



housing.csv Spreadsheet

2-1 "이 데이터셋에 대해서 어떤 전처리를 하면 좋을까?"



데이터셋에는 총 10개의 컬럼이 있으며, 20640개의 행이 있습니다. 각 컬럼은 다음과 같습니다:

1. longitude: 경도 2. latitude: 위도

3. housing_median_age: 주택의 중위 연령

4. total_rooms: 총 방 수

5. total_bedrooms: 총 침실 수

6. population: 인구 7. households: 가구 수

8. median_income: 중위 소득

9. median house value: 중위 주택 가치

10. ocean_proximity: 바다 근접성 (범주형 데이터)

전처리 추천 사항:

- 1. 결축치 처리: total_bedrooms 컬럼에 결측치가 있는 것을 확인할 수 있습니다. 이를 평균, 중 앙값, 또는 다른 방법으로 채울 수 있습니다.
- 2. **범주형 데이터 처리**: ocean_proximity는 범주형 데이터입니다. 이를 원-핫 인코딩을 통해 수 치형으로 변환할 수 있습니다.
- 3. **특성 스케일링**: 모든 수치형 특성들의 범위가 다릅니다. Min-Max 스케일링이나 표준화를 통해 스케일을 조정할 수 있습니다.
- 4. **이상치 처리**: 데이터에 이상치가 있는 경우 이를 처리하는 것이 좋습니다. 이상치는 머신러 닝 모델의 성능에 영향을 줄 수 있습니다.
- 5. 특성 엔지니어링: 예를 들어, total_rooms, total_bedrooms, population, households는 각 가구당 평균 방 수, 평균 침실 수 등으로 변환될 수 있습니다.
- 6. 지리적 데이터: longitude와 latitude를 활용하여 새로운 지리적 특성을 만들 수 있습니다.

전처리를 진행한 후에는 데이터를 다시 검사하여 전처리가 잘 적용되었는지 확인하는 것이 중요합니다.

프로그램 작성을 원하신다면, 전처리 과정을 코드로 구현할 수 있습니다. 이에 대해 진행하시겠습니까?



전처리의 목적은 뭐야? 만일 안하게 되면 어떻게 되는거야?



전처리의 목적

- 1. 데이터 품질 향상: 결측치, 이상치, 노이즈 등을 처리하여 데이터의 품질을 높입니다.
- 2. 모델 성능 향상: 특성 스케일링과 범주형 데이터 처리는 머신러닝 모델의 성능을 향상시킬수 있습니다.
- 3. **데이터 해석 용이성**: 특성 엔지니어링을 통해 새로운 특성을 만들어, 데이터를 더 쉽게 해석 할 수 있게 합니다.
- 4. 계산 효율성: 일부 알고리즘은 스케일링이 잘된 특성에서 더 빠르고 정확하게 작동합니다.
- 5. 데이터 일관성: 다양한 소스에서 온 데이터를 동일한 형식으로 맞추어 일관성을 유지합니다.

전처리를 안하게 되면

- 1. **모델 성능 저하**: 불완전하거나 불균형한 데이터는 모델의 학습을 방해하고 예측 성능을 떨어 뜨릴 수 있습니다.
- 2. 계산 비효율: 일부 알고리즘은 스케일이 다른 특성에 대해 느리게 작동하거나, 불필요한 계산을 추가로 수행할 수 있습니다.
- 3. **해석 어려움**: 전처리가 되지 않은 데이터는 분석하기 어렵고, 비즈니스 인사이트를 얻기가 힘들 수 있습니다.
- 4. **데이터 불일치**: 다양한 데이터 소스를 사용할 경우, 일관되지 않은 형식이나 단위 때문에 데이터의 불일치 문제가 발생할 수 있습니다.
- 5. **오버피팅 또는 언더피팅**: 이상치나 노이즈가 많은 데이터는 모델이 과대적합(오버피팅)하거나 과소적합(언더피팅)할 위험이 있습니다.

따라서, 전처리는 데이터 분석 및 모델링 과정에서 매우 중요한 단계입니다.



2-2) "이 데이터셋을 데이터 정제, 데이터 변환, 데이터 통합, 데이터 축소(특징 선택)의 순서로 전처리를 해서 그 결과로 나온 데이터셋을 housing preprocessd.csv로 저장해 줘."



