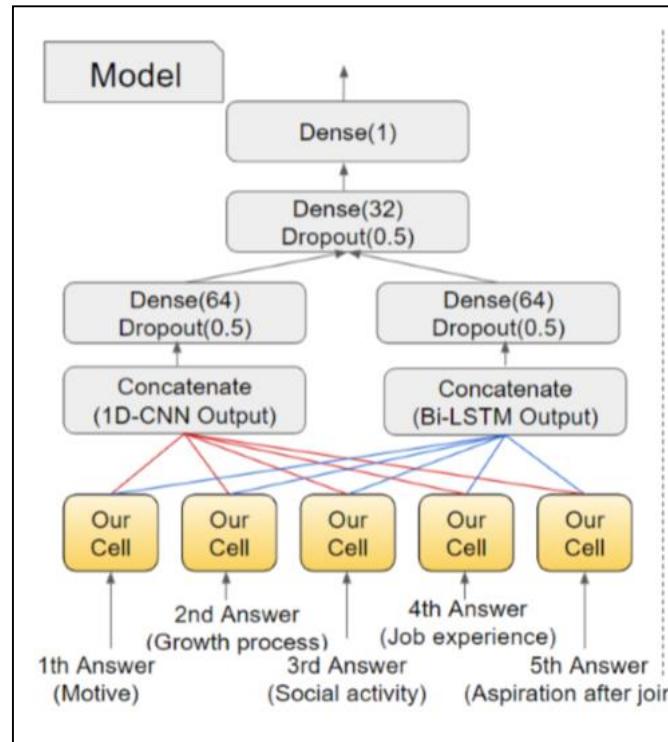


에잇! 이놈도 아녀!
최고의 성능을
내는 놈만
살아남는겨!



1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

- 다음 딥러닝 모델의 Cell은 각각의 답변이 입력된다.
특정 Cell을 제거하면 해당 답변이 입력되지 않은 딥러닝 모델을 얻을 수 있다.
- 지원동기 / 성장과정 / 사회활동 / 직무경험 / 10년계획 답변 중에서 가장 중요한 질문은 '**직무경험**' 이었다.



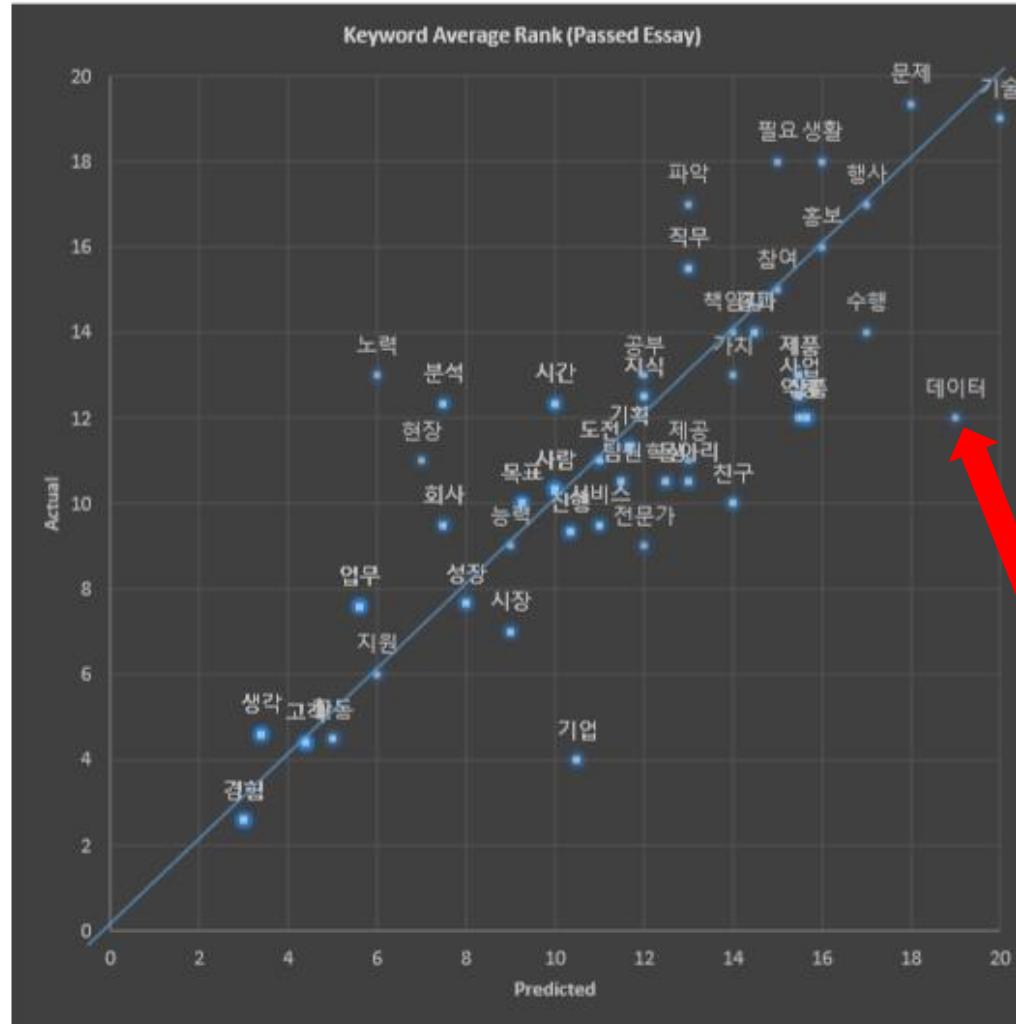
뭣이 중헌디!
어떤 질문이
중허냐고!



1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

Language
Conference

합격 자기소개서 키워드 비교



- 자기소개서를 Komoran 형태소 분석기로 명사 추출 ->
-> 1글자 제거 -> TextRank로 키워드 추출 ->
추출된 키워드 빈도수로 순위 부여 ->
- 정확도가 100% 딥러닝 모델은 키워드가
모두 $y = x$ 직선에 위에 놓인다.
- $y = x$ 직선에서 멀어질수록 실제와 예측간 순위 차이가 크다.
- $y = x$ 직선보다 아래쪽에 있는 키워드는
실제 순위가 예측 순위보다 높다.
- [데이터] 키워드는 채용 담당자는 합격으로 많이 분류했지만,
딥러닝 모델은 그보다 적게 분류했음을 알 수 있다.

