5장 회귀

▼ 01 회귀 소개

회귀는 여러개의 독립변수와 한 개의 종속변수 간의 상관관계를 모델링하는 기법 독립변수는 피처에 해당하며 종속변수는 결정 값 머신 러닝 회귀 예측 핵심은 주어진 피처와 결정 값 데이터 기반에서 학습을 통해 최적의 회귀 계 수를 찾아내는 것

지도학습은 분류와 회귀로 나뉨 분류의 예측값은 이산형 클래스 값이고 회귀는 연속형 숫자 값이라는 것이 가장 큰 차이

선형회귀는 실제 값과 예측값의 차이를 최소화하는 직선형 회귀선을 최적화하는 방식

<대표적인 선형 회귀 모델>

- 일반 선형 회귀: 예측값과 실제 값의 RSS를 최소화 할 수 있도록 회귀 계수를 최적화하며, 규 제를 적용하지 않은 모델
- 릿지: 선형 회귀에 L2 규제를 추가한 회귀 모델(L2는 상대적으로 큰 회귀 계수 값의 예측 영향도를 감소시키기 위해 회귀 계수값을 더 작게 만드는 규제 모델)
- 라쏘: 선형 회귀에 L1 규제를 적용한 방식(L1은 예측 영향력이 작은 피처의 회귀 계수를 0으로 만들어 회귀 예측 시 피처가 선택되지 않게 하는 것)
- 엘라스틱넷: L1,L2 규제를 함께 결합한 모델로 주로 피처가 많은 데이터 세트에서 적용되며 L1 규제로 피처의 개수를 줄임과 동시에 L2 규제로 계수 값의 크기를 조정
- 로지스틱 회귀: 분류에 사용되는 선형 모델

▼ 02 단순 선형 회귀를 통한 회귀 이해

단순 선형 회귀는 독립변수도 하나, 종속변수도 하나인 선형 회귀 실제 값과 회귀 모델의 차이에 따른 오류 값을 남은 오류, 즉 잔차라고 부름 최적의 회귀모델을 만든다는 것은 바로 전체 데이터의 잔차 합이 최소가 되는 모델을 만든다는 의 미

보통 오류 합을 계산 할 때는 절댓값을 취해서 더하거나(Mean Absolute Error) 오류 값의 제곱을 구해서 더하는 방식(RSS, Residual Sum of Square)을 취하고 일반적으로 RSS 방식을 씀

회귀에서 RSS는 비용이며 w변수로 구성되는 RSS를 비용함수 라고 함 머신 러닝 회귀 알고리즘은 데이터를 계속 학습하면서 이 비용함수가 반환하는 값을 지속해서 감 소시키고 최종적으로는 더 이상 감소하지 않는 최소의 오류 값을 구하는 것 비용 함수를 손실 함수라고도 함

▼ 03 비용 최소화하기 - 경사 하강법

경사하강법은 반복적으로 비용 함수의 반환 값, 즉 예측값과 실제 값의 차이가 작아지는 방향성을 가지고 w파라미터를 지속해서 보정해 나감

비용함수가 포물선 형태의 2차 함수라면 경사 하강법은 최초 w에서부터 미분을 적용한 뒤 이 미분 값이 계속 감소하는 방향으로 순차적으로 w를 업데이트함

더 이상 미분된 1차 함수의 기울기가 감소하지 않는 지점을 비용 함수가 최소인 지점으로 간주하고 그때의 w를 반환

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4
5 np.random.seed(0)
6 X = 2 * np.random.rand(100,1)
7 y = 6 +4 * X+ np.random.randn(100,1)
8
9 plt.scatter(X, y)
```

```
1 #비용함수
2 def get_cost(y, y_pred):
3 N=len(y)
4 cost=np.sum(np.square(y-y_pred))/N
5 return cost
```

```
1 #경사하강법
2 def get_weight_updates(w1, w0, X, y, learning_rate=0.01):
       N = Ien(y)
4
       w1_update = np.zeros_like(w1)
5
       w0_update = np.zeros_like(w0)
 6
       y_pred = np.dot(X, w1.T) + w0
7
       diff = y-y\_pred
8
9
       w0_factors = np.ones((N,1))
10
11
       w1\_update = -(2/N)*learning\_rate*(np.dot(X.T, diff))
12
       w0\_update = -(2/N)*learning\_rate*(np.dot(w0\_factors.T, diff))
13
14
       return w1_update, w0_update
 1 def gradient_descent_steps(X, y, iters=10000):
       w0 = np.zeros((1,1))
3
       w1 = np.zeros((1,1))
4
5
       for ind in range(iters):
6
           w1_update, w0_update = get_weight_updates(w1, w0, X, y, learning_rate=0.01)
7
           w1 = w1 - w1_update
8
           w0 = w0 - w0\_update
9
10
       return w1, w0
 1 w1, w0 = gradient_descent_steps(X, y, iters=1000)
2 \text{ print}(\text{"w1:}\{0:.3f\} \text{ w0:}\{1:.3f\}\text{".format}(\text{w1[0,0], w0[0,0]}))
3 \text{ y\_pred} = \text{w1}[0,0] * X + \text{w0}
4 print('Gradient Descent Total Cost:{0:.4f}'.format(get_cost(y, y_pred)))
     w1:4.022 w0:6.162
     Gradient Descent Total Cost:0.9935
 1 plt.scatter(X, y)
2 plt.plot(X,y_pred)
```

실전에서는 대부분 확률적 경사 하강법 이용

확률적 경사 하강법: 전체 입력 데이터로 w가 업데이트되는 값을 계산하는 것이 아니라 일부 데이터만 이용해 w가 업데이트 되는 값을 계산하므로 빠른 속도 보장

대용량 데이터경우 대부분 확률적 경사 하강법이나 미니 배치 확률적 경사 하강법 이용해 최적 비용함수 도출

