

2021. 6. 17

MACRO팀

안미성, Ph.D.

Economist

misung11.ahn@samsung.com

글로벌투자전략팀

김동영, CFA

Quant Analyst

dy76.kim@samsung.com

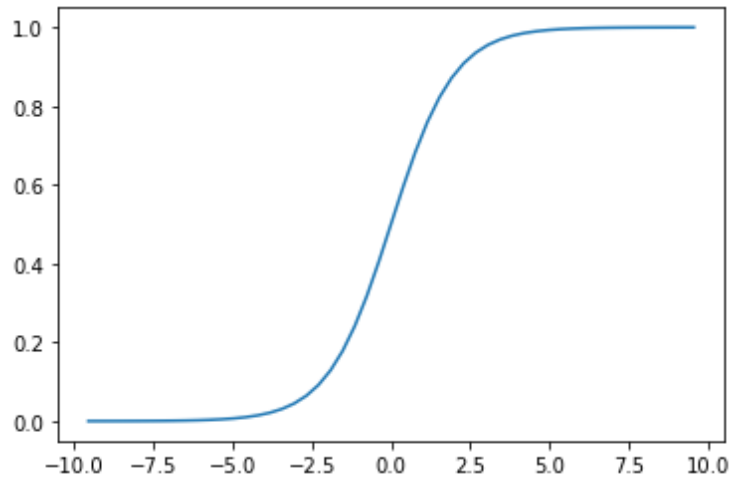
퀀트 모델링 A to Z

(4) 로짓 모델(로지스틱 회귀 모델)

- Logit 모델(Logistic Regression 모델)에 대한 설명 및 사용 방법 수록

이번 편에서는 지난 (1)회귀분석과 (2)HP 필터, (3)Penalized Regression에 이어 Logit 모델(Logistic Regression 모델)을 다룬다. Logit 모델은 특정 이벤트가 발생할 확률을 모델링하고 범주 분류(classification) 작업을 수행하는 기법이다. Logit 모델은 0 또는 1의 값만 취하는 이분적 데이터와 [0,1] 사이인 확률 데이터의 분석과 결과 해석에 용이하다.

이름처럼 Logit 모델은 로지스틱 분포(Logistic distribution, 혹은 Sigmoid function)를 이용하여 모델에서 추정된 이벤트별 확률이 항상 [0,1] 사이에 놓이도록 한다. 분석자는 이 추정된 확률을 바탕으로 임계치에 따라 범주를 분류할 수 있다.



Compliance Note

보고서는 철저히 계량적 분석에 근거한 의견을 제시합니다. 따라서 당사의 대표 투자 의견과 다를 수 있습니다. 본 조사자료는 당사의 저작물로서 모든 저작권은 당사에게 있습니다. 본 조사자료는 당사의 동의없이 어떠한 경우에도 어떠한 형태로든 복제, 배포, 전송, 변경, 대여할 수 없습니다. 본 조사자료에 수록된 내용은 당사 리서치센터가 신뢰할만한 자료 및 정보로부터 얻어진 것이나, 당사는 그 정확성이나 완전성을 보장할 수 없습니다. 따라서 어떠한 경우에도 본 자료는 고객의 주식투자의 결과에 대한 법적 책임 소재에 대한 증빙자료로 사용될 수 없습니다. 본 자료에는 외부의 부당한 압력이나 간섭없이 애널리스트의 의견이 정확하게 반영되었습니다.