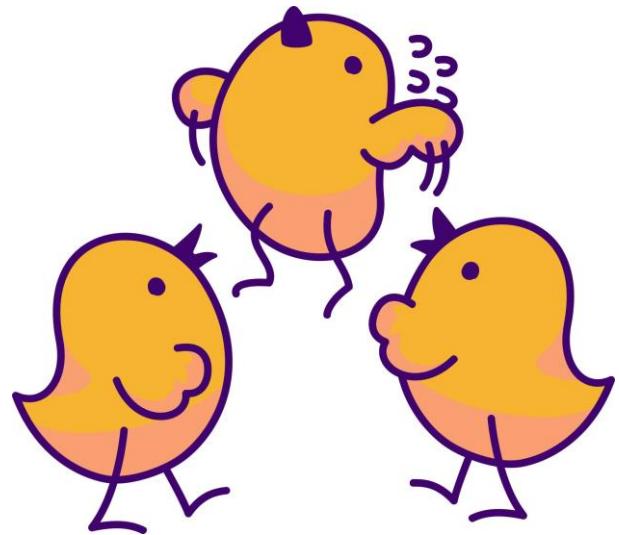
1. 딥러닝으로 우수한 자기소개서를 분류할 수 있을까?

- 혼동행렬을 기준으로 각 영역별 키워드 표현
 - 많이 등장한 키워드일수록 크고, 중앙에 표현
 - [생각], [업무] 등 무난한 의미를 지닌 단어들은 **불용어** 처리
 - 합격자기소개서에 [동아리], [사업] 등의 키워드가 상대적으로 더 많았다.



특정 단어를 쓴다고
혹은 안 쓴다고
합격하거나
불합격하지 않습니다!
오해없기!!

	Pass(Predicted)	Fail(Predicted)
Pass (Actual)	<p>True Positive</p>	<p>False Negative</p>
Fail (Actual)	<p>False Positive</p>	<p>True Negative</p>



2. 자기소개서의 어색한 문장을 찾을 수 있을까?

기술 키워드 : 딥러닝, 유사도, MaLSTM, Transformer

2. 자기소개서의 어색한 문장을 찾을 수 있을까?

Language
Conference

- **Business 요건** : 자기소개서에서 **어색한 문장을** 찾아주세요.
- 외부업체로부터 들어온 요구
- 해당 모델로 자기소개서 중간에 '엄마가 밥 먹으래.'라는 문장이 실제로 검출됨.
- **시연영상** : <https://www.youtube.com/watch?v=klUiScycydc>

The screenshot displays two side-by-side resume analysis results from the Chaeung AI ARAVOM platform.

Left Panel (Screenshot 1):

- Input:** Company: 롯데, Category: 사회활동
- Resume Content:** 성으로 참여한 풀리마켓에서 평생 나눌 일 없을 줄 알았던 일리스트도 생각보다 다를 만했고 배달 전화조차 잘 못 했던 제가 여러 관계자와 수십 번 미팅을 했습니다. 저는 여전히 제 것에 자부심을 느끼고 소중하게 생각합니다. 이후로 어떤 새롭고 어려운 문제에 직면해도 '일단 시작하자'라는 생각을 하게 되었습니다. 또, 누군가와 의견차가 있을 때 이해를 못했던 지난날과 달리 그들의 의견을 듣고 협력할 수 있게 되었습니다. 엄마가 밥 먹으래. 이 경험을 통해 앞으로 새로운 문제를 용기 있게 다룰 수 있고 좋은 협력을 보여줄 수 있으며 맙은 업무를 책임감 있게 수행하고, 저의 프로그램을 소중하게 생각하는 사람이 될 수 있을 것입니다.
- Analysis Results:**
 - Keywords:** 일리스트, 주식, 디자인, 프로그램
 - Main Sentence:** [좋아하는 가수와 함께한 풀리마켓] 휴학 기간에 좋아하는 가수가 풀리마켓을 주최한다는 소식을 접했습니다.
 - Awkward Sentences (Red Box):** 엄마가 밥 먹으래
 - Context-based Stopword:** 비정상적으로 타사명이 언급된 문장은 없습니다.

Right Panel (Screenshot 2):

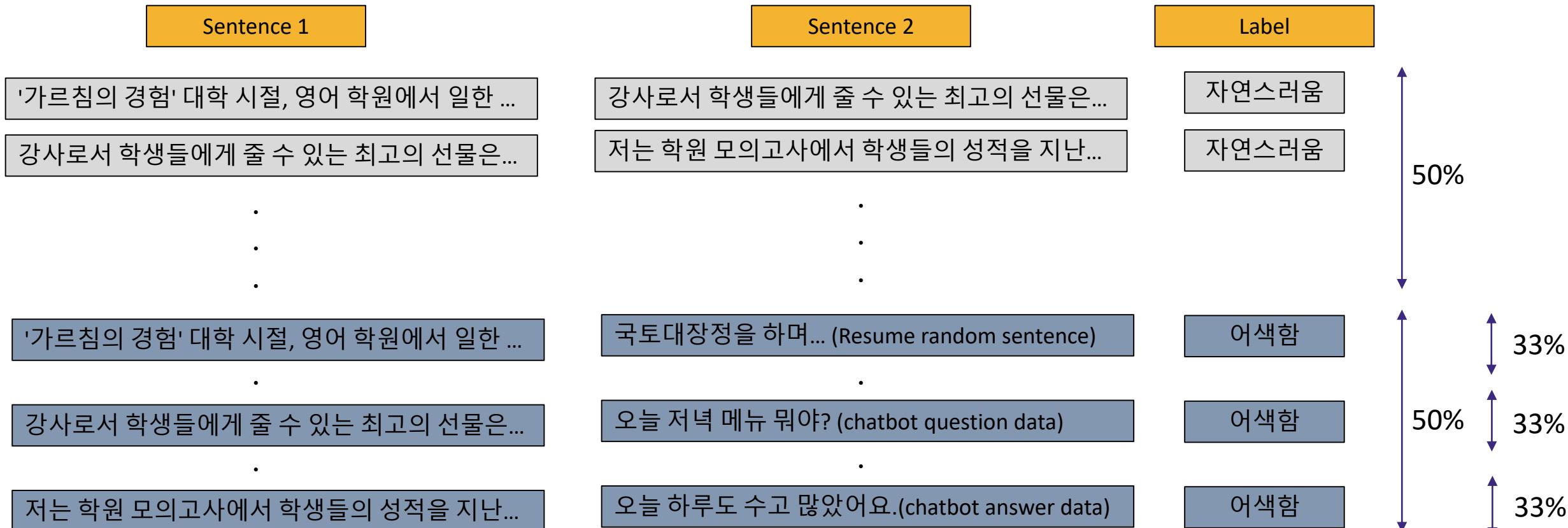
- Input:** Company: 롯데, Category: 사회활동
- Resume Content:** 을 진행하고 wireshark를 사용하여 packet을 분석했습니다. 네트워크 구축에 성공한 이후에는 다른 팀원이 진행 중인 RTOS와 Firmware 기반의 M3 board 코딩에 도움을 주었습니다. IoT 프로젝트 경험이 있었기 때문에 좀 더 적극적으로 참여할 수 있었습니다. 팀원들이 지칠 때마다 먼저 힘을 내어 열심히 하는 모습을 보여주니 참여율이 낮았던 팀원의 참여가 자연스럽게 높아졌고 프로젝트의 완성도가 높아졌습니다. 공익근무요원으로서는 힘들 때도 있지만 그 속에서는 언제나 작은 재미도 있었기에 공무원들도 저희와 일을 하는 것을 즐거워하기 시작하였습니다. 이를 통해 먼저 나서서 적극성을 보여준다면 팀의 사기가 증가한다는 것을 알았습니다. 프로젝트를 통해 1437 bytes
- Analysis Results:**
 - Keywords:** 실패, 감색, 제미, 프로젝트
 - Main Sentence:** 실패 극복을 위한 열정 '스마트 디바이스' 전문가 과정에 참여하여 주차 시스템 IoT 팀 프로젝트를 진행한 경험이 가장 기억에 남는 경험입니다.
 - Awkward Sentences (Red Box):** 공익근무요원으로서는 힘들 때도 있지만 그 속에서는 언제나 작은 재미도 있었기에 공무원들도 저희와 일을 하는 것을 즐거워하기 시작하였습니다.
 - Context-based Stopword:** 비정상적으로 타사명이 언급된 문장은 없습니다.

