



SnS 크롤링을 통한 햄버거 브랜드별 평가 감성분석

김란
김민지
김명섭
김태준
김유진

CONTENTS

01
주제 선정

02
데이터 개요

03
연구방법

04
데이터 수집 / 전처리

05
예측 모델 구현

06
모델링 예측결과

주제 선정

01

주제 선정 이유

리뷰 뒤에 숨겨진 긍정적 또는 부정적인 감정을 예측할 수 있는 모델을 만들 수 있을까? 만들 수 있다면, 이러한 예측 모델을 이용하여 부정적 요소와 긍정적 요소를 마케팅에 활용할 수 있지 않을까 생각하여 주제를 선정하였다.

가장 친숙하고 동네마다 있으며, 대표적인 브랜드들이 생각이 나는 햄버거를 주제로 선정하였고, 2023 버거 프랜차이즈 트렌드 리포트에서 버거 프랜차이즈 이용 순위와 한 번이라도 취식 경험이 있는지 순위를 보았을 때, 다른 점을 발견하여 실제 이용해 본 사람들의 리뷰를 데이터로 선정하였다.



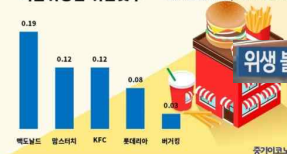
버거 프랜차이즈 브랜드 이용 순위

주 이용		단위: %
1위 버거킹		30.7
2위 맥도날드		22.4
3위 맘스터치		15.1
4위 롯데리아		13.8

한 번이라도 취식 경험		단위: %
1위 맥도날드		95.5
2위 버거킹		95.3
3위 롯데리아		95.0

출처 : 버거 프랜차이즈 트렌드 리포트 2023

주요 햄버거 프랜차이즈 매장 숫자 대비 식품위생법 위반횟수



위생 불량 등 햄버거 프랜차이즈 매장 19곳 적발

중>이코노미

<그림> (해당년 기준) > 순종 (이코노미)

데이터 개요

연구 방법

02



인스타그램

#롯데리아 / #맥도날드 / #버거킹
검색하였을 때 가장 최근 게시물 10000개 수집



네이버 리뷰

각 브랜드 3사별 DT매장
50개 매장의 최근 리뷰 200개씩
총 10000개 수집



유튜브 댓글

각 브랜드 3사별 평가 댓글과
위생 관련 뉴스 댓글
time.sleep(1000) 기준으로
댓글 수집

03

데이터 크롤링

인스타그램
네이버리뷰
유튜브댓글

Target
데이터 선정

형태소 분석

토큰화
정규화
불용어 제거

딥러닝 모델

WRITE HERE

크롤링	라벨링	전처리	딥러닝
BeautifulSoup , Selenium	모든 데이터의 30% 긍정 1 / 부정 0	Konlpy(코엔엘파이) Okt / Mecab	Tensorflow GRU

