

# 범주형 자료분석

## Basics in Logit and Probit Models

김현우, PhD<sup>1</sup>

<sup>1</sup>충북대학교 사회학과 조교수

March 24, 2025



# 진행 순서

- 1 로지스틱 회귀분석의 원리
- 2 회귀계수의 해석
- 3 유의성 검정

## 로지스틱 회귀분석의 원리

# 로지스틱 회귀분석의 원리

종속변수가 가변수라면 회귀분석이 어떻게 달라지는가 상상해보자.

- 단순히 생각하면 가변수를 그대로 종속변수로 한 회귀식에서  $\beta_0$ 와  $\beta_1$ 을 추정할 수 있다.
- 종속변수가 0 또는 1이므로 (이것을 마치 확률처럼 접근하면) “ $X$ 가 한 단위 증가할 때  $Y$ 가 1이 될 확률(probability of being 1)은  $\beta_1$ 만큼 증가한다”고 해석한다(Why?).

$$E(Y|X) = P(Y = 1|X) = \beta_0 + \beta_1 X$$

- 이런 타입의 회귀모형을 **선형확률모형(linear probability model)**이라고 부른다.
- 그러나 이것은 (1)  $\hat{Y}$ 가 0보다 작거나 1보다 큰 값이 나오기도 하고, (2) 회귀분석의 가정 가운데 **등분산성(homoscedasticity)**에도 위배된다(증명 생략).



# 로지스틱 회귀분석의 원리

- 본래 선형회귀식의 종속변수는 양적 변수이므로 당연히  $-\infty < E(Y|X) < \infty$  를 전제한다.
- 그러나 가변수가 종속변수라면 반드시  $E(Y|X) = \{0, 1\}$  여야만 한다.
- 가변수로 부호화(encoding)되었다는 것은 범주가 두 개라는 뜻이다.
- (가변수를 독립변수로 사용할 때와 마찬가지로) 어느 한 쪽이 **기준집단(reference group)** 또는 **기저범주(baseline category)**이 되어 분석에서 제외된다.
- 어느 쪽이 1이고 나머지가 0이 될지는 직접 결정한다.



# 로지스틱 회귀분석의 원리

가변수가 종속변수일 때는 종속변수를 살짝 바꾸면 된다.

- 먼저 **오즈(odds)** 또는 **승산(勝算)**은 다음과 같이 정의된다.

$$odds = \frac{P(Y = 1|X)}{1 - P(Y = 1|X)}$$

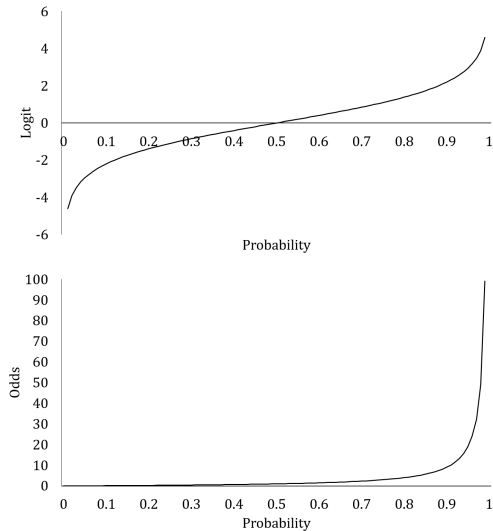
- 오즈에 자연대수를 취한 값을 로그오즈 또는 **로짓(logit)**이라고 부른다.

$$logit = \log_e \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} = \ln \frac{P(Y = 1|X)}{1 - P(Y = 1|X)}$$

- 로짓을 **로지스틱 함수(logistic function)**라고도 부른다.
- 확률은 정의상  $[0, 1]$  사이로 제한되지만, 로짓의 범위는  $[-\infty, \infty]$  이므로 이제  $-\infty < E(Y|X) = P(Y = 1|X) < \infty$  가 성립한다.