I. Logit Model

Logit 모델(Logistic Regression 모델)이 필요한 이유

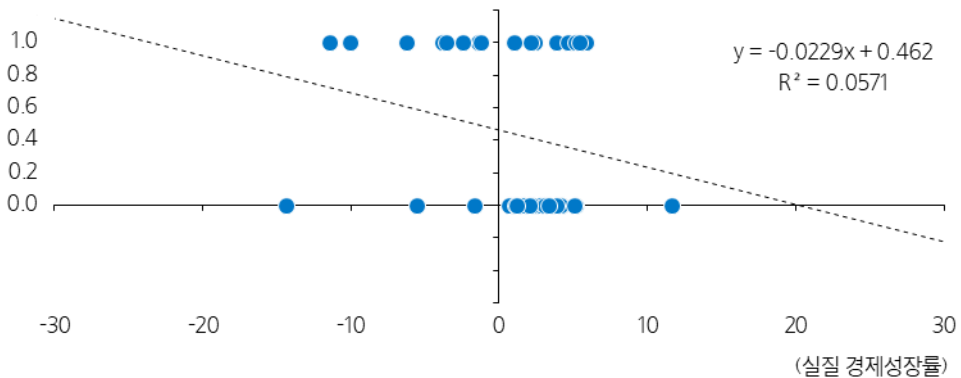
우리는 종종 특정 이벤트가 발생할 확률을 추정하고 싶다. 예컨대, 양당제인 미국 대통령 선거에서 공화당과 민주당 중 공화당이 선거에서 승리할 확률이 얼마나 될지를 알고 싶은 상황이다. 종속변수(y)로 과거 선거에서 공화당이 득표한 비율을, 독립변수(x)로 성장률, 실업률, 물가상승률, 후보가 재선에 도전하는지 등을 생각해볼 수 있다(Fair, 1996). 이 때, 한 가지 방법은 지난 (1)편에서 배운 선형 회귀분석을 사용하는 것이다. 그러나 독립변수에 따라서는 선형 회귀분석의 적합치(fitted value, \hat{y})가 0과 1 범위 밖의 값을 가질 수 있다는 단점이 있다. 이 경우, 자료 해석이 직관적이지 않게 되어 계량분석이 유용하지 않은 상황이 된다.

종속변수가 0 또는 1의 값만 갖는 이분적 데이터(binary data)를 추정하는 일도 마찬가지이다. 가령 어떤 특성을 가진 유권자가 공화당에 투표하느냐를 분석할 수도 있다. 즉, 공화당에 투표하면 1, 공화당에 투표하지 않으면 0의 값을 갖는 이분적 종속변수를 생각할 수 있다. 독립변수로는 소득, 교육수준, 거주 지역, 나이, 성별 등이 적절할 것이다. 이 때, 모델은 어떤 특성(x_i)을 가진 유권자(i)가 공화당에 투표할 확률($y_i = 1$)을 추정하게 된다($P_i = E[y_i|x_i]$). 만약 선형 회귀분석을 사용하면, 위의 예시와 동일하게 적합치가 0과 1 범위 밖의 값을 갖는 경우가 발생할 뿐만 아니라 이번에는 종속변수의 값이 0 또는 1로 제한되어 데이터의 적합도가 크게 떨어질 것이다. 참고로 지난 (1)회귀모형편에서의 예시였던 SK 주가 모델과 보스턴 주택 가격 모델의 종속변수(y)는 주가와 주택가격으로 둘 다 0보다 큰 연속형(continuous) 변수였다.

위와 같은 데이터를 분석할 경우에는 Logit 모델(혹은 Logistic Regression 모델, 이하 Logit 모델)이 유용하다. Logit 모델은 특정 이벤트가 발생할 확률을 모델링하고 추정된 확률에 따라 범주를 분류(classification)하는 작업을 수행하는 기법이다. 위의 예에서는 추정된 확률이 0.5 이상이면 해당 유권자를 공화당 지지자로, 0.5 미만이면 민주당 지지자로 분류하는 작업을 할 수 있다. 한편, Logit 모델은 기술적으로 로지스틱 함수의 누적 확률 분포(Cumulative Distribution Function)를 이용하여 모델의 적합치가 항상 [0,1] 사이에 놓이도록 함으로써 선형 회귀분석의 문제점을 해결한다. 한마디로 Logit 모델은 확률 데이터의 분석과 결과 해석에 용이한 모델이다.

선형 회귀분석 적합시의 문제점: [0,1]을 벗어날 가능성

(공화당 대통령 후보의 다득표 확률)



자료: Fair(1996), 삼성증권

Logit 모형식

Logit 모델은 이름처럼 로지스틱 분포(Logistic Distribution, 혹은 Sigmoid function)인 $F(z) = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}}$ 를 이용한다. 이는 확률값을 계산하기 편리하고 계량적으로 유용한 성질을 지니기 때문이다. Logit 모델의 수학적 모형식은 다음과 같다.