데이터 수집 / 전처리

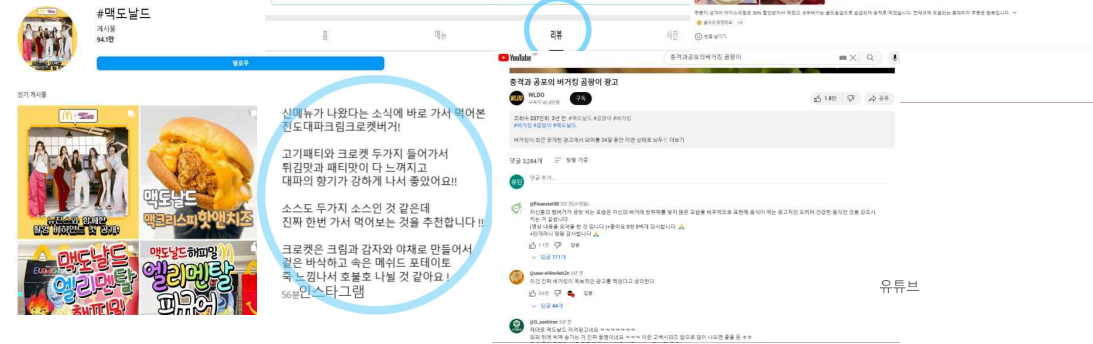
04

데이터 수집 - 인스타그램

<pre># 검색어 조건에 따른 url 생성 def insta_searching(word): url = "https://www.instagram.com/explore/tags/" + str(word) return url</pre>	<pre># 첫 번째 게시물 클릭 후 다음 게시물 클릭 def move_next(driver): right = driver.find_element(By.CSS_SELECTOR, "svg[aria-label='다음']") right.click() time.sleep(3)</pre>	<pre># 검색 결과 페이지 열기 driver.get(url) time.sleep(8)</pre>
<pre># 열린 페이지에서 첫 번째 게시물 클릭 + sleep 여스드 통하여 시차 두기 def select_first(driver): first = driver.find_element(By.CSS_SELECTOR, '._aau') first.click() time.sleep(3)</pre>	<pre># 크롤 브라우저 열기 driver = webdriver.Chrome() driver.get("https://www.instagram.com") time.sleep(3)</pre>	<pre># 첫번째 게시물 클릭 select_first(driver)</pre>
<pre># 본문 내용, 해시태그 가져오기 def get_content(driver): html = driver.page_source soup = BeautifulSoup(html, 'lxml') # 본문 내용 try: content = soup.select('div._adzs')[0].text except: content = '' # 해시태그 tags = re.findall(r'#[\w\#\&\._]+', content) data = [content, tags] return data</pre>	<pre># 인스타그램 로그인 email = '' # 아이디 입력 input_id = driver.find_element(By.XPATH, '//input[@type="text"]') input_id.clear() input_id.send_keys(email) password = '' # 비밀번호 입력 input_pw = driver.find_element(By.XPATH, '//input[@type="password"]') input_pw.clear() input_pw.send_keys(password) input_pw.submit() time.sleep(5)</pre>	<pre># 게시물 수집 시작 results = [] target = 10000 for i in range(target): try: data = get_content(driver) results.append(data) move_next(driver) except: time.sleep(2) move_next(driver)</pre>
	<pre># 게시물 조회를 검색 키워드 입력 요청 word = input("검색어를 입력하세요 ") word = str(word) url = insta_searching(word) 검색어를 입력하세요 : 맥도날드</pre>	

04

데이터 수집 (Crawling)



04

데이터 수집 - 네이버지도 리뷰

<pre>from selenium import webdriver from selenium.webdriver.common.by import By import time import re from bs4 import BeautifulSoup from selenium.webdriver.common.keys import Keys from selenium.common.exceptions import NoSuchElementException driver = webdriver.Chrome()</pre>	<pre>results = [] for i in place : naver_map_search_url = f'https://map.naver.com/v5/search/{i}/place' # 검색 url 만들기 driver.get(naver_map_search_url) # 검색 url 접속 = 검색하기 time.sleep(17) cu = driver.current_url res_code = re.findall(r'place/(#d+)', cu) final_url = f'https://pcmap.place.naver.com/restaurant/{res_code[0]}/review/visitor/' driver.get(final_url) time.sleep(5) for j in range(19) : next = driver.find_element(By.CSS_SELECTOR, ".fH30") next.click() time.sleep(3) html = driver.page_source soup = BeautifulSoup(html, 'lxml') time.sleep(1) one_review = soup.find_all('div', attrs = {'class': 'ZZ40K'}) for k in range(len(one_review)): try: review_content = one_review[k].find('span', attrs = {'class': 'zPv1t'}).text except: # 리뷰가 없다면 pass results.append(review_content)</pre>
--	---

