

[Statistics]

Binary Classification

오영석

고려사이버대학교 AI·데이터과학부 외래교수

Contents

01 Binary Classification	03
1.1 Learning objectives	04
1.2 Parallax Thinking	05
1.3 Binary Classification Overview	06
1.4 Sigmoid Function	14
1.5 Loss Function	21
1.6 Learning Summary	32

Binary Classification

본 단원에서는 이진분류 시스템에 대해 학습합니다.

1.1 Learning objectives

1. 이진분류 시스템에 대해서 설명할 수 있다.
 2. 로지스틱 알고리즘에 대해서 설명할 수 있다.
 3. 시그모이드 함수에 대해서 설명할 수 있다.
 4. 이진 크로스 엔트로피에 대해서 설명할 수 있다.
-

1.2 Parallax Thinking

생각 열기

스팸 메일을 분류하거나 암 조직과 정상조직을
분류하기 위해서는 어떤 데이터를 수집해야 할까?

수집한 데이터는 어떤 알고리즘을 통해
분류할 수 있을까?

본 절에서는 **로지스틱 회귀 알고리즘**을 사용하여
수집한 데이터를 분류하고자 함

1.3 Binary Classification Overview

Idea of Binary Classification

이진 분류는 트레이닝 데이터의 특징을 학습하여, 임의의 테스트 데이터를 사전에 정의된 두 가지 범주 중 하나로 분류하는 예측 모델을 구축하는 과정

- 이메일 스팸 분류, Spam(1) 또는 Ham(0)
- 금융 사기 탐지, 사기 거래(1) 또는 정상 거래(0)
- 의료 진단, 암 조직(1) 또는 정상 조직(0)

1.3 Binary Classification Overview

Idea of Binary Classification

이진 분류는 **트레이닝 데이터**의 특성과 그들 간의 상관관계를 분석하여, 임의의 입력 데이터를 사전에 정의된 두 가지 범주 중 하나로 분류할 수 있는 예측 모델을 만드는 과정

Petal length (x)	Species (t)
1.4	Iris-setosa
1.5	Iris-setosa
4.7	Iris-versicolor
4.5	Iris-versicolor

Training Data

1.3 Binary Classification Overview

Idea of Binary Classification

이진 분류는 **트레이닝 데이터의 특성과 그들 간의 상관관계를 분석하여**, 임의의 입력 데이터를 사전에 정의된 두 가지 범주 중 하나로 분류할 수 있는 예측 모델을 만드는 과정



1.3 Binary Classification Overview

Idea of Binary Classification

이진 분류는 트레이닝 데이터의 특성과 그들 간의 상관관계를 분석하여, **임의의 입력 데이터를** 사전에 정의된 두 가지 범주 중 하나로 분류할 수 있는 예측 모델을 만드는 과정

Petal length (x)	Species (t)
1.4	Iris-setosa
1.5	Iris-setosa
4.7	Iris-versicolor
4.5	Iris-versicolor

