

به نام خداوند بخشنده مهربان

# روش‌های آماری و اقتصادسنجی

در تحلیل و مدل‌سازی داده‌های حمل و نقلی

محمد مهدی بشارتی

---

[besharati@iut.ac.ir](mailto:besharati@iut.ac.ir)

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

مقدمه

پروبیت دوگانه

لوجیت دوگانه (لوجستیک)

معیارهای نیکویی برآزش

ریسک و شانس

لوجیت چندگانه

لوجیت شرطی

لوجیت آشیانه‌ای و ترکیبی



## □ انواع خاصی از متغیرهای گسسته؛

- مقادیر ۰، ۱، ۲، ... داشته باشد.
- مقادیر ۰ و ۱ (دو انتخابی)
- مقادیر محدود (مثلا ۴ گروه: ۰، ۱، ۲، ۳)



متغیر وابسته ممکن است در اختیار کردن مقادیر خود با محدودیت روبرو باشد.

مثال؛

عوامل موثر بر اشتغال افراد (متغیر وابسته؟)،

عوامل موثر بر فوتی بودن یک تصادف (متغیر وابسته؟)،

عوامل موثر بر انتخاب حمل و نقل شخصی (متغیر وابسته؟)،

و ...

به طور کلی، متغیر وابسته می تواند یک متغیر کیفی باشد که نتیجه تصمیم گیری های افراد (یا یک رویداد مشخص) را نشان دهد.

یک انتخاب یا وقوع یک رویداد را به مجموعه ای از متغیرهای توصیفی (توضیحی) مرتبط می کنیم (رگرسیون).

پایه اصلی این موارد در چارچوب کلی مدل های احتمال مورد بررسی قرار می گیرد.



- پایه اصلی این موارد در چارچوب کلی مدل‌های احتمال مورد بررسی قرار می‌گیرد.
- مثلاً وقوع تصادف فوتی را به عنوان متغیر  $Y$  در نظر می‌گیریم: ( وقوع تصادف فوتی:  $Y=1$  )
- احتمال وقوع حادثه مورد نظر (احتمال فوتی بودن تصادف):

$$P(\text{فوتی بودن}) = P(Y = 1)$$

- مقدار احتمال فوتی بودن تصادف می‌تواند تحت تأثیر عوامل مختلفی باشد (سن، سرعت برخورد، نوع خودرو و ...)
- پس می‌توان احتمال بالا را تابعی از مجموعه‌ای از متغیرهای توصیفی دانست.

- در بسیاری از موضوعات، متغیر وابسته ( $Y$ ) بیانگر دو حالت است (مدل‌های دو انتخابی)؛

➤  $Y=0$  عدم وقوع/انتخاب موضوع مورد نظر

➤  $Y=1$  وقوع/انتخاب موضوع مورد نظر



## □ مدل‌های دو انتخابی

○ برای مدل‌های دو انتخابی می‌توان مدل احتمال را به صورت زیر تعریف کرد؛

$$P(Y = 1|x_i) = F(x_i, \beta)$$

$$P(Y = 0|x_i) = 1 - F(x_i, \beta)$$

بردار متغیرهای موثر بر احتمال وقوع  $Y$ :  $x'_i = [1 \ x_{2i} \ \dots \ x_{Ki}]$

سوال: چه رابطه‌ای بین متغیرهای توصیفی و متغیر وابسته وجود دارد؟

پاسخ: بستگی به شکل تابع  $F(x_i, \beta)$  دارد (ممکن است خطی یا غیرخطی باشد).

در چارچوب رگرسیون، مدل  $Y_i$  را به صورت روبرو می‌نویسیم:

$$Y_i = E(Y|x_i) + u_i = F(x_i, \beta) + u_i$$

در واقع، تابع  $F(x_i, \beta)$  همان امیدریاضی شرطی  $Y$  است.

اگر برای  $F(x_i, \beta)$  یک معادله خطی تعریف کنیم، آنگاه معادله روبرو یک مدل رگرسیون خطی چندمتغیره خواهد بود.

به طور کلی برای  $F(x_i, \beta)$  توابع مختلفی معرفی می‌شود که هر کدام با یک نام مخصوص شناخته می‌شود.



□ به طور کلی برای  $F(x_i, \beta)$  توابع مختلفی معرفی می‌شود که هر کدام با یک نام مخصوص شناخته می‌شود؛

- مدل احتمال خطی
- مدل پروبیت
- مدل لوجیت





## □ مدل احتمال خطی (Linear Probability Model (LPM))؛

- در مدل احتمال خطی برای  $F(x_i, \beta)$  یک معادله خطی تعریف می‌شود:

$$F(x_i, \beta) = \mathbf{x}'_i \boldsymbol{\beta} = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$$

$$E(Y|x_i) = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$$

- امید ریاضی شرطی  $Y$ :

$$Y_i = E(Y|x_i) + u_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i$$

- مدل رگرسیون  $Y$ :

- این یک مدل احتمال خطی (LPM) است. مدلی ساده شبیه مدل رگرسیون معمولی و دارای مشکلات زیر؛

۱- ناهمسانی واریانس دارد (واریانس  $u_i$  وابسته به مقدار متغیرهای توصیفی است).

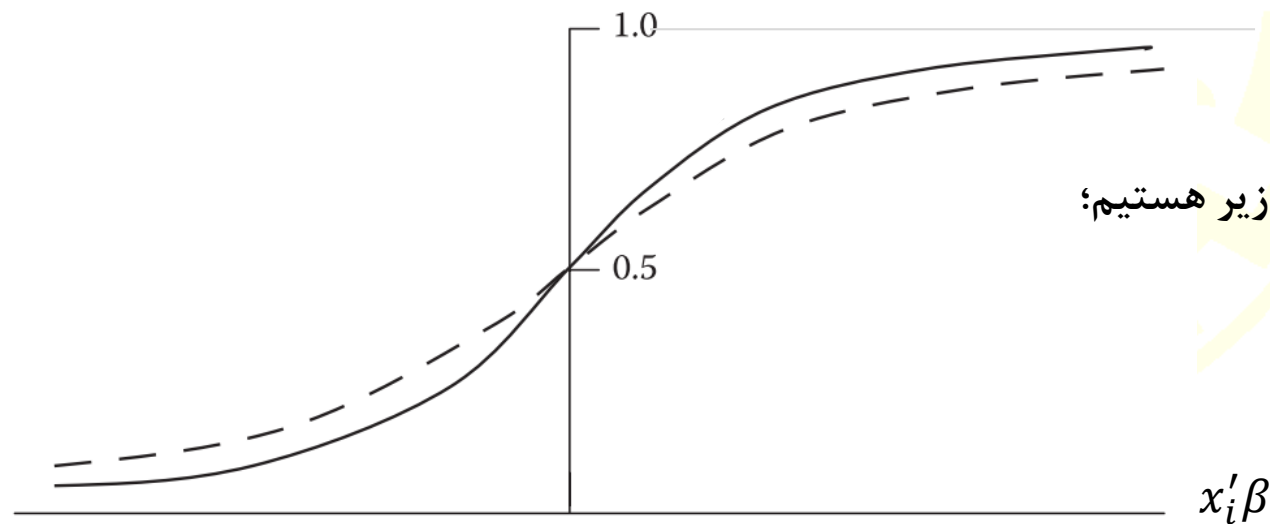
۲- چون  $F(x_i, \beta)$  بیانگر احتمال است، مقادیر آن همواره باید بین صفر و یک قرار بگیرد. در مسائل واقعی ممکن است  $\mathbf{x}'_i \boldsymbol{\beta}$  خارج از محدوده  $[0, 1]$  قرار بگیرد.



### نظریه مطلوبیت تصادفی □

- یکی از مشکلات مدل احتمال خطی (LPM) آن است که مقدار  $F(x_i, \beta)$  که بیانگر یک احتمال است، ممکن است گاهی اوقات مقادیر بیشتر از ۱ و یا کمتر از صفر را اختیار کند.
- برای حل این مشکل، از مدل‌های دیگری استفاده می‌کنیم که **قیدی را روی تابع احتمال  $F(x_i, \beta)$  اعمال** می‌کنند تا مقدار احتمال محدوده به بازه  $[0, 1]$  گردد.
- نمودار چنین تابعی به صورت زیر است.

$$F(x'_i \beta)$$



○ پس برای حل این مشکل، نیازمند موارد زیر هستیم؛

1. یک معادله رگرسیون،
2. یک تابع احتمال



## نظریه مطلوبیت تصادفی □

○ برای حل این مشکل، نیازمند موارد زیر هستیم؛

1. یک معادله رگرسیون،
2. یک تابع احتمال

بر مبنای نظریه مطلوبیت تصادفی می توان این دو موضوع را مرتبط کرده و مشکل مذکور را حل کرد.

مطلوبیت تصادفی از **نظریه انتخاب** استفاده می کند. نظریه انتخاب در حالت دو انتخابی:

فرض کنید هر فرد ۲ انتخاب دارد. او گزینه های خود را به گونه ای انتخاب می کند که مطلوبیتش ( $U_i$ ) حداکثر شود.

برای فرد  $i$ -ام

اگر  $U_{1i} \geq U_{0i}$  و یا  $U_{1i} - U_{0i} \geq 0$  باشد، آنگاه گزینه ۱ انتخاب می شود.

در ادامه، متغیر وابسته  $Y_i^*$  را به صورت روبرو تعریف می کنیم:

$$Y_i^* = U_{1i} - U_{0i}$$

اگر  $Y_i^* \geq 0$  باشد، فرد  $i$ -ام گزینه ۱ و اگر  $Y_i^* < 0$  باشد، فرد  $i$ -ام گزینه ۰ را انتخاب می کند.

$U_{1i}$  : مطلوبیت حاصل از انتخاب گزینه ۱

$U_{0i}$  : مطلوبیت حاصل از انتخاب گزینه ۰



## □ نظریه مطلوبیت تصادفی

می توان ادعا کرد که مطلوبیت، ترجیحات فرد را منعکس می کند، بنابراین، مقداری که  $Y_i^*$  اختیار می کند، می تواند تحت تأثیر مجموعه ای از عوامل و متغیرهای دیگر باشد.

یعنی اینکه کدام گزینه برای فرد  $i$  مطلوب تر است، به عوامل مختلفی بستگی دارد:

$$Y_i^* = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + u_i \quad ; \quad X_{1i} = 1$$

$$Y_i^* = \mathbf{x}'_i \boldsymbol{\beta} + u_i \quad ; \quad \mathbf{x}'_i = [1 \quad X_{2i} \quad \dots \quad X_{ki}] \quad ; \quad \boldsymbol{\beta}' = [\beta_1 \quad \beta_2 \quad \dots \quad \beta_k]$$

برای کاربردهای عملی باید مدل بالا را برآورد کنیم.

**چالش:**  $Y_i^*$  یک متغیر غیرقابل مشاهده است.

برای حل این چالش، متغیر قابل مشاهده  $Y_i$  را معرفی می کنیم (که عبارتست از انتخاب مشاهده شده برای فرد  $i$ -ام).



## □ نظریه مطلوبیت تصادفی

مثلا  $Y_i$  انتخاب مترو برای جابجایی است. اگر  $Y_i = 1$  باشد، فرد  $i$ -ام مترو را برای جابجایی انتخاب کرده است (و بالعکس). رابطه میان  $Y_i$  و  $Y_i^*$  به صورت زیر است؛

۱- اگر  $Y_i = 1$  باشد، بدین معناست که  $Y_i^* \geq 0$  است. (گزینه ۱ برای فرد  $i$ -ام مطلوب تر بوده است).

۲- اگر  $Y_i = 0$  باشد، بدین معناست که  $Y_i^* < 0$  است. (گزینه ۰ برای فرد  $i$ -ام مطلوب تر بوده است).