@MB-wo6ct 3년 전
전 오히려 아주 좋았어요. 공방이가 피니 신선하다는 느낌이 들었습니다.
👍 1 🗨 답글

@user-j2f8ny9i 2년 전
저렇게 광고 만들 수 있으면 재밌겠다 ㅋㅋㅋ
👍 1 🗨 답글

데이터 수집 - 유튜브

```
driver = webdriver.Chrome()
url = 'https://www.youtube.com/watch?v=omxjG5t0I4'
driver.get(url)
last_page_height = driver.execute_script("return document.documentElement.scrollHeight")
while True:
    driver.execute_script("window.scrollTo(0, document.documentElement.scrollHeight);")
    time.sleep(3.0)
    new_page_height = driver.execute_script("return document.documentElement.scrollHeight")
    if new_page_height == last_page_height:
        break
    last_page_height = new_page_height
time.sleep(1000)
html_source = driver.page_source
driver.close()
soup = BeautifulSoup(html_source, 'lxml')

mains = soup.select('div#main')
youtube_user_ids = []
youtube_comments = []
for main in mains:
    youtube_user_ids.append(main.select('div#header div#header-author > h3 > a#author-text'))
    youtube_comments.append(main.select('div#comment-content'))
```



```
str_youtube_userIDs = []
str_youtube_comments = []

for i in range(len(youtube_user_IDs)):
    try:
        str_tmp = str(youtube_user_IDs[i][0].text)
        # print(str_tmp)
        str_tmp = str_tmp.replace('\n', '')
        str_tmp = str_tmp.replace('\t', '')
        str_tmp = str_tmp.replace(' ', '')
        str_youtube_userIDs.append(str_tmp)

        str_tmp = str(youtube_comments[i][0].text)
        str_tmp = str_tmp.replace('\n', '')
        str_tmp = str_tmp.replace('\t', '')
        str_tmp = str_tmp.replace(' ', '')
        str_youtube_comments.append(str_tmp)

    except Exception as e:
        print(e)

for i in range(len(str_youtube_userIDs)):
    print(str_youtube_userIDs[i], str_youtube_comments[i])
```

데이터 라벨링 - 각 브랜드 3사별로 나눈 이유

```
loaded_model = load_model('best_model.h5')
print("### 테스트 정확도: %.4f" % (loaded_model.evaluate(X_test, y_test)[1]))
```

맥도날드(train) 롯데리아, 버거킹 (test)

218/218 [=====] - 22s 100ms/step - loss: 0.6153 - acc: 0.7890
 롯데리아 테스트 정확도: 0.7890

219/219 [=====] - 22s 101ms/step - loss: 0.5888 - acc: 0.7781
버거킹 테스트 정확도: 0.7781

롯데리아(train) 맥도날드, 버거킹 (test)

219/219 [=====] - 13s 57ms/step - loss: 0.4359 - acc: 0.8393
맥도날드 테스트 정확도: 0.8393

219/219 [=====] - 12s 57ms/step - loss: 0.4057 - acc: 0.8592
버거킹 테스트 정확도: 0.8592

버거킹(train) 롯데리아, 맥도날드 (test)

218/218 [] - 23s 103ms/step - loss: 0.4752 - acc: 0.7876

테스트 정확도: 0.7876

219/219 [=====] - 34s 151ms/step - loss: 0.3646 - acc: 0.8574

테스트 정확도: 0.8574

04

데이터 라벨링

인스타그램 + 네이버지도 리뷰 + 유튜브 댓글
각 브랜드 별 파일 합친 뒤

약 22000개의 리뷰 중 30%인 7000개의 리뷰를
target 데이터로 선정하여
긍정을 1, 부정을 0으로 라벨링 진행



04

데이터 전처리

```
stowords = ['도', '내', '나', '대', '가', '자', '한', '한', '해', '하', '고', '말', '들', '아', '못', '고', '대', '해', '하', '못', '지', '아', '일']

train_data['tokenized'] = train_data['content'].apply(lambda words:
train_data['tokenized'] = train_data['tokenized'].apply(lambda x: [item for item in x if item not in stowords])

test_data['tokenized'] = test_data['content'].apply(lambda words:
test_data['tokenized'] = test_data['tokenized'].apply(lambda x: [item for item in x if item not in stowords])

negative_words = np.hstack(train_data[train_data.target == 0]['tokenized'].values)
positive_words = np.hstack(train_data[train_data.target == 1]['tokenized'].values)

negative_word_count = Counter(negative_words)
positive_word_count = Counter(positive_words)

[('한글', 455), ('영자', 493), ('200', 453), ('100', 387), ('1년', 222), ('2년', 206), ('1년', 203), ('1월', 178), ('1년', 175), ('2015년', 174), ('1월', 173), ('1월', 162), ('한글', 161)]

positive_word_count = Counter(positive_words)
[('한글', 203), ('200', 163), ('영자', 162), ('1월', 138), ('1년', 130), ('2015', 104), ('1월', 80), ('2012', 77), ('1월', 74), ('한글', 65), ('1월', 46), ('1월', 44), ('1월', 44)]
```

Mecab을 사용하여 형태소 분석을 한 뒤,
단어의 빈도 수가 높은 20개의 단어를 추출해 보았다.

