프로젝트 최종 발표

인터넷 쇼핑몰 리뷰를 활용한 NLP감성분석 파이프라인

말하는 감자Team



프로젝트 개요

- 1.1 프로젝트 역할 분담
- 1.2 AS-IS, TO-BE
- 1.3 프로젝트 일정

분석 및 설계 구현

- 2.1 프로세스 분석
 - 2.1.1 서비스
 - 2.1.2 기능
 - 2.1.3. 데이터 플로우
- 2.2. 프로세스 설계
 - 2.2.1. 프로세스 아키텍처
- 2.3 프로세스 구현

프로젝트 결론

- 3.1 결과물 활용 방안
- 3.2 문제점 발생 및 해결
- 3.3 향후 개선점
- 3.4 Q&A

프로젝트 개요

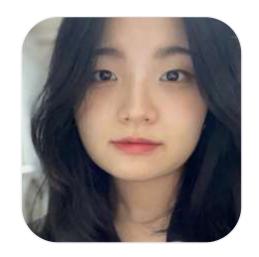
1.1 프로젝트 역할 분담

1.2 AS-IS, TO-BE

1.3 프로젝트 일정



말하는감자 Team



도효주

- 프로젝트 팀장
- 프로젝트 전체 일정 관리
- 문서 작업
- 발표



선우지훈

- 데이터 정제 및 저장
- 대시보드 구현
- 인프라 구축

OpenSearch

Logstash

Kibana

Kafka



오승우

- 프로젝트 기획
- 대시보드 구현
- 데이터 분석 프로그래밍

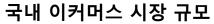


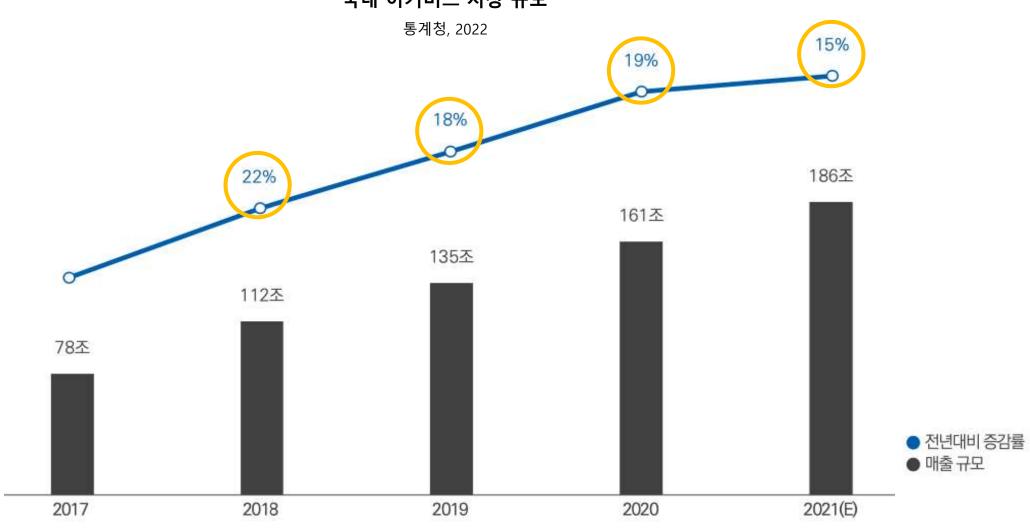
전중석

- 데이터 크롤링
- NLP 처리
- 데이터 분석 프로그래밍
- 인프라 구축
- 대시보드 구현







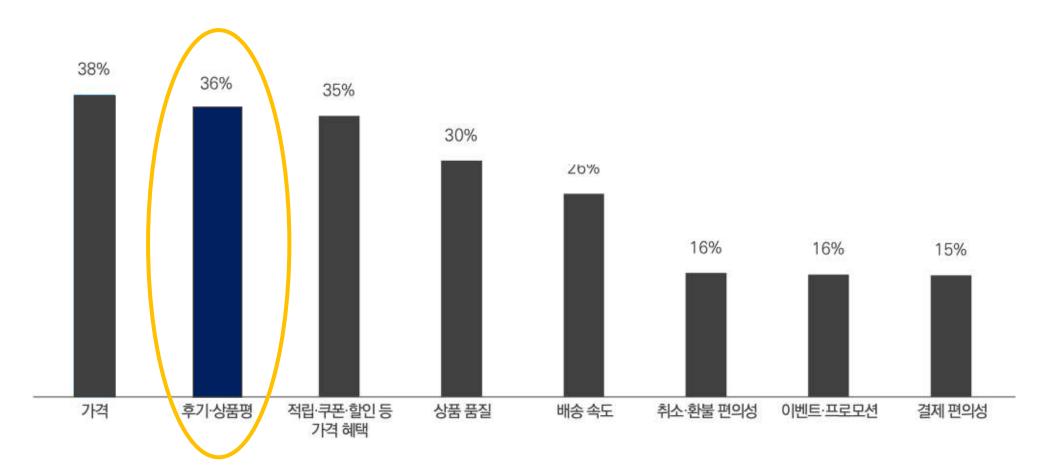


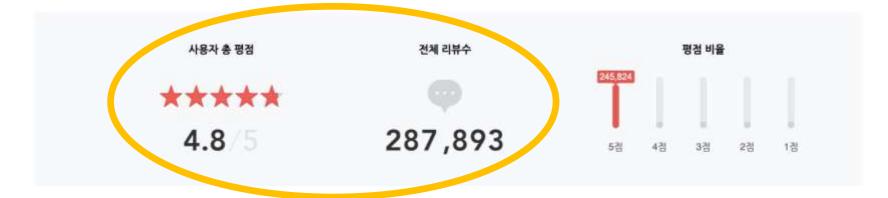




Q. 온라인 쇼핑 구매 고려요인

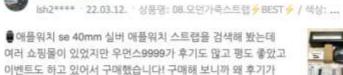
-복수 응답- 통계청, 2022





스토어 PICK 판매자가 직접 선정한 베스트 리뷰입니다.







< 1/15 >

포토&동영상 173,009건







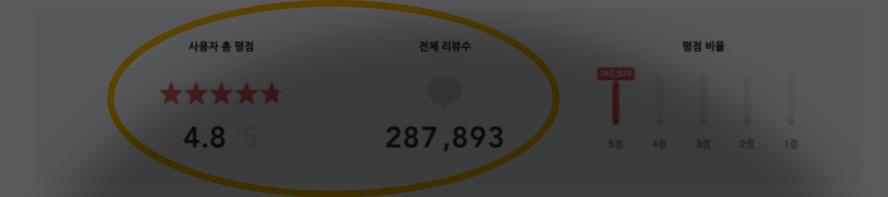












스토어 PICK 판매자가 직접 선정한 베스트 리뷰입니다.



포토&동영상 173,009건



AS-IS, TO-BE



★★★★ 4

zizi****** - 22.08.05. | 신고 [1301반팔] 사이즈: L / 색상: 37.세이프티그린

너무 얇아요ㅠㅜㅠ 구매할 수록 얇아지는 느낌 ㅠㅠ 처음 구매했을때랑도 너무 다른 얇음이네요...



pass****** 22.07.26. | 신고 [1301반팔] 사이즈: XL / 색상: 9 하버블루

18수라고 보기 힘듭니다. 너무 얇아요. 특히 목부분이 너무 쉽게 늘어납니다. 예전에 다른 곳에서 구입한 제품은 3년 이 지난 지금도 두툼하도 목 부분도 늘어나지 않았는데 두번 빨아 입으니 바로 티셔츠가 엉망이 되네요 ㅠㅠ 매우 아쉽습니다.