```
1 #확률적 경사 하강법
2 def stochastic_gradient_descent_steps(X, y, batch_size=10, iters=1000):
      w0 = np.zeros((1,1))
4
      w1 = np.zeros((1,1))
5
       prev_cost = 100000
 6
      iter_index =0
7
8
       for ind in range(iters):
9
           np.random.seed(ind)
           stochastic_random_index = np.random.permutation(X.shape[0])
10
11
           sample_X = X[stochastic_random_index[0:batch_size]]
           sample_y = y[stochastic_random_index[0:batch_size]]
12
13
           w1_update, w0_update = get_weight_updates(w1, w0, sample_X, sample_y, learning_rate=0.0
14
           w1 = w1 - w1\_update
           w0 = w0 - w0 \text{ update}
15
16
17
       return w1, w0
1 w1, w0 = stochastic_gradient_descent_steps(X, y, iters=1000)
2 print("w1:",round(w1[0,0],3),"w0:",round(w0[0,0],3))
3 \text{ y\_pred} = \text{w1}[0,0] * X + \text{w0}
4 print('Stochastic Gradient Descent Total Cost:{0:.4f}'.format(get_cost(y, y_pred)))
     w1: 4.028 w0: 6.156
     Stochastic Gradient Descent Total Cost:0.9937
```

04 사이킷런 LinearRegression을 이용한 보스턴 주택 가격 예측

▼ LinearRegression 클래스 - Ordinary Least Squares

class sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False,copy_X=True,n_jobs=1)

- fit_intercept: 절편을 계산할 것인지 말지 지정(디폴트 true)
- normalize: 회귀 수행전 입력 데이터 세트 정규화(디폴트 false)

- coef_: fit() 메서드 수행했을 때 회귀 계수가 배열 형태로 저장하는 속성
- intercept: intercept 값

ordinary least squares기반의 회귀 계수 계산은 입력 피처의 독립성에 많은 영향 받음 피처 간 상관관계가 매우 높은 경우 분산이 매우 커져 오류에 민감해짐 이러한 현상을 다중공선성 문제라고함 일반적으로 상관관계가 높은 피처가 많은 경우 독립적인 중요한 피처만 남기고 제거하거나 규제 를 적용하거나 PCA통해 차원 축소 수행

<회귀 평가 지표>

- MAE:mean absolute erroer. 실제 값과 예측 값 차이를 절댓값으로 변환해 평균
- MSE:mean squared error, 실제 값과 예측값의 차이를 제곱해 평균
- RMSE: MSE에 루트 씌운 것
- R2: 분산 기반, 실제 값의 분산 대비 예측 값의 분산 비율,1에 가까울 수록 정확도 높음

▼ LinearRegression을 이용해 보스턴 주택 가격 회귀 구현

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 from scipy import stats
6 from sklearn.datasets import load_boston
7 %matplotlib inline
8
9 boston = load_boston()
10
11 bostonDF = pd.DataFrame(boston.data , columns = boston.feature_names)
12
13 bostonDF['PRICE'] = boston.target
14 print('Boston 데이타셋 크기:',bostonDF.shape)
15 bostonDF.head()
```

regplot(): x,y축 값의 산점도와 함께 선형 회귀 직선을 그려줌 matplotlib.subplots():각 ax마다 칼럼과 price 관계 표현

```
fig, axs = plt.subplots(figsize=(16,8) , ncols=4 , nrows=2)
Im_features = ['RM','ZN','INDUS','NOX','AGE','PTRATIO','LSTAT','RAD']
for i , feature in enumerate(Im_features):
    row = int(i/4)
    col = i%4
    sns.regplot(x=feature , y='PRICE',data=bostonDF , ax=axs[row][col])
```