Training Data

Petal length (x)
1.6
4.6

Test Data

1.3 Binary Classification Overview

Idea of Binary Classification

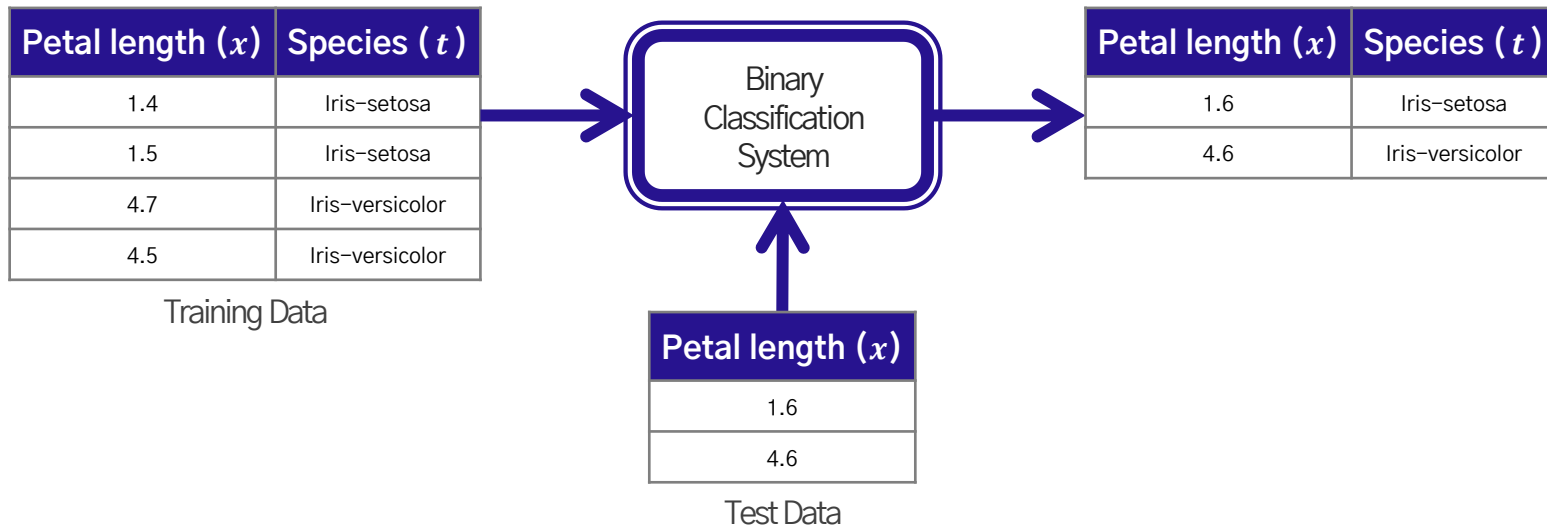
이진 분류는 트레이닝 데이터의 특성과 그들 간의 상관관계를 분석하여, 임의의 입력 데이터를
사전에 정의된 두 가지 범주 중 하나로 분류할 수 있는 예측 모델을 만드는 과정

Petal length (x)	Species (t)
1.6	Iris-setosa
4.6	Iris-versicolor

1.3 Binary Classification Overview

Idea of Binary Classification

이진 분류 과정을 도식화하면 다음과 같음

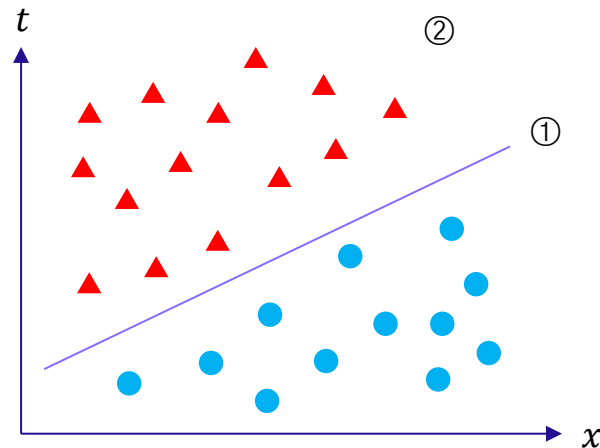


1.3 Binary Classification Overview

Logistic Regression

로지스틱 회귀 알고리즘은 ① 트레이닝 데이터의 특성과 분포를 나타내는 최적의 직선을 찾고, ② 해당 직선을 기준으로 데이터를 위(1)나 아래(0) 또는 왼쪽(1)이나 오른쪽(0) 등으로 분류하는 방법

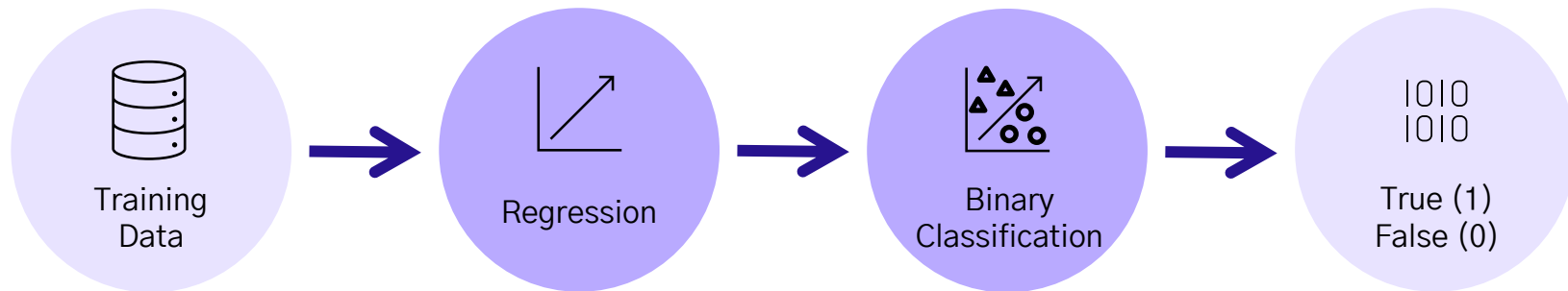
이러한 로지스틱 회귀는 이진 분류 시스템의 알고리즘 중에서도 정확도가 높은 알고리즘으로 알려져 있어서 딥러닝에서도 기본적인 컴포넌트로 사용되고 있음