〈어색한 문장 검출 예시〉

2. 자기소개서의 어색한 문장을 찾을 수 있을까?

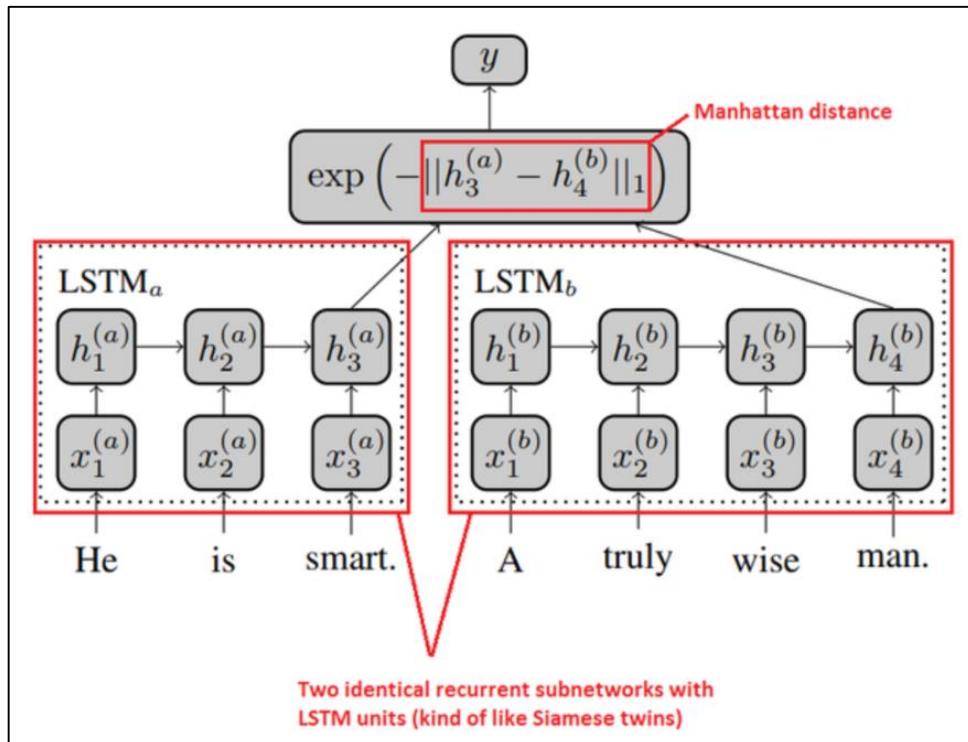
- **Data 구성**

- 우수한 자기소개서를 문장별로 분리하여 연속된 문장을 Pairing



2. 자기소개서의 어색한 문장을 찾을 수 있을까?

- BaseModel : [3]MaLSTM 딥러닝 모델
- Test Data 정확도 : 76%



A screenshot of a software interface titled "제 응 AI ARAVOM". The interface includes fields for "회사명" (Company Name) and "자기 소개 서" (Resume). On the right, there are several sections: "문장 분석 결과" (Text Analysis Result), "키워드(Keywords)", "주요 문장(Main sentence)", "어색한 문장(Awkward sentences)", "문맥 기반 불용어(Context-based stopword)", and "모호한 단어(Ambiguous words)". A large blue arrow points from a thought bubble containing Korean text to the "어색한 문장(Awkward sentences)" section, which is highlighted with a red box.

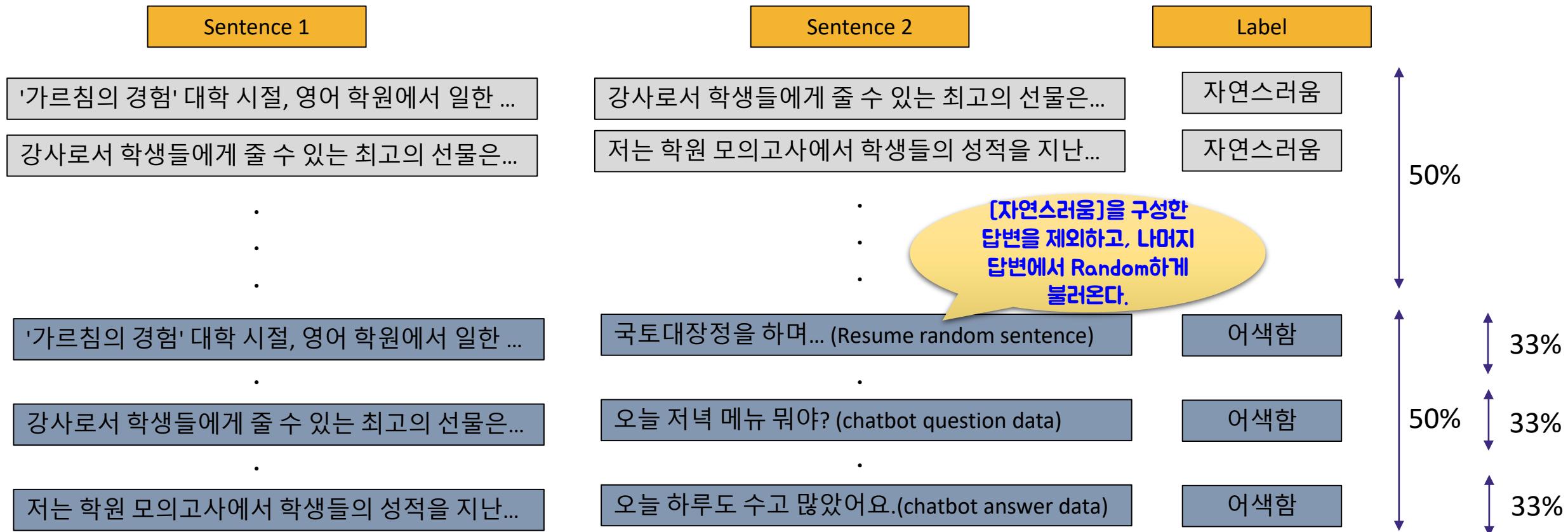
어떻게 하면 정확도를 더 올릴 수 있을까?
응답구분별로 모델을 분리하자!