별점은 5점이지만

<mark>부정적인 내용</mark>의 리뷰들











★★★★★ 5

skdu****** 22.07.24. | 신고

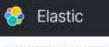
색상: 화이트 / 사이즈: S

예쁘긴한데 티가 너무 길어요 .. 무슨 치마수준이라 바지 안입어도 될 길이에요 ;; 수선해서 입어야할거같아여 .. 재질도 그닥좋은 재질은 아니에요





분석된 데이터를 바탕으로 쇼핑몰 운영의 <mark>편리성,효율성</mark> 증대



감자 모자





분석된 리뷰

전체 리뷰

38,871 analyzed_data - Count

리뷰 갯수

평균 평정

부정 리뷰

88.778% 10.494%

긍정 - Count

4.592

Average Star

82.613% 16.659%

긍정 - Count

긍정 리뷰

♥ 듯사이즈 궁생각 재질 것 가성 감사 비마감 게 상품 머리 제품 개편 구입 품질 여름

Word: Descending - Count

대비 제품 나 불편 世世

Word: Descending - Count

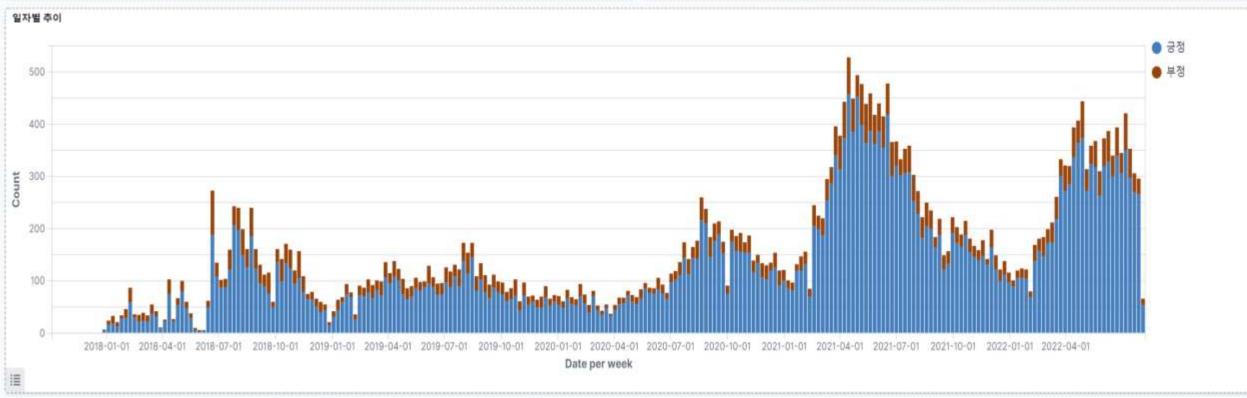
디자인 제품 면 착용감 퀄리티

Word: Descending - Count



다사 키워드		9
가격	긍정	4,208
가격	부정	527
배송	공정	3,788
배송	부정	216
사이즈	공정	2,997
사이즈	부정	903
만족	공정	3,021
만족	부정	75
구매	공정	2,734
구매	부정	242

마사 키워드		(9)
가격	공청	1,186
가격	부정	102
배송	긍정	1,023
배송	부정	32
만축	공정	769
만족	부정	13
사이츠	궁정	613
사이즈	부정	130
구매	공정	636
구매	부정	.47





프로젝트 일정

PHASE		DETAILS	7월 8월																								
				1주차		2주차			3주	다.			4주ㅊ	ŀ	5주차				6주차				7주차				
	PROJECT WEEK: 7 weeks		29 30	1		4	5 6	5 7	8	11 1	.2 13	14	15 18	8 19	20	21 22	25	26 2'	7 28	29	1 2	. 3	4	5	8 9	10	
		이슈 및 적용 사례 조사																									
1	프로젝트 계획 및 분석	이슈 및 적용 사례 분석																									
		데이터 분석 소프트웨어 조사																									
2	ㅠㅋ제ㅌ 서계	프로젝트 기획서 작성																									
2	프로젝트 설계	소프트웨어 아키텍처 설계																									
		ML 프로세스 모델링																									
		데이터 저장소 프로세스																								PR	
		데이터 정제 프로세스																								PROJECT	
3	프로젝트 개발	알림 프로세스																									
		데이터 Queueing 프로세스																								END	
		시각화																									
		추가 구현																									
4		분석 모델 유효성 검증																									
4	프로젝트 테스트	부하 테스트																									
_		보고서 작성																									
5	프로젝트 종료	발표준비																									

분석 및 설계 구현

2.1 프로세스 분석

2.1.1 서비스

2.1.2 기능

2.1.3. 데이터 플로우

2.2. 프로세스 설계

2.2.1. 프로세스 아키텍처

2.3 프로세스 구현

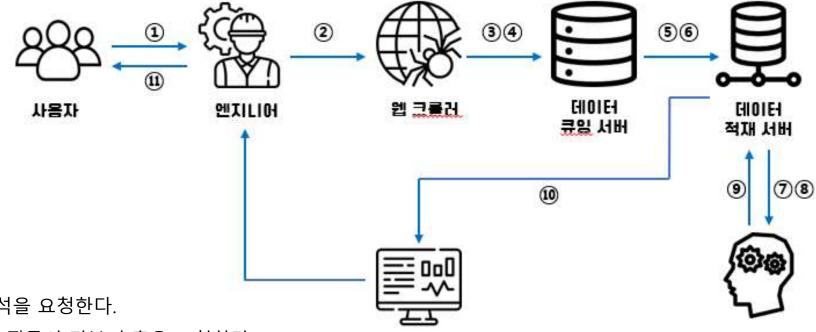


Software	Reason for use				
Python	사용자 쇼핑몰 리뷰 추출 및 자연어 처리에 사용되는 프로그래밍 언어				
AWS CLI	Amazon 서비스 통합 관리 (AWS EKS)				
CUDA	GPU의 가상 명령어셋을 사용할 수 있도록 만들어주는 소프트웨어 레이어				
tensorflow	딥러닝 라이브러리 중 하나이며 Python을 활용하여 연 산처리 작성				
Amazon EKS	AWS 상에서 컨테이너화 된 애플리케이션 관리 자동화				
Anaconda	과학 연구 및 머신러닝 분야에 적합한 Python 및 R 언어의 패키지/의존성 관리 및 배포를 편리하게 해주는 패키지 관리자				
Kafka	쇼핑몰 추출 리뷰 유실 방지				
Opensearch	쇼핑몰 리뷰 데이터, 감성 분석 결과 데이터 적재				

Software	Reason for use
Terraform	클라우드 프로바이더에 IaC 배포 자동화
cuDNN	심층 신경망을 위한 GPU 가속 프리미티브 라이브러리
ChromeDriver	쇼핑몰 댓글 추출 시 렌더링 되는 웹 드라이버
Jupyter Notebook	탐색적 데이터 분석, 데이터 정리 및 변환, 데이터 시각 화, 통계적 모델링, 머신 러닝, 딥러닝 등의 각종 데이 터 사이언스 문서 생성 애플리케이션
Docker	크롤링 파드 이미지 생성
mecab	일본어와 한국어의 유사점으로 한글 분석에도 동작하 는 것을 확인하고 개발한 한국어 형태소 분석기
Logstash	- Kafka에 저장된 데이터 수집 (Consumer) - 쇼핑몰 리뷰 데이터 정제 및 인덱싱
Kibana	리뷰 감성 분석 결과 시각화