# 로지스틱 회귀분석의 원리



# 로지스틱 회귀분석의 원리

- 이제  $E(Y|X) = \beta_0 + \beta_1 X$  대신 다음의 회귀식을 사용한다.

$$\ln \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} = \beta_0 + \beta_1 X$$

- 즉 종속변수로  $\{0, 1\}$ 인 원점수를 그대로 쓰지 않고, 범위 제한이 없는 로짓을 사용하는 편이 선형확률모형보다 낫다!
- 로짓을 종속변수로 하므로 로지스틱 회귀모형 또는 **로짓(logit)** 회귀모형이라고도 부른다.
- 구태여 덧붙이자면, (자료에서 관측되는) 종속변수가  $\{0, 1\}$ 이므로 **이항(binary)** 로지스틱(또는 이항 로짓) 회귀모형이라고도 부른다.
- 교과서에 따라서는 로지스틱 회귀모형과 로짓 회귀모형을 개념상 구분하지만, 실익이 크지 않으므로 너무 집착하지 않아도 된다.





## 회귀계수의 해석

# 회귀계수의 해석

로지스틱 회귀분석은 보통최소제곱이 아닌 다른 알고리즘을 요구한다.

- 선형회귀분석은 오차제곱합(sum of squared error)을 최소화하는  $\beta_0$ 와  $\beta_1$ 를 찾기 위해 보통최소제곱(OLS)이라는 알고리즘을 사용하였다.
- 로지스틱 회귀분석은 비선형모형(nonlinear model)인 탓에 OLS를 사용할 수 없다. 대신 최대우도법(maximum likelihood estimation; MLE)을 주로 사용한다.
- 최대우도법은 우도함수(likelihood function)  $\lambda$ 를 정의한 뒤, 이것의 로그 우도(log-likelihood)  $\ln\lambda$ 를 극대화하는  $\beta_0$ 와  $\beta_1$ 를 찾는 알고리즘이다.

$$\operatorname{argmax}_{\beta_0, \beta_1} \ln \prod_{i=1}^n f_i(Y_i) = \operatorname{argmax}_{\beta_0, \beta_1} \ln \lambda(\beta_0, \beta_1; Y_i)$$



# 회귀계수의 해석

안타깝게도 로그오즈의 해석은 직관적이지 않다.

- 회귀식을 잘 살펴보면 로지스틱 회귀분석 특유의  $\beta_1$ 의 해석법을 깨달을 수 있다.

$$\ln \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} = \beta_0 + \beta_1 X$$

- 가령 종속변수가 낙태권 지지(0=지지하지 않음; 1=지지함)이고 독립변수가 교육연수라고 하자.
- “교육연수  $X$ 가 한 단위 증가하면, 낙태권을 지지할 로그오즈는  $\beta_1$  만큼 증가한다.”
- “로그오즈가 증가한다니... 이게 대체 무슨 소린가?” 나도 모르겠다.



# 회귀계수의 해석

그렇기 때문에 해석상 차라리 오즈비가 선호된다.

- 독립변수  $X$ 가 한 단위 더 증가한다면 예전 상태와 비교하여 **오즈비(odds ratio; OR)**는 다음과 같이 정의된다.

$$\begin{aligned} OR &= \frac{odds_{X+1}}{odds_X} = \frac{\left( \frac{P(Y=1|X+1)}{1 - P(Y=1|X+1)} \right)}{\left( \frac{P(Y=1|X)}{1 - P(Y=1|X)} \right)} \\ &= \frac{e^{\beta_0 + \beta_1(X+1)}}{e^{\beta_0 + \beta_1 X}} = e^{\beta_1(X+1) - \beta_1 X} = e^{\beta_1} \end{aligned}$$

- 즉  $e^{\beta_1}$ 는  $X$ 가  $X+1$ 로 한 단위 증가할 때 ‘오즈가 달라지는 비율’을 보여준다 (Why?)!
- 그러므로  $e^{\beta_1}$ 를 그대로 오즈비라고 부른다.