*참고 논문

[3] Siamese Recurrent Architectures for Learning Sentence Similarity

2. 자기소개서의 어색한 문장을 찾을 수 있을까?

- 예를 들어 '사회활동' 모델의 데이터 구성 시, 해당 답변의 자기소개서로만 [자연스러움] 데이터를 구성한다.
- 더불어 Resume random sentence는 '사회활동'을 제외한 지원동기, 성장과정, 직무경험 등에서 랜덤하게 추출한다.
- 이를 통해 딥러닝 모델이 해당 답변의 의미에 집중할 수 있게 된다. (정확도 : 76% → 84%(8% 향상))

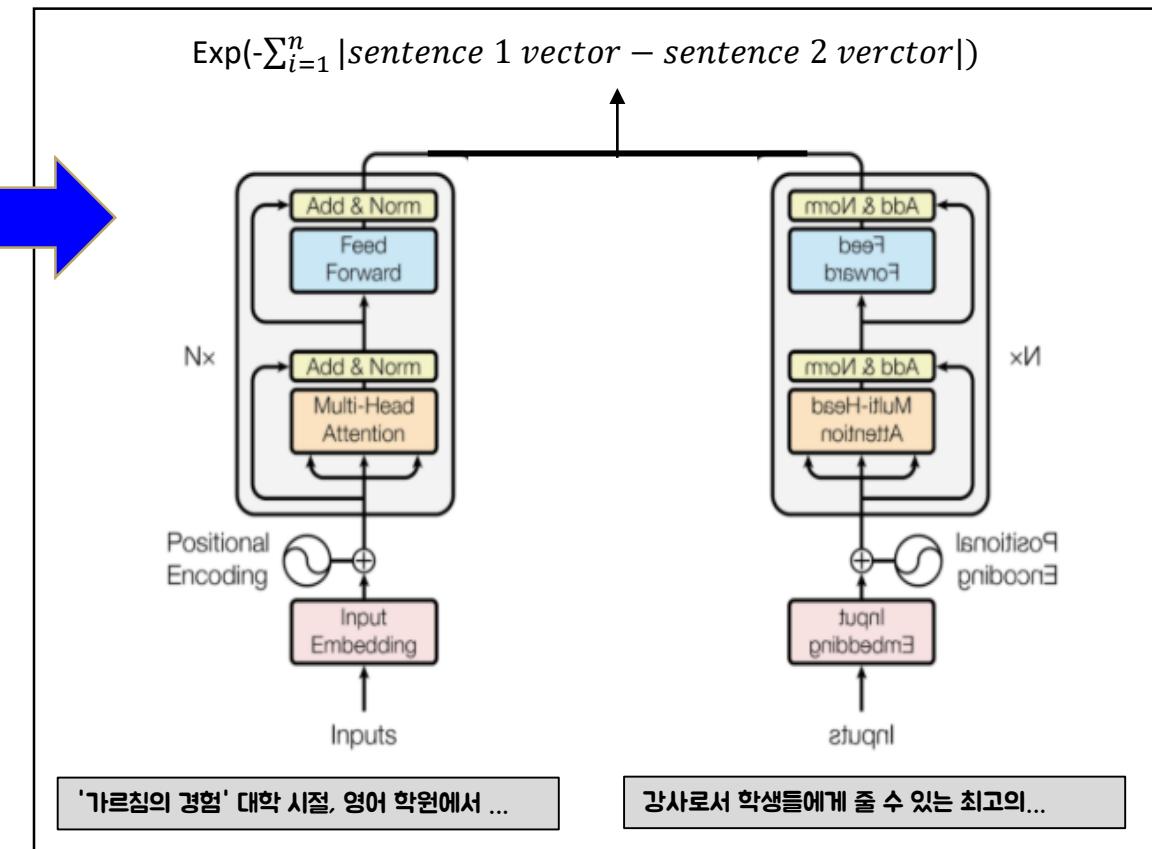
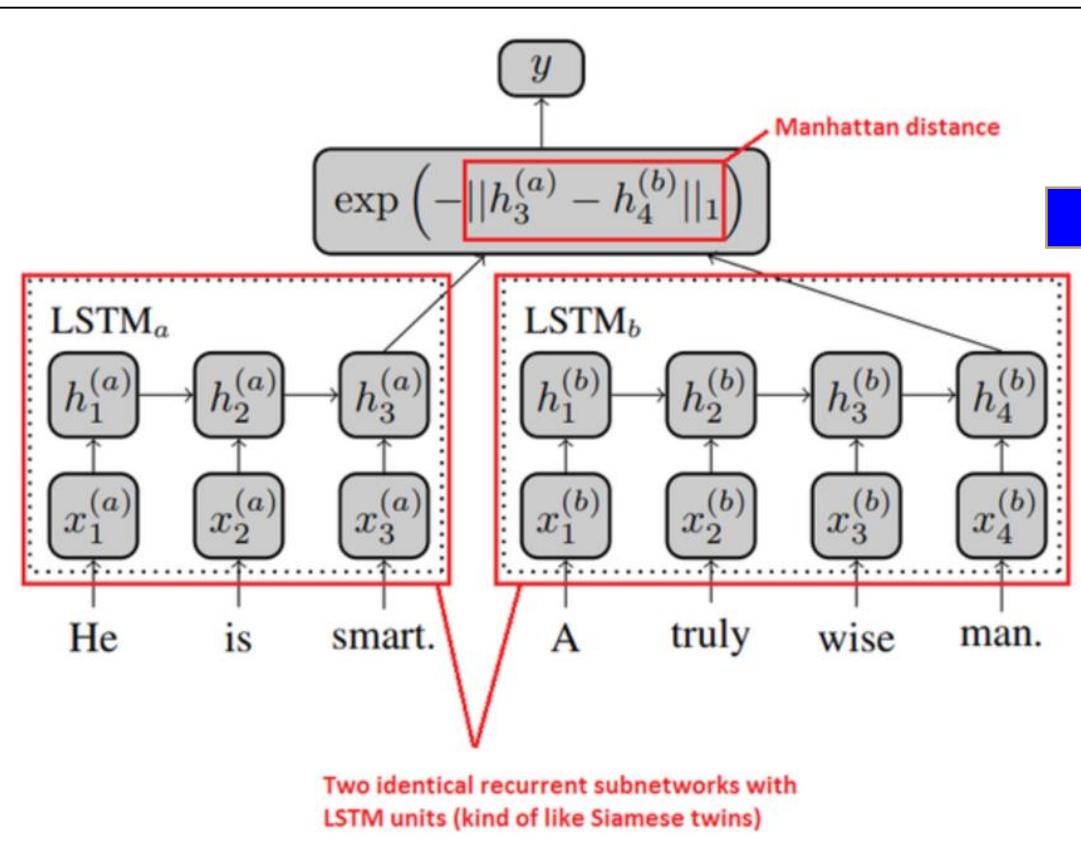