❖ بنابراین، یک مدل آماری داریم که دو جزء دارد؛

۱- یک معادله رگرسیون برای متغیر غیرقابل مشاهده  $Y_i^*$

۲- یک معادله برای مرتبط ساختن متغیر غیرقابل مشاهده  $Y_i^*$  و متغیر قابل مشاهده  $Y_i$  (تابع اتصال: Link function)

✓ برای مرتبط کردن  $Y_i$  و  $Y_i^*$  از مفهوم احتمال استفاده می کنیم.



## نظریه مطلوبیت تصادفی □

برای مرتبط کردن  $Y_i$  و  $Y_i^*$  از مفهوم احتمال استفاده می‌کنیم.

$$P(Y_i = 1) = P(Y_i^* \geq 0) = P(\mathbf{x}_i' \boldsymbol{\beta} + u_i \geq 0) = P(u_i \geq -\mathbf{x}_i' \boldsymbol{\beta})$$

پس نهایتاً احتمال انتخاب مترو برای جابجایی

معادل است با اینکه مترو مطلوبیت بیشتری داشته باشد ( $Y_i^* \geq 0$ )

معادل است با احتمال اینکه متغیر تصادفی  $u_i$  بزرگتر از  $-\mathbf{x}_i' \boldsymbol{\beta}$  باشد.

اکنون لازم است برای  $u_i$  یک تابع احتمال تعریف کنیم و خصوصیات آن را بررسی کنیم.

دو تابع رایج مورد استفاده برای تابع احتمال  $u_i$ : (تابع اتصال: Link function)

۱- تابع توزیع نرمال (مدل پروبیت)

۲- تابع توزیع لوجستیک (مدل لوجیت)

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس

لوجیت چندگانه

لوجیت شرطی

لوجیت آشیانه‌ای و ترکیبی

مقدمه

پروبیت دوگانه

لوجیت دوگانه (لوجستیک)

معیارهای نیکویی برآزش



در رابطه با مدل پروبیت، در اسلایدهای بعدی ۳ موضوع را دنبال می‌کنیم؛

1. ساختن تابع توزیع خطای مدل (تابعی برای مرتبط کردن  $Y_i$  و  $Y_i^*$ )
2. برآورد پارامترهای مدل
3. تفسیر مدل



## 1. ساختن تابع توزیع خطای مدل

فرض کنید  $u_j$  دارای توزیع نرمال باشد. برای هر متغیری مانند  $Z$  که تابع چگالی احتمال آن نرمال استاندارد باشد، می توان تابع توزیع یا تابع احتمال تجمعی را به صورت زیر معرفی کرد؛

$$P(Z \leq z) = \int_{-\infty}^z \varphi(z) dz = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz = \Phi(z)$$

$\varphi(z)$  تابع چگالی و  $\Phi(z)$  تابع توزیع  $Z$  است.

اگر  $u_j$  دارای توزیع نرمال باشد (تابع اتصال: نرمال)؛

$$\begin{aligned} P(Y_i = 1) &= P(Y_i^* \geq 0) = P(u_j \geq -\mathbf{x}'_i \boldsymbol{\beta}) = 1 - P(u_j < -\mathbf{x}'_i \boldsymbol{\beta}) = 1 - \Phi(-\mathbf{x}'_i \boldsymbol{\beta}) \\ &= \Phi(\mathbf{x}'_i \boldsymbol{\beta}) \end{aligned}$$

از سوی دیگر،

$$P(Y_i = 0) = P(Y_i^* < 0) = 1 - P(Y_i = 1) = 1 - \Phi(\mathbf{x}'_i \boldsymbol{\beta}) = \Phi(-\mathbf{x}'_i \boldsymbol{\beta})$$



## 2. برآورد پارامترهای مدل

برای برآورد پارامترهای مدل از روش MLE استفاده می‌کنیم. می‌دانیم متغیر تصادفی  $Y_i$  دارای توزیع برنولی است.

$Y_i$	0	1
$p_i$	$1-\Phi_i$	$\Phi_i$

که در آن،  $\Phi_i = \Phi(\mathbf{x}'_i \boldsymbol{\beta})$

این توزیع را به صورت زیر نیز می‌توان نشان داد؛

$$p_i = P(Y_i | \mathbf{x}_i) = \Phi_i^{Y_i} (1 - \Phi_i)^{1 - Y_i}, \quad Y_i = 0, 1$$

تابع درستنمایی؛

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n P(Y_i | \mathbf{x}_i) = \prod_{i=1}^n \Phi_i^{Y_i} (1 - \Phi_i)^{1 - Y_i}$$



## 2. برآورد پارامترهای مدل

تابع درستنمایی؛

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n P(Y_i | \mathbf{x}_i) = \prod_{i=1}^n \Phi_i^{Y_i} (1 - \Phi_i)^{1 - Y_i}$$

لگاریتم تابع درستنمایی را محاسبه کرده و نسبت به  $\boldsymbol{\beta}$  مشتق می‌گیریم (سپس برابر با صفر قرار داده و معادله را حل می‌کنیم).

بدین ترتیب، برآوردگر  $\hat{\boldsymbol{\beta}}_{ML}$  به دست می‌آید.



### 3. تفسیر نتایج مدل پروبیت

در هر مدل رگرسیون،  $\beta$  بیانگر اثرات نهایی متغیرهای توصیفی بر متغیر وابسته است. در اینجا، متغیر وابسته، بیانگر مطلوبیت است (که غیرقابل مشاهده است).

اما در مدل پروبیت، دیدیم که احتمال انتخاب گزینه ۱ برابر با  $P(Y_i = 1) = \Phi(\mathbf{x}'_i\boldsymbol{\beta})$  است. احتمال برآورد شده در مدل پروبیت (بر مبنای  $\hat{\beta}$  برآورد شده) به صورت زیر است؛

$$p_i = P(Y_i = 1|x_i) = \Phi(\mathbf{x}'_i\hat{\boldsymbol{\beta}}) = \Phi(\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki})$$



## 3. تفسیر نتایج مدل پروبیت

مثال؛ فرض کنید یک مدل تک متغیره داریم (و پارامترها به صورت  $-0.8$  و  $0.01$  برآورد شده است)؛

$Y_i$ : انتخاب حمل و نقل شخصی برای جابجایی ( $Y_i = 1$  به معنی انتخاب خودرو شخصی و بالعکس)

$$p_i = P(Y_i = 1|x_i) = \Phi(-0.8 + 0.01 X_i)$$

$X_i$ : سطح درآمد فرد  $i$ -ام

احتمال انتخاب خودرو شخصی به ازای سطوح مختلف درآمدی؛

$$P(Y_i = 1|x_i = 40) = \Phi(-0.8 + 0.01 (40)) = 0.345$$

سطح درآمدی  $X_i=40$

$$P(Y_i = 1|x_i = 80) = \Phi(-0.8 + 0.01 (80)) = 0.5$$

سطح درآمدی  $X_i=80$

$$P(Y_i = 1|x_i = 150) = \Phi(-0.8 + 0.01 (150)) = 0.758$$

سطح درآمدی  $X_i=150$

$$P(Y_i = 1|x_i = 200) = \Phi(-0.8 + 0.01 (200)) = 0.885$$

سطح درآمدی  $X_i=200$

پس با افزایش سطوح درآمدی، احتمال انتخاب خودرو شخصی افزایش می یابد.



### 3. تفسیر نتایج مدل پروبیت

برای بررسی اثرات نهایی، اثر تغییر در  $X_j$  ها را بر  $Y$  اندازه گیری می کنیم.

در این مدل ها که  $Y_i^*$  یک متغیر کیفی و غیرقابل مشاهده است،  $\beta$  اثر تغییرات  $X_j$  ها را بر  $Y_i$  اندازه گیری می کند. اگر  $\beta$  مثبت باشد، مطلوبیت انتخاب گزینه ۱ همراه با افزایش  $X_j$ ، افزایش می یابد.

اما سوال اینست: در واکنش به تغییر  $X_j$ ، احتمال انتخاب گزینه ۱ چقدر افزایش می یابد؟

اثر تغییر  $X_j$  بر احتمال انتخاب گزینه ۱ به صورت روبرو محاسبه می شود؛  
$$\frac{d P(Y_i = 1)}{dX_j} = \varphi(\mathbf{x}'_i \hat{\beta}) \hat{\beta}_j$$

مثلاً اثر تغییر در  $X_{ki}$  بر  $P(Y_i = 1)$  به صورت روبرو محاسبه می شود؛

$$\frac{d P(Y_i = 1)}{dX_{ki}} = \varphi(\mathbf{x}'_i \hat{\beta}) \hat{\beta}_k = \varphi(\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki}) \hat{\beta}_k$$

در مثال قبلی، اثر تغییر در درآمد بر احتمال انتخاب خودرو شخصی برای جابجایی:

$$\frac{d P(Y_i = 1)}{dX_i} = \varphi(-0.8 + 0.01 X_j) \times 0.01$$



### 3. تفسیر نتایج مدل پروبیت

در مثال قبلی، اثر تغییر در درآمد بر احتمال انتخاب خودرو شخصی برای جابجایی:

$$\frac{d P(Y_i = 1)}{dX_i} = \varphi(-0.8 + 0.01 X_i) \times 0.01$$

این بدان معناست که به ازای مقادیر مختلف  $X$ ، اثر نهایی  $X$  بر  $Y$  نیز تغییر می‌کند.

می‌توانیم اثر نهایی  $X$  بر  $Y$  را براساس مقدار متوسط  $X$  برآورد و گزارش کنیم.

$$\frac{d P(Y_i = 1)}{dX_{ki}} = \varphi(\hat{\beta}_1 + \hat{\beta}_2 \bar{X}_2 + \dots + \hat{\beta}_k \bar{X}_k) \hat{\beta}_k$$

در مثال قبلی، اگر مقدار متوسط درآمد برابر با ۱۵۰ باشد، احتمال انتخاب خودرو شخصی = 0.758 بود.

اکنون، به طور متوسط اثر تغییر در درآمد بر احتمال انتخاب خودرو شخصی برای جابجایی برابر است با:

$$\frac{d P(Y_i = 1)}{dX_i} = \varphi(-0.8 + 0.01 \bar{X}_i) \times 0.01 = \varphi(0.7) \times 0.01 = 0.3123 \times 0.01 = 0.031$$

پس برای یک فرد با درآمد متوسط، اگر درآمد یک واحد افزایش یابد، احتمال انتخاب خودرو شخصی ۳/۱٪ افزایش می‌یابد.



### 3. تفسیر نتایج مدل پروبیت

اگر درآمد فردی کمتر از متوسط باشد (مثلا ۴۰ باشد)، احتمال انتخاب خودرو شخصی توسط او با افزایش درآمد، چقدر است؟

$$\frac{d P(Y_i = 1)}{dX_i} = \varphi(-0.8 + 0.01 \times 40) \times 0.01 = 0.03683$$

حدود ۳.۷٪ افزایش می‌یابد.

پس با افزایش درآمد برای این فرد احتمال انتخاب خودرو شخصی بیشتر از متوسط افزایش می‌یابد (حدود ۳.۷٪).

**سوال:** بیشترین تاثیر درآمد بر احتمال انتخاب خودرو شخصی به ازای چه شرایطی رخ می‌دهد؟

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس

لوجیت چندگانه

لوجیت شرطی

لوجیت آشیانه‌ای و ترکیبی

مقدمه

پروبیت دوگانه

لوجیت دوگانه (لوجستیک)

معیارهای نیکویی برآزش



## 1. ساختن تابع توزیع خطای مدل

به جای تابع توزیع نرمال استاندارد (مدل پروبیت)، توابع توزیع دیگری نیز می‌توان استفاده کرد. یکی از رایج‌ترین توابع، تابع توزیع لوجستیک است.