# 회귀계수의 해석

- 오즈비는 정의상 비율(ratio)이므로, 가령 오즈비가 1.5라면 분모의 오즈보다 분자의 오즈가 50% 크다는 것을 뜻한다.
- 그러므로 ‘오즈가 달라지는 비율’을 다음과 같이 백분율로 나타낼 수 있다(Why?).

$$\Delta\% = 100 \cdot (e^{\beta_1} - 1)$$

- “ $X$ 가 한 단위 증가하면,  $Y = 1$ 의 오즈가  $100 \times (e^{\beta_1} - 1)$  퍼센트 증가한다.”
- 로지스틱 회귀계수의 오즈비 해석은 혼동하기 쉬우므로 많이 연습해야 한다.



# 회귀계수의 해석

예제 1. lbw.dta를 이용하여 저출생체중아 여부(low)를 종속변수로, 임신전 마지막 월경 당시 산모체중(lwt)을 독립변수로 하는 회귀식을 추정하시오. 두 변수의 연관성에 관해 해석하시오.



# 회귀계수의 해석

- 종속변수인 저출생체중아 여부(low)를 잘 살펴보면 0 (저출생체중 아님) 또는 1 (저출생체중)로 부호화되어 있으므로 (OLS를 사용한 선형회귀분석보다) 로지스틱 회귀분석이 바람직하다.
- Stata를 사용하여 다음과 같이 회귀식을 추정한다.

$$\ln \frac{P(\text{low} = 1|X)}{1 - P(\text{low} = 1|X)} = 0.996 - 0.014\text{wt}$$

- “산모의 마지막 월경 당시 체중이 1 파운드 증가하면 산아가 저출생체중일 로그오즈는 0.014 만큼 감소한다.”
- “산모의 마지막 월경 당시 체중이 1 파운드 증가하면 산아가 저출생체중일 오즈는 1.4% ( $=100 \times (e^{0.014} - 1)$ ) 감소한다.”



# 회귀계수의 해석

어쩌면 오즈비조차도 별로 직관적이지 않다.

- 오즈보다는 확률이야말로 우리에게 가장 직관적으로 와닿는 해석을 제공한다!
- 가령 “내가 이길 확률은 80% 정도 된다구” 라고 말하면 모를까, “내가 이길 오즈는 4 (=0.8/0.2) 정도 된다구” 라고 말하는 것은 몹시 이상하다.
- 그러므로 (로그)오즈를 아예 확률로 변환할 수 있다면 해석에 직관성을 더할 수 있다.
- 로그오즈와 확률의 관계는 다음과 같다(증명).

$$\text{logit} = \ln \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} = \beta_0 + \beta_1 X$$
$$P(Y = 1|X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}}$$

- 그러면 추정된  $\hat{\beta}_0$ 와  $\hat{\beta}_1$ 를 대입하고  $X$ 의 증가에 따라 예측확률(predicted probability)  $P(Y = 1|X)$ 이 어떻게 변화하는지 그림을 그릴 수 있다.





# 회귀계수의 해석

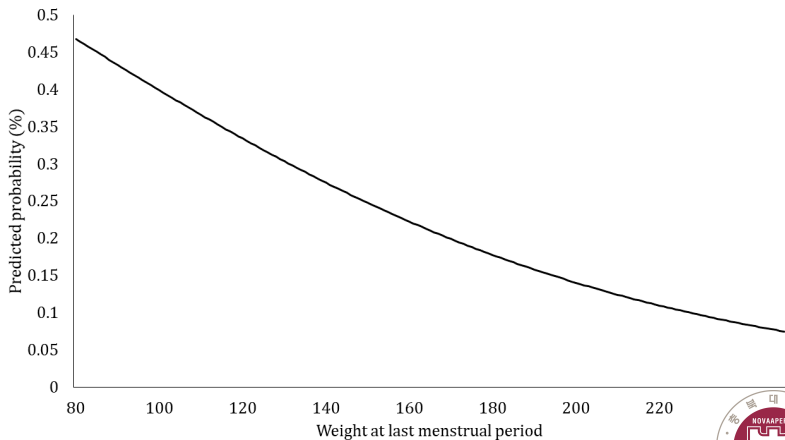
- Stata를 사용하여 예측확률  $P(\widehat{Y=1|X})$ 을 계산할 수 있다. 엑셀에서도 직접 계산해보자(이것은 사실 훌륭한 연습이 된다).
- 임신전 체중(lwt)은 [80, 250] 파운드 사이에 놓이므로 각각의 예측확률을 계산한다.