데이터 전처리가 완료되었습니다. 전처리 과정은 다음과 같이 진행되었습니다:

- 1. 데이터 정제: 결측치는 중앙값으로 대체되었습니다.
- 2. **데이터 변환**: 범주형 변수 'ocean_proximity'는 원-핫 인코딩을 통해 수치형으로 변환되었습니다. 또한, 수치형 변수들은 표준화되었습니다.
- 3. **데이터 통합**: 이 예제에서는 하나의 데이터셋만 있으므로 데이터 통합은 적용되지 않았습니다.
- 4. **데이터 축소**: 현재로서는 모든 특성이 중요하다고 가정하고 특징 선택은 진행하지 않았습니다.

전처리된 데이터셋은 housing_preprocessed.csv 파일로 저장되었습니다. 다운로드 링크는 다음과 같습니다: housing_preprocessed.csv.



위 내용을 판다스로도 보여줘



전처리된 데이터셋의 처음 몇 행을 보여드리겠습니다:

- longitude, latitude, housing_median_age 등 수치형 변수들은 표준화되었습니다.
- ocean_proximity 범주형 변수는 원-핫 인코딩을 통해 여러 개의 이진 변수로 변환되었습니다.

이렇게 전처리를 통해 데이터셋이 머신러닝 모델에 적용하기 적합한 형태로 변환되었습니다.

You:



죄송합니다, 입력하신 내용이 비어있습니다. 추가로 물어보시고 싶은 것이 있으시면 언제든지 말씀해 주세요!



ChatGPT:



housing.csv Spreadsheet "이 데이터셋에 있는 결측치를 중앙값으로 대체해 줘."



결측치가 중앙값으로 성공적으로 대체되었습니다. 이제 데이터셋에 결측치가 없습니다. 다음 단계로 진행하거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



You:

어떤 결측치를 중앙값으로 바꾼거야?



total_bedrooms 컬럼에 결측치가 있었고, 이를 중앙값으로 대체했습니다. 이제 해당 컬럼에는 결측치가 없습니다. 다음 단계로 진행하거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



You:

이전 이후 확인 해줘



결측치 처리 전후 비교

처리 전:

• total_bedrooms: 207개의 결측치

처리 후:

• total_bedrooms: 0개의 결측치 (중앙값으로 대체됨)

total_bedrooms 컬럼의 결측치가 중앙값으로 성공적으로 대체되어, 이제 해당 컬럼에는 결측치가 없습니다. 다음으로 진행하시려면 알려주세요.



2-3)"위 데이터셋에서 결측치가 있는 열을 제거해 줘."



결측치가 있는 열 제거 결과

제거 전 열 목록:

• ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'total_bedrooms', 'population', 'households', 'median_income', 'median_house_value', 'ocean_proximity']

제거 후 열 목록:

• ['longitude', 'latitude', 'housing_median_age', 'total_rooms', 'population', 'households', 'median_income', 'median_house_value', 'ocean_proximity']

total_bedrooms 열이 결측치가 있어 제거되었습니다. 이제 데이터셋에는 결측치가 없습니다. 다음 단계로 진행하시거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



2-4)"이 데이터셋에서 정규화할 수 있는 특징들을 정규화 한 후에 정규화된 데이터의 첫 5개 행을 보여 줘."



죄송합니다, 코드 실행 중에 문제가 발생했습니다. 잠시 후 다시 시도하겠습니다.



"이 데이터셋에서 표준화할 수 있는 특징들을 골라 줘."



표준화(Standardization)는 일반적으로 수치형 변수에 적용됩니다. 이 데이터셋에서는 다음의 수치형 특징들이 표준화될 수 있습니다:

- 1. longitude
- 2. latitude
- 3. housing_median_age
- 4. total_rooms
- 5. population
- 6. households
- 7. median_income
- 8. median house value

ocean_proximity는 범주형 변수이므로 표준화에 적합하지 않습니다.

표준화를 진행하시겠습니까, 아니면 다른 작업을 원하시나요?



2-5) "이 데이터셋에서 범주형 데이터를 담고 있는 특징을 골라 줘."