[Logit 모델]

특정 이벤트가 발생할 확률 P_i 를 i 번째 관측치의 특성을 나타내는 독립변수 x_i 간 선형결합, $z_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + u_i$ 를 이용하여 다음과 같이 표기할 수 있음

$$P_i = \frac{1}{1 + e^{-z_i}} = \frac{e^{z_i}}{1 + e^{z_i}}, \text{ 여기서 } z_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + u_i$$

$$\frac{P_i}{1 - P_i} = e^{z_i} = e^{\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + u_i}$$

Logit 모델은 추정에 용이하도록 아래와 같이 선형화한 모델임

$$\log\left(\frac{P_i}{1 - P_i}\right) = \log(e^{z_i}) = z_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + u_i$$

위 세 개 식은 수학적으로 동일함

첫 번째 식에서 알 수 있듯이, 로지스틱 분포를 이용하면 $(-\infty, +\infty)$ 값을 가진 $z_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki}$ 가 추정된 확률(P_i)이 $[0,1]$ 범위 내에 놓이게 된다. 또한, z_i 의 값이 커질수록, 이벤트가 발생할 확률인 P_i 도 커진다.

그러나 첫 번째 식을 통한 β 계수의 추정과 해석은 쉽지 않다. 그래서 두 번째와 같은 변형을 거쳐 세 번째 식을 이용하여 추정한다. 두 번째 식의 좌변 $\frac{P_i}{1-P_i}$ 는 승산비(odds ratio)로 알려져 있다. 이벤트가 발생하지 않을 확률 대비 이벤트가 발생할 확률, 혹은 실패 대비 성공 확률로 해석된다. 공화당이 많은 표를 얻지 못할 확률 대비 공화당이 많은 표를 얻을 확률의 비율이다. 만약 $P_i = 0.8$ 라면, 공화당이 집권할 가능성은 그렇지 않을 가능성의 4배가 된다.

두 번째 식에 로그를 취하면, 우리가 익숙한 선형모델이 된다. 세 번째 식의 좌변 $\log\left(\frac{P_i}{1-P_i}\right)$ 는 승산비에 로그를 취한 값으로, 로그 승산비라고 부른다. z_i 의 값이 커질수록, P_i 는 물론, 승산비 $\frac{P_i}{1-P_i}$ 와 로그 승산비 $\log\left(\frac{P_i}{1-P_i}\right)$ 도 커진다. 물론 승산비는 $[0, +\infty)$, 로그 승산비는 $(-\infty, +\infty)$ 범위의 값을 취한다.

Logit 모델 해석

세 번째 식은 Logit 모델에서 로그 승산비가 독립변수에 대해 선형으로 나타낼 수 있음을 의미한다. 즉, 종속 변수가 확률 데이터인 경우, 로그 승산비 식을 통해 추정하면 된다. 이 식은 선형모델이기 때문에 (1)편 회귀 분석에서와 같은 해석을 적용할 수 있다. 즉, 다른 변수의 값은 일정하고 x_k 만 한 단위 커질 때, 로그 승산비는 β_k 만큼 커진다. 예컨대, 소득이 한 단위(1,000불) 증가할 때 공화당에 투표할 확률에 대한 로그 승산비는 $\beta_{\text{소득}}$ 만큼 증가한다고 해석할 수 있다. 종속변수가 확률 데이터가 아닌 경우도 있다. 이 때에는 파이썬 (Python)이나 다른 통계 패키지에 내장된 모듈을 이용하여 Logit 모델을 추정한다. III장에서 이 방법을 소개한다. 해석은 같은 방식으로 로그 승산비의 변화로 설명한다.

다른 방식으로 로그 승산비에 지수함수를 취해서 해석할 수 있다. 다른 독립변수의 값은 일정하고 x_k 만 한 단위 커질 때, 승산비가 e^{β_k} 만큼 커진다고 해석하는 것이다. 위의 예에 적용하면, 소득이 한 단위(1,000불) 증가할 때, 공화당에 투표할 확률은 그렇지 않을 확률에 대비해서 $e^{\beta_{\text{소득}}}$ 이 커진다.