1.3 Binary Classification Overview

Logistic Regression

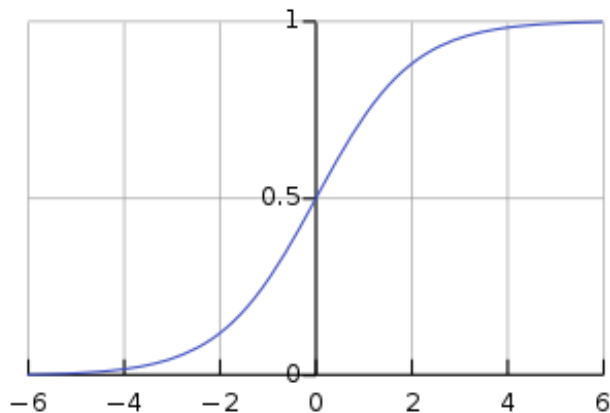
로지스틱 회귀 알고리즘을 도식화하면 다음과 같음



1.4 Sigmoid function

Sigmoid function

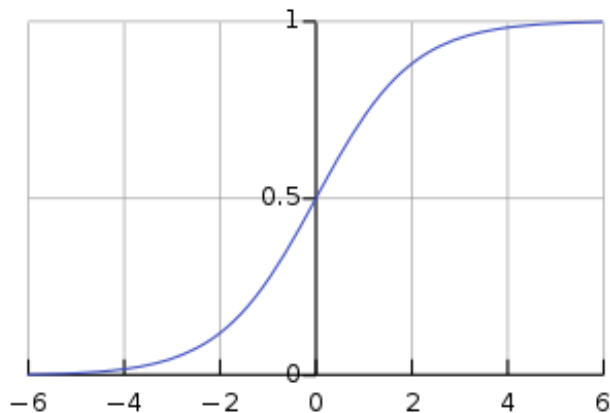
아래의 그림은 시그모이드 함수에 대한 그래프 표현



1.4 Sigmoid function

Sigmoid function

시그모이드 함수의 함수식 표현은 아래와 같으며 (0, 1) 범위의 값으로 변환함



$$y = \text{Sigmoid}(z) = \frac{1}{1 + \exp(-z)}$$

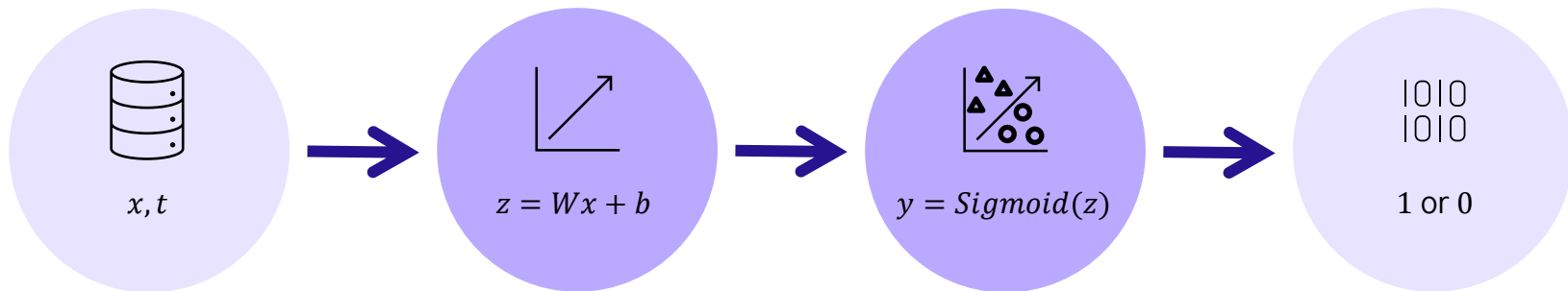
[2] https://en.wikipedia.org/wiki/Sigmoid_function

[3] 김동근 (2024), Step by step 파이토치 딥러닝 프로그래밍, 서울: 도서출판 가메.