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error , r2_score
5 y_target = bostonDF['PRICE']
6 X_data = bostonDF.drop(['PRICE'],axis=1,inplace=False)
8 X_train , X_test , y_train , y_test = train_test_split(X_data , y_target ,test_size=0.3, randor
10 # Linear Regression OLS로 학습/예측/평가 수행.
11 Ir = LinearRegression()
12 Ir.fit(X_train ,y_train )
13 y_preds = Ir.predict(X_test)
14 mse = mean_squared_error(y_test, y_preds)
15 rmse = np.sqrt(mse)
16
17 print('MSE : {0:.3f} , RMSE : {1:.3F}'.format(mse , rmse))
18 print('Variance score : {0:.3f}'.format(r2_score(y_test, y_preds)))
     MSE : 17.297 , RMSE : 4.159
     Variance score: 0.757
 1 print('절편 값:', Ir. intercept_)
2 print('회귀 계수값:', np.round(Ir.coef_, 1))
     절편 값: 40.995595172164755
     회귀 계수값: [-0.1 0.1 0. 3. -19.8 3.4 0. -1.7 0.4 -0. -0.9
      -0.61
 1 coeff = pd.Series(data=np.round(lr.coef_, 1), index=X_data.columns )
2 coeff.sort_values(ascending=False)
     RM
                3.4
     CHAS
                3.0
     RAD
                0.4
     ΖN
                0.1
     INDUS
                0.0
     AGE
                0.0
```

```
-0.0
     TAX
                0.0
                -0.1
     CRIM
                -0.6
     LSTAT
     PTRAT I O
                -0.9
     DIS
                -1.7
     NOX
               -19.8
     dtype: float64
1 from sklearn.model_selection import cross_val_score
3 y_target = bostonDF['PRICE']
4 X_data = bostonDF.drop(['PRICE'],axis=1,inplace=False)
5 Ir = LinearRegression()
6
7 neg_mse_scores = cross_val_score(lr, X_data, y_target, scoring="neg_mean_squared_error", cv = {
8 \text{ rmse\_scores} = \text{np.sqrt}(-1 * \text{neg\_mse\_scores})
9 avg_rmse = np.mean(rmse_scores)
11 print(' 5 folds 의 개별 Negative MSE scores: ', np.round(neg_mse_scores, 2))
12 print(' 5 folds 의 개별 RMSE scores : ', np.round(rmse_scores, 2))
13 print(' 5 folds 의 평균 RMSE : {0:.3f} '.format(avg_rmse))
      5 folds 의 개별 Negative MSE scores: [-12.46 -26.05 -33.07 -80.76 -33.31]
      5 folds 의 개별 RMSE scores : [3.53 5.1 5.75 8.99 5.77]
      5 folds 의 평균 RMSE : 5.829
```

▼ 05 다항 회귀와 과적합/과소적합 이해

▼ 다항 회귀 이해

회귀가 독립변수의 단항식이 아닌 2차, 3차 방정식과 같은 다항식으로 표현되는 것을 다항 회귀라고 함

다항 회귀는 선형 회귀임

```
1 from sklearn.preprocessing import PolynomialFeatures 2 import numpy as np 3 
4 X = np.arange(4).reshape(2,2) 
5 print('일차 단항식 계수 feature:\n',X) 6 
7 poly = PolynomialFeatures(degree=2) 
8 poly.fit(X) 
9 poly_ftr = poly.transform(X) 
10 print('변환된 2차 다항식 계수 feature:\n', poly_ftr) 
일차 단항식 계수 feature:
[[0 1] [2 3]]
```

```
변환된 2차 다항식 계수 feature:
     [[1. 0. 1. 0. 0. 1.]
      [1. 2. 3. 4. 6. 9.]]
 1 def polynomial_func(X):
      y = 1 + 2*X[:,0] + 3*X[:,0]**2 + 4*X[:,1]**3
3
      print(X[:, 0])
4
      print(X[:, 1])
5
      return y
6
7 X = \text{np.arange}(0,4).\text{reshape}(2,2)
9 print('일차 단항식 계수 feature: ₩n',X)
10 y = polynomial_func(X)
11 print('삼차 다항식 결정값: ₩n', y)
     일차 단항식 계수 feature:
     [[0 1]
     [2 3]]
     [0\ 2]
     [1 3]
     삼차 다항식 결정값:
      [ 5 125]
 1 poly_ftr = PolynomialFeatures(degree=3).fit_transform(X)
2 print('3차 다항식 계수 feature: ₩n',poly_ftr)
3
4 model = LinearRegression()
5 model.fit(poly_ftr,y)
6 print('Polynomial 회귀 계수₩n', np.round(model.coef_, 2))
7 print('Polynomial 회귀 Shape :', model.coef_.shape)
     3차 다항식 계수 feature:
      [[ 1. 0. 1. 0. 0. 1. 0. 0. 0. 1.]
      [ 1. 2. 3. 4. 6. 9. 8. 12. 18. 27.]]
     Polynomial 회귀 계수
      [0. 0.18 0.18 0.36 0.54 0.72 0.72 1.08 1.62 2.34]
     Polynomial 회귀 Shape : (10,)
1 from sklearn.preprocessing import PolynomialFeatures
2 from sklearn.linear_model import LinearRegression
3 from sklearn.pipeline import Pipeline
4 import numpy as np
5
6 def polynomial_func(X):
7
      y = 1 + 2*X[:,0] + 3*X[:,0]**2 + 4*X[:,1]**3
8
      return y
10 model = Pipeline([('poly', PolynomialFeatures(degree=3)),
11
                    ('linear', LinearRegression())])
12 X = np.arange(4).reshape(2,2)
13 y = polynomial_func(X)
14
15 model = model.fit(X, y)
```