또한, 라벨링한 긍정 리뷰와 부정 리뷰의 평균 길이 수를 구해보니 길이는 비슷하였다.

content

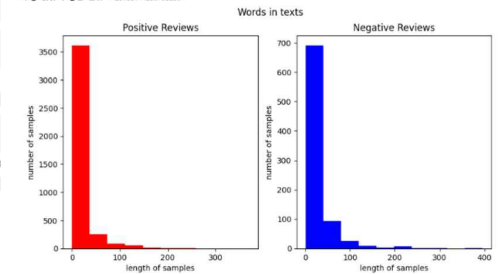
target

13396	새벽에 드라이브 쓰루로 방문하기 좋아요!	1
13397	안정적인 맛	1
13398	궁금한게 있는데 맥도날드 알바 왜 인상쓰고 일함? 가는 곳마다 잘 맞출것엔 한데 처 맞출거 같아	0
13399	굿굿다 다음에도갈게요~	1
13400	협팔인 굿	1
13401	파티를 레어로 먹혔네 ㅋㅋ	0
13402	조아올 마시어요	1
13403	잘 먹었습니다	1
13404	맥스파이스 크림이 너무나 맛있어요. 상하이보다 풍미가 더 있었요. 크림이하면서 더 매콤해요. 소스랑 베이컨 더 들어간것같아요. 그냥 맥스파이스 상하이보다 친원 이상 더 주고 먹을거 같진 않아요. 맥도날드도 이제 버찌네요.	0
13405	오늘의 점심#맥도날드#상하이스파이스버거세트	0
13406	나 호주 갈때 햄버거 하나 시켰는데 세입만엔 다먹어서 여자가 없었다 매뉴판에서는 개크게 그려놓고 나오니까 계산할때 사기 같다..사이나 버거집만 가자..	0
13407	주자공간이 항상 차요.그래도 드라이브 스루는 가능	1
13408	여러가지 주문후에 함께 담아달라했더니 종이가방을 던져주듯이 주더라.. 바빠서 그러면 이해하겠지만 자기들끼리 떠들면서 말이 다..패스트푸드라서 다들 그렇고, 실망이었으면 다시는 안가고 싶다는ㅜ	0
13409	맥도날드는 역시 지즈버거!	1
13410	리뉴얼 되고 처음 방문했는데너무 깨끗해요는 지난 커피 마시기 막 종류요^^	1

```
fig=plt.subplot(2,1,figsize=(8,6))
text_len=train_data['text_data'][target==0]['tokenized'].max(axis=1).len(x)
ax=plt.text(0.5, 0.5, 'Positive Texts')
ax.set_xlabel('Positive Texts')
fig=plt.subplot(2,1,figsize=(8,6))
text_len=train_data['text_data'][target==1]['tokenized'].max(axis=1).len(x)
ax=plt.text(0.5, 0.5, 'Negative Texts')
ax.set_xlabel('Negative Texts')
fig=plt.subplot(2,1,figsize=(8,6))
word_len=train_data['text_data'][target==0]['tokenized'].max(axis=1).len(x)
ax=plt.text(0.5, 0.5, 'Words in Texts')
ax.set_xlabel('length of samples')
ax.set_ylabel('Number of samples')
print '평균 길이의 평균 : %g'%np.mean(text_len)

plt.show()

평균 길이의 평균 : 14.99991030696967
길이의 평균 : 25.65142379269237
```



예측 모델 - GRU

05

리뷰의 길이가 390 이하인 샘플만 가져와
서로 다른 길이의 샘플들의 길이를 동일하게 맞춰주는 패딩 진행

모델은 다대일 구조의 LSTM을 개선한 모델인 GRU를 사용하여
두 개의 선택지 (긍정 / 부정) 중 하나를 예측하는 이진 분류 문제 수행하는 모델 구현

