



SAMSUNG CARD IDEA CHALLENGE



Track 1 모델 기획서

상담원 및 삼성카드 서비스에 대한 만족/불만족 피드백 분류 모델 개발

Team 람지



알고리즘에 입력될 토큰 & Feature의 구성 방법 (모델 입력까지 전처리)

ㅁ : 14 ㅅ : 13 ㅊ : 7 ㄱ : 5
 ㅂ : 14 ㅈ : 13 ㅌ : 7 ㄴ : 5 ...
 ㅍ : 15 ㅊ : 12 ㅍ : 6.5 ㅋ : 5.5

① 문장 내 불필요한 요소 제거 (. , ; ! 띄어쓰기 등)

안녕하세요, 삼성카드입니다!

② 문장을 2음절씩 자르기

통화대기시간이너무길다
 -> 통화 화대 대기 기시 시간
 간이 이너 너무 무길 길다

③ 쪼개진 2음절을 자음과 모음으로 나누기

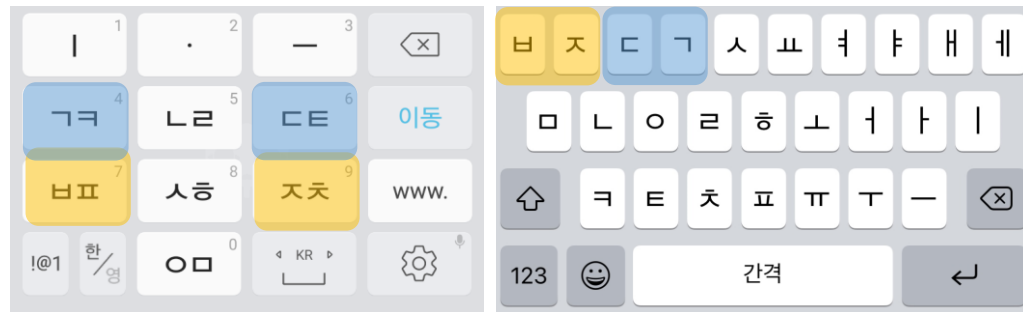
삼성 -> ㅅ ㅌ ㅁ ㅅ ㅌ ㅇ

유니코드로 자음 모음 분리

ord 함수로 텍스트를 유니코드를 확인 후,
 아래 식에 따라 초성, 중성, 종성으로 분리

글자의 자모를 1차원 리스트에 저장 후,
 2차원 리스트에 이를 순차적으로 추가

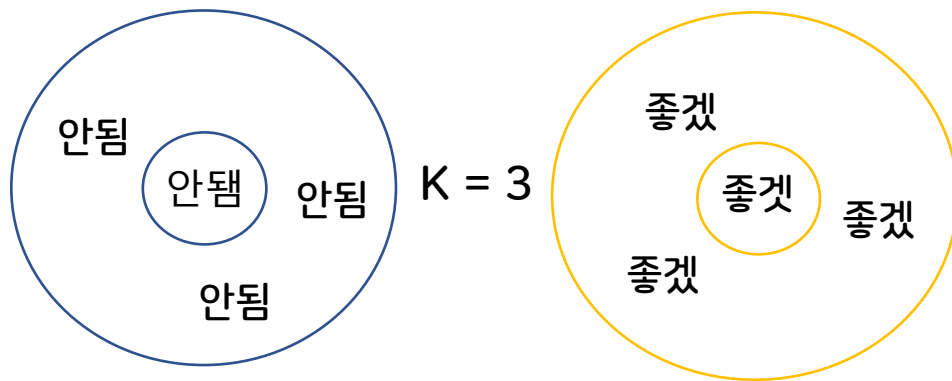
한글 코드 : (초성*21+중성)*28+종성+0xAC00
 한글 코드 범위 : 44032 ~ 55199 (외 값은 예외처리)



오타 발생시 키보드의 근처에 위치한
 자모음들로 바뀌었을 가능성이 높으므로,
 쿼티와 천지인 키보드의 위치를
 모두 고려한 좌표로 설정

2음절이므로 각 글자의 초성 중성 종성을
 좌표로 설정하여 6차원의 좌표가 생성됨
 ex. 상담 : [1, 9, 6, 7, 9, 10]

④ 키보드 위치를 고려하여 자음 모음의 좌표 설정하기



비슷한 단어들을 묶는 것만으로
 실제 오타를 모두 발견하여
 고치기에는 한계가 있음

추후 적용할 classification에서
 오타를 다른 feature로 착각하는
 것을 방지하며, 같은 feature로
 동일 하는 것에 의미를 둠

⑤ 각 단어들의 좌표를 기준으로 k-Nearest-Neighbor(kNN) 알고리즘을 적용하기

현재음절로부터 거리가 가까운 3개의 음절 중
 빈도수가 높은 음절로 수정함

전체 데이터 셋에서의 등장 횟수가 2 번 이하일 경우,
 오타이거나 오타를 바로잡을 음절로 적합하지 않다고
 판단하여, 해당 음절로 수정하지 않음