برای هر متغیر تصادفی  $Z$  تابع توزیع لاجستیک به صورت زیر است،

$$P(Z \leq z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$

احتمال آنکه  $Y_i = 1$  باشد، برابر است با؛

$$p_i = P(Y_i = 1 | \mathbf{x}'_i) = G(\mathbf{x}'_i \boldsymbol{\beta}) = \frac{1}{1 + e^{-\mathbf{x}'_i \boldsymbol{\beta}}} = \frac{e^{\mathbf{x}'_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}'_i \boldsymbol{\beta}}}$$

احتمال آنکه  $Y_i = 0$  باشد، برابر است با؛

$$P(Y_i = 0 | \mathbf{x}'_i) = 1 - P(Y_i = 1 | \mathbf{x}'_i) = \frac{1}{1 + e^{\mathbf{x}'_i \boldsymbol{\beta}}}$$



## 2. برآورد پارامترهای مدل

در واقع شکل مدل رگرسیون لوجستیک به صورت روبروست؛

$$Y_i = \text{logit}(p_i) = \ln\left(\frac{p_i}{1-p_i}\right) = \mathbf{x}'_i \boldsymbol{\beta} = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$$

$$p_i = P(Y_i = 1 | \mathbf{x}'_i) = G(\mathbf{x}'_i \boldsymbol{\beta}) = \frac{1}{1 + e^{-\mathbf{x}'_i \boldsymbol{\beta}}} = \frac{e^{\mathbf{x}'_i \boldsymbol{\beta}}}{1 + e^{\mathbf{x}'_i \boldsymbol{\beta}}}$$

تابع اتصال لوجیت:  $\ln\left(\frac{p_i}{1-p_i}\right) = \mathbf{x}'_i \boldsymbol{\beta}$

برای برآورد پارامترهای مدل از روش MLE استفاده می‌کنیم.

می‌دانیم متغیر تصادفی  $Y_i$  دارای توزیع برنولی است.

تابع درستنمایی؛

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n G_i^{Y_i} (1 - G_i)^{1-Y_i} \quad ; \quad G_i = G(\mathbf{x}'_i \boldsymbol{\beta})$$

برای به دست آوردن برآوردگرهای MLE،

لگاریتم تابع درستنمایی را به دست آورده، از آن مشتق گرفته و برابر با صفر قرار می‌دهیم.



## 3. تفسیر نتایج مدل لوجیت

احتمال برآورد شده در مدل لوجیت؛

$$p_i = P(Y_i = 1|x_i) = G(\mathbf{x}'_i\boldsymbol{\beta}) = G(\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki})$$

❖ به عبارت دیگر، احتمال اینکه  $Y_i = 1$  باشد، بستگی به پارامترهای برآورد شده و مقدار متغیرهای توصیفی دارد.

❖ پارامترهای مدل ثابت هستند، اما مقدار متغیرهای توصیفی مدل تغییر می کند.



## 3. تفسیر نتایج مدل لوجیت

مثال؛ فرض کنید یک مدل تک متغیره داریم (و پارامترها به صورت  $-0.8$  و  $0.01$  برآورد شده است)؛

$$P(Y_i = 1|x_i) = G(-0.8 + 0.01 X_i)$$

$Y_i$ : انتخاب حمل و نقل شخصی برای جابجایی

$X_i$ : سطح درآمد فرد  $i$ -ام

$$P(Y_i = 1|x_i = 40) = G(-0.8 + 0.01 (40)) = \frac{1}{1 + e^{-(-0.4)}} = 0.401$$

سطح درآمدی  $X_i=40$

$$P(Y_i = 1|x_i = 80) = G(-0.8 + 0.01 (80)) = \frac{1}{1 + e^{-0}} = 0.5$$

سطح درآمدی  $X_i=80$

$$P(Y_i = 1|x_i = 150) = G(-0.8 + 0.01 (150)) = \frac{1}{1 + e^{-0.7}} = 0.668$$

سطح درآمدی  $X_i=150$

$$P(Y_i = 1|x_i = 200) = G(-0.8 + 0.01 (200)) = \frac{1}{1 + e^{-1.2}} = 0.7685$$

سطح درآمدی  $X_i=200$

○ بنابراین، با افزایش درآمد، احتمال انتخاب حمل و نقل شخصی بیشتر می شود.



# لو جیت دو گانه (لو جستیک)

### 3. تفسیر نتایج مدل لو جیت

اثر تغییر  $X_i$  بر احتمال انتخاب گزینه ۱ به صورت روبرو محاسبه می شود؛  
$$\frac{d P(Y_i = 1)}{d X_i} = g(\mathbf{x}'_i \hat{\beta}) \hat{\beta}$$

مثلاً اثر تغییر در  $X_{ki}$  بر  $P(Y_i = 1)$  به صورت روبرو محاسبه می شود؛

$$\frac{d P(Y_i = 1)}{d X_{ki}} = g(\mathbf{x}'_i \hat{\beta}) \hat{\beta} = g(\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki}) \hat{\beta}_k$$

پس مقدار اثر حاشیه ای به مقدار  $X$  بستگی دارد.

$$g(\mathbf{x}'_i \beta) = \frac{e^{\mathbf{x}'_i \beta}}{(1 + e^{\mathbf{x}'_i \beta})^2}$$

در مثال قبلی، اثر تغییر در درآمد بر احتمال انتخاب خودرو شخصی برای جابجایی:

$$\frac{d P(Y_i = 1)}{d X_i} = g(-0.8 + 0.01 X_i) \times 0.01$$



## 3. تفسیر نتایج مدل لوجیت

معمولا اثر نهایی  $X$  بر  $Y$  را براساس مقدار متوسط  $X$  برآورد و گزارش می کنیم (EAM) (“Effect At the Means”).

$$\frac{d P(Y_i = 1)}{dX_{ki}} = g(\hat{\beta}_1 + \hat{\beta}_2 \bar{X}_2 + \dots + \hat{\beta}_k \bar{X}_k) \hat{\beta}_k \quad ; \quad k = 1, 2, \dots, K$$

در مثال قبلی، اگر مقدار متوسط درآمد برابر با ۱۵۰ باشد، احتمال انتخاب خودرو شخصی = 0.668 بود.

اکنون، **به طور متوسط** اثر تغییر درآمد بر احتمال انتخاب خودرو شخصی برای جابجایی برابر است با:

$$\frac{d P(Y_i = 1)}{dX_i} = g(-0.8 + 0.01 \bar{X}_i) \times 0.01 = g(0.7) \times 0.01 = \frac{e^{0.7}}{(1 + e^{0.7})^2} \times 0.01 = 0.00222$$

پس برای یک فرد با درآمد متوسط، اگر درآمد یک واحد افزایش یابد، احتمال انتخاب خودرو شخصی ۰.۲۲٪ افزایش می یابد.



## 3. تفسیر نتایج مدل لوجیت

نوع جدیدتری از برآورد اثر حاشیه‌ای  $X$  بر  $Y$  با نام متوسط اثر حاشیه‌ای (Average Marginal Effect (AME)) شناخته می‌شود.

تفاوت این روش با روش قبلی: در این روش، اثرات حاشیه‌ای برای تک تک مشاهدات نمونه محاسبه گردیده و سپس، میانگین این اثرات، به عنوان «متوسط اثر حاشیه‌ای» گزارش می‌شود.

$$AME_k = \frac{1}{N} \sum_{i=1}^N (\hat{\beta}_k \times \frac{e^{x_i' \hat{\beta}}}{(1 + e^{x_i' \hat{\beta}})^2})$$

برای یک متغیر پیوسته  $X_k$  متوسط اثر حاشیه‌ای برابر است با:

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس	مقدمه
لوجیت چندگانه	پروبیت دوگانه
لوجیت شرطی	لوجیت دوگانه (لوجستیک)
لوجیت آشیانه‌ای و ترکیبی	معیارهای نیکویی برآزش



## ❖ نسبت درستنمایی (Likelihood Ratio)

برای برآورد پارامترهای مدل‌های پروبیت و لوجیت از روش MLE استفاده می‌شود. بنابراین، با استفاده از نسبت درستنمایی می‌توان هر محدودیتی را مورد آزمون قرار داد. منظور از محدودیت (restriction): حذف یک یا چند متغیر توصیفی از مدل نسبت درستنمایی برای مقایسه دو مدل رگرسیون غیرمقید و مقید استفاده می‌شود؛

$$LR = 2 [\ln L_{UR} - \ln L_0]$$

$L_{UR}$ : مقدار تابع درستنمایی غیرمقید (مدل ساخته شده)

$L_0$ : مقدار تابع درستنمایی مدل ساده (فقط شامل عرض از مبدأ و سایر متغیرها کنار گذاشته شده است)

نسبت درستنمایی از توزیع  $\chi^2_m$  پیروی می‌کند که درجات آزادی =  $m$  تعداد محدودیت‌ها (تعداد متغیرهای کنارگذاشته شده) را نشان می‌دهد.



❖ شاخص نسبت درستنمایی (Likelihood ratio index) یا  $\rho^2$  مک فادن

$$\rho_M^2 = LRI = 1 - \frac{\ln L_{UR}}{\ln L_0}$$

$L_{UR}$ : مقدار تابع درستنمایی غیرمقید (مدل ساخته شده)

$L_0$ : مقدار تابع درستنمایی مدل ساده (فقط شامل عرض از مبدأ)

شاخص  $\rho^2$  مک فادن **میزان بهبود در نسبت لگاریتم درستنمایی** را نسبت به مدل ساده، اندازه گیری می کند.

$LRI$  بین صفر و یک است.

○  $LRI = 0$ : اگر  $L_{UR} = L_0$  یعنی قدرت توضیح مدل ساخته شده با یک مدل صفر، یکسان است.

○  $LRI = 1$ : به معنای **برازش کامل** است. (یعنی به ازای  $Y_i = 1$  همواره  $G(x_i'\beta) = 1$  و به ازای  $Y_i = 0$  همواره

$G(x_i'\beta) = 0$  باشد) بدین ترتیب لگاریتم تابع درستنمایی برابر با 0 می شود و  $LRI = 1$  - این یک **حالت خاص**

است.



## ❖ نکته:

1. برای مدل پروبیت و لوجیت دوگانه، مقادیر **log-Likelihood همواره منفی (یا نهایتاً صفر)** است.
2. مقادیر **AIC و BIC** نه همیشه **مثبت** است و نه همیشه **منفی**.
3. در واقع، هنگامی که لگاریتم درست‌نمایی مقدار منفی دارد و جمله جریمه مقدار کوچکی دارد، مقادیر **AIC و BIC** منفی به دست می‌آید.
4. اما همیشه مدلی که مقدار **AIC و BIC** **کمتری** دارد، **بهتر** است



## ❖ درصد پیش‌بینی صحیح (Percent correctly predicted)

به ازای هر  $i$ ، احتمال تخمینی برای وضعیت  $Y_i = 1$  را حساب می‌کنیم. این احتمال برابر است با  $\Phi(\mathbf{x}'_i \hat{\beta})$  یا  $G(\mathbf{x}'_i \hat{\beta})$ . در حالت کلی به صورت  $F(\mathbf{x}'_i \hat{\beta})$  نمایش می‌دهیم.

برای مقادیر  $F(\mathbf{x}'_i \hat{\beta}) > 0.5$  پیش‌بینی می‌شود که  $Y_i = 1$  است.

برای مقادیر  $F(\mathbf{x}'_i \hat{\beta}) \leq 0.5$  پیش‌بینی می‌شود که  $Y_i = 0$  است.

درصد مواردی که  $Y_i$  پیش‌بینی شده با  $Y_i$  مشاهده شده مطابقت دارد = «درصد پیش‌بینی صحیح»



## ❖ درصد پیش‌بینی صحیح (Percent correctly predicted)

مثال؛ در یک نمونه ۲۰۰ تایی، نتایج به دست آمده از مدل لجوجیت به صورت جدول زیر است. مطلوبست محاسبه درصد پیش‌بینی صحیح؟

		پیش‌بینی شده		
		$Y_i = 0$	$Y_i = 1$	جمع
مشاهده شده	$Y_i = 0$	۵۰	۳۰	۸۰
	$Y_i = 1$	۲۰	۱۰۰	۱۲۰
	جمع	۷۰	۱۳۰	۲۰۰

✓ وقتی  $Y_i = 0$  در ۵۰ مورد پیش‌بینی با مشاهدات مطابقت دارد.