```
16 print('Polynomial 회귀 계수₩n', np.round(model.named_steps['linear'].coef_, 2))
17
Polynomial 회귀 계수
  [0. 0.18 0.18 0.36 0.54 0.72 0.72 1.08 1.62 2.34]
```

▼ 다항회귀를 이용한 과소적합 및 과적합 이해

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.pipeline import Pipeline
4 from sklearn.preprocessing import PolynomialFeatures
5 from sklearn.linear_model import LinearRegression
6 from sklearn.model_selection import cross_val_score
7 %matplotlib inline
8
9 def true_fun(X):
10
       return np.cos(1.5 * np.pi * X)
11
12 np.random.seed(0)
13 \text{ n\_samples} = 30
14 X = np.sort(np.random.rand(n_samples))
16 y = true_fun(X) + np.random.randn(n_samples) * 0.1
17
 1 plt.figure(figsize=(14, 5))
2 \text{ degrees} = [1, 4, 15]
4 for i in range(len(degrees)):
5
      ax = plt.subplot(1, len(degrees), i + 1)
      plt.setp(ax, xticks=(), yticks=())
6
7
8
       polynomial_features = PolynomialFeatures(degree=degrees[i], include_bias=False)
9
       linear_regression = LinearRegression()
       pipeline = Pipeline([("polynomial_features", polynomial_features),
10
11
                            ("linear_regression", linear_regression)])
12
       pipeline.fit(X.reshape(-1, 1), y)
13
       scores = cross_val_score(pipeline, X.reshape(-1,1), y,scoring="neg_mean_squared_error", cv=
14
       coefficients = pipeline.named_steps['linear_regression'].coef_
15
       print('₩nDegree {0} 회귀 계수는 {1} 입니다.'.format(degrees[i], np.round(coefficients),2))
16
17
       print('Degree {0} MSE 는 {1:.2f} 입니다.'.format(degrees[i] , -1*np.mean(scores)))
18
19
      X_{\text{test}} = \text{np.linspace}(0, 1, 100)
20
21
       plt.plot(X_test, pipeline.predict(X_test[:, np.newaxis]), label="Model")
22
       plt.plot(X_test, true_fun(X_test), '--', label="True function")
23
24
       plt.scatter(X, y, edgecolor='b', s=20, label="Samples")
25
26
       plt.xlabel("x"); plt.ylabel("y"); plt.xlim((0, 1)); plt.ylim((-2, 2)); plt.legend(loc="bes
```

```
plt.title("Degree {}\WnMSE = {:.2e}(+/- {:.2e})".format(degrees[i], -scores.mean(), scores.general 28 29 plt.show()
```

첫번째는 과소적합모델, 고편향성 가짐 세번째는 과적합이 심한 모델, 고분산성 가짐

▼ 편향 분산 트레이드오프

일반적으로 편향과 분산은 한 쪽이 높으면 한 쪽이 낮아지는 경향 높은 편향/낮은 분산에서 과소적합되기 쉬우며 낮은 편향/높은 분산에서 과적합되기 쉬움

- ▼ 06 규제 선형 모델 릿지, 라쏘, 엘라스틱넷
- ▼ 규제 선형 모델의 개요

비용함수는 RSS 최소화 방법과 과적합을 방지하기 위해 회귀 계수 값이 커지지 않도록 하는 방법 이 서로 균형을 이뤄야함

비용함수 목표 =Min(RSS(W)+alpha*||W||2) alpha는 학습 데이터 적합 정도와 회귀 계수 값의 크기 제어를 수행하는 튜닝 파라미터 alpha=0인 경우 W가 커도 비용함수는 Min(RSS(W)) alpha=무한대인경우 비용함수는 W를 0에 가깝게 최소화 해야함

alpha를 0에서 부터 지속적으로 값을 증가시키면 회귀 계수 값의 크기를 감소시킬 수 있음 규제: 비용함수에 alpha값으로 페널티를 부여해 회귀 계수 값의 크기를 감소시켜 과적합을 개선하 는 방식