[Logit 모델의 해석]

$$\exp \left[\log \left(\frac{P_i}{1 - P_i} \right) \right] = \frac{P_i}{1 - P_i} = e^{\beta_0} \cdot e^{\beta_1 x_{1i}} \dots e^{\beta_k x_{ki}}$$

(x_1, x_2, \dots, x_k) 에서 $(x'_1, x'_2, \dots, x_k + 1)$ 로의 승산비 변화는 다음과 같음

$$\text{승산비의 변화} = \frac{P'/(1 - P')}{P/(1 - P)} = e^{\beta_k}$$

Logit 모델을 해석하는 또 다른 방법은 특정 독립변수 $x^o = (x_1^o, x_2^o, \dots, x_k^o)$ 수준에서 공화당에 투표할 확률에 대해 이야기하는 것이다. 예를 들어, 소득이 30,000불인 유권자가 공화당에 투표할 확률이다. 평균 소득(mean income) 혹은 중앙 소득(median income)을 가진 유권자가 공화당에 투표할 확률을 구하면, 이를 경제의 대표적 유권자가 공화당에 투표할 확률이라고 해석할 수도 있다. 다음 II장에서 실제 경제 데이터와 추정치를 통해 Logit 모델을 해석하고 분류작업을 시행한다.

II. 예제

미국의 공화당이 어떤 경제적 상황에서 많은 표(popular vote)를 얻었는지를 Logit 모델로 추정해보자. (참고로, 미국의 투표 제도상, 많은 표를 받았다고 해서 대통령에 당선되는 것은 아니다. 2016년 트럼프 대통령은 적은 표를 받았지만 대통령에 당선됐다.)

종속변수는 더 많은 표를 받았는지의 여부, 즉, 이분적 데이터이다. 미국 공화당이 민주당보다 더 많은 표를 얻었으면 1, 더 적은 표를 얻었으면 0의 값을 취한다. 독립변수로는 앞서 말한 Fair(1996)를 따라서 선거 연도의 첫 3분기 실질 경제성장률, 임기중인 대통령의 첫 15개월 간 물가상승률, 현재 집권당인지 여부 등을 고려할 수 있다. 여기에서는 1880년부터 2016년까지 데이터로 한정하고, 성장률, 물가상승률, 공화당의 현재 집권 여부의 세 가지 독립변수만 고려한다. Logit 모델을 다음과 같이 세울 수 있다.

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 [\text{경제 성장률}] + \beta_2 [\text{물가 상승률}] + \beta_3 [\text{현재 집권 여부}] + u$$

여기서 추정된 P 는 공화당이 더 많은 표를 얻을 이벤트의 확률이다.

P_i 는 경제 성장률, 물가 상승률, 현재 집권당 여부와 같은 특성을 가진 i 연도에 공화당 대통령 후보가 더 많은 표를 얻을 확률이다. 데이터에서 직접 관찰되지 않으나 Python을 비롯한 여러 통계 패키지들에서는 주로 MLE(maximum likelihood estimation)을 이용하여 이를 추정한다. 표본 데이터의 일부와 추정된 Logit 모델은 다음과 같다.

표본 데이터

연도	y :공화당 다득표 여부	x_1 :경제 성장률	x_2 :물가 상승률	x_3 :현재 집권 여부
1880	1	3.879	1.974	1
1884	0	1.589	1.055	1
		⋮		
2016	0	1.208	1.411	0

$$\log\left(\frac{P_{\text{공화당}}}{1-P_{\text{공화당}}}\right) = -0.93 - 0.09 [\text{경제 성장률}] + 0.08 [\text{물가 상승률}] + 1.07 [\text{현재 집권 여부}]$$

$$P_{\text{공화당}} = \frac{1}{1 + e^{-(-0.93 - 0.09[\text{경제 성장률}] + 0.08[\text{물가 상승률}] + 1.07[\text{현재 집권 여부}])}}$$

여기서 $P_{\text{공화당}}$ 은 공화당 표가 더 많은 이벤트의 확률이다.

간단한 위 모델에 의하면 미국의 공화당 대통령 후보는 경제성장률이 낮고 물가상승률이 높으며, 현재 집권을 한 상태일 때 유권자들의 지지를 얻을 확률이 높아진다.

대표적으로 공화당이 투표연도에 집권당일 때의 효과를 살펴보자. 공화당이 집권당일 때는 민주당이 집권당일 때보다 로그 승산비는 1.07만큼 커지고, 승산비는 $e^{1.07} = 2.92$ 만큼 상승한다. 즉, 공화당이 집권한 연도에는 민주당이 집권한 연도보다 민주당 후보 대비 공화당 후보가 우세할 확률이 2.9배 가량 커진다는 뜻이다. 물론 이 해석은 경제 성장률과 물가 상승률이 변하지 않았다는 가정하에서 성립한다.