2. 트랜스포머를 활용한 자기소개서 문장 생성

- [4] Transformer Encoder를 활용한 딥러닝 모델로 고도화
- 정확도 : 74% -> 84%-> 88%(4% 향상)



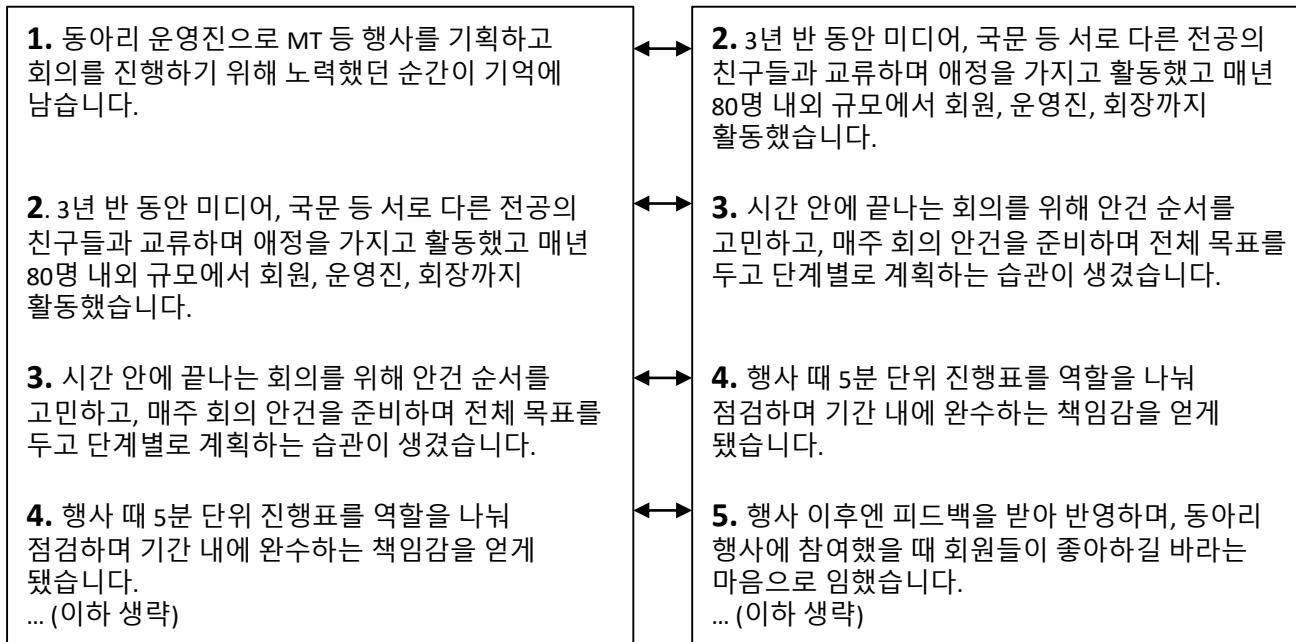
이제
추론 속도를
빠르게 할
방법은 없을까?