L2규제: W의 제곱에 페널티를 부여하는 방식 -> 릿지 회귀 L1규제: W의 절댓값에 대해 페널티 부여 -> 라쏘 회귀

▼ 릿지 회귀

```
1 from sklearn.linear_model import Ridge
2 from sklearn.model_selection import cross_val_score
3
4 ridge = Ridge(alpha = 10)
5 neg_mse_scores = cross_val_score(ridge, X_data, y_target, scoring="neg_mean_squared_error", cv
6 rmse_scores = np.sqrt(-1 * neg_mse_scores)
7 avg_rmse = np.mean(rmse_scores)
8 print(' 5 folds 의 개별 Negative MSE scores: ', np.round(neg_mse_scores, 3))
9 print(' 5 folds 의 개별 RMSE scores : ', np.round(rmse_scores,3))
10 print(' 5 folds 의 평균 RMSE : {0:.3f} '.format(avg_rmse))
     5 folds 의 개별 Negative MSE scores: [-11.422 -24.294 -28.144 -74.599 -28.517]
     5 folds 의 개별 RMSE scores : [3.38 4.929 5.305 8.637 5.34 ]
     5 folds 의 평균 RMSE : 5.518
1 \text{ alphas} = [0, 0.1, 1, 10, 100]
3 for alpha in alphas:
     ridge = Ridge(alpha = alpha)
    neg_mse_scores = cross_val_score(ridge, X_data, y_target, scoring="neg_mean_squared_error"
7
      avg_rmse = np.mean(np.sqrt(-1 * neg_mse_scores))
      print('alpha {0} 일 때 5 folds 의 평균 RMSE : {1:.3f} '.format(alpha,avg_rmse))
    alpha 0 일 때 5 folds 의 평균 RMSE : 5.829
    alpha 0.1 일 때 5 folds 의 평균 RMSE : 5.788
    alpha 1 일 때 5 folds 의 평균 RMSE : 5.653
    alpha 10 일 때 5 folds 의 평균 RMSE : 5.518
    alpha 100 일 때 5 folds 의 평균 RMSE : 5.330
```

```
1 fig , axs = plt.subplots(figsize=(18,6) , nrows=1 , ncols=5)
3 coeff_df = pd.DataFrame()
5 for pos , alpha in enumerate(alphas) :
       ridge = Ridge(alpha = alpha)
6
7
       ridge.fit(X_data , y_target)
8
      coeff = pd.Series(data=ridge.coef_ , index=X_data.columns )
9
10
      colname='alpha:'+str(alpha)
      coeff_df[colname] = coeff
11
12
13
      coeff = coeff.sort_values(ascending=False)
14
      axs[pos].set_title(colname)
      axs[pos].set_xlim(-3,6)
15
       sns.barplot(x=coeff.values , y=coeff.index, ax=axs[pos])
16
17
18 plt.show()
19
```

```
1 ridge_alphas = [0 , 0.1 , 1 , 10 , 100]
2 sort_column = 'alpha:'+str(ridge_alphas[0])
3 coeff_df.sort_values(by=sort_column, ascending=False)
```

▼ 라쏘 회귀

```
1 from sklearn.linear_model import Lasso, ElasticNet
3 def get_linear_reg_eval(model_name, params=None, X_data_n=None, y_target_n=None,
                         verbose=True, return_coeff=True):
4
5
      coeff_df = pd.DataFrame()
      if verbose : print('####### ', model_name , '#######')
6
7
      for param in params:
          if model_name == 'Ridge': model = Ridge(alpha=param)
8
          elif model_name == 'Lasso': model = Lasso(alpha=param)
9
10
          elif model_name == 'ElasticNet': model = ElasticNet(alpha=param, l1_ratio=0.7)
11
          neg_mse_scores = cross_val_score(model, X_data_n,
12
                                             y_target_n, scoring="neg_mean_squared_error", cv =
          avg_rmse = np.mean(np.sqrt(-1 * neg_mse_scores))
13
          print('alpha {0}일 때 5 폴드 세트의 평균 RMSE: {1:.3f} '.format(param, avg_rmse))
14
          # cross_val_score는 evaluation metric만 반환하므로 모델을 다시 학습하여 회귀 계수 추출
15
16
17
          model.fit(X_data_n , y_target_n)
18
          if return_coeff:
              # alpha에 따른 피처별 회귀 계수를 Series로 변환하고 이를 DataFrame의 컬럼으로 추가
19
              coeff = pd.Series(data=model.coef_ , index=X_data_n.columns )
20
21
              colname='alpha:'+str(param)
22
              coeff_df[colname] = coeff
23
24
      return coeff_df
25
 1 lasso_alphas = [0.07, 0.1, 0.5, 1, 3]
2 coeff_lasso_df =get_linear_reg_eval('Lasso', params=lasso_alphas, X_data_n=X_data, y_target_n=y
     ###### Lasso ######
     alpha 0.07일 때 5 폴드 세트의 평균 RMSE: 5.612
     alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 5.615
     alpha 0.5일 때 5 폴드 세트의 평균 RMSE: 5.669
     alpha 1일 때 5 폴드 세트의 평균 RMSE: 5.776
     alpha 3일 때 5 폴드 세트의 평균 RMSE: 6.189
```

```
1 sort_column = 'alpha:'+str(lasso_alphas[0])
2 coeff_lasso_df.sort_values(by=sort_column, ascending=False)
```

▼ 엘라스틱넷 회귀

▼ 선형 회귀 모델을 위한 데이터 변환

1.standardscaler 클래스를 이용해 평균이 0, 분산이 1인 표준 정규 분포를 가진 데이터 세트로 변환하거나 minmaxscaler클래스를 이용해 최솟값이 0이고 최댓값이 1인 값으로 정규화 수행 2.스케일링/정규화 수행한 데이터 세트에 다시 다항 특성을 적용하여 변환하는 방법 3.로그변환