서비스 플로우 -1



대시보드

자연어 처리 (감성 분석 모델)

쇼핑몰 리뷰 분석 요청하기

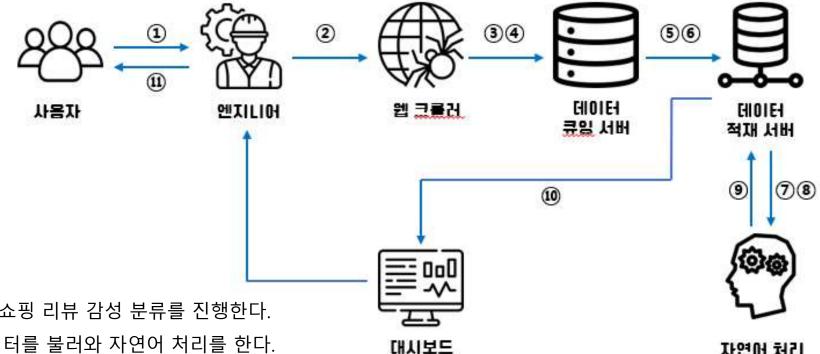
- 1. 사용자가 온라인 쇼핑몰 리뷰 분석을 요청한다.
- 2. 엔지니어는 크롤러에게 사용자 쇼핑몰의 리뷰 추출을 요청한다.
- 3. 크롤러는 사용자의 쇼핑몰에 접근한여 리뷰, 별점, 작성일자를 추출한다.
- 4. 크롤러는 추출한 데이터를 카프카 브로커의 토픽에 저장한다.

쇼핑몰 리뷰 데이터 정제 및 저장하기

- 5. 토픽에 저장된 데이터를 불러와 특수문자를 제거하고 JSON 포맷으로 변환한다.
- 6. 정제된 데이터를 인덱싱하여 데이터 적재 서버(OpenSearch)에 저장한다



서비스 플로우 -2



자연어 처리 (감성 분석 모델)

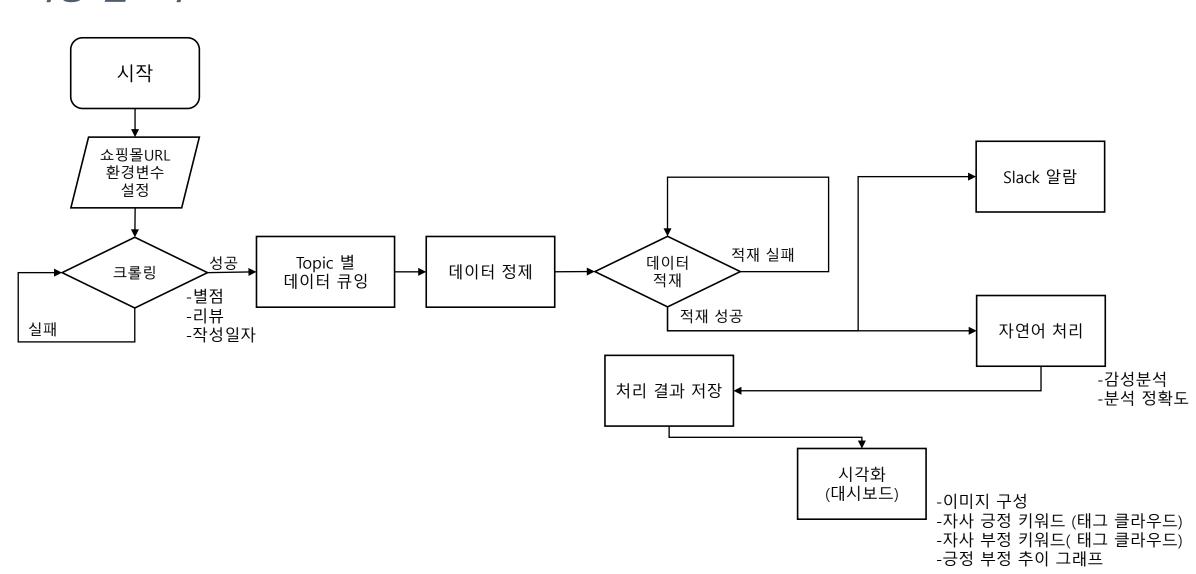
쇼핑몰 리뷰 NLP 처리

- 7. GRU를 사용하여 자사·타사·유사 쇼핑 리뷰 감성 분류를 진행한다.
- 8. 데이터 적재 서버에 저장 된 데이터를 불러와 자연어 처리를 한다.
- 9. 처리 결과를 다시 데이터 적재 서버에 저장한다.

분석 결과 시각화

- 10. 저장 된 데이터의 인덱스를 기반으로 대시보드 형태로 시각화 한다.
- 11. 사용자에게 해당 대시보드의 주소를 제공한다.

기능 플로우







***** 5

pr**** · 22.04.18. | 신고

선택1: (스트랩)핵빅사이즈 / 선택2: 베이지

머리 크신분들은 여기 무조건 추천합니다.



★★★★★ 5

pr**** · 22.04.18. | 신고

선택1: (스트랩)핵빅사이즈 / 선택2: 블랙

머리 크면 제약이 많은데 여기 제품은 너무 알잘딱이네여







Logstash

```
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "머리 크신분들은 여기 무조건 추천합니다.", "date": "22.04.18."}
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "머리 크면 제약이 많은데 여기 제품은 너무 알잘딱이네여", "date": "22.04.18."}
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "너무 좋아요 담에 또 이용하께요", "date": "22.04.18."}
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "흐물거리지 않고 빳빳 하니 정말 좋아요 정사이즈입니다", "date": "22.04.18."}
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "흐물거리지 않고 빳빳 하니 정말 좋아요 정사이즈입니다", "date": "22.04.18."}
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "머리 크면 제약이 많은데 여기 제품은 너무 알잘딱이네여", "date": "22.04.18."}
{"topic": "smartstore.drstyle.review", "star": "5", "comment": "너무 좋아요 담에 또 이용하께요", "date": "22.04.18."}
```



