

산학교육센터 T-SUM 주차별 활동 보고서

팀명	TSUANMI	장소	과학기술 2관 225호	일시 (주차)	4주차, 2024. 04. 30.(화), 19시 ~ 21시 30분 (150분)
----	---------	----	--------------	---------	---

활동 주제	이진 분류
-------	-------

[로지스틱 회귀 학습]

4주차에서는 이진 분류 수업을 진행했다. 이전 회차까지 선형 회귀에 대한 수업을 진행했는데, 선형 회귀는 한 종속 변수와 여러 독립 변수 간의 관계를 파악하는 것으로, 집값 예측, 당뇨 수치 예측 등이 있다. 이러한 회귀 모델을 통해 특정 변수의 값을 분류할 수 있다. 예를 들어, 당뇨 수치를 통해 당뇨병이 있는지/없는지 집값을 예측한 값을 통해 집을 구매할 것인지/아닌지 등으로 분류할 수 있다. 이러한 분류 문제는 단순히 두 개로 분류하는 이진 분류 모델이 아니라, 여러 카테고리로 분류하는 다중 분류로도 이어질 수 있다. 당뇨병 수치를 더 잘게 나누어 매우 높음, 높음, 정상, 낮음, 매우 낮음과 같이 5개의 카테고리로 나눌 수 있다. 이러한 분류 모델은 이전에 배운 선형 회귀 모델에서 조금의 확장으로 분류 모델을 설계할 수 있다. 이전에 만든 선형 회귀 모델의 선형인 값을 비선형으로 바꿔주는 활성화 함수 하나만을 추가하면 분류 모델이 만들어진다. backpropagation 과정도 선형 회귀와 동일하고, 미분된 값도 동일하기 때문이다. 즉, 로지스틱 회귀는 새로운 개념이 아닌, 이전에 배운 개념의 추가한 것임을 멘티들에게 강조했다. 이론 수업이 끝난 후 실제 코드를 작성해보므로써 이론 내용을 한 번 더 정리하고 코드로 분류 모델을 어떻게 구축하는지를 학습했다.

활동 내용

1. Logistic Regression
분류(classification)

- 데이터의 레이블로 학습하고 레이블이 없는 데이터를 예측
- 이진 분류(binary Classification) 외 다중 분류(multiclass Classification) 로 구분됨
- 결정 경계를 기준으로 분류를 진행

1. Logistic Regression
구조

- 이진 분류 알고리즘
- $[-\infty, \infty]$ 의 범위를 가지는 z의 값을 시그모이드 함수를 통해 $[0, 1]$ 로 변환 (학습처럼 해석이 가능)
- 임계 값을 통해 0 또는 1로 구분
- 손실 함수: 크로스 엔트로피 함수

[그림 1] (좌) 분류 정의 (우) 로지스틱 회귀 구조

```

1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score
3
4 model = LogisticRegression()
5 model.fit(X_train, y_train)
6
7 pred = model.predict(X_test)
8 print(f'Acc : {accuracy_score(y_test, pred)}')

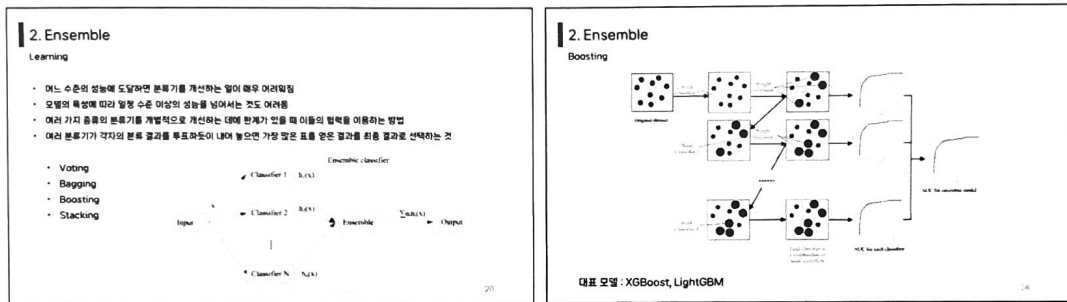
```

executed in 58ms, finished 20:28:48 2024-04-30
Acc : 0.956140350877193

[그림 2] 로지스틱 회귀 코드

[앙상블 기법 학습]