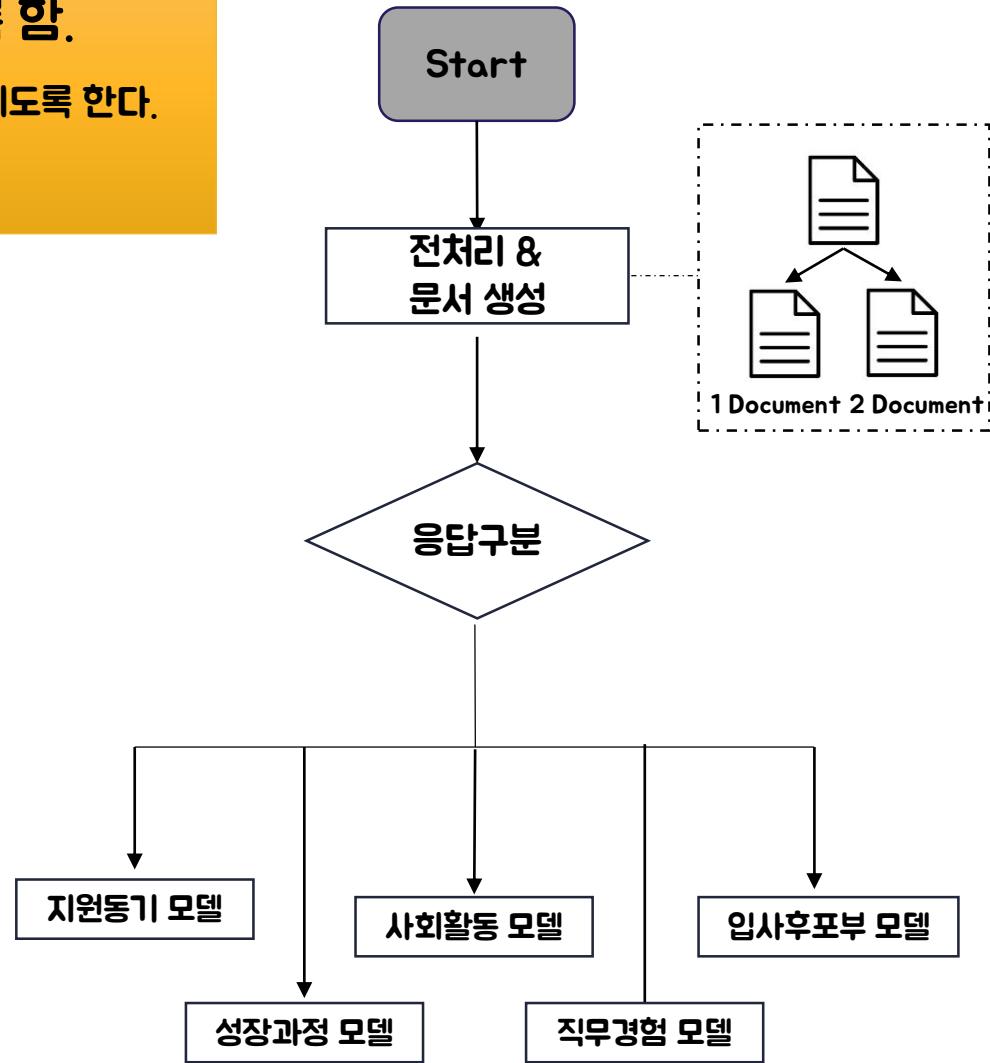
*참고 논문

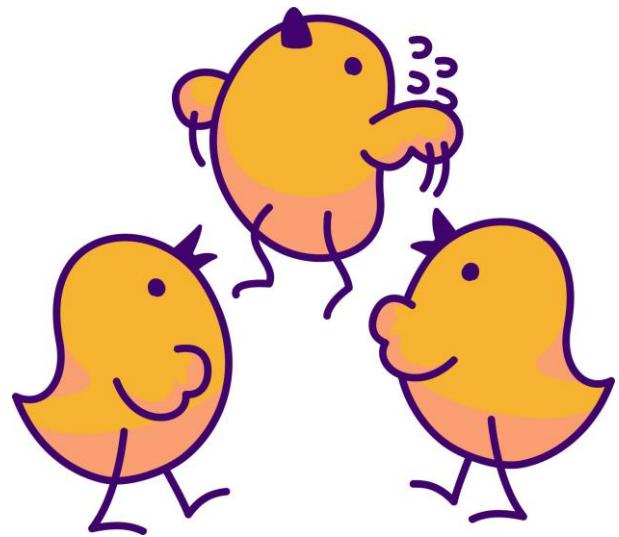
2. 트랜스포머를 활용한 자기소개서 문장 생성

- 딥러닝 추론 속도 향상을 위해 문장이 아닌 문서 단위 배치로 추론 되도록 함.
- [마지막 문장이 제거된 문서]와 [첫번째 문장이 제거된 문서]를 만들어 연속된 문장이 Paring 되도록 한다.
- 자기소개서 1건 기준 추론 속도 32% 향상



2 Document
(remove first sentence)





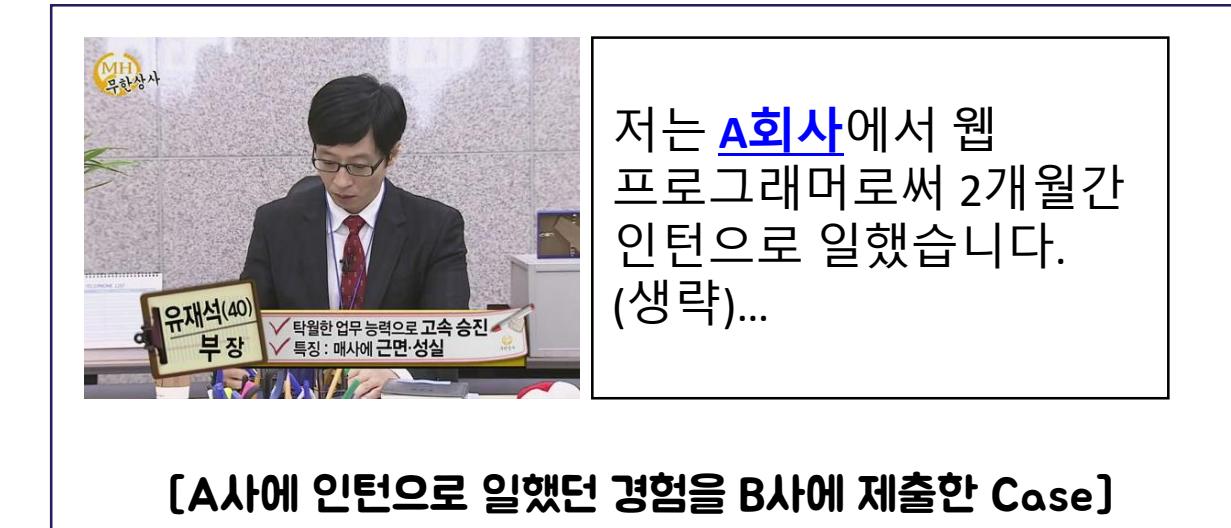
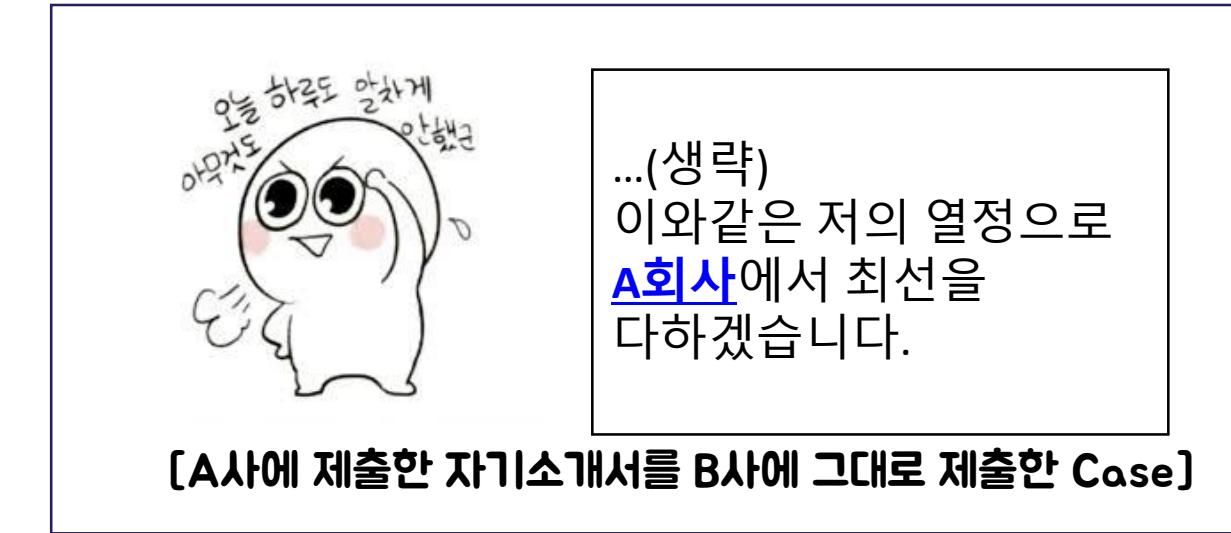
잡았다 요놈!



3. [Ctrl+C, Ctrl+V] 자기소개서 어떻게 찾을까? 기술 키워드 : 딥러닝, BERT

3. [Ctrl+C, Ctrl+V] 자기소개서 어떻게 찾을까?

Language
Conference



3. [Ctrl+C, Ctrl+V] 자기소개서 어떻게 찾을까?