타깃 값의 경우는 일반적으로 로그 변환 적용

```
1 from sklearn.preprocessing import StandardScaler, MinMaxScaler, PolynomialFeatures
3 def get_scaled_data(method='None', p_degree=None, input_data=None):
      if method == 'Standard':
5
           scaled_data = StandardScaler().fit_transform(input_data)
      elif method == 'MinMax':
6
7
          scaled_data = MinMaxScaler().fit_transform(input_data)
      elif method == 'Log':
8
9
          scaled_data = np.log1p(input_data)
10
      else:
          scaled_data = input_data
11
12
      if p_degree != None:
13
14
           scaled_data = PolynomialFeatures(degree=p_degree,
                                             include_bias=False).fit_transform(scaled_data)
15
16
17
      return scaled_data
1 \text{ alphas} = [0.1, 1, 10, 100]
3 scale_methods=[(None, None), ('Standard', None), ('Standard', 2),
                  ('MinMax', None), ('MinMax', 2), ('Log', None)]
5 for scale_method in scale_methods:
      X_data_scaled = get_scaled_data(method=scale_method[0], p_degree=scale_method[1],
6
7
                                       input_data=X_data)
```

- 8 print(X_data_scaled.shape, X_data.shape)
- 9 print('\m## 변환 유형:{0}, Polynomial Degree:{1}'.format(scale_method[0], scale_method[1])
- 10 get_linear_reg_eval('Ridge', params=alphas, X_data_n=X_data_scaled,
- 11 y_target_n=y_target, verbose=False, return_coeff=False)

(506, 13) (506, 13)

변환 유형:None, Polynomial Degree:None alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 5.788 alpha 1일 때 5 폴드 세트의 평균 RMSE: 5.653 alpha 10일 때 5 폴드 세트의 평균 RMSE: 5.518 alpha 100일 때 5 폴드 세트의 평균 RMSE: 5.330 (506, 13) (506, 13)

변환 유형:Standard, Polynomial Degree:None alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 5.826 alpha 1일 때 5 폴드 세트의 평균 RMSE: 5.803 alpha 10일 때 5 폴드 세트의 평균 RMSE: 5.637 alpha 100일 때 5 폴드 세트의 평균 RMSE: 5.421 (506, 104) (506, 13)

변환 유형:Standard, Polynomial Degree:2 alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 8.827 alpha 1일 때 5 폴드 세트의 평균 RMSE: 6.871 alpha 10일 때 5 폴드 세트의 평균 RMSE: 5.485 alpha 100일 때 5 폴드 세트의 평균 RMSE: 4.634 (506, 13) (506, 13)

변환 유형:MinMax, Polynomial Degree:None alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 5.764 alpha 1일 때 5 폴드 세트의 평균 RMSE: 5.465 alpha 10일 때 5 폴드 세트의 평균 RMSE: 5.754 alpha 100일 때 5 폴드 세트의 평균 RMSE: 7.635 (506, 104) (506, 13)

변환 유형:MinMax, Polynomial Degree:2 alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 5.298 alpha 1일 때 5 폴드 세트의 평균 RMSE: 4.323 alpha 10일 때 5 폴드 세트의 평균 RMSE: 5.185 alpha 100일 때 5 폴드 세트의 평균 RMSE: 6.538 (506, 13) (506, 13)

변환 유형:Log, Polynomial Degree:None alpha 0.1일 때 5 폴드 세트의 평균 RMSE: 4.770 alpha 1일 때 5 폴드 세트의 평균 RMSE: 4.676 alpha 10일 때 5 폴드 세트의 평균 RMSE: 4.836 alpha 100일 때 5 폴드 세트의 평균 RMSE: 6.241

▼ 07 로지스틱 회귀

시그모이드 함수 y=1/(1+e^(-x))