미국 경제가 평균적인 경제 상태일 때, 공화당이 더 많은 표를 얻을 확률도 구할 수 있다. 1880년부터 2016년동안의 경제 성장률, 물가 상승률, 현재 집권당 여부를 각각 산술평균한 값인 0.65%, 2.56%, 0.54를 대입하면, $P_{\text{공화당}} = 0.45$ 가 추정된다. 즉, 지난 36년간의 평균적인 경제 상황에서는 미국의 공화당이 대통령 선거에서 더 많은 표를 받을 확률은 45%이다.

여기서 분류작업을 수행할 수 있다. 추정된 공화당 득표율이 특정 임계값 이상인 연도를 “이벤트 1(공화당 표가 더 많음)”으로, 임계값 이하인 연도를 “이벤트 0(공화당 표가 더 많지 않음)”으로 분류하는 일이다. 1880년부터 2016년을 포함하는 기간 동안 평균적인 경제 상황에서는 해당 확률이 45%였는데, 50%를 임계값으로 삼아 이를 “이벤트 0”으로 분류할 수 있다. 2016년의 관측치에 이를 적용하면 $P_{\text{공화당}} = 0.28$ 이 추정된다. 2016년 역시 “이벤트 0”으로 분류하게 된다.

III. 사용법

1. 파이썬(Python) 사용법

엑셀 기본 프로그램에서는 로지스틱 회귀를 제공하지 않고 있다. 대신, 파이썬에서는 여러 머신러닝 및 통계 패키지에서 로지스틱 회귀를 사용할 수 있다. 대표적인 sklearn 패키지에서의 사용법은 다음과 같다.

```
# -*- coding: utf-8 -*-

#모듈 가져오기 # 파이썬에서는 샵(#)기호로 주석을 작성
import pandas as pd # pandas를 가져와 pd로 사용
from sklearn.datasets import load_breast_cancer #load_breast_cancer 로드
from sklearn.linear_model import LogisticRegression #로지스틱 회귀 가져오기

#유방암 진단자료(파이썬 내재 데이터) 샘플 가져오기
# y에는 유방암이 양성(benign)이면 1, 악성(malignant)이면 0으로 기록, 2종만 존재
# X는 유방암에 영향을 주는 radius, texture, perimeter 등의 30개 독립변수임
X, y = load_breast_cancer(return_X_y=True)

print(y[:30]) # y 샘플의 첫 30개를 출력 (전체는 569개임)

[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 0 0]

print(X[:2]) # X 샘플의 첫 2개를 출력

[[1.799e+01 1.038e+01 1.228e+02 1.001e+03 1.184e-01 2.776e-01 3.001e-01
 1.471e-01 2.419e-01 7.871e-02 1.095e+00 9.053e-01 8.589e+00 1.534e+02
 6.399e-03 4.904e-02 5.373e-02 1.587e-02 3.003e-02 6.193e-03 2.538e+01
 1.733e+01 1.846e+02 2.019e+03 1.622e-01 6.656e-01 7.119e-01 2.654e-01
 4.601e-01 1.189e-01]
 [2.057e+01 1.777e+01 1.329e+02 1.326e+03 8.474e-02 7.864e-02 8.690e-02
 7.017e-02 1.812e-01 5.667e-02 5.435e-01 7.339e-01 3.398e+00 7.408e+01
 5.225e-03 1.308e-02 1.860e-02 1.340e-02 1.389e-02 3.532e-03 2.499e+01
 2.341e+01 1.588e+02 1.956e+03 1.238e-01 1.866e-01 2.416e-01 1.860e-01
 2.750e-01 8.902e-02]]

#로지스틱 회귀 실행
logistic = LogisticRegression(max_iter=10000) #로지스틱 회귀 모델 하나 생성
logistic.fit(X, y) #X와 y 데이터로 로지스틱 회귀 실시

print('logistic coef:\n', logistic.coef_) #독립변수의 각 회귀계수 출력
print('\nlogistic intercept:\n', logistic.intercept_) #절편값 출력

logistic coef:
[[ 0.94685226  0.18294513 -0.26494659  0.02261067 -0.1866965  -0.20411948
 -0.54240207 -0.30569951 -0.27516362 -0.02961991 -0.07743169  1.27259801
  0.13674665 -0.11071874 -0.02827259  0.08230945 -0.03278114 -0.0401502
 -0.03189051  0.01647542  0.15817157 -0.43910338 -0.11162044 -0.01340686
 -0.37764545 -0.64213555 -1.45347724 -0.62513887 -0.7271355  -0.09074315]]

logistic intercept:
[28.20058428]
```