Language
Conference

- 본 세션에서는 딥러닝을 활용하여 타 회사에 지원한 자기소개서를 그대로 복사, 붙여넣기 했을 때 어떻게 찾아낼 수 있는지 소개해드리겠습니다.(표절에 대한 기술은 제외)
- 시연영상 : <https://www.youtube.com/watch?v=klUiScycydc>



채용 AI ARAVOM

회사명

롯데

용답구분

사회활동

자기소개서

시각화한 자료를 보여주면 결정에 도움이 될 것으로 판단해 Onewapp으로 간단한 프로토타입을 만들었습니다. 삼성에서 진행하는 기업 실무형 교육에 참가하였습니다. 화면 구성은 바꿔가며 토론하고 주요 기능을 2일 만에 합의했습니다. 높은 우선순위부터 개발했으며 문제시 프로토타입을 수정해 의견했습니다. 구글 문서로 화면별 기능과 담당할 사람을 정리했고, 각자 자신이 맡은 일을 해결하면 점검표를 표시하는 방법을 제안했습니다. 한 팀원이 구글 문서에 회의이 필요하 내용을 기록 옆에 정리해두어 문제 상황 공유가 빨라졌습니다. 삼성에서도 목표를 두고 우선순위를 정해 개발하도록 하겠습니다!

1405 bytes

자기소개서 분석하기

랜덤 불러오기

문장 분석 결과

키워드(Keywords)

행사,안건,회원,점검

주요 문장(Main sentence)

동아리 운영진으로 mt 등 행사를 기획하고 회의를 진행하기 위해 노력했던 순간이 기억에 남았습니다.

어색한 문장(Awkward sentences)

추출된 첨삭대상 문장이 없습니다.

문맥 기반 불용어(Context-based stopword)

삼성에서도 목표를 두고 우선순위를 정해 개발하도록 하겠습니다. - 삼성

<롯데에 지원한 자기소개서>

채용 AI ARAVOM

회사명

삼성

용답구분

사회활동

자기소개서

간단한 프로토타입을 만들었습니다. 삼성에서 진행하는 기업 실무형 교육에 참가하였습니다. 화면 구성은 바꿔가며 토론하고 주요 기능을 2일 만에 합의했습니다. 높은 우선순위부터 개발했으며 문제시 프로토타입을 수정해 의견했습니다. 구글 문서로 화면별 기능과 담당할 사람을 정리했고, 각자 자신이 맡은 일을 해결하면 점검표를 표시하는 방법을 제안했습니다. 한 팀원이 구글 문서에 확인이 필요한 내용을 기록 옆에 정리해두어 문제 상황 공유가 빨라졌습니다. 삼성에서도 목표를 두고 우선순위를 정해 개발하도록 하겠습니다!

1405 bytes

자기소개서 분석하기

랜덤 불러오기

문장 분석 결과

키워드(Keywords)

행사,안건,회원,점검

주요 문장(Main sentence)

동아리 운영진으로 mt 등 행사를 기획하고 회의를 진행하기 위해 노력했던 순간이 기억에 남았습니다.

어색한 문장(Awkward sentences)

추출된 첨삭대상 문장이 없습니다.

문맥 기반 불용어(Context-based stopword)

비정상적으로 티사명이 언급된 문장은 없습니다.

<삼성에 지원한 자기소개서>

3. [Ctrl+C, Ctrl+V] 자기소개서 어떻게 찾을까?

Language
Conference

데이터를 어떻게 만들지?
손으로 한땀한땀 적어야될까?



- 단순히 자기소개서에 타회사명이 언급된다고 문제가 되지 않습니다. 인턴, 공모전 등 타사에서 경험한 사례가 자기소개서에 등장할 수 있기 때문입니다.
- 그래서 아래와 같은 방법을 고안했습니다!

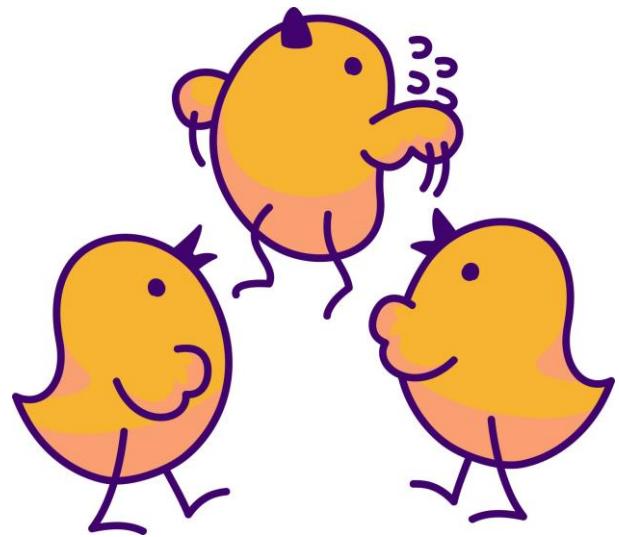
- 1) 회사명이 등장한 모든 문장 수집!(회사명 사전 및 정규표현식 활용)
- 2) 회사명을 학습하지 않게 하기 위해서 문장에서 회사명 제거
(자사뿐만 아니라 다양한 회사에서 사용할 수 있도록 하기 위함.)
- 3) 반드시 지원하는 회사명이 있어야 하는 문장인지 혹은
타사명이 있어도 괜찮은 문장인지 레이블링
- 4) *basemodel [5] 1D - CNN 모델 학습(정확도 : 88%)
- 5) BERT의 Finetunning 하여 transfer learning 한 경우(정확도 : 94%)
- 6) 다음 두 조건을 동시에 만족하면 문장이 있다면 복불한 자기소개서로 판별
* (타사명이 존재) AND (반드시 지원하는 회사명이 와야하는 문장)
⇒ 실제 타사 PoC 결과 100% 정확도 보였음.

TRAIN DATA 예시>

에서 열심히 일하겠습니다.(0 - 반드시 지원하는 회사명이 와야함)
에서 주최한 공모전에서 최우수상을 수상했습니다.(1 - 타 회사명이 와도 괜찮음)

*참고 논문

[5] Y. Kim. (2014). Convolutional neural networks for sentence classification. In Proceedings of EMNLP.



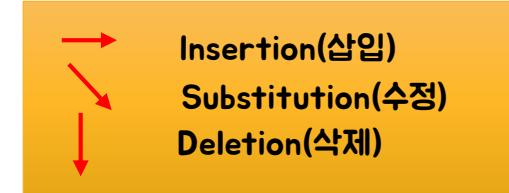
4. 오타 어떻게 찾을까?

기술 키워드 : Levenshtein distance 알고리즘

4. 오타 어떻게 찾을까?

- Insertion(삽입), Deletion(삭제), Substitution(수정)은 비용을 계산
Ex) [사랑, 사랑해]는 한 글자가 삽입되어 거리 1
[나무, 나]는 한 글자가 삭제되어 거리 1
[자연어, 자연언]는 연의 초성 수정(0.3), 어의 종성 수정(0.3)되어 거리 0.6
- 그렇다면 [자연어처리짱]과 [자여어철리굿]사이의 거리는?