- 1 import pandas as pd
- 2 import matplotlib.pyplot as plt
- 3 %matplotlib inline

```
22. 4. 4. 오후 8:01
                                               Week5_예습과제.ipynb - Colaboratory
     4
    5 from sklearn.datasets import load_breast_cancer
    6 from sklearn.linear_model import LogisticRegression
    7
    8 cancer = load_breast_cancer()
     1 from sklearn.preprocessing import StandardScaler
    2 from sklearn.model_selection import train_test_split
    4 scaler = StandardScaler()
    5 data_scaled = scaler.fit_transform(cancer.data)
    7 X_train , X_test, y_train , y_test = train_test_split(data_scaled, cancer.target, test_size=0.0
     1 from sklearn.metrics import accuracy_score, roc_auc_score
    3 Ir_clf = LogisticRegression()
    4 Ir_clf.fit(X_train, y_train)
    5 Ir_preds = Ir_clf.predict(X_test)
    7 print('accuracy: {:0.3f}'.format(accuracy_score(y_test, Ir_preds)))
    8 print('roc_auc: {:0.3f}'.format(roc_auc_score(y_test , Ir_preds)))
         accuracy: 0.977
         roc_auc: 0.972
     1 from sklearn.model_selection import GridSearchCV
    3 params={'penalty':['12', '11'],
    4
              'C':[0.01, 0.1, 1, 1, 5, 10]}
    5
    6 grid_clf = GridSearchCV(Ir_clf, param_grid=params, scoring='accuracy', cv=3)
    7 grid_clf.fit(data_scaled, cancer.target)
    8 print('최적 하이퍼 파라미터:{0}, 최적 평균 정확도:{1:.3f}'.format(grid_clf.best_params_,
    9
                                                        grid_clf.best_score_))
         최적 하이퍼 파라미터:{'C': 1, 'penalty': 'I2'}, 최적 평균 정확도:0.975
         /usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_validation.py:372: FitFailedM
         18 fits failed out of a total of 36.
         The score on these train-test partitions for these parameters will be set to nan.
         If these failures are not expected, you can try to debug them by setting error_score='raise'.
         Below are more details about the failures:
         18 fits failed with the following error:
         Traceback (most recent call last):
           File "/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_validation.py", line
             estimator.fit(X_train, y_train, **fit_params)
           File "/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py", line 1461,
```

```
18 fits failed with the following error:
Traceback (most recent call last):
    File "/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_validation.py", line
    estimator.fit(X_train, y_train, **fit_params)
    File "/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py", line 1461
    solver = _check_solver(self.solver, self.penalty, self.dual)
    File "/usr/local/lib/python3.7/dist-packages/sklearn/linear_model/_logistic.py", line 449,
    % (solver, penalty)
ValueError: Solver lbfgs supports only '12' or 'none' penalties, got 11 penalty.

warnings.warn(some_fits_failed_message, FitFailedWarning)
```

/usr/local/lib/python3.7/dist-packages/sklearn/model_selection/_search.py:972: UserWarning: C 0.97539218 nan 0.97011974 nan 0.96661097 nan] category=UserWarning,



▼ 08 회귀 트리

회귀를 위한 트리를 생성하고 이를 기반으로 회귀 예측을 하는 것

분류 트리가 특정 클래스 레이블을 결정하는 것과는 달리 회귀 트리는 리프 노드에 속한 데이터 값의 평균값을 구해 회귀 예측 값을 계산

```
1 from sklearn.datasets import load_boston
2 from sklearn.model_selection import cross_val_score
3 from sklearn.ensemble import RandomForestRegressor
4 import pandas as pd
5 import numpy as np
7 boston = load boston()
8 bostonDF = pd.DataFrame(boston.data, columns = boston.feature_names)
10 bostonDF['PRICE'] = boston.target
11 y_target = bostonDF['PRICE']
12 X_data = bostonDF.drop(['PRICE'], axis=1,inplace=False)
13
14 rf = RandomForestRegressor(random_state=0, n_estimators=1000)
15 neg_mse_scores = cross_val_score(rf, X_data, y_target, scoring="neg_mean_squared_error", cv = {
16 rmse_scores = np.sqrt(-1 * neg_mse_scores)
17 avg rmse = np.mean(rmse scores)
18
19 print(' 5 교차 검증의 개별 Negative MSE scores: ', np.round(neg_mse_scores, 2))
20 print(' 5 교차 검증의 개별 RMSE scores : ', np.round(rmse_scores, 2))
21 print(' 5 교차 검증의 평균 RMSE : {0:.3f} '.format(avg_rmse))
22
```

/usr/local/lib/python3.7/dist-packages/sklearn/utils/deprecation.py:87: FutureWarning: Functi

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::

import pandas as pd import numpy as np

```
data_url = "http://lib.stat.cmu.edu/datasets/boston"
             raw_df = pd.read_csv(data_url, sep="\st", skiprows=22, header=None)
             data = np.hstack([raw_df.values[::2, :], raw_df.values[1::2, :2]])
             target = raw_df.values[1::2, 2]
         Alternative datasets include the California housing dataset (i.e.
         :func:`~sklearn.datasets.fetch_california_housing`) and the Ames housing
         dataset. You can load the datasets as follows::
             from sklearn.datasets import fetch_california_housing
             housing = fetch_california_housing()
         for the California housing dataset and::
             from sklearn.datasets import fetch_openml
             housing = fetch_openml(name="house_prices", as_frame=True)
         for the Ames housing dataset.
       warnings.warn(msg, category=FutureWarning)
      5 교차 검증의 개별 Negative MSE scores: [ -7.88 -13.14 -20.57 -46.23 -18.88]
      5 교차 검증의 개별 RMSE scores : [2.81 3.63 4.54 6.8 4.34]
      5 교차 검증의 평균 RMSE : 4.423
 1 def get_model_cv_prediction(model, X_data, y_target):
2
      neg_mse_scores = cross_val_score(model, X_data, y_target, scoring="neg_mean_squared_error"
3
      rmse_scores = np.sqrt(-1 * neg_mse_scores)
 4
      avg_rmse = np.mean(rmse_scores)
      print('##### ',model.__class__.__name__ , ' #####')
5
      print(' 5 교차 검증의 평균 RMSE : {0:.3f} '.format(avg_rmse))
 1 from sklearn.tree import DecisionTreeRegressor
2 from sklearn.ensemble import GradientBoostingRegressor
3 from xgboost import XGBRegressor
4 from lightgbm import LGBMRegressor
6 dt_reg = DecisionTreeRegressor(random_state=0, max_depth=4)
7 rf_reg = RandomForestRegressor(random_state=0, n_estimators=1000)
8 gb_reg = GradientBoostingRegressor(random_state=0, n_estimators=1000)
9 xgb_reg = XGBRegressor(n_estimators=1000)
10 lgb_reg = LGBMRegressor(n_estimators=1000)
11
12
13 models = [dt_reg, rf_reg, gb_reg, xgb_reg, lgb_reg]
14 for model in models:
15
      get_model_cv_prediction(model, X_data, y_target)
     ##### DecisionTreeRegressor #####
      5 교차 검증의 평균 RMSE : 5.978
     ##### RandomForestRegressor #####
      5 교차 검증의 평균 RMSE : 4.423
     ##### GradientBoostingRegressor
                                      #####
      5 교차 검증의 평균 RMSE : 4.269
     [10:58:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
     [10:58:58] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
```