요구사항 정의서

구분	ID	요구사항명	기능	세부사항						
				소비자의 쇼핑몰 리뷰에서 댓글, 별점, 작성일자를 추출						
	SFR-001	리뷰 크롤링을 통한 쇼핑몰 데이터 추출	데이터 추출	추출하고자 하는 URL에 따라 추출 진행						
				추출이 완료된 후 성공/실패 여부와 추출에 걸린 시간을 kafka에 전송						
	SFR-002	리뷰 데이터 적재 알림 전송	알림 전송	Kafka에 크롤링 한 데이터가 적재 완료될 경우 Slack을 통해 알림 전송						
			상위 키워드 추출	빈도 수에 따라 상위 키워드 추출하여 시각화						
기능 요구사항			감성 분석	20만개의 리뷰와 별점 데이터를 크롤링한 후 test데이터와 train 데이터로 나눔						
	SFR-003	온라인 쇼핑몰 리뷰을 분석하기 위한 NLP		train 데이터로 리뷰에 대한 긍정/부정 분석 모델링						
			분석 정확도 계산	test 데이터로 모델 검증						
			단어 간 유사도 계산	키워드를 지정 후 키워드와 유사한 단어를 추출						
			월별 추이	분석하고자 하는 쇼핑몰 리뷰의 긍정/부정 추이를 간트 차트로 시각화						
	SFR-004	시각화를 통한 대시보드 생성	대시보드 생성	분석 결과를 바탕으로 데이터 프레임 생성 데이터 프레임에 맞게 대시보드에서 시각화						
	DD 004	크롤링 데이터 유실 방지를 위한 큐잉 서비		크롤링 한 데이터를 받아 오기 위해 상품 별 토픽 생성						
	DR-001	스	네이터 유실 망시	토픽에 쌓인 데이터를 차례로 전달						
	DR-002		수신 데이터 특수 문자 제거,	원하는 상품에 해당하는 토픽의 데이터 수신						
데이터			무단 데이터 국무 군시 세기, 인덱싱 및 형식 변환	수신 받은 데이터의 특수 문자 제거 및 인덱싱						
요구사항				정제 된 데이터를 JSON 형태로 변환						
	DR-003	온라인 쇼핑몰 리뷰 데이터 NLP를 위한 데 이터 적재	크롤링 후 정제된 데이터 저장	자연어 처리를 하기 위한 데이터를 받아와 저장						
	DR-004	온라인 쇼핑몰 리뷰 데이터를 NLP한 결과 데이터 적재	감성 분석 결과 데이터 저장	분석 결과를 데이터 프레임에 맞추어 저장						
				시스템의 모든 기능은 화면에 출력할 때 브라우저의 영향을 받지 않고 보여주지 않아야 함						
성능 요구사항	PER-001	질의 · 응답 시간	질의응답 시간 및 오류메시지 응답시간	구축 시스템의 사용자 서비스 페이지는 평균 5초 이내에 처리되어야 함						
				사용자 요청 작업 관련 평균 시간 초과 응답 시 성능향상 방안을 강구하여야 함						

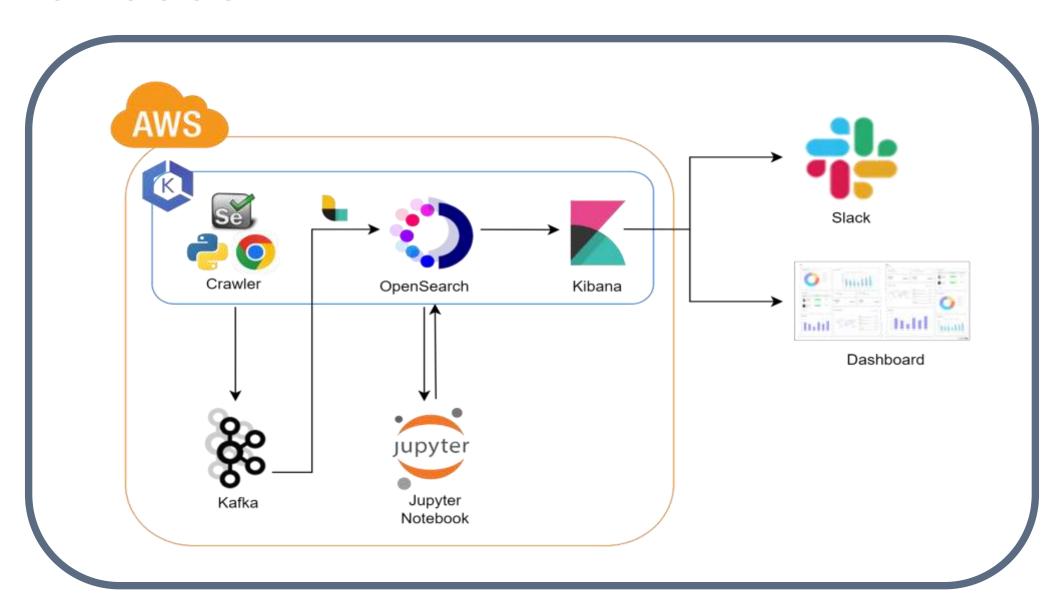


요구사항 정의서

구분	ID	요구사항명	기능	세부사항
				Crawler: t3.medium * 2 EKS 내 2개의 노드에 지정
시스템 장비구성 요구사항	ECR-001	시스템 장비구성 요구사항	시스템 장비구성 요구사항	ElasticSearch-master: t3.large 1~3 ElasticSearch-data: t3.large 1~3 ElasticSearch-client: t3.large 1~3 Kibana: t3.medium 1~3 Logstash: t3.large 1~3 EKS 내 1~3개의 노드로 가용성 확보
				Jupyter Notebook: g4dn.xlarge 1 Kafka Cluster: t3.medium 3 EC2 내 1 혹은 3개의 노드로 구성
			단위 테스트	프로그램 개발 일정 및 테스트 일정에 따라 개발된 프로그램에 대해 단위 테스 트 실시 방안을 수립
테스트 요구사항	TER-001	테스트 요구사항	통합 테스트	통합 테스트는 최소 1회 이상 실시해야 하며, 테스트 일정에 따라 구체적인 실시 방안을 수립하여 수행 후 결과 보고
			시험 운영	계획된 일정에 따라 시험운영 방안을 수립하여 제시
품질 요구사항	QR-001	품질 보증 활동	기능 구현의 정확성 향상 방안 수립 및 준수	개발 시스템은 제공되기로 한 기능 요구사항을 모두 제공해야하며, 초기 협의한 요구사항에서 변경이 필요한 경우 주관기관 담당자와 협의하여 요구사항을 변경 협의 및 과업 변경 가능
	PMR-001	프로젝트 일정관리	프로젝트 일정관리에 관한 사항	일정계획 제시 - 개발의 완성도를 높이기 위해서 프로젝트 착수에서 종료까지 체계적으로 프로 젝트를 관리해야 함
프로젝트 관리		. 2021	. 2021.22.10	단계별 일정관리 - 각 업무단위별 단계별 일정 계획을 수립하여 제시하고, 일정 지연이 발생하는 경우 별도 계획수립 등 만회 계획을 작성
요구사항	PMR-002	투입인력 및 역할 분담	투입인력에 관한 사항	투입인력 및 역할 분담 - 모든 투입인력은 프로젝트 인력으로 구성하여야 하며, 업무분담(역할)을 관리 및 작성
	PMR-003	형상 관리	형상 관리 방안 수립 및 준수	형상 관리 어플리케이션 개발부터 소멸까지의 소스 코드를 포함한 각종 산출물은 중앙에서 통합적으로 변경관리가 가능하도록 관리