-> [1,1]은 [자, 자] 같기 때문에 0으로 초기화한다.(다르면 1)
-> [1,j]는 j가 1씩 증가함에 따라 한 글자씩 추가되었기 때문에 1씩 더한다.
-> [2,1]은 [자연, 자]가 한 글자가 추가되었기 때문에 1이다.
-> [2,2]은 [자연, 자여]의 종성이 다르기 때문에 0.334(약0.3)이다.
-> [2,3]은 [자연, 자여어]는 [여]와 [연]의 종성이 다르고, 한 글자가 추가되었기 때문에 1.3이다.
-> 위와 같은 방식으로 매트릭스를 계산한다.
-> [6,6]의 값 1.6이 [자연어처리짱]과 [자여어철리굿]의 거리이다.



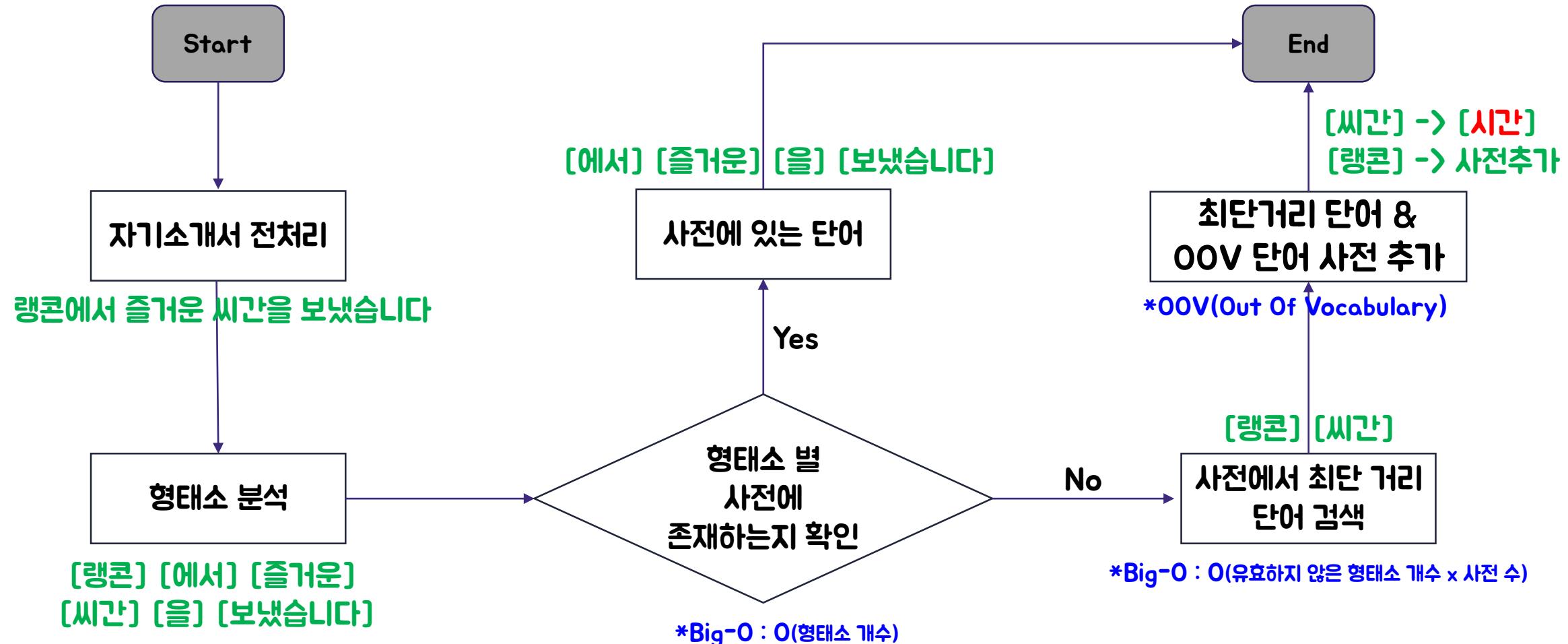
자	여	어	철	리	굿	
자	0	1	2	3	4	5
연	1	0.3	1.3	2.3	3.3	4.3
어	2	1.3	0.3	1.3	2.3	3.3
처	3	2.3	1.3	0.6	1.6	2.6
리	4	3.3	2.3	1.6	0.6	1.6
짱	5	4.3	3.3	2.6	1.6	1.6

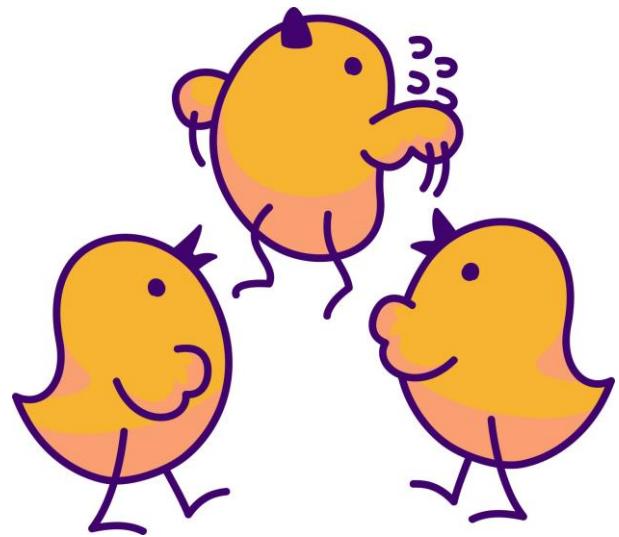
4. 오타 어떻게 찾을까?

- 오타 검정 순서도(simple version)

예시 문장

[랭콘에서 즐거운 **씨간**을 보냈습니다.★]





5. 롯데그룹 채용 안내

5. 롯데그룹 채용 안내

Language
Conference

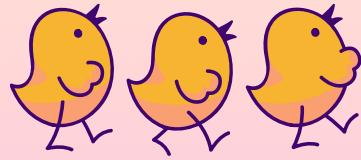


*롯데그룹 채용 URL
<http://job.lotte.co.kr/>

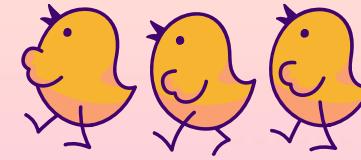
- 신입 채용
 - 2020년 하반기 공채 예정
 - 스펙태클 채용 예정

: SPEC태클(스펙태클) 채용은 2015년에 시작된 롯데 고유의 블라인드 채용 전형입니다. 학벌이나 스펙 중심의 서류 전형에서 벗어나 지원자 직무 수행 능력과 역량 중심으로 인재를 선발합니다.

화려한 볼거리(Spectacle)처럼 다양한 지원자의 이력에 집중하고, 무분별한 스펙 쌓기에 태클을 건다(spec-tackle)는 중의적인 의미를 담고 있습니다.
- 경력 채용
 - 롯데 채용 사이트내 수시 공고



LangCon 2020



QnA
감사합니다.



김영준

data.juny@gmail.com