1.4 Sigmoid function

Sigmoid function

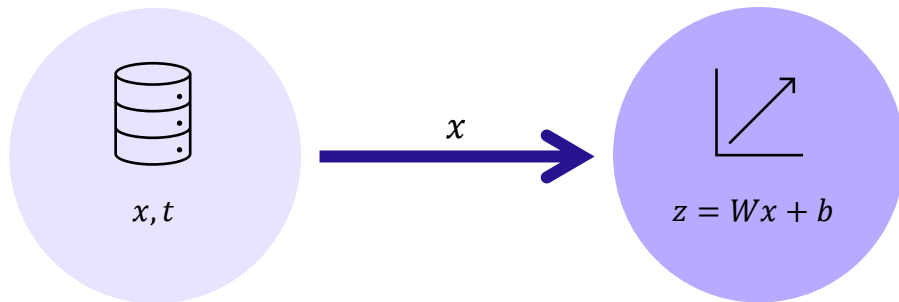
시그모이드 함수를 사용해서 이진 분류 시스템을 구현하는 방법은 다음과 같음



1.4 Sigmoid function

Sigmoid function

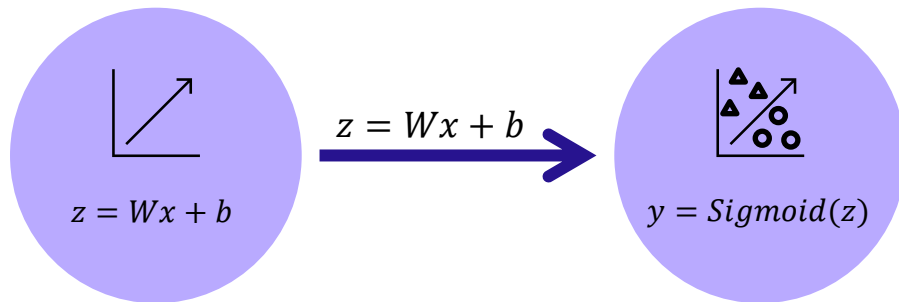
① 트레이닝 데이터의 특징 변수 x 값이 회귀에 입력으로 들어가서 $z = Wx + b$ 값으로 계산됨



1.4 Sigmoid function

Sigmoid function

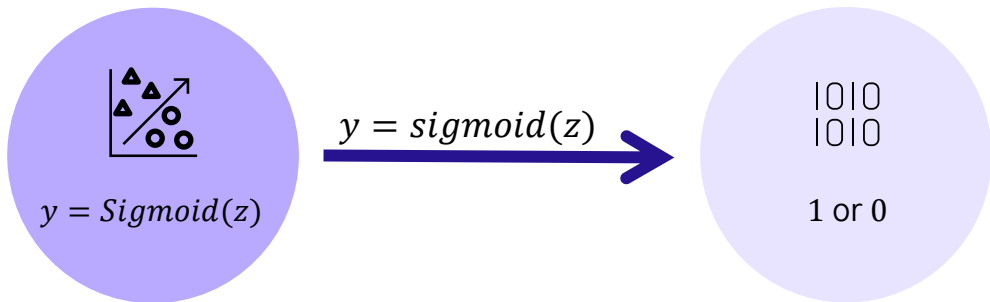
② z 값이 이진 분류에 입력에 입력으로 들어가서 $y = \text{Sigmoid}(z)$ 값으로 계산됨



1.4 Sigmoid function

Sigmoid function

- ③ 출력된 $y = \text{Sigmoid}(z)$ 값이 0.5 이상이면 논리적인 결과 값을 True (1) 상태로 정의하고,
 $y = \text{Sigmoid}(z)$ 값이 0.5 미만이면 논리적인 결과 값을 False (0) 상태로 정의하여 이진 분류 시스템을 구현함