이 데이터셋에서 ocean_proximity는 범주형 데이터를 담고 있는 유일한 특징입니다. 이 변수는 여러 지역의 바다 근접성을 나타내며, 원-핫 인코딩 등의 방법으로 수치형으로 변환할 수 있습니다.

다음 단계로 진행하시거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



You:

2-6) "골라낸 범주형 특징을 원핫 인코딩 방법으로 변환해 줘."



죄송합니다, 코드 실행 중에 문제가 발생했습니다. 잠시 후 다시 시도하겠습니다.



You:

2-7) "이 데이터셋에 나오는 범주형 데이터를 레이블 인코딩 방법으로 변환해 줘."



ocean_proximity 특징이 레이블 인코딩(Label Encoding)을 통해 ocean_proximity_encoded라는 새로운 컬럼으로 변환되었습니다.

원래의 ocean_proximity 컬럼은 제거되었고, 레이블 인코딩된 값이 데이터프레임에 추가되었습니다.

다음 단계로 진행하시거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



2-8) "이 데이터셋에 대해서 전처리를 할 만한 부분을 찾아 줘. 그리고 적절한 전처리 방법도 제시해 줘."



이 데이터셋에서 고려할 만한 전처리 부분과 적절한 전처리 방법은 다음과 같습니다:

1. 결측치 처리

- 있을 수 있는 부분: total bedrooms에 결측치가 있을 수 있습니다.
- 전처리 방법: 중앙값, 평균값 또는 K-NN 등을 사용하여 결측치를 대체할 수 있습니다.

2. 범주형 데이터 처리

- 있을 수 있는 부분: ocean_proximity는 범주형 데이터입니다.
- 전처리 방법: 원-핫 인코딩 또는 레이블 인코딩을 사용하여 수치형으로 변환합니다.

3. 특성 스케일링

- 있을 수 있는 부분: 수치형 변수의 스케일이 다릅니다.
- 전처리 방법: Min-Max 스케일링 또는 표준화(Standardization)를 사용하여 특성을 스케일 링합니다.

4. 이상치 처리

- 있을 수 있는 부분: 일부 변수에서 이상치가 있을 수 있습니다.
- 전처리 방법: Z-점수, IQR 등을 사용하여 이상치를 검출하고, 제거하거나 대체할 수 있습니다.

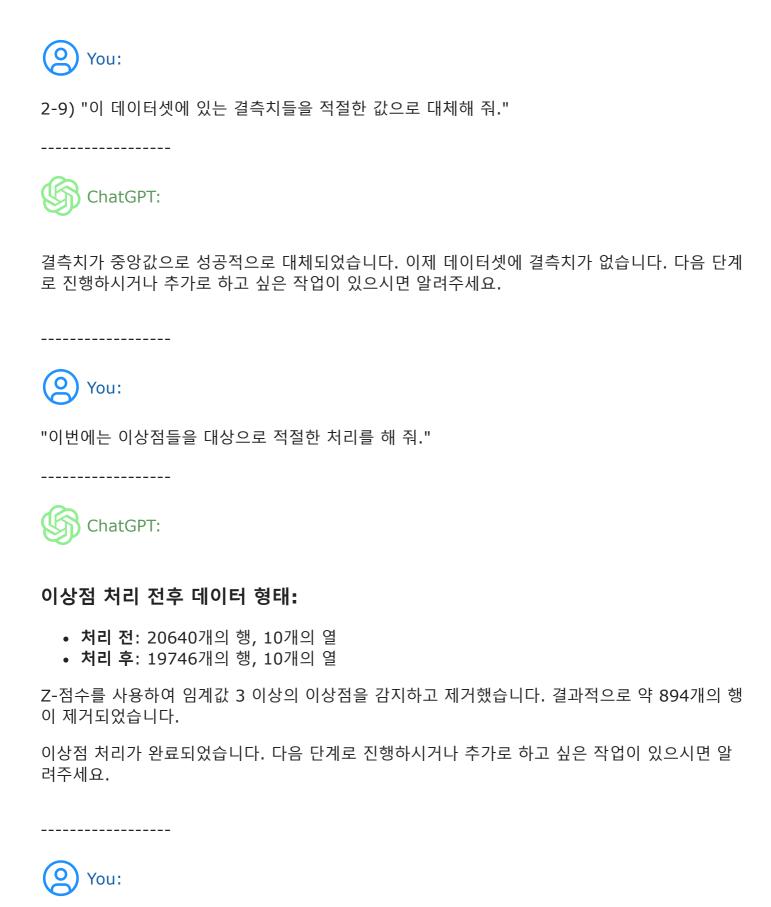
5. 특성 엔지니어링

- **있을 수 있는 부분:** total_rooms, total_bedrooms, population, households 등은 가구당 평균으로 변환될 수 있습니다.
- 전처리 방법: 새로운 특성을 생성하여 모델의 성능을 향상시킵니다.