✓ وقتی  $Y_i = 1$  در ۱۰۰ مورد پیش‌بینی با مشاهدات مطابقت دارد.

✓ در مجموع ۱۵۰ مورد از ۲۰۰ تا مطابقت داشته. پس **درصد پیش‌بینی صحیح = ۰.۷۵**



## ❖ شاخص انحراف (Deviance)

در این روش دو مدل را با هم مقایسه می‌کنیم؛

۱- **مدل کامل (Saturated model)**: مدلی که تعداد پارامترهای آن برابر با تعداد مشاهدات است.

در این حالت، مدل به صورت کامل برازش می‌شود.

مثلا اگر یک مدل رگرسیون با یک متغیر توصیفی داشته باشیم، و فقط ۲ مشاهده داشته باشیم، آنگاه هیچ خطایی وجود ندارد، و پارامترهای مدل به صورت دقیق به دست می‌آید.

در حالت کلی، در مدل کامل،  $n$  پارامتر و  $n$  مشاهده داریم.

۲- **مدل موردنظر (Proposed model)**: مدل ساخته شده که تعداد پارامترهای آن از تعداد مشاهدات **کمتر** است.

در حالت کلی، به تعداد  $p$  پارامتر و  $n$  مشاهده داریم ( $p < n$ ).

$$E(Y_i|x_i) = \beta_1 + \beta_2 X_{2i} + \dots + \beta_p X_{pi}$$



## ❖ شاخص انحراف (Deviance)

برای مقایسه این دو مدل از تابع درستنمایی استفاده می‌کنیم.

$$D = 2[\ln L(\mu_i) - \ln L(\beta)]$$

$L(\mu_i)$  بیشترین مقدار برای تابع درستنمایی است.

هرچه  $D$  کمتر ← نشانه‌ای از برازش بهتر

کوچک‌تر بودن مقدار  $D$  بدین معناست که مدل رگرسیون موردنظر، **برازش بهتری بر روی داده‌ها داشته است.** زیرا انحراف آن از مدل کامل، کمتر است.

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس	مقدمه
لوجیت چندگانه	پروبیت دوگانه
لوجیت شرطی	لوجیت دوگانه (لوجستیک)
لوجیت آشیانه‌ای و ترکیبی	معیارهای نیکویی برآزش



❖ آنچه در مدل رگرسیون لوجستیک بر آورد می شود؟

- ریسک؟
- نسبت ریسک؟
- شانس؟
- نسبت شانس؟





## ❖ ریسک (Risk)

- ریسک از مبانی مدل سازی داده‌های شمارشی است.
- ریسک عبارت است از مواجهه با احتمال یک پیامد.
- مثلا ریسک تصادف: احتمال وقوع تصادف
- ریسک به معنای تجربه یک پیامد معین توسط یک فرد در شرایط مشخص است.
- ریسک بیانگر احتمال آن است که فرد واقعا آن پیامد را تجربه کند.
- از طرف دیگر، این احتمال، بیانگر نسبت افرادی است که این پیامد را تجربه می‌کنند.



## ❖ نسبت ریسک (Risk Ratio (RR))

- عامل ریسک: شرط یا شرایطی که ریسک تحت آن اندازه‌گیری می‌شود.
- ریسک: معیاری از رابطه میان پیامد و عامل (عوامل) ریسک خاص.
- ریسک نسبی (نسبت ریسک (RR)): نسبت احتمال وقوع پیامد برای افرادی که یک عامل ریسک معین را دارند در مقایسه با احتمال وقوع پیامد برای آن‌هایی که آن عامل ریسک را ندارند.
- پس ریسک نسبی، نسبت دو نسبت است.
- وقتی که با داده‌های شمارشی کار می‌کنیم، به آن «نسبت نرخ وقوع (Incidence Rate Ratio (IRR))» می‌گویند.



## ❖ نسبت ریسک (RR)

○ مثال؛

(X)				
جمع	۱	۰		
A+B	B	A	۰	(Y)
C+D	D	C	۱	
n	B+D	A+C	جمع	

ریسک  $Y = 1$  به ازای  $X = 1$  برابر است با؛

$$\frac{D}{B + D}$$

ریسک  $Y = 1$  به ازای  $X = 0$  برابر است با؛

$$\frac{C}{A + C}$$

نسبت ریسک (ریسک نسبی) برای  $Y = 1$  به ازای  $X = 1$  در مقایسه با  $X = 0$  برابر است با؛

$$IRR = \frac{\frac{D}{B + D}}{\frac{C}{A + C}} = \frac{DA + CD}{CB + CD}$$



## ❖ نسبت ریسک (RR)

○ مثال ریسک فوت/مصدومیت در تصادف؛

ریسک فوت/مصدومیت برای بزرگسالان:

$$\frac{445}{1200} = 0.371$$

ریسک فوت/مصدومیت برای خردسالان

$$\frac{55}{100} = 0.55$$

نسبت ریسک (ریسک نسبی) فوت/مصدومیت بزرگسالان در مقایسه با خردسالان؛

$$\frac{0.371}{0.55} = 0.6742$$

❖ تفسیر: احتمال فوت/مصدومیت برای بزرگسالان نسبت به

خردسالان ۳۲.۶٪ کمتر است.

سن (X)

جمع	بزرگسال: ۱	خردسال: ۰	
۸۰۰	۷۵۵	۴۵	بدون مصدومیت: ۰
۵۰۰	۴۴۵	۵۵	مصدوم/فوت شده: ۱
۱۳۰۰	۱۲۰۰	۱۰۰	جمع

وضعیت  
مصدومیت  
(Y)

(X)

جمع	۱	۰	
A+B	B	A	۰
C+D	D	C	۱
n	B+D	A+C	جمع

(Y)



## ❖ فاصله اطمینان برای نسبت ریسک

محاسبه فاصله اطمینان نیازمند خطای استاندارد می باشد.

توجه: ابتدا محاسبات را برای لگاریتم نسبت ریسک انجام می دهیم و سپس آن را به نسبت ریسک تبدیل می کنیم

سن (X)			
جمع	بزرگسال: ۱	خردسال: ۰	
۸۰۰	۷۵۵	۴۵	بدون مصدومیت: ۰
۵۰۰	۴۴۵	۵۵	مصدوم/فوت شده: ۱
۱۳۰۰	۱۲۰۰	۱۰۰	جمع

وضعیت  
مصدومیت (Y)

در جدول ۲\*۲ خطای استاندارد لگاریتم نسبت ریسک را می توان به صورت زیر حساب کرد:

$$SE(\ln(RR)) = \sqrt{\frac{1}{D} - \frac{1}{B+D} + \frac{1}{C} - \frac{1}{A+C}} = \sqrt{\frac{1}{445} - \frac{1}{1200} + \frac{1}{55} - \frac{1}{100}} = 0.098$$

قبلا نسبت ریسک را برابر با 0.6742 به دست آوردیم؛ لگاریتم نسبت ریسک:  $\ln(0.6742) = -0.3942$

فاصله اطمینان با فرض توزیع نرمال برای  $\ln(RR)$ :  $\ln(RR) \pm Z_{\alpha/2} SE(\ln(RR))$

فاصله اطمینان ۹۵ درصدی برای لگاریتم نسبت ریسک:  $-0.3942 \pm 1.96 * 0.098 \rightarrow (-0.5863, -0.2021)$

فاصله اطمینان ۹۵ درصدی برای نسبت ریسک (RR):  $(e^{-0.5863}, e^{-0.2021}) = (0.556, 0.817)$

سن (X)			
جمع	بزرگسال: ۱	خردسال: ۰	
۸۰۰	۷۵۵	۴۵	بدون مصدومیت: ۰
۵۰۰	۴۴۵	۵۵	مصدوم/فوت شده: ۱
۱۳۰۰	۱۲۰۰	۱۰۰	جمع

## ❖ تفاضل ریسک

تفاضل ریسک بیانگر کاهش مطلق ریسک است.

این تفاضل را به عنوان «معیار کاهش ریسک» در نظر می‌گیریم.

ریسک فوت/مصدومیت برای بزرگسالان:

$$\frac{445}{1200} = 0.371$$

ریسک فوت/مصدومیت برای خردسالان

$$\frac{55}{100} = 0.55$$

تفاضل ریسک (اختلاف ریسک برای افرادی که در معرض خطر هستند با افرادی که در معرض نیستند):

$$risk(exposed) - risk(unexposed) = 0.371 - 0.55 = -0.179$$

❖ تفسیر: نرخ فوت/مصدومیت بزرگسالان (به صورت مطلق) حدود ۱۸٪ کمتر از خردسالان است.



## ❖ نسبت ریسک (RR)

○ مثال ریسک فوت/مصدومیت در تصادف؛

اکنون یک **متغیر سه سطحی** را در نظر بگیرید.

جمع	سن (X)			وضعیت مصدومیت (Y)
	بزرگسال: ۲	میانسال: ۱	خردسال: ۰	
۸۰۰	۵۲۰	۱۶۰	۱۲۰	بدون مصدومیت: 0
۵۰۰	۱۸۰	۱۲۰	۲۰۰	مصدوم/فوت شده: 1
۱۳۰۰	۷۰۰	۲۸۰	۳۲۰	جمع

- در اینجا باید یک ستون را به عنوان گروه مرجع انتخاب کنیم.
- ستون ۳: **بزرگسال** را به عنوان **مرجع** انتخاب می‌کنیم.

$$\frac{\frac{120}{280}}{\frac{180}{700}} = 1.6667$$

نسبت ریسک فوت/مصدومیت برای گروه ۲ (میانسال)

$$\frac{\frac{200}{320}}{\frac{180}{700}} = 2.4305$$

نسبت ریسک فوت/مصدومیت برای گروه ۱ (خردسال)

مقدار ۲.۴۳۰۵ یعنی ریسک فوت/مصدومیت خردسالان نسبت به بزرگسالان (گروه مرجع) ۲.۴۳ برابر یا ۱۴۳٪ بیشتر است.



## ❖ نسبت شانس (بخت) – Odds Ratio (OR)

شانس وقوع یک واقعه برابر است با احتمال وقوع آن تقسیم بر احتمال عدم وقوع آن؛  
شانس وقوع =  $\frac{p}{1-p}$

$P$  احتمال وقوع است.