1.4 Sigmoid function

Sigmoid function

시그모이드 함수는 0과 1사이의 값으로 계산되므로 시그모이드 함수의 결과를 확률로 해석할 수도 있음

- 시그모이드 함수 값이 0.72라면, 결과가 나올 확률이 72%임을 의미함

1.5 Loss Function

Binary Cross Entropy

이진 분류 시스템에서 최종 출력 값 y 는 시그모이드 함수의 계산 값이므로 논리적으로 True (1) 또는 False (0) 값을 가지기 때문에, 선형 회귀 때와는 다른 손실함수가 필요함

1.5 Loss Function

Binary Cross Entropy

손실 함수는 $y = \frac{1}{1+\exp(-z)}$ 시그모이드 함수에 의해서 0과 1사이로 계산된 값 y 와 정답 $t_i = 0$ or 1 와의 오차를

토대로 다음과 같은 함수식으로 표현함

$$E(W, b) = - \sum_{i=1}^n \{t_i \log y_i + (1 - t_i) \log(1 - y_i)\}$$

1.5 Loss Function

Binary Cross Entropy

그리고 이진 크로스 엔트로피를 사용하면, 선형회귀 강의에서 배웠던 것과 같이

가중치는 $W^* = W - \alpha \frac{\partial E(W,b)}{\partial W}$ 와 같이 계산하여 손실함수 $E(W,b)$ 의 최소값을 가지는 W 를 구할 수 있고,

편향 또한 $b^* = b - \alpha \frac{\partial E(W,b)}{\partial b}$ 와 같이 계산하여 손실함수 $E(W,b)$ 의 최소값을 가지는 b 를 구할 수 있음

학습퀴즈

1

이진 분류 시스템은 여러가지 범주 중 하나로 분류할 수 있는 예측모델을 만드는 과정이다.

O

X

학습퀴즈

1 이진 분류 시스템은 여러가지 범주 중 하나로 분류할 수 있는 예측모델을 만드는 과정이다.

O X

해설

이진 분류는 두 가지 범주 중 하나로 분류할 수 있는 예측 모델을 만드는 과정이다.

학습퀴즈

2

로지스틱 회귀는 이진 분류 알고리즘 중에서도
정확도가 높지 않은 알고리즘으로 알려져 있다.

O

X

학습퀴즈

2

로지스틱 회귀는 이진 분류 알고리즘 중에서도
정확도가 높지 않은 알고리즘으로 알려져 있다.

O

X

해설

로지스틱 회귀는 이진 분류 알고리즘 중에서도 정확도가 높은 알고리즘이다.

학습퀴즈

3

시그모이드 함수는 0과 1사이의 값으로 계산되어
시그모이드 함수의 결과를 확률로 해석할 수도 있다.

0

X

학습퀴즈

3

시그모이드 함수는 0과 1사이의 값으로 계산되어
시그모이드 함수의 결과를 확률로 해석할 수도 있다.

0

X

해설

시그모이드 함수의 결과는 확률로 해석할 수 있다.

학습퀴즈

4

이진 크로스 엔트로피는 분류 시스템에서의 손실함수의 한 종류이다.

O

X

학습퀴즈

4

이진 크로스 엔트로피는 분류 시스템에서의 손실함수의 한 종류이다.

0

X

해설

이진 크로스 엔트로피는 이진 분류 시스템의 손실함수이다.

1.6 Learning Summary

학습 정리

이진 분류는 트레이닝 데이터의 특성과 그들 간의 상관관계를 분석하여, 임의의 입력 데이터를 사전에 정의된 두 가지 범주 중 하나로 분류할 수 있는 예측 모델을 만드는 과정이다.

로지스틱 회귀 알고리즘은 ① 트레이닝 데이터의 특성과 분포를 나타내는 최적의 직선을 찾고, ② 해당 직선을 기준으로 데이터를 위(1)나 아래(0) 또는 왼쪽(1)이나 오른쪽(0) 등으로 분류하는 방법이다.

시그모이드 함수는 0과 1사이의 값으로 계산되므로 시그모이드 함수의 결과를 확률로 해석할 수도 있다.

이진 크로스 엔트로피를 사용하면, 손실함수 $E(W, b)$ 의 최소값을 가지는 가중치와 편향을 구할 수 있다.

End of Document

Thank you

본 콘텐츠를 어떠한 경로로든 외부로 일부 또는 전부를 복사, 복제, 판매, 재판매, 공유 등을 할 수 없습니다.
이를 위반하는 경우, 지식재산권 침해 및 관련 법률에 따라 책임을 질 수 있습니다.