프로세스 아키텍처

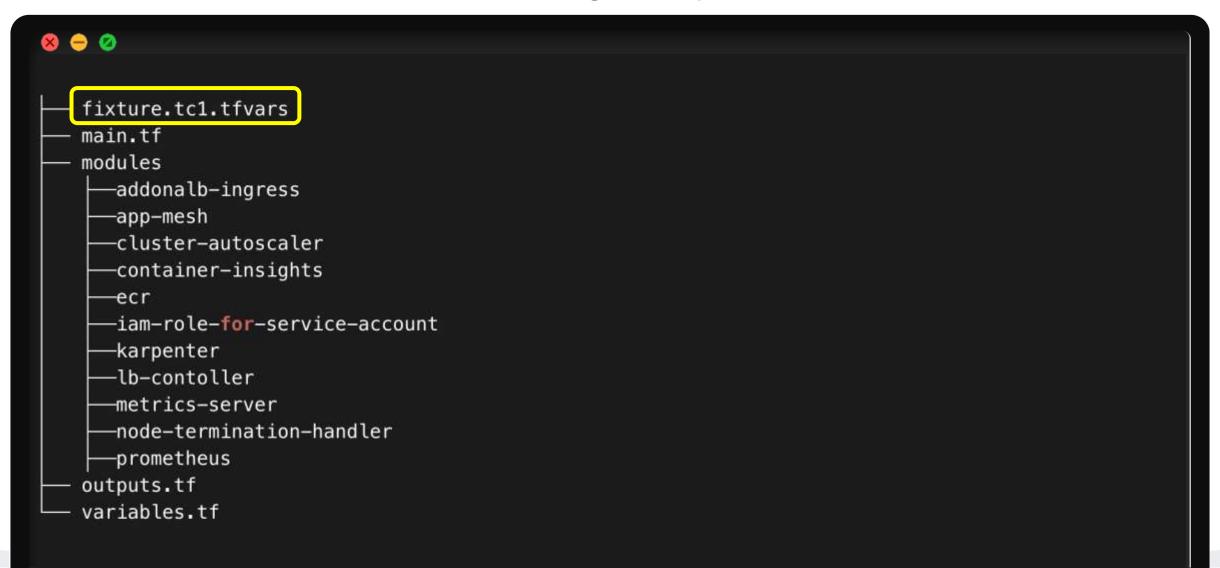






인프라 구현 – EKS

프로젝트의 소프트웨어 중 파이썬 크롤러, Logstash, Opensearch, Kibana는 EKS를 통해 배포



```
8 😑 🕢
fixture.tc1.tfvars
      aws_region
                    = "ap-northeast-2"
                    = ["ap-northeast-2a", "ap-northeast-2b", "ap-northeast-2c"]
     azs
      cidr
                = "10.1.0.0/16"
     enable_igw = true
     enable_ngw = true
     single_ngw = true
      name
                    = "eks-autoscaling-tc1"
      tags = {
       env = "dev"
       test = "tc1"
11
12
      kubernetes_version = "1.21"
      enable_ssm
13
                        = true
     managed_node_groups = [
15
                      = "crawler"
         name
17
         min_size = 1
         max_size
                  = 6
         desired_size = 1
19
20
         instance_type = "t3.medium"
21
22
23
                      = "ElasticSearch-master"
         name
24
         min_size = 1
25
         max_size
                  = 3
         desired_size = 1
27
         instance_type = "t3.large"
       1,
```



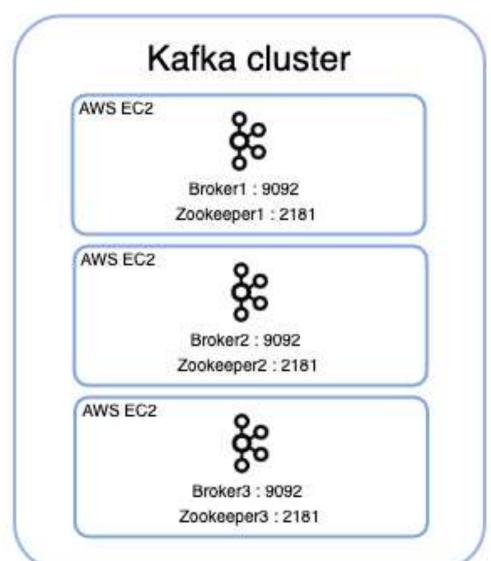
인프라 구현 - Python Crawler

Crawler파드 생성을 위한 YAML파일

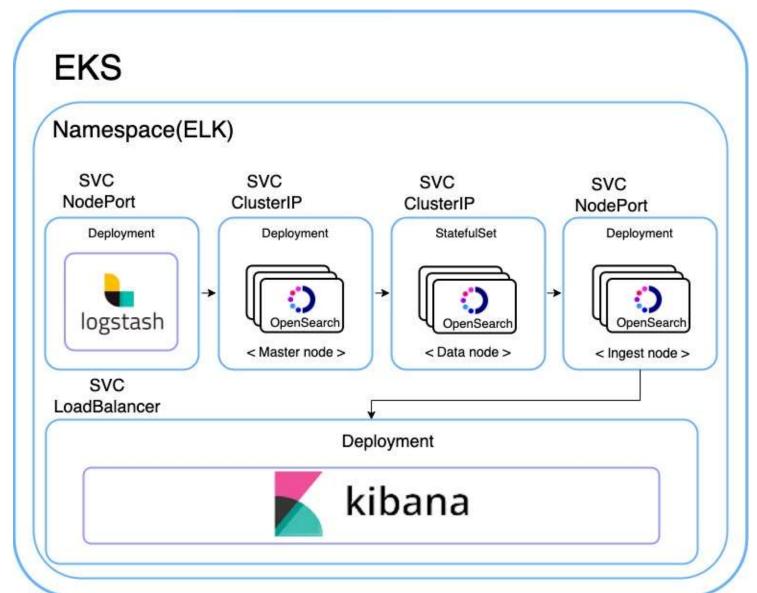
```
🛞 😑 🙆
               atch/v1
  kind: CronJob
                                                                                                memory: "3000Mi"
                                                                                                cpu: "1500m"
    namespace: crawler
    name: crawler-1
                                                                                            - name: url
                                                                                              value: 'https://smartstore.naver.com/goodnara/products/37118c9
    labels:
      app: crawler
                                                                                             - name: topic
                                                                                              value: smartstore.goodnara.review
    schedule: "00 00 * * *"
                                                                                             - name: server
                                                                                              value: "3.38.10.106:9092,3.34.18.190:9092,13.209.146.71:9092"
    jou remptate:
                                                                                             command: ["/bin/sh", "-c"]
                                                                                            args: ["cd /Datapipeline_Project/crawler; ./kafka_producer.py"]
        template:
                                                                                          restartPolicy: Never
          metadata:
            labels:
              app: crawler
            annotations:
              "cluster-autoscaler.kubernetes.io/safe-to-evict": "false"
            nodeSelector:
              Name: Crawler
            containers:
            - name: crawler-1
              image: ddung1203/kafkacrawler:10
              resources:
                requests:
                  memory: "3000Mi"
```