توجه: در مدل رگرسیون لجستیک در واقع لگاریتم بخت (log odds) را مدل می‌کنیم:

$$Y_i = \text{logit}(p_i) = \ln\left(\frac{p_i}{1-p_i}\right) = \ln(\text{Odds}_i) = \mathbf{x}'_i \boldsymbol{\beta} = \beta_1 + \beta_2 X_{2i} + \dots + \beta_k X_{ki}$$

وقتی شانس وقوع را در هر سطحی از عامل ریسک ( $X$ ) مقایسه می‌کنیم، نسبت بین شانس وقوع  $X = 1$  با شانس

عدم وقوع  $X = 0$  معروف به «نسبت شانس **Odds Ratio (OR)**» است؛

$$OR = \frac{\frac{p_1}{1-p_1}}{\frac{p_0}{1-p_0}} = \frac{p_1(1-p_0)}{p_0(1-p_1)}$$



## ❖ نسبت شانس (بخت)

(X)				
جمع	۱	۰		
A+B	B	A	۰	(Y)
C+D	D	C	۱	
n	B+D	A+C	جمع	

شانس وقوع (Y = 1) برای X = 1 برابر است با؛

$$\frac{p_1}{1 - p_1} = \frac{\frac{D}{B + D}}{\frac{B}{B + D}} = \frac{D}{B}$$

شانس وقوع (Y = 1) برای X = 0 برابر است با؛

$$\frac{p_0}{1 - p_0} = \frac{\frac{C}{A + C}}{\frac{A}{A + C}} = \frac{C}{A}$$

$$OR = \frac{\frac{D}{B}}{\frac{C}{A}} = \frac{DA}{CB}$$

**نسبت شانس** وقوع (Y = 1) برای X = 1 در مقایسه با X = 0 برابر است با؛



## ❖ نسبت شانس (بخت)

(X)			
جمع	۱	۰	
A+B	B	A	۰
C+D	D	C	۱
n	B+D	A+C	جمع

نسبت شانس وقوع (Y = 1) برای X = 1 در مقایسه با X = 0

$$OR = \frac{\frac{D}{B}}{\frac{C}{A}} = \frac{DA}{CB}$$

**صورت کسر:** حاصلضرب (تعداد مواردی که X = 1 و Y = 1) \* (تعداد مواردی که X = 0 و Y = 0)

**مخرج کسر:** حاصلضرب (تعداد مواردی که X = 0 و Y = 1) \* (تعداد مواردی که X = 1 و Y = 0)

$$RR = \frac{DA + CD}{CB + CD} = \text{نسبت ریسک}$$

صورت و مخرج  
تقسیم بر CB

$$RR = \frac{\frac{DA}{CB} + \frac{CD}{CB}}{\frac{CB}{CB} + \frac{CD}{CB}} = \frac{OR + \frac{D}{B}}{1 + \frac{D}{B}}$$



## ❖ نسبت شانس (بخت)

(X)				
جمع	۱	۰		
A+B	B	A	۰	(Y)
C+D	D	C	۱	
n	B+D	A+C	جمع	

نسبت ریسک =  $RR = \frac{DA + CD}{CB + CD}$

صورت و مخرج  
تقسیم بر CB

$$RR = \frac{\frac{DA}{CB} + \frac{CD}{CB}}{\frac{CB}{CB} + \frac{CD}{CB}} = \frac{OR + \frac{D}{B}}{1 + \frac{D}{B}}$$

پس تفاوت نسبت ریسک و نسبت شانس، به مقدار  $\frac{D}{B}$  بستگی دارد.

◀ نسبت  $\frac{D}{B}$  شانس رویداد موردنظر ( $Y=1$ ) در گروه دارای مواجهه ( $X=1$ ) را بیان می کند. پس؛

1. هنگامی که رویداد موردنظر در گروه دارای مواجهه **نادر** باشد؛  $\frac{D}{B}$  کوچک (نزدیک به صفر) است. بدین ترتیب،  $RR \cong OR$ . در این حالت **می توان** نسبت شانس را مانند نسبت ریسک تفسیر نمود.

2. هنگامی که رویداد موردنظر در گروه دارای مواجهه **رایج** باشد؛ مقدار  $\frac{D}{B}$  بزرگ است. بدین ترتیب،  $RR$  کوچکتر از  $OR$  می شود (در این حالت **نمی توان** نسبت شانس را مانند نسبت ریسک تفسیر نمود).



## ❖ نسبت شانس

مثال ۱: شانس فوت/مصدومیت در تصادف؛

سن (X)			
جمع	بزرگسال:	خردسال:	
۸۰۰	۷۵۵	۴۵	بدون مصدومیت: 0
۵۰۰	۴۴۵	۵۵	مصدوم/فوت شده: 1
۱۳۰۰	۱۲۰۰	۱۰۰	جمع

وضعیت  
مصدومیت  
(Y)

شانس فوت/مصدومیت برای بزرگسالان:  $\frac{445}{775} = 0.5894$

شانس فوت/مصدومیت برای خردسالان  $\frac{55}{45} = 1.222$

نسبت ریسک: ریسک فوت/مصدومیت خردسالان نسبت به ریسک فوت/مصدومیت بزرگسالان؛

$$\frac{1}{0.6742} = 1.483$$

نسبت شانس فوت/مصدومیت بزرگسالان در مقایسه با خردسالان؛

$$OR = \frac{0.5894}{1.222} = 0.48224$$

نسبت شانس فوت/مصدومیت خردسالان در مقایسه با بزرگسالان؛  $OR = \frac{1.222}{0.5894} = \frac{1}{0.48224} = 2.07$

تفسیر: شانس فوت/مصدومیت خردسالان بیش از ۲ برابر بالاتر از شانس فوت/مصدومیت بزرگسالان است.

توجه: در این حالت چون هر دو ریسک، مقدار بالایی دارد، اگر به دنبال ریسک نسبی هستیم، نباید از مدل رگرسیون لجستیک و OR استفاده کرد، چون در برآورد ریسک موردنظر اغراق می کند (باید RR را گزارش کرد).



## ❖ نسبت شانس

مثال ۲: شانس فوت/مصدومیت در تصادف؛

شانس فوت/مصدومیت برای بزرگسالان:

$$\frac{6}{994} = 0.00604$$

شانس فوت/مصدومیت برای خردسالان

$$\frac{8}{992} = 0.00806$$

نسبت شانس فوت/مصدومیت خردسالان در مقایسه با بزرگسالان؛

$$OR = \frac{0.00806}{0.00604} = 1.336$$

همانطور که ملاحظه می‌شود، برای این حالت (که رویداد فوت/مصدومیت نادر است)، مقدار OR بسیار نزدیک به مقدار RR است. بنابراین، در این حالت می‌توان مقدار OR را همانند RR تفسیر کرد.

سن (X)			
جمع	بزرگسال:	خردسال:	
۱۹۸۶	۹۹۴	۹۹۲	بدون مصدومیت: 0
۱۴	۶	۸	مصدوم/فوت شده: 1
۲۰۰۰	۱۰۰۰	۱۰۰۰	جمع

وضعیت  
مصدومیت  
(Y)

نسبت ریسک: ریسک فوت/مصدومیت خردسالان نسبت به ریسک فوت/مصدومیت بزرگسالان؛

$$\frac{0.008}{0.006} = 1.333$$

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس	مقدمه
لوجیت چندگانه	پروبیت دوگانه
لوجیت شرطی	لوجیت دوگانه (لوجستیک)
لوجیت آشیانه‌ای و ترکیبی	معیارهای نیکویی برآزش



❖ شرایطی که پدیده مورد بررسی، بیش از دو حالت دارد. مثلاً؛

○ مسیر انتخاب شده از بین ۳ مسیر آلترناتیو،

○ نوع تصادف (صرفاً خسارتی، جرحی خفیف، جرحی شدید، فوتی)

○ شیوه حمل‌ونقل انتخابی (تاکسی، اتوبوس، مترو، خودرو شخصی)

○ ... ؟

○ در این شرایط یکی از سطوح ( $Y=0$ ) را به عنوان مرجع در نظر گرفته، و بقیه سطوح را با آن مقایسه می‌کنیم.

○ از مدل لوجیت چندگانه برای بررسی این پدیده‌ها استفاده می‌کنیم.



## ❖ مدل لوجیت سه گانه

○ فرض کنید  $p$  متغیر توصیفی داریم.

$$\mathbf{x}'_i = [ X_{1i} \quad X_{2i} \quad \dots \quad X_{pi} ] \quad ; \quad X_{1i} = 1$$

تابع لوجیت برای  $Y=1$  را به صورت زیر تعریف می‌کنیم.

$$g_1(\mathbf{x}'_i) = \ln \frac{P(Y = 1 | \mathbf{x}'_i)}{P(Y = 0 | \mathbf{x}'_i)} = \beta_{11} + \beta_{12}X_{2i} + \dots + \beta_{1p}X_{pi} = \mathbf{x}'_i \boldsymbol{\beta}_1$$

برای  $Y=2$  را به صورت زیر تعریف می‌کنیم.

$$g_2(\mathbf{x}'_i) = \ln \frac{P(Y = 2 | \mathbf{x}'_i)}{P(Y = 0 | \mathbf{x}'_i)} = \beta_{21} + \beta_{22}X_{2i} + \dots + \beta_{2p}X_{pi} = \mathbf{x}'_i \boldsymbol{\beta}_2$$



## ❖ مدل لوجیت سه گانه

$$P(Y = 0|\mathbf{x}'_i) = \frac{1}{1 + e^{g_1(\mathbf{x}'_i)} + e^{g_2(\mathbf{x}'_i)}}$$

○ احتمال‌های شرطی برای هر پیامد عبارت است از؛

$$P(Y = 1|\mathbf{x}'_i) = \frac{e^{g_1(\mathbf{x}'_i)}}{1 + e^{g_1(\mathbf{x}'_i)} + e^{g_2(\mathbf{x}'_i)}}$$

$$P(Y = 2|\mathbf{x}'_i) = \frac{e^{g_2(\mathbf{x}'_i)}}{1 + e^{g_1(\mathbf{x}'_i)} + e^{g_2(\mathbf{x}'_i)}}$$

برای سادگی و اختصار، از نام‌گذاری‌های زیر استفاده می‌کنیم؛

$$\pi_j(\mathbf{x}'_i) = P(Y = j|\mathbf{x}'_i) = \frac{e^{g_j(\mathbf{x}'_i)}}{\sum_{k=0}^2 e^{g_k(\mathbf{x}'_i)}} \quad ; \quad j = 0,1,2 \quad ; \quad g_0(\mathbf{x}'_i) = 0$$

بنابراین، هر احتمال تابعی از بردار  $\boldsymbol{\beta}' = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)$  است؛ که تعداد ضرایب برابر با  $2p$  است.

ضرایب این مدل با روش حداکثر درست‌نمایی برآورد می‌شود.



# مدل‌های لوجیت – لوجیت چندگانه

## ❖ مدل لوجیت سه گانه

○ برای تشکیل تابع درست‌نمایی، سه متغیر دوتایی (باینری) می‌سازیم که هر یک شامل 0 و 1 می‌شود.

○ توجه: این متغیرها را فقط برای تشکیل تابع درست‌نمایی استفاده می‌کنیم و در تفسیر مدل لوجیت کاربردی ندارد.

○  $Y = 0$                        $Y_0 = 1, Y_1 = 0, Y_2 = 0$                       متغیرها به صورت روبرو تعریف می‌شوند؛

$Y = 1$                        $Y_0 = 0, Y_1 = 1, Y_2 = 0$

$Y = 2$                        $Y_0 = 0, Y_1 = 0, Y_2 = 1$

○ جمع این متغیرها باید برابر با 1 باشد؛  
$$\sum_{j=0}^2 Y_j = 1$$

○ تابع درست‌نمایی (برای n مشاهده مستقل)؛

$$L(\beta) = \prod_{i=1}^n [\pi_0(\mathbf{x}'_i)^{Y_{0i}} \pi_1(\mathbf{x}'_i)^{Y_{1i}} \pi_2(\mathbf{x}'_i)^{Y_{2i}}]$$



## ❖ نسبت شانس

○ فرض کنید ۳ حالت برای متغیر تصادفی  $Y$  متصور است؛ که عبارتند از،  $Y = j$  ;  $j = 0,1,2$

$Y = j = 0$  را به عنوان **حالت مرجع** در نظر گرفته و دو حالت دیگر را با آن مقایسه می‌کنیم.

اگر  $X$  تغییر کند موجب تغییر در  $P(Y = j)$  می‌شود.

فرض کنید  $X=a$  باشد. آنگاه؛

$$X = a \Rightarrow \begin{cases} P(Y = 0 | X = a) \\ P(Y = j | X = a) \quad ; j = 1,2 \end{cases}$$

بنابراین، تغییر در احتمال به معنای تغییر از  $P(Y = 0 | X = a)$  به  $P(Y = j | X = a)$  است.

حال اگر  $X$  تغییر کند (مثلا از  $a$  به  $b$ ) در این صورت  $P(Y = j | X = b)$  را داریم.