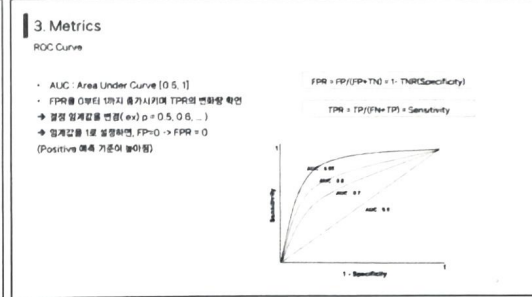
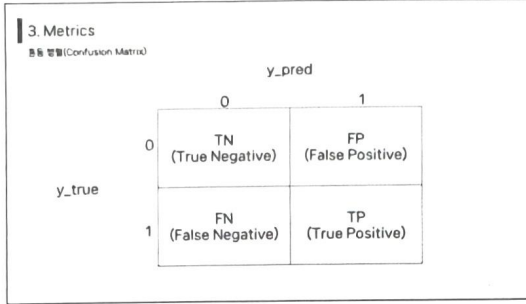
다음으로 앙상블 모델에 대해 학습했다. 우리가 의사 결정을 할 때 한 명의 전문가의 의견보다 여러 명의 전문가들의 의견을 모아서 결정하는 것이 더 나은 경우가 많다. 마찬가지로, 앙상블은 여러 모델의 의견을 모아서 더 나은 예측을 할 수 있도록 한다. 즉, 앙상블은 여러 개의 다양한 기본 모델을 조합하여 더 강력한 모델을 만드는 기법이다. 주로 사용되는 앙상블 기법으로는 Voting, Bagging, Boosting, Stacking 등이 있다. 이러한 다양한 방식으로 다수의 모델을 조합하여 예측을 수행하며, 그 특징에 따라 다양한 상황에 적용될 수 있다. 멘티들에게 앙상블 기법을 강조하여 수업을 진행했는데, 앞서 말했듯이, 단순히 로지스틱 회귀 모델 혹은 SVM 모델 등의 단일 모델만 사용해서 데이터 분석을 진행하지 않는다. 이는 좋은 결과를 보이지 못할 뿐만 아니라, 모델의 특성 중요도를 파악할 수 없는 등의 문제가 있다. 그래서, 이번 경진대회에서도 앙상블 모델을 사용하여 예측을 수행할 것이기 때문에 멘티들에게 앙상블 기법에 대해 잘 이해할 수 있도록 최대한 자세하고 반복적으로 설명해주었다.



[그림 3] (좌) 앙상블 기법 (우) Bagging 과정

[평가 지표 학습]

머신러닝 모델을 만든 뒤, 그 모델의 성능을 평가하는 것도 중요하다. 모델의 정확도와 같은 평가 지표를 확인하고 그 모델을 그대로 사용할 것인지 아니면 조금 더 학습시키고 사용할 것인지 정해야 한다. 특히, 분류 모델에서는 이 과정이 매우 중요하다. 단순히 정확도만을 보고 모델을 평가하면 안 된다. 예를 들어, 정확도가 90%라고 해서 반드시 좋은 모델은 아니다. 예를 들어, 100명의 환자 중 98명이 양성, 2명이 음성이라고 하자. 만약 이 모델을 정확도만을 보고 환자의 양성/음성 여부를 예측한다면, 단순히 100명 모두 양성이라고 해도 정확도는 98%가 되기 때문에 좋은 모델이라고 착각할 수 있다. 그렇기에 단순히 정확도만 보는 것이 아닌, 다른 지표 또는 혼동 행렬(confusion matrix)을 확인함으로써 모델을 평가하는 것이 좋다. 그렇다고 해서, 계속해서 모델을 복잡하게 만들면 이 모델은 오로지 훈련 데이터에서만 좋은 성능을 보이고 실제 예측 데이터에서는 좋은 성능을 보이지 못하는 과적합이 발생할 수 있다. 그렇기에 모든 일반적인 성능을 내는 모델을 만드는 것이 중요하다. 이 내용은 인공지능 분야에서 매우 중요한 내용이고, 새로운 모델을 만들 때, 자신이 만든 모델의 성능을 설명해야 하므로 이 내용을 정확하게 알고 있는 것이 중요하다. 멘티들에게 이 내용을 강조하여 단순히 수치만을 보고 모델의 성능이 좋다고 판단하지 않도록 했다. 이론 수업 후 간단한 퀴즈를 직접 손으로 계산해보면서 평가 지표를 이해하도록 했고, 그 후 코드를 통해 그 계산한 값이 맞는지 확인해보았다.