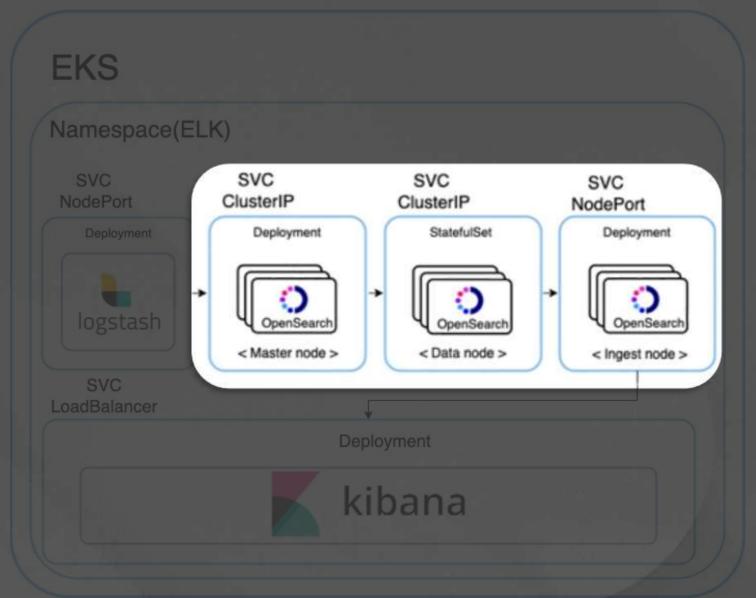


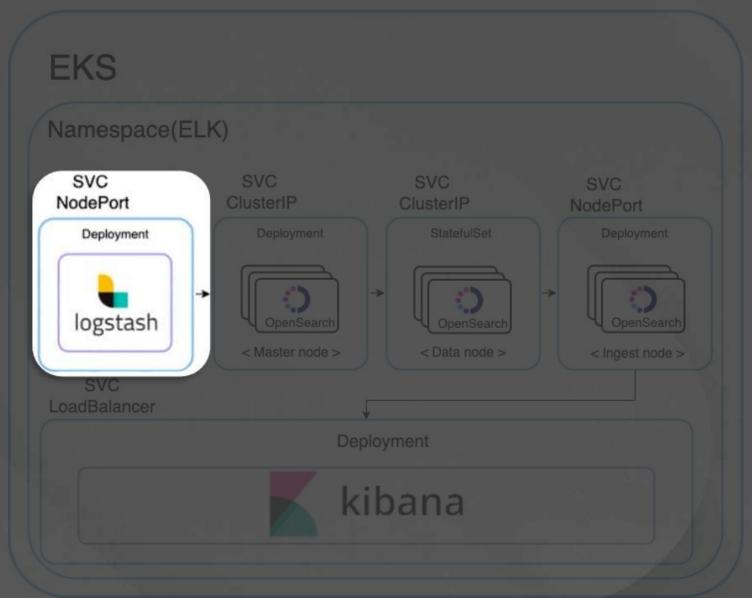
인프라 구현 - Kafka cluster

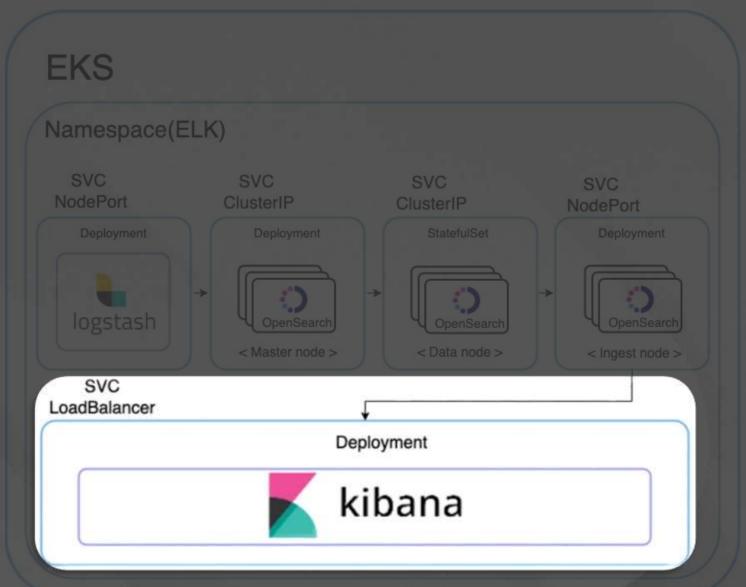














기능 구현 – 데이터 크롤링

```
⊗ ⊝ ⊘
 2 producer = KafkaProducer(acks=1, compression_type='gzip',bootstrap_servers=[server],
                             value serializer=lambda x: json.dumps(x, ensure_ascii=False).encode('utf-8'))
 6 for cmt in range(1,21):
       time.sleep(1)
       list comnt = driver.find element(By.XPATH, '경로'.format(cmt))
       if list comnt.text ==
11
12
13
14
           list_comnt = driver.find_element(By.XPATH, '경로[2]'.format(cmt))
       if list comnt.text ==
           list comnt = driver.find element(By.XPATH, '경로[3]'.format(cmt))
       elif list comnt.text ==
           list comnt = driver.find element(By.XPATH, '경로[2]'.format(cmt))
       comment.append(list comnt.text)
17
18
19
20
21
22
23
24
25
       list_star = driver.find_element(By.XPATH, '경로'.format(cmt))
       star.append(list star.text)
       list_date = driver.find_element(By.XPATH, '경로'.format(cmt))
       date.append(list_date.text)
27 tmp={'star':star.pop(), 'comment':comment.pop(), 'date':date.pop()}
28 producer.send(topic, value=tmp)
30 producer.flush()
31
32
33 |
34
```



로그 데이터의 파이프 라인 설정 파일

```
🔕 😑 🙆
 1 logstash.conf/
 2 >> input
       kafka
             topics => "smartstore.goodnara.review",
                      "smartstore.drstyle.review",
                      "smartstore.thecheaper.review",
                      "smartstore.180store.review",
                      "smartstore cloony review",
                      "smartstore.theshopsw.review"]
11
             consumer_tirreaus -> 3
             isolation_level => "read_committed"
             value_deserializer_class => "org.apache.kafka.common.serialization.StringDeserializer"
             auto_offset_reset => "earliest"
             # 처음 브로커에 진입했을때, 데이터를 가지고오는 시작점을 지정한다.
             #추후 컨슈머를 재시작 하더라도 데이터 중복 이슈를 해결할 수 있다.
21
             group_id => "smartstore" # 컨슈머 그룹이름을 지정한다.
22
23
```



filter -timestamp 설정

```
8 😑 🙆
1 filter
         mutate (
          add field =>
            "timestamp" => "" # timestamp 필드 생성(새로 생성된 필드의 기본 테이터 타입은 String이다.)
         # ruby 코드로 "@timestamp" 필드의 UTC 기준 현재 시간에 9시간을 더한 값을 timestamp 필드에 저장한다.
         ruby
          code => "event.set('timestamp', event.get('@timestamp').time.localtime('+09:00').strftime('%Y-%m-%d %H:%M:%S'))"
         date
          match => ["timestamp", "ISO8601", "YYYY-MM-dd HH:mm:ss"]
          target => "timestamp" # date 필터가 적용될 필드 지정
         grok
          match => {
            "timestamp" => "\d\d%{INT:yy}-%{MONTHNUM:mm}-%{MONTHDAY:dd}%{GREEDYDATA}"
          add_field => <
            "[@metadata][yymmdd]" => "%{yy}%{mm}%{dd}"
```



기능 구현 - 데이터 정제

filter – message 필드 정제

```
8 😑 🙆
           mutate {
             gsub => ["message", "[\"/{}]", ""]
           kv {
             field_split => ","
             value_split => ":"
11
12
          mutate {
13
             remove_field => [ "port","@version","host","message","@timestamp", "yy", "mm", "dd" ]
             rename => {" comment" => "comment"}
             rename => {" date" => "date"}
             rename => { " star" => "star" }
          mutate
20
             convert => \{
21
               "star" => "integer"
23
24
```

