



“뛰어난 팀은 서로 감추지 않습니다. 치부를 드러내길 꺼리지 않습니다. 비난을 두려워하지 않고 자신의 실수, 약점, 걱정을 인정합니다.” — 미국 작가 Patrick Lencioni

프로젝트 주요 일정

2023.6/12 ~ 6/21

1. 6/12 ~ 6/14 : 매일 2시간

2. 6/15 ~ 6/21 프로젝트 몰입

컨플루언스에서 제공하는 로드맵 위젯이 일 단위 VIEW 지원이 안되니, 일자를 표기한 라인을 참조할 것!

extension

미션

해당 프로젝트의 궁극적인 목적이 무엇인지, 무얼 위해 하는지 항상 상기해야 방향을 잃지 않을 수 있습니다.

1. 이후 취업을 위한 면접 과정에서 본인이 맡은 프로젝트 내용을 어필할 수 있어야 합니다.

목표

1. 프로젝트의 성과 및 결과를 구체적으로 정의하고 Task를 정의합니다.
- a. 정의한 문제와 해결하고자 하는 방향을 최대한 구체적으로 정리해야 합니다. (이 정도까지 정리한다고?) 할 정도로.. 그렇지 않으면 진행 과정에서 방향을 잃기 쉽습니다. 구체적인 성능기준이 필요하고 해당 기준이 곧 평가 시의 체크 항목으로 활용해서 객관적으로 프로젝트 결과를 평가해야 합니다.

b. 반드시 성공할 필요는 없다고 생각합니다. 오히려 기한 내에 실패했다면 해당 원인에 대해 명확하게 분석하고 문서화가 필요합니다. 현업에서는 오히려 실패에 대한 회고를 자주 했던 기억이 있네요.
2. 정의한 Task 기준으로 일정에 맞추어 프로젝트를 진행합니다.
- a. 계획하는 과정에서 이슈가 생기거나 생길 것 같다면 팀원에게 미리 공유합니다.

b. 문제를 들고 있는 것 보다 신속하게 그리고 정확하게 공유하는 게 프로젝트 진행에 상당히 중요한 태도입니다.
3. 2023.06.21(수) 프로젝트 마감 기한 전날까지 모든 과정을 완료하고, 발표할 수 있어야 합니다.
- a. 발표 자료는 곧 팀원들이 각 블로그(깃합)에 정리할 아티클이 됩니다.

주요 인원

팀원	자기소개	MBTI
김지현	<div><div>• 야구 경기(하이라이트 직관 노상관) 보는거 좋아합니다</div><div>• 나중에 한 줄 더 써 볼게요..!</div></div>	ISFJ
이상훈	<div><div>• 2년차 유부남입니다.</div><div>• 경험이나 지식이 없지만 할 수 있는한 최대한 해보겠습니다!</div></div>	ISFP or ENFP
정은아	<div><div>• 도파민 중독 스포츠 팬! 여자배구랑 야구 그리고 농구</div><div>• 추구미는 재미와 적당한 시련입니다.</div></div>	INTP
이영직	<div><div>• 취미는 농구 테니스 small talk</div><div>• 빅데이터는 한글 처리 작업 해봤습니다</div><div>• 24일 빅분기 시험을 앞두고 있습니다</div></div>	ENFP
김범기	<div><div>• 유부남 경력 1년 차입니다.</div><div>• 취미로 커피 하는 걸 좋아합니다.)</div></div>	INFJ

액션 아이템

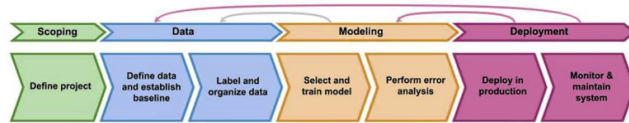
모든 회의 이후에는 회의록 및 액션 아이템이 정의되어야 합니다.

extension

- b. 발표 목차를 정하고 내용을 정리해야 합니다.
- c. 마지막 날에 발표내용을 정리하는 것 보다는 조금씩이라도 매일 진행 과정을 정리해야 합니다.
- d. 발표자를 정합시다.
  - i. 발표자는 꼭 1명이 아니어도 괜찮지 않을까?

## 머신러닝 프로젝트 프로세스 [🔗](#)

▼ 여기를 클릭하여 펼치기...



### **(2) 머신러닝 프로젝트: 기획부터 배포까지**

#### 1. 프로젝트 정의(문제 정의)

- a. 어떤 문제를 해결하고자 하는 것인가?
- b. 프로젝트를 통해 얻고자 하는 결과값은 무엇인가? (Y값)
- c. 현재 보유하고 있는 데이터가 무엇인가?(X 값)
- d. **성과지표를 설정한다.**
  - Accuracy: 음성 인식 정확도가 어느정도여야 하는지?
  - Latency: 음성 인식을 하는데 시간이 얼마나 걸리는지?
  - Throughput: 초당 쿼리 처리수는 얼마나 되는지?