GRU 모델

```
def below_threshold(len(max_len, nested_list):
    count = 0
    for sentence in nested_list:
        if (len(sentence) <= max_len):
            count = count + 1
    print('전체 샘플 중 길이가 %s 이하인 샘플의 비율: %s' % (max_len, (count / len(nested_list) * 100)))

max_len = 390
below_threshold(len(max_len, X_train))

전체 샘플 중 길이가 390 이하인 샘플의 비율: 99.97948297086582

X_train = pad_sequences(X_train, maxlen=max_len)
X_test = pad_sequences(X_test, maxlen=max_len)

from tensorflow.keras.layers import Embedding, Dense, GRU
from tensorflow.keras.models import Sequential
from tensorflow.keras.models import load_model
from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint

embedding_dim = 100
hidden_units = 128

model = Sequential()
model.add(Embedding(vocab_size, embedding_dim))
model.add(GRU(hidden_units))
model.add(Dense(1, activation='sigmoid'))

es = EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience=4)
mc = ModelCheckpoint('best_model.h5', monitor='val_acc', mode='max', verbose=1, save_best_only=True)

model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['acc'])
history = model.fit(X_train, y_train, epochs=15, callbacks=[es, mc], batch_size=64, validation_split=0.2)

loaded_model = load_model('best_model.h5')
print("테스트 정확도: %.4f" % (loaded_model.evaluate(X_test, y_test)[1]))
```

05

Train set / Test set

```
train_data, test_data = train_test_split(data, test_size = 0.3, random_state = 42)
print('훈련용 리뷰의 개수 : ', len(train_data))
print('테스트용 리뷰의 개수 : ', len(test_data))

훈련용 리뷰의 개수 : 4887
테스트용 리뷰의 개수 : 2095

train_data['target'].value_counts().plot(kind = 'bar')

train_data['content'] = train_data['content'].str.replace("[^ㄱ-ㅎㅏ-ㅣ가-힣]*","")
train_data['content'].replace('', np.nan, inplace=True)
print(train_data.isnull().sum())
train_data = train_data.dropna(how='any') # Null 값 제거
print(train_data['content'].head(20))
print('전처리 후 훈련용 샘플의 개수 : ', len(train_data))

test_data['content'] = test_data['content'].str.replace("[^ㄱ-ㅎㅏ-ㅣ가-힣]*","") # 정규 표현식 수행
test_data['content'].replace('', np.nan, inplace=True) # 공백은 Null 값으로 변경
print(train_data.isnull().sum())
test_data = test_data.dropna(how='any') # Null 값 제거
print('전처리 후 테스트용 샘플의 개수 : ', len(test_data))

content    13
target      0
dtype: int64
11835      롯데리아 햄버거 너무 비싸
4721        굿
1554        좋아요
13087      드라이브스루라 좋아요
5606      허쉬 핫초코 파우치하고 적당한 양 좋네요
13455      롯데리아에서 귀여운 포켓몬 하우스 데려왔어요피카츄는 없네요롯데리아포켓몬하우스하우스고라파덕
9408        빠름
81         매번 먹는 단골 롯데리아ㅎㅎㅎㅎ조아요
7450      미치겠다 아직도 롯데잇츠 서버쪽발했나애티도 들어가지 못했다니니 롯데잇츠 롯데리아 대기...
2659      햄버거 먹으러 속속 맛집 롯데리아 롯데속초리조트 왔어요 새우버거 완전 맛있어
4437        굿
810      마라가들여간 햄버거라고 새로 출시되었다고해서 함사서 먹어요먹팔해요 멕스타맛팔 멕스타...
3369      롯데리아 비빔한마당 정범빅전주비빔라이스버거
7355        굿
```

05

긍정 / 부정 판단 - 예측해보기

```
def sentiment_predict(new_sentence):
    new_sentence = re.sub(r"[^ㄱ-ㅎㅏ-ㅣ가-힣]*", '', new_sentence)
    new_sentence = Mecab.morphs(new_sentence)
    new_sentence = [word for word in new_sentence if not word in stopwords]
    encoded = tokenizer.texts_to_sequences([new_sentence])
    pad_new = pad_sequences(encoded, maxlen = max_len)

    score = float(loaded_model.predict(pad_new))
    if (score > 0.5):
        print("{: .2f}% 확률로 긍정 리뷰입니다.".format(score * 100))
    else:
        print("{: .2f}% 확률로 부정 리뷰입니다.".format((1 - score) * 100))
```

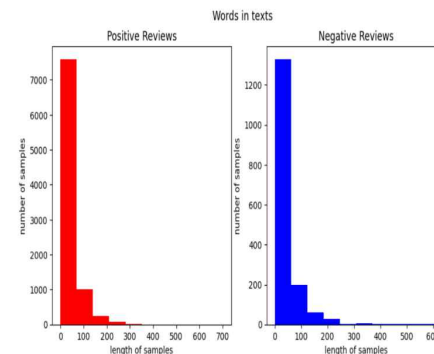
sentiment_predict('친절합니다. ~~맛은기분으로 만나요~~')

1/1 [=====] - 1s 525ms/step
99.81% 확률로 긍정 리뷰입니다.