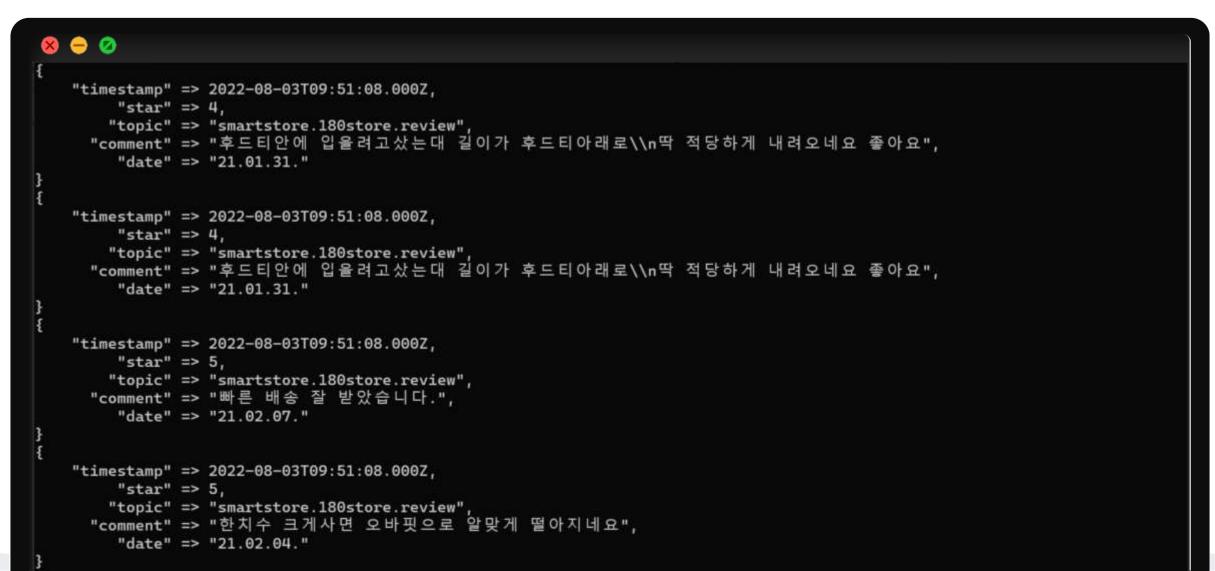


output

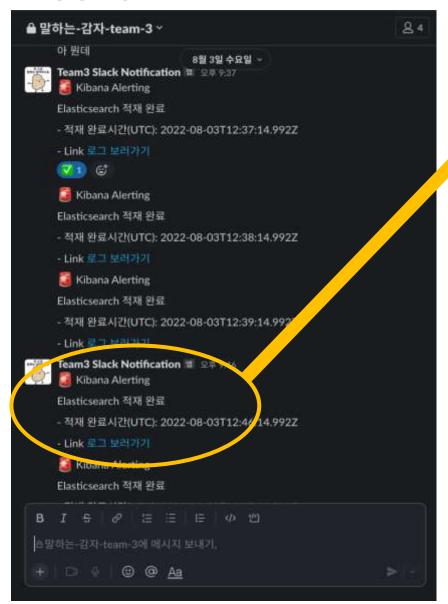
```
🛞 😑 🙆
1 output
       stdout ( codec => rubydebug )
                                                                                            else if |topic| =~ "smartstore.180store.review" {
                                                                                               elasticsearch {
                                                                                                 hosts => "http://elasticsearch-client-http.elk.svc.cluster.local:9200"
       if |topic| = "smartstore.goodnara.review" {
                                                                                                 index => "smartstore.180store.review-%{[@metadata][yymmdd]}"
         elasticsearch (
           hosts => "http://elasticsearch-client-http.elk.svc.cluster.local:9200"
                                                                                                 codec => "json"
           index => "smartstore.goodnara.review-%{[@metadata][yymmdd]}"
                                                                                                 timeout => 120
           codec => "ison"
           timeout => 128
                                                                                             else if [topic] =~ "smartstore.cloony.review" {
                                                                                               elasticsearch (
       else if [topic] =~ "smartstore.drstyle.review" [
                                                                                                 hosts => "http://elasticsearch-client-http.elk.svc.cluster.local:9200"
         elasticsearch
                                                                                                 index => "smartstore.cloony.review-%{[@metadata][yymmdd]}"
           hosts => "http://elasticsearch-client-http.elk.svc.cluster.local:9200"
                                                                                                 codec => "json"
           index => "smartstore.drstyle.review-%{[@metadata][yymmdd]}"
                                                                                                 timeout => 120
           codec => "ison"
           timeout => 128
                                                                                             else if |topic| =~ "smartstore.theshopsw.review" |
       else if [topic] = "smartstore.thecheaper.review" {
                                                                                               elasticsearch
         elasticsearch
                                                                                                 hosts => "http://elasticsearch-client-http.elk.svc.cluster.local:9200"
           hosts => "http://elasticsearch-client-http.elk.svc.cluster.local:9200"
                                                                                                 index => "smartstore.theshopsw.review-%{[@metadata][yymmdd]}"
           index => "smartstore.thecheaper.review-%{[@metadata][yymmdd]}"
                                                                                                 codec => "json"
           codec => "ison"
                                                                                                 timeout => 120
           timeout => 128
```



쇼핑몰 리뷰 데이터가 성공적으로 적재된 것을 알 수 있다.



기능 구현 - Slack Alert







Team3 Slack Notification 2 오후 9:46

Kibana Alerting

Elasticsearch 적재 완료

- 적재 완료시간(UTC): 2022-08-03T12:46:14.992Z
- Link 로그 보러가기



기능 구현 - Slack Alert

```
Define extraction query
   1 - {
            "size": 0,
            "query": {
                "bool": {
    4 +
    5 +
                    "filter": [
    6 +
                            "range": {
    7 *
                                "timestamp": {
    8 +
                                    "from": "{{period_end}}||+9h-2m",
    9
                                    "to": "{{period_end}}||+9h",
   10
   11
                                    "include_lower": true,
   12
                                    "include upper": true,
                                    "format": "epoch millis",
   13
   14
                                    "boost": 1
   15
   16
   17
   18 *
   19 *
                            "match": {
                                " status": {
   20 *
                                    "query": "Success",
                                    "operator": "OR",
   22
                                    "prefix_length": 0,
   23
   24
                                    "max_expansions": 50,
                                    "fuzzy_transpositions": true,
   25
   26
                                    "lenient": false,
                                    "zero_terms_query": "NONE",
   27
                                    "auto_generate_synonyms_phrase_quer
   28
                                    "boost": 1
   29
   30
```

Extraction query response

```
1 + {
         "_shards": {
            "total": 6,
            "failed": 0,
            "successful": 6,
             "skipped": 0
         "hits": {
 8 +
            "hits": [],
 9
            "total": {
10 *
                 "value": 0,
11
12
                 "relation": "eq"
13
             "max score": null
14
15
         "took": 9,
16
17
         "timed out": false
18
```



기능 구현 - 감성분석 (NLP)

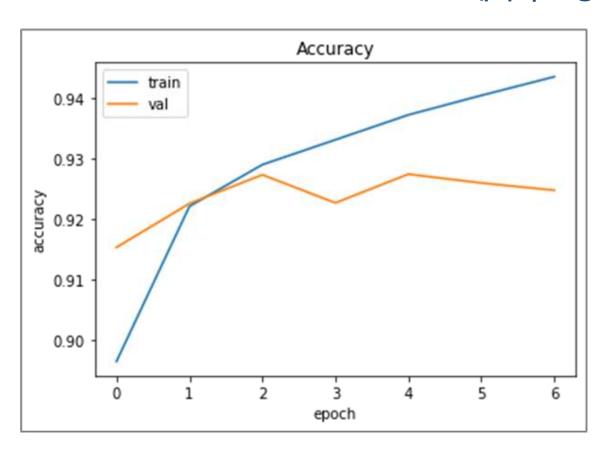
네이버 쇼핑 리뷰 감성 분류

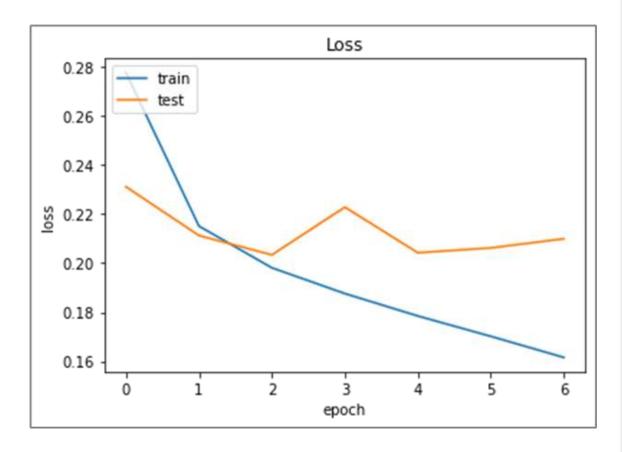
```
8 😑 9
1 from tensorflow.keras.layers import Embedding, Dense, GRU
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.models import load model
4 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
      embedding_dim = 100
      hidden_units = 128
      model = Sequential()
10
      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)
16
      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)
```