#### 2. 데이터 정의

- a. 보유 데이터(X)의 데이터 타입?
- b. 예측을 위한 데이터의 양은 충분한지?
- c. 데이터를 신뢰할 수 있는가?
- d. **데이터 레이블이 일정한가?**
- e. 데이터 표준화에 대한 방법 고려
- f. 라벨링을 한다면 어떻게 할 것인가?

#### 3. 베이스라인 설정

- a. 기존방식은 성능이 어떤가?
- b. 기존보다 대략 얼마정도 성능을 개선할 것인가?

#### 4. 라벨링 & 데이터 정제하는 단계

- a. 수많은 데이터 중, 중요한 데이터는 무엇인가?
- b. 중요한 데이터가 있다면 그 이유는 무엇인가?
- c. 데이터 수집 가이드라인을 어떻게 세울것인가?
- d. 데이터에 라벨링 할 필요가 있다면, 어떻게 할 것인가?

#### 5. 모델링

- a. Train 모델 선정
  - i. 어떤 모델을 선정할 것인가?
  - ii. 선정한 이유는 무엇인가?
- b. 오류 원인 분석
  - i. 성능이 베이스라인(Baseline) 혹은 목표치(Goal)를 달성할 만큼 우수한가?
  - ii. 데이터 수가 충분하여 신뢰할만 한가? (부족하다고 판단될경우 전단계로 돌아간다)

#### 6. 배포

- a. 제품 배포
- b. **모니터링 & 시스템 유지**

회의록

날짜	참석자	논의 사항	결정 사항	액션 아이템	비고
2023년 6월 12일	1. 김지현 2. 정은아 3. 이영직 4. 김범기 5. 이상훈	1. <b>조정 결정 - 이상훈</b> a. 발표는 별도 인원 지정 가능 2. <b>최종 결과물에 대한 주요 형식 정의</b> a. 참고 프로젝트 i. <a href="#">nlp_hate_speech/README.md at master · hayoon/nlp_hate_speech</a> ii. <a href="#">R. 이우오, 오토콜리오</a> 3. <b>주제 논의 회의 6/13</b> a. 각자 하고 싶은 주제(약 2~3개) 선정 i. 정은아 : 교육 플랫폼 관련 주제 혹은 금융 쪽 공모전을 참고하는 방향 제시 ii. 이영직 : 금융 쪽 주제에 선호 iii. 김범기 : GCP나 AWS와 같은 클라우드 플랫폼 활용하여 데이터 파이프라인 구축 선호 b. 데이터 수집할지? i. 어디서 수집할 수 있는지? c. 사례 i. 실제 프로젝트 사례가 있는지? 4. <b>기타</b> a. 5명에서 분석 프로젝트를... 어떻게 분배할지? i. 시각화 ii. 담당자별 담당 모델링, 분석 결과 정리	<b>프로젝트 주제 후보</b> 1. <b>교육플랫폼 강의 추천 시스템(인프런,유데미 등)</b> a. 구글트렌드 등 이용 b. 깃(언어 점유율), 유튜브 강의 등/사람인 등 구축사이트 분석/현업자 기술블로그에서 키워드 크롤링 2. <b>환자정보와 기침소리를 이용한 covid-19 진단 모델</b> a. <b>▲[토이프로젝트] 환자정보와 기침소리를 이용한 코로나 감염여부 판단 모델 (정형데이터 + 오디오 멀티모달 실습 코드, 머신러닝 딥러닝 프로젝트 주제)</b>	<input checked="" type="checkbox"/> 6/13_프로젝트 주제 후보에 대한 예상 진행 과정 및 결과에 대한 조사  <input checked="" type="checkbox"/> 6/13_프로젝트 주제 후보 중 주제 선정	
2023년 6월 13일	1. 김지현 2. 정은아 3. 이영직 4. 김범기 5. 이상훈	주제 브레인스토밍 1. 상훈 a. 검색을 토대로 한 사실확률에 대한 예측 모델 b. 장바구니 담겨있는 상품 중, 구매확률 예측 모델 2. 상품 묶음에 따른 판매확률 예측? a. 연관 분석	선정 주제 1. Coupon Purchase Prediction  • Predict which coupons a customer will buy  • <a href="#">k Coupon Purchase Prediction</a>	<input checked="" type="checkbox"/> 테이블 스키마 및 컬럼 분석 <input checked="" type="checkbox"/> 6/13_프로젝트 과정에 대한 역할 분담 <input checked="" type="checkbox"/> 6/13_프로젝트 마일스톤 정리	<a href="#">k An EDA of coupon's transaction (in Vietnamese)</a>
2023년 6월 14일	1. 김지현 2. 정은아 3. 이영직 4. 김범기 5. 이상훈	• EDA 목록 정하기 ◦ 이상치, 결측치 확인 ◦ 데이터 특징 및 패턴 파악 ◦ 데이터 시각화 ▪ 히트맵 ▪ 산점도 ▪ 히스토그램 • 전처리 ◦ 결함 내의 다중 값 처리하기? ◦ NULL 값 변환? , 평균값? 중앙값? ◦ 전처리 할 목록 정하기- 일본어 전처리 - 이해할 수 있는 값으로 (영어 or 한국어) ◦ hash-id 필터링 ▪ hash 로 들지? 이진값으로 변환할지? ◦ 테이블 별 변수 타입별 분류 ▪ 타입별 전처리 단계에서 방식을 구분하는 접근? ◦ 테이블 조인(어떤 테이블을 어떻게 조인할지) • 피처링 ◦ 모델에 넣을 피쳐, 라벨 정하기 ◦ 타겟(=라벨) 선정? ▪ 쿠폰 사용 카테고리 ▪ 쿠폰 할인 가격 ▪ 쿠폰 할인율 ▪ 쿠폰 사용 가능 일수 ▪ 홀리데이 사용 가능 여부 • 머신러닝 모델 선정 ◦ (딥러닝) 사용해서 모델 만들어라 ▪ 예)합성곱... ◦ • 모델 평가 • 모델 튜닝 및 개선 • 결과 해석	1. 피쳐 선정 진행 중 2. ML 모델 검토 진행 중 3. 빅쿼리 환경 세팅 완료  1. 전처리 수행 2. 훈련, 테스트 데이터 분리 3. 훈련 시작 .....	<input checked="" type="checkbox"/> 6/13_프로젝트 작업 리스트업 @BK <input checked="" type="checkbox"/> 6/14_VertexAI 환경 세팅 @BK <input checked="" type="checkbox"/> 6/14_프로젝트 과정에 대한 역할 분담 <input checked="" type="checkbox"/> 6/14_프로젝트 마일스톤 정리	
2023년 6월 15일		<b>TO_DO_List</b>  1. <b>분석 주제(가설) 구체적으로 설정</b>  분석의 목표/구체적으로 적용해볼 모델 종류/해당 모델에 넣은 피쳐(input)데이터 설정/모델을 통해 도출할 예측값(target)설정/모델 성능 평가 지표 설정(ex. AUC, confusion matrix etc...)  • 캐글 컴페티션에서 요구 모델1 - 코사인 유사도 - 이론 및 코드 해석이 어려워 후순위로  • 캐글에서 요구 모델2 - 로지스틱으로 구현 + 하이퍼 파라미터 튜닝 추가 - For 구매여부 예측		<input type="checkbox"/> coupon_list 테이블 usable_date ~ (na.0.1.2) 값 의미 파악_ @정은아  <input type="checkbox"/> 각 다른 테이블별 EDA 결과 6/16 공유  <input type="checkbox"/> @지현 - Cupon_area, Cupon_list_ @정은아  <input checked="" type="checkbox"/> @훈씨네마 - Coupon_detail  <input checked="" type="checkbox"/> @이영직 - Coupon_visit	