문장에서 앞 뒤에 위치한 음절들 중
 빈도수가 더 높은 단어를 대상으로만 변환하여
 해당구역의 의미가 바뀌지 않도록 함

ex. '~해택만족~' 에서 '택만'을 '객만'으로 수정할 때,
 '객만'보다 '만족'의 빈도수가 높으므로 '택만' 수정 X

Feature Importance & Tokenizer

Tokenizer를 활용한 2음절 상위 n개 단어

상답	니다	친절	감사	합니
으면	절한	답답	하고	화면
면중	기시	전화	좋겠	불편
네요	세요	원과	까지	해서
안내	보이	터치	화대	음것
성카	드사	지알	가너	는것
했습	주시	사하	절상	되었
무이	기다	결되	었으	질문
문의	어려	필요	했으	다는
되어	결시	지만	는상	절했

Tokenizer를 활용한 3음절 상위 n개 단어

합니다	감사할	사할니	친절하	습니다
기시간	셔셔감	주셔셔	해주셔	하게상
상답직	담직원	원연결	사드립	너무길
성카드	카드사	가너무	상답에	주세요
무이자	었으면	화연결	게해주	다감사
친절감	절감사	연과외	연결되	연결시
불편함	상답을	통화하	는상답	고객의
하게해	는화면	사했습	요감사	직원과
의통화	다른카	른카드	중아요	설명해
드상답	었는데	화면으	면으로	해서감

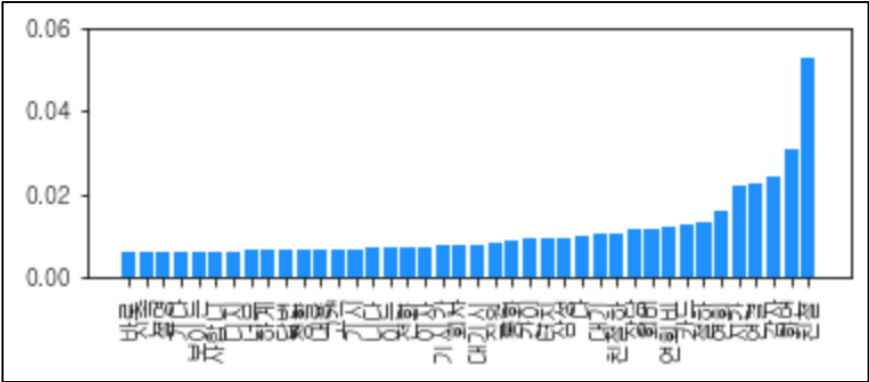
Tokenizer를 활용해 빈도수에 따른
상위 n개 음절을 선별하여 1차 feature 설정



	이 는 a	다 감 사	하 기	절 한 상	연 결 하	합 니 다	이 길	상 답 하	rs	...	것 갈	대 해	무 전	원 이	말 투	드 답	연 결 대	가 너 무	이 분 이
발화																			
직원상답시간너무지연될까중유발	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
상답특이부족	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
상답사와연결부분이어려웠어요	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
특상답인원을늘려주세요요즘같은때에는전화상답보다특상답이더신속하고일처리 가빠른데인원이너무적은듯매번기다리다일처리못하고그냥지나간적이한두번이아 닙니다	0	0	0	0	0	0	0	0	0	...	0	0	0	1	0	0	0	0	0
ARS지동응답불요한안내속소상답직원과빠른연결을원합니다추석연휴잘보내세 요	0	0	0	0	0	0	1	0	0	...	0	0	0	0	0	0	0	0	0
1상답원과의대화까지소요되는시간이길어서힘들었음2상성카드사의직원과면대면 하여질문+답변이되길수있으면좋겠습니	0	0	0	0	0	0	1	0	0	...	0	0	0	0	0	1	0	0	0
문자상답이나온라인상답이가능한부분이많아졌으면좋겠습니	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
상답직원과연결이좀신속했으면좋겠네요	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
상답번호를남겼는데너무늦게주심이넘불만입니	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
상답직원대기시간이짧았으면한다	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0

10 rows × 400 columns

각 문장에서 400개의 1차 feature의 포함유무를 표기하는
One-Hot Encoding을 진행 후 Data set으로 저장



중요도 상위 40개에 대한 feature의 수치 그래프
(최상위 5개 음절 : 친절, 화면, 감사, 연결, 시간)



