

Song popularity prediction - (Base ML model)

목표 · Dataset 이해 / clean up (필요시)

- song popularity 예측 가능 regression 모델 생성
- 모델 평가 (R^2 , RMSE, ...)

순서

1. Data exploration → 2. EDA → 3. 전처리 → 4. Data manipulation → 5. Feature selection / extraction.
6. predictive modeling → 7. outcome / conclusion.

□ Data exploration

→ 필요한 library / dataset 로드.

- song_name 은 지명, target 변수: "song_popularity"
- `df.info()` 를 통해 null 값 확인 → null 값 삭제
- `df.nunique().sort_values()` 이용하여 변수 확인
- numerical / categorical 변수 나누기.

□ EDA

- target 변수의 분포를 확인 (점점대로 줄어듦 확인)
⇒ 전반적으로 점점대로 줄어듦 모양을 가진다.
- categorical 변수의 분포를 plot을 통해 확인
 { audio mode: 0/1 (0이 2배 많다)
 time signature (0/1/2/3/4/5 - 4가 제일 많이 많다)
 key (0/11) 글자 분포 확인.
- numerical 변수의 분포를 plot을 통해 확인
⇒ skewed 된 모습이 많이 보인다.
- pair-plot 을 이용하여 모든 변수 간 상관관계를 파악.

□ 데이터 전처리

- Duplicate row 삭제. (총 311개 삭제)
- categorical 변수를 numerical 로 전환.
 { one-hot encoding
 dummy encoding
- outlier 지기 (IQR을 이용)
- 삭제된 data의 비율을 살펴봄
 { 삭제: 62.5%
 4등: 49.5% } Cleanup 까지 완료.

④ Data manipulation.

- train set / test set 으로 data 구분.
- Feature scaling (standardization)

⑤ Feature selection / Extraction.

- 변수 간의 correlation (heatmap) 표현.
- ⇒ 다양한 correlation 이 보인 것으로 판단. → 여러 변수들 간의 상관성이 높음을 판단.
- OLS를 통해 OLS Regression Result를 알아낸다. (AIC, BIC 등..)
- ⇒ 다중공선성 (multicollinearity)이 나타난다.

- 1) minimal method - Variance Inflation Factor (VIF)
- 2) Automatic method - Recursive Feature Elimination (RFE)
- 3) Feature elimination using PCA decomposition.

1) VIF

- 다중공선성을 판단할 때 이용
- 예측변수들이 상관성이 높으면 평균 해 계수의 신뢰 크기를 측정하고.
- 신뢰가 저절수록 신뢰도 ↓ (VIF 값이 10이면 다중공선성이 있어 모형은 신뢰, VIF 10 → 변수 선택의)

2) RFE

- 직관적인 feature selection 방식.
- 강하게 feature에 대해 훈련 후, 증명하지 않은 feature 제거.

3) PCA decomposition. (주성분분석)

⑥ Predictive modeling

1) Multiple Linear Regression (MLR)

- Actual model, prediction model의 분포도 그래프로 표현.
- Training set / Testing set의 R2-score, MSE, RMSE 구함.
- residual plot까지 그려기.

2) Ridge regression model.

- 0과 동일한 값 구한다.

3) Lasso regression model

4) Elastic-net

5) polynomial regression model

⑦ 위의 5개의 모델을 비교.

- 비교 결과, model 별로 R2 값이 다양하게 나온다.

→ 5개의 model에 대한 비교정렬 (plot)

→ 결과중, polynomial model의 R2 값이 가장 크다.

→ model 별 RMSE 값도 비교

↳ RMSE 값이 작을수록, training / test 의 값이 비슷할 수록 좋은 model.

→ RMSE 값으로는 MLR model 이 가장 좋다. PNR은 overfitting 된 것이다.

VIF / RFE / PCA

(Feature selection 의 방식)

1. Filter method (Feature 간 상관성) - VIF (변량팽창)
2. Wrapper method (Feature를 조정하여 최적 예측 → 성능파악, Feature 선택)
3. Embedded method (머신러닝 최적화, 리지회귀 특징 강제에서 각 Feature를 선택)

• VIF: 10이 넘으면 다중공선성이 있는 것으로 판단. 5로 넘겨라로 취급.

→ VIF값이 무조건 10이 넘는다고 해서는 안되며, 상황에따라 제거하거나 혹은.

• RFE: Backward 방식 중 하나. 모든 변수 다 포함 → 반복적으로 제거 → 중요도가 낮은 것부터 삭제.

• PCA: 주성분이 여러 개일 때.