**참조** 

10

#### 2019년도 구축 빅데이터 플랫폼

- o 통신 빅데이터 플랫폼: [bdp.kt.co.kr](http://bdp.kt.co.kr)
- o 교통 빅데이터 플랫폼: [국가고통 데이터 오픈마켓](#)
- o 문화 빅데이터 플랫폼: [init-page](#)
- o 환경 빅데이터 플랫폼: [환경 빅데이터 플랫폼](#)
- o 중소기업 빅데이터 플랫폼: [datastore.wehago.com](http://datastore.wehago.com)
- o 지역경제 빅데이터 플랫폼: [ggdata.kr](http://ggdata.kr)
- o 금융 빅데이터 플랫폼: [fnbigdata.com](http://fnbigdata.com)
- o 헬스케어 빅데이터 플랫폼: [cancerportal.kr](http://cancerportal.kr)
- o 유통소비 빅데이터 플랫폼: [kdx.kr](http://kdx.kr)
- o 산림 빅데이터 플랫폼: [forestdata.kr](http://forestdata.kr)

#### 2020년도 구축 빅데이터 플랫폼

- o 소방안전 빅데이터 플랫폼: [소방안전 빅데이터 플랫폼](#)
- o 스마트치안 빅데이터 플랫폼: [스마트 치안 빅데이터 플랫폼](#)
- o 해양수산 빅데이터 플랫폼: [해양수산빅데이터 거래소](#)
- o 농식품 빅데이터 플랫폼: <https://kadx.co.kr/>
- o 라이프로그 빅데이터 플랫폼: <https://www.bigdata-lifelog.kr/portal>
- o 디지털 산업혁신 빅데이터 플랫폼: [디지털산업혁신 거래소](#)

#### GCP 관련 링크

1. [딥러닝 프로젝트를 위한 클라우드 GPU 자원, Google Cloud Platform](#)
2. [모듈 4: GCP를 활용한 머신러닝](#)
3. [GCP AI-Platform에서 딥러닝 학습하기](#)
4. [Google Cloud Big Data and Machine Learning Fundamentals 한국어](#)
5. [Google Cloud에서 머신러닝 구현을 위한 권장사항 | 클라우드 아키텍처 센터](#)
  - a. [GCP](#)