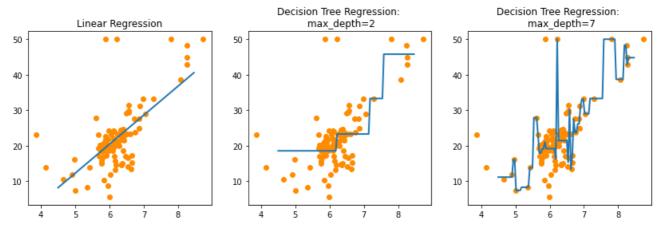
```
[10:58:58] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca [10:58:59] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca [10:58:59] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca ##### XGBRegressor ##### 5 교차 검증의 평균 RMSE : 4.089 ##### LGBMRegressor ##### 5 교차 검증의 평균 RMSE : 4.646
```

```
1 import seaborn as sns
2 %matplotlib inline
3
4 rf_reg = RandomForestRegressor(n_estimators=1000)
5
6 rf_reg.fit(X_data, y_target)
7
8 feature_series = pd.Series(data=rf_reg.feature_importances_, index=X_data.columns)
9 feature_series = feature_series.sort_values(ascending=False)
10 sns.barplot(x= feature_series, y=feature_series.index)
```

```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3
4 bostonDF_sample = bostonDF[['RM','PRICE']]
5 bostonDF_sample = bostonDF_sample.sample(n=100,random_state=0)
6 print(bostonDF_sample.shape)
7 plt.figure()
8 plt.scatter(bostonDF_sample.RM , bostonDF_sample.PRICE,c="darkorange")
```

```
1
     import numpy as np
2
     from sklearn.linear_model import LinearRegression
3
4
     Ir_reg = LinearRegression()
     rf_reg2 = DecisionTreeRegressor(max_depth=2)
5
6
     rf_reg7 = DecisionTreeRegressor(max_depth=7)
7
8
     X_{\text{test}} = \text{np.arange}(4.5, 8.5, 0.04).\text{reshape}(-1, 1)
9
     X_feature = bostonDF_sample['RM'].values.reshape(-1,1)
10
     y_target = bostonDF_sample['PRICE'].values.reshape(-1,1)
11
12
13
     Ir_reg.fit(X_feature, y_target)
14
     rf_reg2.fit(X_feature, y_target)
15
     rf_reg7.fit(X_feature, y_target)
16
     pred_Ir = Ir_reg.predict(X_test)
17
18
     pred_rf2 = rf_reg2.predict(X_test)
19
     pred_rf7 = rf_reg7.predict(X_test)
1
     fig , (ax1, ax2, ax3) = plt.subplots(figsize=(14,4), ncols=3)
2
3
     ax1.set_title('Linear Regression')
4
     ax1.scatter(bostonDF_sample.RM, bostonDF_sample.PRICE, c="darkorange")
     ax1.plot(X_test, pred_Ir, label="linear", linewidth=2)
5
6
7
     ax2.set_title('Decision Tree Regression: \max_depth=2')
     ax2.scatter(bostonDF_sample.RM, bostonDF_sample.PRICE, c="darkorange")
8
     ax2.plot(X_test, pred_rf2, label="max_depth:3", linewidth=2 )
9
10
11
     ax3.set_title('Decision Tree Regression: \max_depth=7')
12
     ax3.scatter(bostonDF_sample.RM, bostonDF_sample.PRICE, c="darkorange")
13
     ax3.plot(X_test, pred_rf7, label="max_depth:7", linewidth=2)
\Box
```

[<matplotlib.lines.Line2D at 0x7f7c95347550>]



✓ 1초 오후 8:00에 완료됨

×