여기서는 분류 문제에서 사용할 수 있는 `breast_cancer` 데이터를 가져왔다. `load_breast_cancer` 안의 `y`(조속변수) 데이터는 1과 0의 값만 가지며, 유방암이 양성(benign)이면 1, 악성(malignant)이면 0으로 기록된다. `load_breast_cancer` 안의 `X`(독립변수) 데이터는 `radius`(종양 반경)부터 `fractal_dimension`까지 30개의 종양 관련 지표로 구성되어 있다. `load_breast_cancer`는 총 569개 샘플로 구성되어 있다. 즉, `y` 변수는 569x1의 행렬 변수이며, `X` 변수는 569x30의 행렬 변수다.

로지스틱 모델을 하나 생성해서, `fit` 메서드를 통해서 분석(학습)을 하면, 로지스틱의 회귀계수가 결정된다. `coef_` 속성과 `intercept_` 속성을 출력하면, 결정된 회귀분석식의 계수들을 볼 수 있다.

(“[[0.94685226 0.18294513 ... “으로 출력된 부분)

여기서 결정된 로지스틱 회귀분석식을, 다시 풀어서 쓰면 다음과 같다.

$$\log\left(\frac{p_{y=1}}{1-p_{y=1}}\right) = 28.20 + 0.95radius + 0.18texture + \dots - 0.09fractal_dimension$$

$$p_{y=1} = \frac{1}{1 + e^{-(28.20+0.95radius+0.18texture+\dots-0.09fractal_dimension)}}$$

각 독립변수 수치가 주어지면, 위의 공식에 의해서 유방암이 양성($y=1$)일 확률 결과가 계산으로 도출된다.

*predict*는 input에 따른, 분류 결과나 회귀분석 결과를 반환하는 메서드다. 여기서는 X를 input으로 하여 0 혹은 1의 분류 결과를 반환한다.

*predict_proba*는 예측과정에서의 각 분류별 확률을 보여주는 메서드다. 여기서는 y가 0이 될 확률, 1이 될 확률 값을 순서대로 보여준다. Input 데이터가 현재 569개이므로, 0이 될 확률과 1이 될 확률 2개로 열을 구성하여 총 569x2 사이즈의 수치를 반환한다.

1이 될 확률 수치는, 앞의 로지스틱 회귀에서 본 수식 $p_{y=1} = \frac{1}{1+e^{-(\beta_0+\beta_1x_1+\dots+\beta_kx_k)}}$ 에 각 수치를 넣어서 나온 값이다.

*predict_log_proba*는 각 분류별 확률의 log값을 반환하는 메서드다. 즉, 로지스틱 회귀식에서 $\log(p_{y=1})$ 와 $\log(p_{y=0})$ 를 반환한다고 볼 수 있다. 여기서는 y가 0이 될 확률 log 값, 1이 될 확률 log 값을 순서대로 보여준다. 그리고, Input 데이터가 현재 569개이므로 출력값도 569x2 사이즈의 수치다.

예제 상, *predict_log_proba* 반환값의 첫째 열이 $\log(p_{y=0}) = \log(1 - p_{y=1})$ 이고 둘째 열이 $\log(p_{y=1})$ 다. 따라서 둘째 열에서 첫째 열을 뺀 값이, 로지스틱 회귀 공식 상에서의 선형 회귀 부분식에 해당한다.

$$\text{둘째 열} - \text{첫째 열} = \log(p_{y=1}) - \log(1 - p_{y=1}) = \log\left(\frac{p_{y=1}}{1 - p_{y=1}}\right) = \beta_0 + \beta_1x_1 + \dots + \beta_kx_k$$

<참고문헌>

Fair, Ray C. 1996. "Econometrics and Presidential Elections." *Journal of Economic Perspectives*, 10 (3): 89-102.

신뢰에 가치로 답하다

삼성증권



삼성증권주식회사

서울특별시 서초구 서초대로74길 11(삼성전자빌딩)
Tel: 02 2020 8000 / www.samsungpop.com

삼성증권 지점 대표번호: 1588 2323 / 1544 1544

고객 불편사항 접수: 080 911 0900



MEMBER OF
**Dow Jones
Sustainability Indices**
In Collaboration with RobecoSAM