[그림 4] (좌) 혼동 행렬 (우) ROC Curve

```

1 y_test = [0, 0, 1, 1, 0, 1, 1, 0, 0, 1]
2 pred = [0, 1, 1, 1, 1, 0, 0, 1, 1, 1]
3
4 print(f'accuracy_score : {accuracy_score(y_test, pred)}')
5 print(f'precision_score : {precision_score(y_test, pred)}')
6 print(f'recall_score : {recall_score(y_test, pred)}')
7 print(f'f1_score : {f1_score(y_test, pred)}')

```

executed in 25ms, finished 20:16:15 2024-04-30

accuracy_score : 0.4
precision_score : 0.42857142857142855
recall_score : 0.6
f1_score : 0.5

[그림 5] 다양한 평가 지표 코드

[3주차 경진대회 활동 피드백]

머신러닝 수업이 끝나고 멘티들이 약 3주 동안 준비한 경진대회 내용에 대해 피드백을 해주었다. 멘티들에게 이번 4주차 활동 전까지 경진대회 주제를 생각해오라고 공지했으며, 이번 회차에 그에 대한 피드백을 해주었다. 멘티들은 멘토링 활동 후 남아서 회의를 진행했으며, 활동이 없는 날에도 지속적으로 회의를 진행했다. 그 결과, 각자 5개의 주제를 정했으며, 그중 가장 좋은 주제 5개를 골라 그 주제에 대한 데이터를 탐색하여 멘토에게 보고서 형식으로 제출했다. 각 주제에 대한 피드백은 [그림 6]과 같다. 요약하자면, 아직 멘티들이 ‘머신러닝을 활용한 경진대회’에 대해 정확하게 파악하지 못했다. 인공지능으로 어떠한 target 값을 예측을 할 수 있지만, 단순히 예측하는 것에 초점을 맞추다 보니 데이터 분석의 의미를 잃는 경우가 많다. 왜 인공지능 방법을 사용해야 되는지, 기존의 방법 말고 새로운 방법을 왜 사용해야되는지, 이것을 예측함으로써 얻을 수 있는 의미는 무엇인지 등을 생각하지 않다 보니, 주제와 데이터가 부실해지는 경우가 많다. 현재 멘티들도 이러한 상황에 있어, 이러한 피드백을 해주며, 단순히 인공지능을 사용하는 것이 아니라, 현실 세계의 문제점을 해결하는 것에 초점을 맞추어 경진대회 주제를 선정하도록 조언해주었다. 또한, 멘티들이 데이터를 선정할 때, 단순히 데이터의 ‘크기’만을 보고 선정하는 것이 아니라, 데이터의 ‘특성(feature)’을 잘 보고 결정하도록 피드백 해주었다. 결국 데이터의 품질이 좋아야 좋은 결과가 나오는데, 데이터가 어떤 정보를 담고 있는지 명확하게 알아야 유의미한 분석을 진행할 수 있기 때문에, 데이터를 유심히 살펴보도록 강조했다.