بنابراین، نسبت شانس برابر است با؛

$$OR_j(a, b) = \frac{P(Y = j | X = b) / P(Y = 0 | X = b)}{P(Y = j | X = a) / P(Y = 0 | X = a)} \quad ; j = 1,2$$



## ❖ نسبت شانس

$$OR_j(a, b) = \frac{P(Y = j | X = b) / P(Y = 0 | X = b)}{P(Y = j | X = a) / P(Y = 0 | X = a)} ; j = 1, 2$$

نسبت شانس برابر است با؛

○ بنابراین، پیامد  $Y=j$  را در مقایسه با  $Y=0$  به ازای  $X=b$  در مقابل  $X=a$  می‌سنجیم.

اگر متغیر توصیفی **باینتری** باشد (مقادیر دوگانه 0 و 1 را داشته باشد)؛ نسبت شانس به صورت زیر محاسبه خواهد شد؛

$$OR_j(0,1) = \frac{P(Y = j | X = 1) / P(Y = 0 | X = 1)}{P(Y = j | X = 0) / P(Y = 0 | X = 0)} ; j = 1, 2$$



# مدل‌های لوجیت – لوجیت چندگانه

❖ مثال: فرض کنید جدول فراوانی سطوح دو متغیر X و Y را داریم؛

وقتی سه پیامد وجود داشته باشد، دو مدل لوجیت داریم.

این توابع را به گونه‌ای تعریف می‌کنیم که دو ضریب برآوردشده (یکی برای هر تابع لوجیت) برابر با لگاریتم نسبت شانس برای هر جدول 2\*2 باشند که از طبقه‌بندی Y=j و Y=0 به دست می‌آید.

Y=0 را به عنوان مرجع در نظر می‌گیریم.

$$OR_j(0,1) = \frac{P(Y = j | X = 1)/P(Y = 0 | X = 1)}{P(Y = j | X = 0)/P(Y = 0 | X = 0)}$$

OR	(X)			
	جمع	1	0	
OR <sub>0</sub> = 1	۲۵۱	۱۷۱	۸۰	0
OR <sub>1</sub> = 1.965	۱۳۰	۱۰۵	۲۵	1
OR <sub>2</sub> = 3.244	۱۱۹	۱۰۴	۱۵	2
	۵۰۰	۳۸۰	۱۲۰	جمع

$$OR_1 = \frac{105}{\frac{171}{25}} = 1.965$$

$$OR_2 = \frac{104}{\frac{171}{15}} = 3.244$$



❖ مثال؛

$$P(Y = 0) = \frac{1}{1 + e^{g_1(X)} + e^{g_2(X)}}$$

$$P(Y = 1) = \frac{e^{g_1(X)}}{1 + e^{g_1(X)} + e^{g_2(X)}}$$

$$P(Y = 2) = \frac{e^{g_2(X)}}{1 + e^{g_1(X)} + e^{g_2(X)}}$$

از طرف دیگر، احتمال‌های روبرو را داریم؛

نسبت شانس برای  $Y=1$  که با  $OR_1$  نشان می‌دهیم، به ازای تغییر  $X$  از  $a$  به  $b$ ؛

$$OR_1(a, b) = \frac{P(Y = 1 | X = b) / P(Y = 0 | X = b)}{P(Y = 1 | X = a) / P(Y = 0 | X = a)} = \frac{e^{g_1(X=b)}}{e^{g_1(X=a)}} = e^{g_1(X=b) - g_1(X=a)} ; j = 1, 2$$

اگر  $g_1(X)$  برابر با  $g_1(X) = \alpha_1 + \beta_1 X$  باشد، آنگاه، نسبت شانس برابر است با؛

$$OR_1(a, b) = e^{(\alpha_1 + \beta_1 b) - (\alpha_1 + \beta_1 a)} = e^{\beta_1(b-a)}$$

حالت خاص: متغیر  $X$  فقط دو مقدار داشته باشد:  $X=0, 1$

$$OR_1(0, 1) = e^{\beta_1}$$



در حالت کلی؛

اگر  $g_j(X)$  برابر با  $g_j(X) = \alpha_j + \beta_j X$  باشد (معادله تک‌متغیره باشد)، آنگاه، نسبت شانس برابر است با؛

$$OR_1(a, b) = e^{(\alpha_1 + \beta_1 b) - (\alpha_1 + \beta_1 a)} = e^{\beta_1(b-a)}$$

$$OR_2(a, b) = e^{(\alpha_2 + \beta_2 b) - (\alpha_2 + \beta_2 a)} = e^{\beta_2(b-a)}$$

حالت خاص: متغیر  $X$  فقط دو مقدار داشته باشد:  $X=0, 1$

$$OR_1(0, 1) = e^{\beta_1}$$

$$OR_2(0, 1) = e^{\beta_2}$$

بنابراین، نسبت شانس در حالت کلی برابر است با؛

$$OR_j(a, b) = e^{g_j(X=b) - g_j(X=a)}$$

در حالتی که  $X=0, 1$  باشد ( $\Delta X = b - a = 1$ ) برآورد نسبت شانس ( $\widehat{OR}_j$ ) برابر است با

$$\widehat{OR}_j(a, b) = e^{\widehat{\beta}_j}$$

لگاریتم گیری

$$\ln \widehat{OR}_j = \widehat{\beta}_j$$

○ یعنی در این حالت، ضریب شیب برآورد شده (برای مدل لوجیت  $j$ ) برابر با لگاریتم نسبت شانس است.



بر آورد فاصله اطمینان؛

فاصله اطمینان ۹۵ درصدی برای  $\hat{\beta}_j$

$$\hat{\beta}_j \pm 1.96 SE(\hat{\beta}_j)$$

فاصله اطمینان ۹۵ درصدی برای  $\widehat{OR}_j$

$$Exp(\hat{\beta}_j \pm 1.96 SE(\hat{\beta}_j))$$



# مدل‌های لوجیت – لوجیت چندگانه

مثال: فرض کنید از داده‌های جدول روبرو، مدل‌های زیر برآورد شده است؛

مقادیر OR با رابطه  $OR_j = e^{\beta_j}$  به دست می‌آید.

خطای استاندارد ضرایب برابر است با؛

جزر مجموع معکوس فراوانی‌ها

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{80} + \frac{1}{171} + \frac{1}{25} + \frac{1}{104}} = 0.261$$

OR	(X)			
	جمع	1	0	
$OR_0 = 1$	۲۵۱	۱۷۱	۸۰	0
$OR_1 = 1.965$	۱۳۰	۱۰۵	۲۵	1
$OR_2 = 3.244$	۱۱۹	۱۰۴	۱۵	2
	۵۰۰	۳۸۰	۱۲۰	جمع

(Y)

فاصله اطمینان ۹۵ برای نسبت شانس	نسبت شانس	خطای استاندارد ضرایب	برآورد ضرایب	ضرایب	متغیر	مدل لوجیت
	1.965	0.225	-1.125	$\alpha_1$	ثابت	۱
1.178 , 3.276		0.261	0.675	$\beta_1$	X	
	3.244	0.25	-1.545	$\alpha_2$	ثابت	۲
1.774 , 5.934		0.308	1.177	$\beta_2$	X	

$$OR_1 : Exp(\hat{\beta}_1 \pm 1.96 SE(\hat{\beta}_1)) = (1.178 , 3.276)$$

فاصله اطمینان ۹۵ درصدی برای  $\widehat{OR}_j$  :



# مدل‌های لوجیت – لوجیت چندگانه

مثال: فرض کنید از داده‌های جدول روبرو، مدل‌های زیر برآورد شده است؛

مقادیر OR با رابطه  $OR_j = e^{\beta_j}$  به دست می‌آید.

خطای استاندارد ضرایب برابر است با؛

جزر مجموع معکوس فراوانی‌ها

$$SE(\hat{\beta}_1) = \sqrt{\frac{1}{80} + \frac{1}{171} + \frac{1}{25} + \frac{1}{104}} = 0.261$$

OR	(X)			
	جمع	1	0	
$OR_0 = 1$	۲۵۱	۱۷۱	۸۰	0
$OR_1 = 1.965$	۱۳۰	۱۰۵	۲۵	1
$OR_2 = 3.244$	۱۱۹	۱۰۴	۱۵	2
	۵۰۰	۳۸۰	۱۲۰	جمع

(Y)

فاصله اطمینان ۹۵ برای نسبت شانس	نسبت شانس	خطای استاندارد ضرایب	برآورد ضرایب	ضرایب	متغیر	مدل لوجیت
	1.965	0.225	-1.125	$\alpha_1$	ثابت	۱
1.178 , 3.276		0.261	0.675	$\beta_1$	X	
	3.244	0.25	-1.545	$\alpha_2$	ثابت	۲
1.774 , 5.934		0.308	1.177	$\beta_2$	X	

$$OR_1 : Exp(\hat{\beta}_1 \pm 1.96 SE(\hat{\beta}_1)) = (1.178 , 3.276)$$

فاصله اطمینان ۹۵ درصدی برای  $\widehat{OR}_j$  :



❖ حالتی که X بیش از دو وضعیت داشته باشد

○ مثال: فرض کنید جدول فراوانی سطوح دو متغیر X و Y را داریم؛

حالت X=0 و Y=0 را به عنوان مرجع در نظر می‌گیریم.

نسبت شانس‌ها به صورت زیر به دست می‌آید؛

	(X)					
جمع	۳	۲	۱	۰		
۲۴۰	۱۰۰	۶۰	۵۰	۳۰	۰	(Y)
۱۳۵	۴۰	۵۰	۳۰	۱۵	۱	
۱۲۵	۶۰	۴۰	۲۰	۵	۲	
۵۰۰	۲۰۰	۱۵۰	۱۰۰	۵۰	جمع	

$$OR_2(1, 0) = \frac{\frac{20}{50}}{\frac{5}{30}} = 2.4$$

$$OR_1(1, 0) = \frac{\frac{30}{50}}{\frac{15}{30}} = 1.2$$

$$OR_2(2, 0) = \frac{\frac{40}{60}}{\frac{5}{30}} = 4$$

$$OR_1(2, 0) = \frac{\frac{50}{60}}{\frac{15}{30}} = 1.67$$

$$OR_2(3, 0) = \frac{\frac{60}{100}}{\frac{5}{30}} = 3.6$$

$$OR_1(3, 0) = \frac{\frac{40}{100}}{\frac{15}{30}} = 0.8$$



# مدل‌های لوجیت – لوجیت چندگانه

❖ حالتی که  $X$  بیش از دو وضعیت داشته باشد

○ مثال: فرض کنید جدول فراوانی سطوح دو متغیر  $X$  و  $Y$  را داریم؛

		(X)					
		۳	۲	۱	۰		
جمع		۳	۲	۱	۰		
۲۴۰	۱۰۰	۶۰	۵۰	۳۰	۰	(Y)	
۱۳۵	۴۰	۵۰	۳۰	۱۵	۱		
۱۲۵	۶۰	۴۰	۲۰	۵	۲		
۵۰۰	۲۰۰	۱۵۰	۱۰۰	۵۰	جمع		

از طرف دیگر، نسبت شانس برابر است با  $e$  به توان ضریب رگرسیون لوجیت

$$OR_j(i, o) = e^{\beta_{ji}} \quad ; \quad j = 1, 2 \quad ; \quad i = 1, 2, 3$$

$j$  حالت‌های  $Y$  و  $i$  حالت‌های  $X$  را نشان می‌دهد.

به ازای یک واحد افزایش در متغیر  $X$  (و ثابت ماندن سایر متغیرها) بخت وقوع پیامد  $j$  نسبت به حالت پایه، با ضریب  $e^{\beta_{ji}}$  تغییر می‌یابد.