	A	B	C	D	E	F
1	Variables in the Equation					
2			B	S.E.	Wald	df
3	Step 1 <sup>a</sup>	Weight at last	-0.014	0.006	5.178	
4		Constant	0.996	0.785	1.608	
5	a. Variable(s) entered on step 1: Weight at last menstrual period.					
6						
7	X	b_0+b_1X	P			
8	80	=C\$4+C\$3*A8				
9	81	-0.14124	0.464748			
10	82	-0.15528	0.461258			
11	83	-0.16932	0.457771			



# 회귀계수의 해석

- 마지막으로 꺾은선 그래프(line chart)로 그린다.



## 유의성 검정

# 유의성 검정

로지스틱 회귀모형에서는 색다른 유의성 검정이 사용된다.

- OLS에서는 회귀계수의 통계적 유의성을 확인하기 위해 다음과 같은  $t$  검정( $t$  test)을 사용하였다.

$$t = \frac{b - \beta}{SE_b} = \frac{b}{SE_b}$$

- $t$  분포의 꼬트머리를 그리고 그 면적을 통해  $p$  값을 구하는 방식으로 귀무가설 ( $H_0 : \beta = 0$ )을 기각할 수 있다.



# 유의성 검정

- 로지스틱 회귀모형과 같은 비선형모형에서는 자유도(degree of freedom: df)가 1인  $\chi^2$  검정( $\chi^2$  test)을 사용한다(증명 생략).

$$Wald = \left( \frac{b - \beta}{SE_b} \right)^2 = \left( \frac{b}{SE_b} \right)^2 \sim \chi_1^2$$

- 이  $\chi^2$  검정을 고안한 Abraham Wald의 이름을 따 특별히 **왈드 검정(Wald test)**이라고 부른다.
- $\chi_1^2$  값이 충분히 크면 통계적으로 유의하게 귀무가설( $H_0 : \beta = 0$ )을 기각할 수 있다.



예제 2. lowbwt.sav는 저출생체중 여부(low) 뿐 아니라 산아의 체중 원점수(bwt)를 양적 변수로도 제공하고 있다(참고로 저출생체중아의 기준은 2,500g 이다). 모두 똑같이 smoke, age, lwt, ht를 독립변수로 하되, (1) 종속변수는 체중 원점수(bwt)인 선형회귀모형, (2) 종속변수는 저출생체중 여부(low)인 선형확률모형, 그리고 (3) 종속변수는 저출생체중 여부(low)인 로지스틱 회귀모형을 각각 추정하시오. 특히 임신중 흡연 여부에 초점을 두고 그 결과를 비교하시오.



# 유의성 검정

- 결과표를 엑셀로 옮겨 정리하자. 어떤 식으로 회귀식을 추정하던지 유의성 검정 결과는 거의 똑같다.
  - (1) “(다른 변수의 영향력을 통제할 때) 산모가 임신중 흡연할 경우 산아의 체중은 그렇지 않은 산모보다 평균적으로 261.9g 적다.”
  - (2) “(다른 변수의 영향력을 통제할 때) 산모가 임신중 흡연할 경우 산아가 저출생체중아가 될 확률은 13.9% 증가한다.”
  - (3a) “(다른 변수의 영향력을 통제할 때) 산모가 임신중 흡연할 경우 산아가 저출생체중아가 될 오즈는 97.2% 증가한다.”
- 예측확률을 직접 계산해 보면 다음과 같이 해석할 수도 있다(Why?).
  - (3b) “(다른 변수의 영향력을 통제할 때) 산모가 임신중 흡연할 경우 산아가 저출생체중아가 될 확률은 25.7%에서 39.3%로 약 13.6% 증가한다.”



# 유의성 검정

- 선형확률모형과 로지스틱 회귀모형의 추정값은 대체로 비슷하다(늘 그런 것은 아니다). 직접 OLS와 logit의 결과물 차이를 비교해보자.