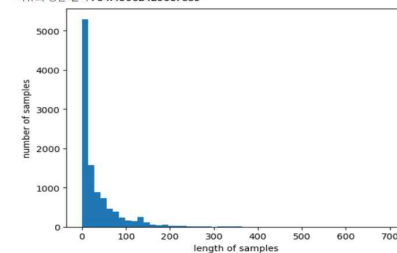
모델링 예측 결과

06

모델로 예측된 라벨링 결과 - 롯데리아



리뷰의 최대 길이 : 702
리뷰의 평균 길이 : 34.45062429057889



전체 리뷰의 최대 / 평균 길이

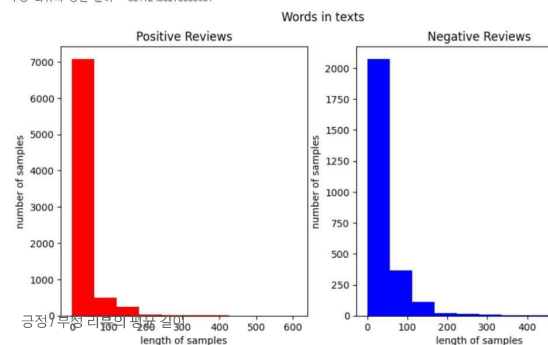
143/143 [=====] - 15s 102ms/step - loss: 0.1582 - acc: 0.9439
테스트 정확도: 0.9439

모델의 정확도

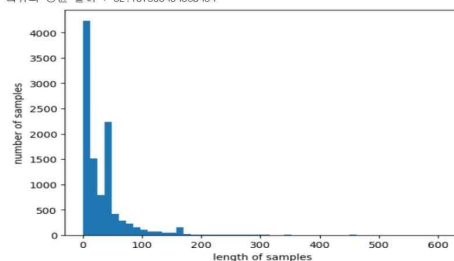
06

모델로 예측된 라벨링 결과 - 맥도날드

긍정 리뷰의 평균 길이 : 29.8975932675483
부정 리뷰의 평균 길이 : 38.72493276898091



리뷰의 최대 길이 : 608
리뷰의 평균 길이 : 32.187956434508454



전체 리뷰의 최대 / 평균 길이

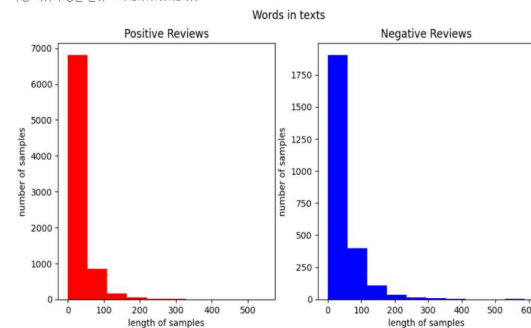
141/141 [=====] - 9s 58ms/step - loss: 0.1809 - acc: 0.9410
테스트 정확도: 0.9410

모델의 정확도

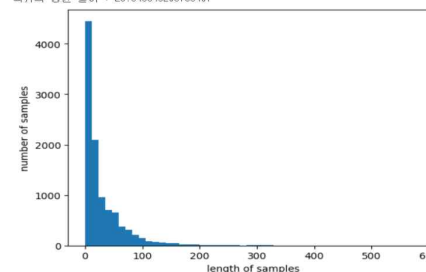
06

모델로 예측된 라벨링 결과 - 버거킹

긍정 리뷰의 평균 길이 : 25.37641045087598
부정 리뷰의 평균 길이 : 43.278137951821965



리뷰의 최대 길이 : 588
리뷰의 평균 길이 : 29.6439546208759407



전체 리뷰의 최대 / 평균 길이

139/139 [=====] - 17s 119ms/step - loss: 0.1629 - acc: 0.9327

테스트 정확도: 0.9327

모델의 정확도

Thank you

A decorative graphic on the left side of the slide. It features a solid yellow rectangular background. Overlaid on this background are two overlapping circles. The circle on the left is a medium orange color, and the circle on the right is a darker, reddish-orange color. They overlap in the center, with the darker circle partially obscuring the lighter one.