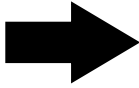
	답 은 면	요 감	는 상 답	절 한 상	감 사 합	있 습 니	비 가	발 부	요 감 사	...	보 다	대 기	정 확	보 이	친 절	디 지	화 면	시 간	연 회 비	무 이
발화																				
직원상답시간너무지연될까중유발	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	1	0	0
상답특이부족	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
상답사와연결부분이어려웠어요	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
특상답인원을늘려주세요요즘같은때에는전화상답보다특상답이더신속하고일처리 가빠른데인원이너무적은듯매번기다리다일처리못하고그냥지나간적이한두번이아 닙니다	0	0	0	0	0	0	0	0	0	...	1	0	0	0	0	0	0	0	0	0
ARS지동응답불요한안내속소상답직원과빠른연결을원합니다추석연휴잘보내세 요	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
1상답원과의대화까지소요되는시간이길어서힘들었음2상성카드사의직원과면대면 하여질문+답변이되길수있으면좋겠습니	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	1	0	0
문자상답이나온라인상답이가능한부분이많아졌으면좋겠습니	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
상답직원과연결이좀신속했으면좋겠네요	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
상답번호를남겼는데너무늦게주심이넘불만입니	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
상답직원대기시간이짧았으면한다	0	1	0	0	0	0	0	0	0	...	0	1	0	0	0	0	0	1	0	0

10 rows × 242 columns

완성된 1차 feature와 label을
RandomForest와 XGB에 fit 후, feature importance를
상위 n개씩 뽑아 겹치는 feature 242개 최종 선별

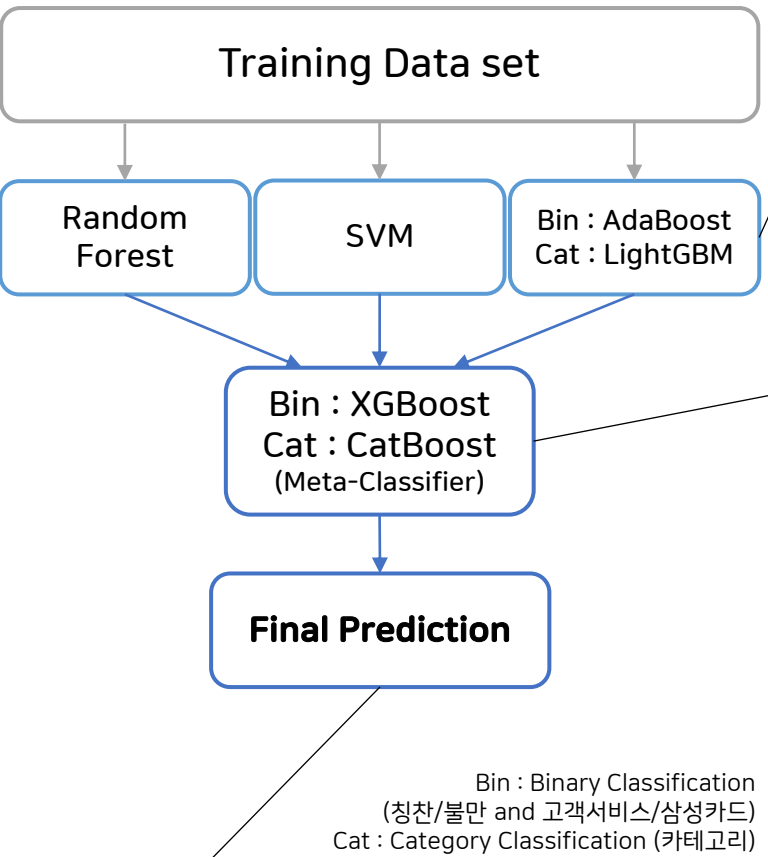
선택한 알고리즘 - Stacking Ensemble

여러 개 모델의 예측 데이터를 통해
최종 메타 모델이 데이터셋을 재생성



장점

단일 모델 대비 **잠재적인 결과**를 이끌어 내어 성능 향상
데이터 특성상 적당한 overfit은 이로울 것이라 예측



사용된 개별 모델 & 활용 이유

RandomForest → 예측 변동성이 낮아 overfitting 방지
SVM(SupportVectorMachine) → 여러 카테고리를 다룰 때 좋은 성능
AdaBoost → 분류 쉬운 feature의 가중치를 낮춰 성능 상승
LightGBM → 카테고리 feature의 자동 변환 및 최적 분할 기능

사용된 최종 모델

① 칭찬/불만 & ②고객서비스/삼성카드
XGBoost
Greedy를 활용하여 overfitting 방지
타 알고리즘과의 연계성 우수

③ 하위 카테고리
CatBoost
Categorical Boost로 카테고리
feature가 많을 때 성능 우수

각각의 모델에서 예측된 칭찬/불만, 고객센터/삼성카드, 하위 카테고리
총 3가지 Label을 모두 합쳐 하나의 최종 Label이 예측됨

튜닝 예정 파라미터

Overfitting 예방을 위한 파라미터

learning_rate
학습 속도와 Stacking을 감안한 overfit을 고려하여 조정 ←
max_depth
overfit 방지를 위해 최소한의 분류 가능 수준에서 각 모델마다 조정 ←

여러 feature를 활용함에 따른 파라미터

reg_lambda
많은 feature를 활용함에 있어 가중치의 L2 reg 미세한 조정 필수 ←
reg_alpha
마찬가지로 많은 feature에 대응하기 위한 L1 reg 미세한 조정 필수 ←