기능 구현 – 감성분석 (NLP)

네이버 쇼핑 리뷰 감성 분류





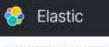
기능 구현 - 감성분석 (NLP)

```
# 불러온 데이터 전처리
     # 중복 제거
     data.drop_duplicates(subset = ['comment'], inplace=True, ignore_index=True)
     # 정규 표현식 수행
     data['comment'] = data['comment'].str.replace("[^¬-っト-|가-헿]","")
     # 공백은 Null 값으로 변경
     data['comment'].replace('', np.nan, inplace=True)
     # Null 값 제거
     data = data.dropna(how='any')
11
     stopwords = ['도', '는', '다', '의', '가', '이', '은', '한', '에', '하', '고', '을',
12
                 '플', '인', '돗', '과', '와', '네', '들', '돗', '지', '임', '게']
13
     data=data.reset_index(drop=False)
     data=data.loc[:, ['comment', 'date', 'star', 'topic']].dropna()
17
     for i in range(len(data)):
         # 한국어 맞춤법 검사
18
         sent = data['comment'][i]
         spelled_sent = spell_checker.check(sent)
20
         hanspell_sent = spelled_sent.checked
21
         data['comment'][i] = hanspell_sent
22
23
     for i in range(len(data)):
25
         # 감성분석
         sp = sentiment_predict(data['comment'][i])
26
         if sp == 1:
27
             data['new_label'][i] = '긍정'
28
         elif sp == 0:
30
             data['new_label'][i] = '부정'
         # 댓글 토론화
         data['tokenized_comment'][i] = mecab.nouns(data['comment'][i])
```



기능 구현 - 분석 데이터 인덱스

```
8 😑 9
2 PUT analyzed_data
    "settings":
     "number_of_shards": 5,
     "number_of_replicas": 3
    "mappings": |
      "properties": {
       "Name" : {
         "type": "keyword"
       1.
       "Star"
        "type": "long"
       1.
        "Date"
         "type" : "date"
        1.
        "Word" : {
         "type" "keyword"
        "Sentiment" : {
        "type" : "keyword"
        "Adjustment-sentiment" : {
         "type" : "keyword"
```



감자 모자





분석된 리뷰

전체 리뷰

38,871 analyzed_data - Count

리뷰 갯수

평균 평정

부정 리뷰

88.778% 10.494%

긍정 - Count

4.592

Average Star

82.613% 16.659%

긍정 - Count

긍정 리뷰

♥ 듯사이즈 궁생각 재질 것 가성 감사 비마감 게 상품 머리 제품 개편 구입 품질 여름

Word: Descending - Count

대비 제품 나 불편 世世

Word: Descending - Count

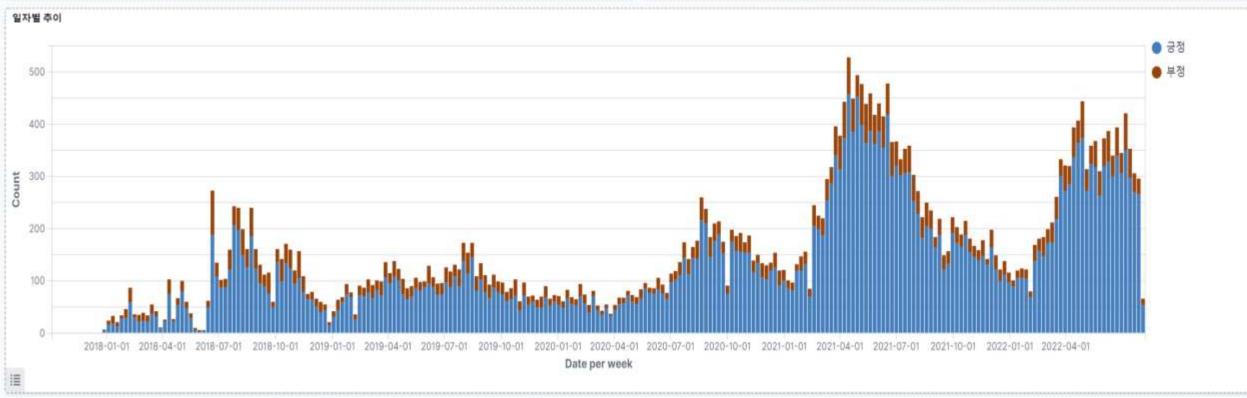
디자인 제품 면 착용감 퀄리티

Word: Descending - Count



다사 키워드		9
가격	긍정	4,208
가격	부정	527
배송	공정	3,788
배송	부정	216
사이즈	공정	2,997
사이즈	부정	903
만족	공정	3,021
만족	부정	75
구매	공정	2,734
구매	부정	242

ł사 키워드		0	
가격	공청	1,186	
가격	부정	102	
배송	긍정	1,023	
배송	부정	32	
만축	공정	769	
만족	부정	13	
사이츠	궁정	613	
사이즈	부정	130	
구매	공정	636	
구매	부정	.47	



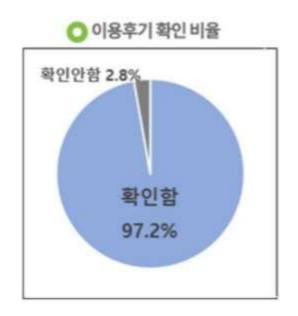
프로젝트 결론

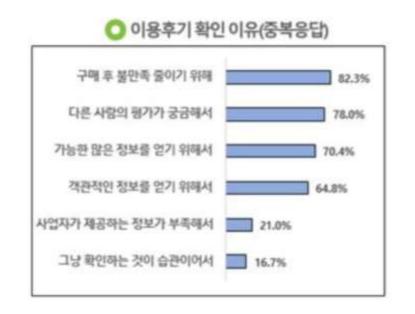
- 3.1 결과물 활용 방안
- 3.2 문제점 발생 및 해결
- 3.3 향후 개선점
- 3.4 Q&A

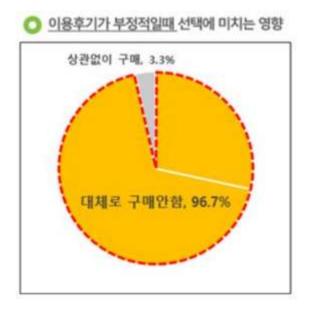




[표=한국소비자연맹]







조사결과 소비자의 대부분에 해당하는 97.2%가 구매 전 이용후기를 확인 특히 <mark>부정적 이용후기가</mark> 소비자의 구매결정에 더 많은 영향을 미치는 것으로 조사됨

문제점 발생 및 해결

- 크롤링 코드 작동 오류
- OOMKilled 오류로 인한 크롤러 파드 강제 종료
- 크롤링 진행 중 파이썬 크롤러 파드 재부팅
- Replication Factor 에러
- Zookeeper 클러스터 실행 오류
- Logstash Timestamp UTC 설정

향후 개선점

- EKS에서 Kafka 배포
- Jupyter Notebook 애플리케이션화
- 데이터 전송 단의 로그 분석 대시보드 생성
- 효율적인 협업 툴 사용

Thank you