6. 지리적 데이터 활용

- 있을 수 있는 부분: longitude와 latitude를 사용하여 새로운 특성을 만들 수 있습니다.
- 전처리 방법: 지리적 클러스터링, 거리 계산 등을 통해 새로운 특성을 생성합니다.

이러한 전처리 단계를 거치면 모델의 성능을 향상시킬 수 있고, 데이터의 해석력을 높일 수 있습니다. 전처리를 진행하시겠습니까?



2-10) "이번에는 적절한 정규화와 표준화를 진행해 줘."

ChatGPT:

정규화와 표준화가 완료되었습니다. 다음은 처리된 데이터의 첫 5개 행입니다:

- 수치형 변수 (Iongitude, latitude, housing_median_age, 등)는 표준화(Standardization)가 적용되었습니다.
- 범주형 변수 (ocean_proximity)는 원-핫 인코딩을 통해 변환되었습니다.

이로써 모든 전처리 단계가 완료되었습니다. 다음 단계로 진행하시거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



2-11) "이번에는 회귀 모형을 개발하기에 알맞게 범주형 데이터를 적절한 방법으로 인코딩을 해줘."



범주형 변수 ocean_proximity는 원-핫 인코딩을 통해 다음과 같이 변환되었습니다:

- ocean_proximity_<1H OCEAN
- ocean_proximity_INLAND
- ocean_proximity_ISLAND
- ocean_proximity_NEAR BAY
- ocean_proximity_NEAR OCEAN

이러한 변환은 회귀 모델에 적용하기 적합합니다, 왜냐하면 원-핫 인코딩은 변수 간의 순서나 크기를 부여하지 않기 때문입니다.

이제 데이터셋은 회귀 모형 개발에 적합한 상태로 전처리되었습니다. 다음 단계로 진행하시거나 추가로 하고 싶은 작업이 있으시면 알려주세요.



2-12) "전처리가 완료된 데이터셋을 housing_preprocessed.csv라는 이름으로 저장해 주고 내려받을 수 있게 링크를 제시해 줘."



전처리가 완료된 데이터셋이 housing_preprocessed.csv 파일로 저장되었습니다. 다음 링크를 통해 데이터셋을 내려받을 수 있습니다: housing_preprocessed.csv.

다른 작업이 필요하시면 언제든지 알려주세요!



2-13) "이 데이터셋을 대상으로 특징 공학 기법을 동원한 특징량 추가, 데이터 척도 재구성(data scaling), VIF 방법을 동원한 다중 공선성 해소, 주성분 분석을 차례로 진행해 줘."



전처리 및 특징 공학 결과

- 1. 특징 공학: 'rooms_per_household'와 'population_per_household'라는 두 개의 새로운 특징이 추가되었습니다.
- 2. 데이터 척도 재구성: 새로운 특징들은 표준화되었습니다.
- 3. **다중 공선성 해소**: VIF (Variance Inflation Factor)를 사용하여 다음의 특징들이 고 VIF 값을 가지고 있어 제거되었습니다.
- 'longitude', 'latitude', 'total_rooms', 'total_bedrooms', 'households', 'ocean_proximity_<1H OCEAN', 'ocean_proximity_INLAND', 'ocean_proximity_NEAR BAY', 'ocean_proximity_NEAR OCEAN'
- 4. 주성분 분석 (PCA): 95%의 분산을 유지하도록 차원을 축소하였습니다.

주성분 분석 결과

• 설명된 분산 비율: 약 99.995%

이러한 전처리와 특징 공학을 통해 데이터셋이 더욱 효율적으로 사용될 수 있습니다. 다음 단계로 진행하시거나 추가로 하고 싶은 작업이 있으시면 알려주세요.