<p>1. 태양광 발전소 용역 발전소 발전량 예측 인증지능 (서지연)</p> <p>가. 설명 : 태양광 발전소 및 풍력 발전소의 발전량을 예측하여 발전소를 어디에, 얼마나 더 지어야 하는지 예측</p> <p>나. 데이터 : 태양광, 풍력, 발전량 : https://www.data.go.kr/data/15065269/fileData.do</p> <p>피드백</p> <p>1. 발전량을 예측하는 발전소를 어디에 얼마나 더 지어야 하는지 예측을 왜하는지 모르겠다</p> <p>2. 발전량을 예측하는 과정이 단순 시계열 데이터이기 때문에 머신러닝으로 풀리는게 직감인데 모르겠다 (예측 과정이 너무 쉽고 보여줄만한게 없을 것 같음)</p> <p>2. 체중 관리 인증지능 (이소호)</p> <p>가. 설명 : 운동량 섭취칼로리, 체중 데이터들을 학습하여 목표 체중이 되기 위한 운동 및 식단관리 계획에 도움(혹은운동량, 섭취칼로리에 따른 체중 변화예측)</p> <p>나. 데이터 : https://www.kaggle.com/datasets/vechoo/diet-plan-recommendation</p> <p>피드백</p> <p>1. 데이터 관점에서만 현재 본 후보 중 그나마 가장 best</p> <p>2. 주제와 데이터가 맞지 않음. 또한, 인증지능을 굳이 왜 해야되는지 모르겠음 (그냥 체중계 쓰면 되는데 어딴가)</p> <p>5. 신제품 감가 예상 인증지능(보유) (최연우)</p> <p>가. 설명 : 신제품이 나왔을때 나중에 얼마나 감가될지를 예측해서 소비자가 합리적인 소비를 할 수 있도록 도움</p> <p>나. 데이터 : 다나와 : 출시일, 평점, 가격변동표, 후기 수를 크롤링하여 csv데이터 셋 제작 https://danawalab.github.io/</p> <p>피드백</p> <p>1. 목적은 가장 좋음</p> <p>2. 나중에가 언제이고 평점, 후기 수 등이 감가에 얼마나 영향을 주는지 모르겠음</p> <p>3. 여러 구매처가 있는데, 이것들을 어떻게 활용할 수 있는지 모르겠음</p>	<p>3. 연료 소비량, 전력 사용량에 따른 지구온도 예측 인증지능 (서지연)</p> <p>가. 설명 : 연료 소비량, 전력 사용량과 지구온난화의 상관관계를 연료소비량을 얼마나 줄여야 지구온난화를 낮출 수 있을지 예측하여 기후위기 문제 해결에 도움</p> <p>나. 데이터 : 1) 연료소비량, 전력 2) 전력사용량, 전력 3) 지구온도 : https://www.kaggle.com/datasets/vechoo/parac/temperature-chaan 이어서 시계열 분석에서 사용</p> <p>피드백</p> <p>1. 연료 소비량과 지구온도, 전력 사용량 데이터가 보편화되어있지 않음</p> <p>2. 난수적 연도에 따른 값이모순 같아지니 있는데 연료, 전력이 시계열 값을 갖지 않으면 안되는 것 같습니다</p> <p>3. 지구온난화로 인해 온도가 예측을 어떻게 해야 할지 모르겠다</p> <p>4. 나이 및 소득별 지출의한 통계 인증지능(보유) (최정민)</p> <p>가. 설명 : 나이와 소득에 따라 들어 어느곳에 소비를 할지 후 불필요한 곳에 쓰이는 돈을 아낄 수 있도록 도움(이와 소득별로 어디에 지출을 더 많이 할것인지 예측하는 인증지능)</p> <p>나. 데이터 : https://www.bigdata-culture.kr/bigdata/user/data_market/detail.do?id=2d4cb9f0-4386-11ec-a3e9-d36274dc304b#</p> <p>피드백</p> <p>1. 예측 목적을 명확하게 모르겠음</p> <p>2. 어디에 지출을 할지에 대한 요건 그 사실에 대한 것이 편입시 같이 더 보충을 주지 않았기 때문이다</p>
---	---

[그림 6] 3주차 경진대회 활동 피드백

[4주차 과제]

4주차 과제로 개인 과제, 경진대회 과제를 부여했다. 개인 과제로는 다양한 분류 모델을 각자 2개씩 조사해서 ppt로 만드는 것이다. 이번 주차에서 로지스틱 회귀 모델을 배웠는데 분류 모델은 매우 많아서 이를 활동 시간에 모두 다루는 것은 무리가 있다. 그래서 멘티들이 각자 공부하고 자료를 만듦으로써 자신의 학습 능력을 향상시키고, 다양한 분류 모델에 대한 깊은 이해도를 갖출 수 있다. 또한, 이러한 과정을 통해 멘티들은 자신이 조사한 모델의 장단점, 주요 특성, 성능 등을 교류하며, 여러 분류 모델 간의 비교 분석을 통해 문제 해결에 적합한 모델을 선택하는 능력을 개발할 수 있을 것으로 기대된다. 경진대회 과제로는 이번 피드백을 통해 주제 구체화, 데이터 선정, 주제 목적 등을 다시 생각해보는 것이다. 이번 과제를 통해 경진대회 주제 및 목적이 매우 구체화될 것으로 기대되고, 이는 전반적인 데이터 분석의 과정을 설계하고 분석하는데 큰 도움이 될 것으로 기대된다.