# مدل‌های لوجیت – لوجیت چندگانه

## ❖ آزمون برابری دو نسبت شانس

فرضیه صفر برابری دو نسبت شانس:

$$H_0: OR_1 = OR_2 \implies H_0: OR_1 - OR_2 = 0$$

این فرضیه معادل آزمون فرضیه زیر است:

$$H_0: \frac{OR_1}{OR_2} = 1 \quad \text{یا} \quad H_0: \ln(OR_1) = \ln(OR_2)$$

از رابطه تفاضل بین دو شیب معادلات رگرسیون لوجیت استفاده می‌کنیم؛

$$\ln(\widehat{OR}_2) - \ln(\widehat{OR}_1) = \hat{\beta}_2 - \hat{\beta}_1$$

اگر فرضیه صفر درست باشد، آنگاه این تفاضل از توزیع نرمال با میانگین و واریانس زیر پیروی می‌کند؛

$$E(\hat{\beta}_2 - \hat{\beta}_1) = \beta_2 - \beta_1$$

$$Var(\hat{\beta}_2 - \hat{\beta}_1) = Var(\hat{\beta}_2) - Var(\hat{\beta}_1) - 2Cov(\hat{\beta}_1, \hat{\beta}_2)$$

آماره t را به صورت روبرو تعریف می‌کنیم

$$t = \frac{\hat{\beta}_2 - \hat{\beta}_1}{\sqrt{var(\hat{\beta}_2 - \hat{\beta}_1)}}$$



## ❖ آزمون برابری دو نسبت شانس

$$(\hat{\beta}_2 - \hat{\beta}_1) \pm 1.96 \sqrt{\text{Var}(\hat{\beta}_2 - \hat{\beta}_1)}$$

فاصله اطمینان ۹۵ درصدی برای  $(\beta_2 - \beta_1)$

اگر حد بالا و پایین این فاصله اطمینان برابر با  $L_1$  و  $L_2$  باشد، آنگاه

$$(e^{L_1}, e^{L_2})$$

فاصله اطمینان ۹۵ درصدی برای نسبت  $\frac{OR_2}{OR_1}$  برابر است با؛

سوال:

در چه صورتی می‌گوییم  $OR_2$  و  $OR_1$  تفاوت معناداری ندارند؟



## ❖ آزمون معنادار بودن ضرایب مدل (شیب رگرسیون لوجیت)

- برای آزمون معنادار بودن شیب رگرسیون لوجیت (ضریب  $X$ ) از نسبت درست‌نمایی استفاده می‌کنیم.
- این آزمون معادل با مقایسه دو رگرسیون است؛
  - رگرسیونی که شامل  $X$  است (رگرسیون غیرمقید)
  - رگرسیونی که فاقد  $X$  است (رگرسیون مقید)
- این کار را برای هر یک از مدل‌های لوجیت انجام می‌دهیم.
  - برای رگرسیون غیرمقید، نسبت درست‌نمایی را با  $L_1$  نشان می‌دهیم.
  - برای رگرسیون مقید، نسبت درست‌نمایی را با  $L_0$  نشان می‌دهیم.

فرضیه صفر: ضریب  $X$  صفر است. مثلاً برای رگرسیون لوجیت ۱ فرضیه صفر:  $H_0: \beta_1 = 0$

در صورت درست بودن فرضیه صفر، نسبت درست‌نمایی  $\lambda$  (رابطه زیر) از توزیع  $\chi^2$  با ۲ درجه آزادی پیروی می‌کند.

$$\lambda = 2(L_1 - L_0)$$

در حالت کلی، اگر تعداد حالت‌های متغیر وابسته برابر با  $s$  و تعداد گروه‌های متغیر  $X$  برابر با  $k$  باشد، درجه آزادی برابر با  $(s - 1)(k - 1)$  است.

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس	مقدمه
لوجیت چندگانه	پروبیت دوگانه
لوجیت شرطی	لوجیت دوگانه (لوجستیک)
لوجیت آشیانه‌ای و ترکیبی	معیارهای نیکویی برآزش



❖ تصور کنید افراد برای رفت و آمد روزانه می‌توانند از چهار وسیله‌نقلیه استفاده کنند؛

- خودرو شخصی
- تاکسی
- اتوبوس
- مترو

$$Y_{i1} = 1 \quad , \quad Y_{i2} = Y_{i3} = Y_{i4} = 0$$

$$Y_{i2} = 1 \quad , \quad Y_{i1} = Y_{i3} = Y_{i4} = 0$$

$$Y_{i3} = 1 \quad , \quad Y_{i1} = Y_{i2} = Y_{i4} = 0$$

$$Y_{i4} = 1 \quad , \quad Y_{i1} = Y_{i2} = Y_{i3} = 0$$

بنابراین فرد  $i$  دارای چهار گزینه است  $j=1,2,3,4$

اگر  $Y_i$  بیانگر انتخاب فرد  $i$  باشد، آنگاه؛

$Y_i$  می‌تواند شامل چهار مقدار ۱، ۲، ۳، ۴ باشد:

احتمال انتخاب گزینه  $j$  بستگی به خصوصیات فرد  $i$  دارد.

اگر بردار متغیرهای توصیفی برای انتخاب فرد  $i$  را با  $x$  نشان دهیم، می‌تواند شامل متغیرهایی مانند درآمد،

سن، شغل، جنسیت و غیره باشد.



اگر بردار متغیرهای توصیفی برای انتخاب فرد  $i$  را با  $\mathbf{x}$  نشان دهیم، می‌تواند شامل متغیرهایی مانند درآمد، سن، شغل، جنسیت و غیره باشد.

در این حالت، یک مدل لوجیت چندگانه را به صورت زیر تعریف می‌کنیم؛

$$\pi_j(\mathbf{x}'_i) = P(Y = j | \mathbf{x}'_i) = \frac{e^{\mathbf{x}'_i \beta_j}}{\sum_{j=1}^4 e^{\mathbf{x}'_i \beta_j}} \quad ; \quad i = 0, 1, \dots, n$$

توجه شود که در اینجا فقط ویژگی‌های فردی است که اهمیت دارند، و هیچ توجهی به ویژگی‌های گزینه‌ها (در اینجا ویژگی‌های شیوه حمل‌ونقلی) نداریم.



اکنون فرض کنید بخواهیم ویژگی‌های خاص گزینه‌ها نیز درون مدل قرار دهیم. بنابراین هر فرد  $i$  دارای چهار گزینه است  $j=1,2,3,4$  و هر گزینه  $j$  دارای ویژگی‌های خاص خود می‌باشد. مثلاً فرض کنید دو متغیر هزینه  $(C_{ij})$  و زمان سفر  $(T_{ij})$  برای گزینه‌های موردنظر، متفاوت است. پس هر فرد  $i$ ، هزینه  $(C_{ij})$  را می‌پردازد و زمان سفر  $(T_{ij})$  را صرف می‌کند تا به مقصد خود برسد. بدین ترتیب، تصمیم‌گیری و انتخاب هر گزینه برای فرد  $i$  فقط براساس ویژگی‌های گزینه‌هاست (و هیچ وابستگی به ویژگی‌های فردی ندارد).

❖ چنین مدلی معروف به **مدل لوجیت شرطی (Conditional Logit)** است.

**توجه:** ویژگی‌های گزینه  $j$  (مثلاً هزینه گزینه  $j$ ) می‌تواند برای فرد  $i$  متفاوت با سایر افراد باشد (مثلاً قیمت متفاوت اتوبوس در مناطق مختلف شهری و یا تخفیف برای دانش‌آموزان و دانشجویان)

اگر قیمت برای همه افراد یکسان باشد، آنگاه  $C_{ij} = C_j$

برای سایر ویژگی‌ها و متغیرها هم به همین شکل.



تصور کنید که  $Y_i$  بیانگر انتخاب بین  $m$  گزینه باشد.

از طرف دیگر،  $U_{ij}$  را مطلوبیت انتخاب گزینه  $j$  برای فرد  $i$  در نظر می‌گیریم.

$U_{ij}$  یک متغیر تصادفی است که شامل دو جزء است:

$V_{ij}$  جزء غیرتصادفی

$\varepsilon_{ij}$  جزء تصادفی

$$U_{ij} = V_{ij} + \varepsilon_{ij}$$

فرض: فرد عاقلانه عمل می‌کند و گزینه‌ای را انتخاب می‌کند که بیشترین مطلوبیت را برای او دارد.

پس گزینه  $j$  در صورتی انتخاب می‌شود که  $U_{ij}$  برابر با بزرگترین مقدار از بین  $(U_{i1}, \dots, U_{im})$  باشد.

$$p_{ij} = P(Y_i = j) = P(U_{ij} = \max(U_{i1}, \dots, U_{im}))$$



بر این مبنا، ثابت می‌شود که متغیر تصادفی  $\varepsilon_{ij}$  دارای تابع چگالی گامبل است، که عبارتست از:

$$f(\varepsilon_{ij}) = e^{-\varepsilon_{ij}} e^{-e^{-\varepsilon_{ij}}} = e^{(-\varepsilon_{ij} - e^{-\varepsilon_{ij}})}$$

از طرف دیگر، برای مقایسه دو گزینه  $j$  و  $k$  توسط فرد  $i$  عبارتست از:

$$\begin{aligned}
 p_{ij} &= P(Y_i = j) = P(U_{ij} \geq U_{ik}) \\
 &= P(V_{ij} + \varepsilon_{ij} \geq V_{ik} + \varepsilon_{ik}) \\
 &= P(\varepsilon_{ij} - \varepsilon_{ik} \geq V_{ik} - V_{ij}) = P(\varepsilon_{ik} - \varepsilon_{ij} \leq V_{ij} - V_{ik})
 \end{aligned}$$

یعنی احتمال اینکه تفاضل  $\varepsilon_{ik} - \varepsilon_{ij}$  کوچکتر از  $V_{ij} - V_{ik}$  باشد.

چون  $\varepsilon_{ik}$  و  $\varepsilon_{ij}$  دو متغیر تصادفی هستند که از توزیع گامبل پیروی می‌کنند، ثابت می‌شود که تفاضل آنها (یعنی  $\varepsilon_{ij} - \varepsilon_{ik}$ ) از توزیع لوجستیک پیروی می‌کند (و از این طریق، فرمول احتمال لوجیت به دست می‌آید).

$$p_{ij} = \frac{e^{V_{ij}-V_{ik}}}{1 + e^{V_{ij}-V_{ik}}} = \frac{e^{V_{ij}}}{e^{V_{ij}} + e^{V_{ik}}} = \frac{e^{V_{ij}}}{\sum_{j=1}^2 e^{V_{ij}}}$$

حالت کلی:  
مقایسه  $U_{ij}$  با  $m-1$   
گزینه دیگر  
( $U_{i1}, \dots, U_{im}$ )

$$p_{ij} = \frac{e^{V_{ij}}}{\sum_{k=1}^m e^{V_{ik}}}$$



اگر  $V_{ij}$  را بر حسب ویژگی‌های گزینه‌ها تعریف کنیم، خواهیم داشت (تابع مطلوبیت خطی):

$$V_{ij} = \mathbf{x}'_i \boldsymbol{\beta} = \alpha_j + \beta_1 C_{ij} + \beta_2 T_{ij}$$

$$p_{ij} = \frac{e^{\alpha_j + \beta_1 C_{ij} + \beta_2 T_{ij}}}{\sum_{j=1}^4 e^{\alpha_j + \beta_1 C_{ij} + \beta_2 T_{ij}}}$$

## نکته:

- از آنجا که پایه مدل لوجیت شرطی، نظریه مطلوبیت تصادفی و مقایسه مطلوبیت حاصل از گزینه‌هاست؛ و  $p_{ij}$  بر حسب تفاضل مطلوبیت گزینه‌ها بیان شده است.
- بنابراین، فقط می‌توانیم «تفاضل  $\alpha_j$  ها» را حساب کنیم؛ و امکان برآورد هر یک از  $\alpha_j$  ها به صورت جداگانه نیست.
- بر این اساس، اگر گزینه ۱ را به عنوان مبنا در نظر بگیریم، می‌توان  $\alpha_1$  را نرمال‌سازی نمود و **برابر با صفر** در نظر گرفت؛ سپس، بقیه  $\alpha_j$  ها برآورد نمود.