<p>과제 3</p> <ul style="list-style-type: none"> 개인 과제 (-5/6(목) 23:59까지) 분류 모델 공부 <p>[포함 내용]</p> <ol style="list-style-type: none"> 분류 알고리즘 (머신러닝으로 분류를 하는지) 모델의 장단점 sklearn 모델의 5요 하이퍼파라미터 설명 breast cancer 데이터를 가지고 학습 후 분류 (로지스틱 회귀 모델처럼) 분류 결과 출력 <ol style="list-style-type: none"> confusion matrix Accuracy, Precision, Recall, F1-score ppt로 만들기 	<p>과제 3</p> <ul style="list-style-type: none"> 경진대회 과제 (-5/6(목) 23:59까지) <ol style="list-style-type: none"> 주제 설명 기본 방법 문제점 탐색 baseline 탐색 <ul style="list-style-type: none"> 각자 역할 표시 팀장이 회의를 날 달하며 표시 (필요한지) 각자 맡은 역할 정리해서 경진대회 출제에 각자 업로드 팀장은 최종 정리본 하나 업로드 <p>>> 필수개의 파일이 있어야 함(= 사용 data)</p>
--	---

[그림 7] 4주차 과제

2. 간단하게 오늘 수업에 대한 자기 이해도 피드백
응답 4개

Metrics 퀴즈랑 코드 돌려보면서 조금 감 잡을 수 있었던 것 같습니다

따라오기 괜찮았가

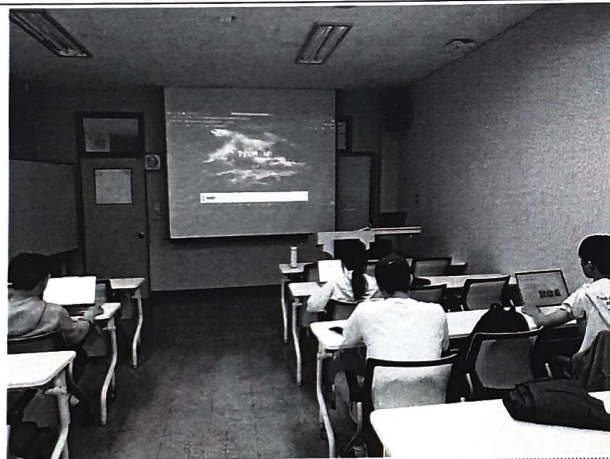
로지스틱 회귀를 잘 이해했고, 혼동행렬의 재현률, 정밀도 등의 계산방법을 잘 숙지했다.

예시 데이터로 예측해보고 heatmap으로 결과 확인

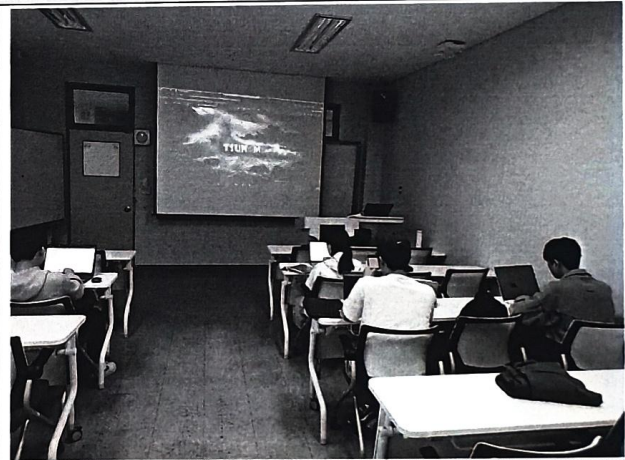
[그림 8] 멘티 셀프 피드백

멘티별 활동내용	멘티 이름	활동 내용 및 역량 증가 정도
	황성백	수학적 이해력이 매우 높고, 수업 내용을 잘 이해하고 있는 것으로 보인다.
	이소호	질문을 자주 하고, 대답도 열심히 하면서 수업에 매우 적극적으로 임한다.
	서지연	큰 어려움 없이 수업 내용을 잘 따라오는 것으로 보인다.
	최연우	프로젝트에 적극적으로 임하며, 자신의 아이디어를 적극적으로 잘 표현한다.

증빙 사진



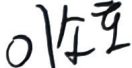

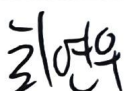


[시작]



[끝]

참여자 서명부

구분	소속	학번	성명	서명
멘토	컴퓨터융합소프트웨어학과	2019270641	박재현	
멘티	컴퓨터융합소프트웨어학과	2022270608	황성백	
멘티	컴퓨터융합소프트웨어학과	2022270611	이소호	
멘티	컴퓨터융합소프트웨어학과	2022271323	서지연	
멘티	컴퓨터융합소프트웨어학과	2023270639	최연우	
지도교수		노경숙 교수님		