## ❖ داده‌های مورد نیاز

در مدل لوجیت شرطی با مسأله تجربه یا آزمایش انتخاب گسسته (Discrete Choice Experiment) مواجه هستیم. داده‌های موردنیاز در این حوزه مربوط به ترجیحات افراد است.

○ بر این مبنا، دو نوع داده قابل جمع‌آوری و استفاده است؛

۱- داده‌هایی که به طور مستقیم جمع‌آوری می‌شود؛ و بنابراین، نشان دهنده انتخاب‌ها (رفتارها)ی واقعی افراد است که در شرایط واقعی از خود نشان داده‌اند.

در این حالت از «**ترجیحات آشکارشده**» (Revealed preferences) استفاده می‌شود. مثلاً: گردآوری داده‌های واقعی استفاده از مترو، اتوبوس، تاکسی و غیره.

۲- داده‌هایی که از طریق پرسشنامه گردآوری می‌شوند. برای این منظور، فرد را در موقعیت موردنظر قرار داده و سپس از او در خصوص ترجیحاتش سوال می‌گردد.

در این حالت از «**ترجیحات اظهارشده یا بیان شده**» (Stated preferences) زیر سناریوهای مختلف استفاده می‌شود. مثلاً فرد را در موقعیت انتخاب بین خودرو شخصی و اتوبوس قرار می‌دهیم و در مورد تمایل به پرداخت هزینه در نسبت با کاهش زمان جابجایی سوالات مختلفی می‌پرسیم.



## ❖ مفروضات در مورد $\varepsilon_{ij}$

فرض می‌شود  $\varepsilon_{ij}$  ها هم در میان گزینه‌ها ( $j$ ) هم در میان مشاهدات ( $i$ ) مستقل هستند. بنابراین،  $\varepsilon_{ij}$  ها iid بوده و فرض می‌شود دارای توزیع گامبل هستند.



## ❖ برآورد مدل لوجیت شرطی

برآورد مدل لوجیت شرطی با استفاده از روش MLE احتمال انتخاب گزینه‌ها برای فرد  $i$  برابر است با؛

$$L_i = P(Y_{i1} = 1, Y_{i2} = 0, \dots, Y_{im} = 0) \\ \times P(Y_{i1} = 0, Y_{i2} = 1, \dots, Y_{im} = 0) \\ \times \dots \times P(Y_{i1} = 0, Y_{i2} = 0, \dots, Y_{im} = 1) = \prod_{j=1}^m p_{ij}^{Y_{ij}}$$

چون  $n$  فرد و  $m$  گزینه داریم، تابع درستنمایی عبارتست از ؛

$$L = L_1 \times L_2 \times \dots \times L_n = \prod_{i=1}^n L_i = \prod_{i=1}^n \prod_{j=1}^m p_{ij}^{Y_{ij}}$$

$$l = \ln L = \sum_{i=1}^n \sum_{j=1}^m Y_{ij} \ln p_{ij}$$

لگاریتم تابع درستنمایی عبارتست از ؛

با مشتق‌گیری نسبت به تمام پارامترهای مدل، برآوردهایی برای پارامترهای مدل به دست می‌آید.



## ❖ تفسیر نتایج مدل لوجیت شرطی

برای محاسبه اثرات نهایی تغییر در  $X_{ij}$  (بر احتمال انتخاب گزینه  $j$ )، از  $p_{ij}$  نسبت به  $X_{ij}$  مشتق می‌گیریم؛

$$\frac{\partial P(Y_i = j)}{\partial X_{ij}} = \frac{\partial p_{ij}}{\partial X_{ij}} = p_{ij}(1 - p_{ij})\beta$$

$\beta$  ضریب متغیر  $X$  است.

برای محاسبه اثرات نهایی تغییر در  $X_{ik}$  (بر احتمال انتخاب گزینه  $j$ )، از  $p_{ij}$  نسبت به  $X_{ik}$  مشتق می‌گیریم؛

$$\frac{\partial P(Y_i = j)}{\partial X_{ik}} = \frac{\partial p_{ij}}{\partial X_{ik}} = -p_{ij}p_{ik}\beta$$



## ❖ تفسیر نتایج مدل لوجیت شرطی

مثال؛ اثر تغییر در هزینه (قیمت) وسیله نقلیه شخصی، بر احتمال استفاده از وسیله نقلیه شخصی عبارتست از؛

$$\frac{\partial P(Y_i = 1)}{\partial C_{i1}} = p_{i1}(1 - p_{i1})\beta_1$$

$\beta_1$  ضریب متغیر  $C_{ij}$  است.

اثر تغییر در هزینه (قیمت) تاکسی، بر احتمال استفاده از وسیله نقلیه شخصی عبارتست از؛

$$\frac{\partial P(Y_i = 1)}{\partial C_{i2}} = -p_{i1}p_{i2}\beta_1$$



## ❖ تفسیر نتایج مدل لوجیت شرطی

مثال؛ اثر تغییر در زمان سفر وسیله نقلیه شخصی، بر احتمال استفاده از وسیله نقلیه شخصی عبارتست از؛

$$\frac{\partial P(Y_i = 1)}{\partial T_{i1}} = p_{i1}(1 - p_{i1})\beta_2$$

$\beta_2$  ضریب متغیر  $T_{ij}$  است.

اثر تغییر در زمان سفر تاکسی، بر احتمال استفاده از وسیله نقلیه شخصی عبارتست از؛

$$\frac{\partial P(Y_i = 1)}{\partial T_{i2}} = -p_{i1}p_{i2}\beta_2$$



## ❖ مثال،

مدل مطلوبیت شیوه‌های اصلی حمل‌ونقل در یک شهر به صورت زیر برآورد شده است. مطلوبست؛

الف) برآورد سهم هر شیوه جابجایی با استفاده از داده‌های موجود در جدول (احتمال انتخاب هر شیوه).

ب) اگر هزینه پارکینگ ۱۰ هزار ریال برای هر سفر افزایش یابد، سهم هر شیوه حمل‌ونقلی چه تغییری می‌کند؟

$$V_1 = -0.3 - 0.002C_1 - 0.05T_1$$

$$V_2 = -0.35 - 0.002C_2 - 0.05T_2$$

$$V_3 = -0.4 - 0.002C_3 - 0.05T_3$$

$C_j$ : هزینه شیوه حمل‌ونقلی  $j$  (بر حسب صد ریال)

$T_j$ : زمان سفر با شیوه حمل‌ونقلی  $j$  (بر حسب دقیقه)

شیوه حمل‌ونقلی	هزینه (هزار ریال)	زمان سفر (دقیقه)
خودرو شخصی	۱۳	۲۵
اتوبوس	۷.۵	۳۵
مترو	۹	۴۰



❖ مثال

$$V_1 = -0.3 - 0.002C_1 - 0.05T_1$$

$$V_2 = -0.35 - 0.002C_2 - 0.05T_2$$

$$V_3 = -0.4 - 0.002C_3 - 0.05T_3$$

شیوه حمل‌ونقلی	هزینه (هزار ریال)	زمان سفر (دقیقه)
خودرو شخصی	۱۳	۲۵
اتوبوس	۷.۵	۳۵
مترو	۹	۴۰

حل الف) برآورد سهم هر شیوه جابجایی با استفاده از داده‌های موجود در جدول (احتمال انتخاب هر شیوه).

$$V_1 = -0.3 - 0.002 * 130 - 0.05 * 25 = -1.81$$

$$V_2 = -0.35 - 0.002 * 75 - 0.05 * 35 = -2.25$$

$$V_3 = -0.4 - 0.002 * 90 - 0.05 * 40 = -2.58$$

$$e^{V_1} = 0.1636 \quad ; \quad e^{V_2} = 0.1054 \quad ; \quad e^{V_3} = 0.0758$$

$$p_{car} = \frac{e^{V_1}}{\sum_{k=1}^3 e^{V_{ik}}} = \frac{0.1636}{0.3448} = 0.4745 \text{ (47.45\%)}$$

$$p_{bus} = \frac{e^{V_2}}{\sum_{k=1}^3 e^{V_{ik}}} = \frac{0.1054}{0.3448} = 0.3057 \text{ (30.57\%)}$$

$$p_{metro} = \frac{e^{V_3}}{\sum_{k=1}^3 e^{V_{ik}}} = \frac{0.0758}{0.3448} = 0.2198 \text{ (21.98\%)}$$



❖ مثال

$$V_1 = -0.3 - 0.002C_1 - 0.05T_1$$

$$V_2 = -0.35 - 0.002C_2 - 0.05T_2$$

$$V_3 = -0.4 - 0.002C_3 - 0.05T_3$$

شیوه حمل‌ونقلی	هزینه (هزار ریال)	زمان سفر (دقیقه)
خودرو شخصی	۱۳	۲۵
اتوبوس	۷.۵	۳۵
مترو	۹	۴۰

حل ب) اگر هزینه پارکینگ ۱۰ هزار ریال برای هر سفر افزایش یابد، سهم هر شیوه حمل‌ونقلی چه تغییری می‌کند؟  
✓ براساس مقیاس‌ها، ۱۰۰ واحد در متغیر قیمت افزایش ایجاد شده است.

$$V_{1-new} = -0.3 - 0.002 * 230 - 0.05 * 25 = -2.01$$

$$V_2 = -0.35 - 0.002 * 75 - 0.05 * 35 = -2.25$$

$$V_3 = -0.4 - 0.002 * 90 - 0.05 * 40 = -2.58$$

$$e^{V_{1-new}} = 0.1339 \quad ; \quad e^{V_2} = 0.1054 \quad ; \quad e^{V_3} = 0.0758$$

$$p_{car} = \frac{e^{V_1}}{\sum_{k=1}^3 e^{V_{ik}}} = \frac{0.1339}{0.3151} = 0.4250 \quad (42.5\%)$$

$$p_{bus} = \frac{e^{V_2}}{\sum_{k=1}^3 e^{V_{ik}}} = \frac{0.1054}{0.3151} = 0.3345 \quad (33.45\%)$$

$$p_{metro} = \frac{e^{V_3}}{\sum_{k=1}^3 e^{V_{ik}}} = \frac{0.0758}{0.3151} = 0.2405 \quad (24.05\%)$$

✓ بدین ترتیب، سهم حمل‌ونقل شخصی از ۴۷/۴۵٪ به ۴۲/۵۰٪ کاهش می‌یابد و سهم اتوبوس ۲/۸۸٪ و مترو ۲/۰.۷٪ افزایش می‌یابد.



## ❖ نسبت احتمال (Probability Ratio)

نسبت احتمال (احتمال نسبی) انتخاب گزینه  $j$  نسبت به گزینه  $k$  برابر است با؛

$$PR = \frac{P(Y_i = j)}{P(Y_i = k)} = \frac{p_{ij}}{p_{ik}} = \frac{e^{x'_{ij}\beta}}{e^{x'_{ik}\beta}} = e^{(x'_{ij} - x'_{ik})\beta}$$

با لگاریتم‌گیری، لگاریتم نسبت احتمال برابر است با؛

$$\ln PR = (x'_{ij} - x'_{ik})\beta$$

برای مثال در مسأله حمل‌ونقل، اگر فقط  $C_{ij}$  تغییر کند، آنگاه داریم؛

$$\ln PR = \ln \frac{e^{\alpha_1 + \beta_1 C_{i1} + \beta_2 T_{i1}}}{e^{\alpha_2 + \beta_1 C_{i2} + \beta_2 T_{i2}}} = (\alpha_1 - \alpha_2) + \beta_1 (C_{i1} - C_{i2}) + \beta_2 (T_{i1} - T_{i2})$$

**تفسیر:** به ازای یک واحد افزایش در مقدار یک متغیر برای گزینه  $j$  نسبت به  $k$  (مثلا هزینه خودرو شخصی یک واحد افزایش یابد و تاکسی تغییر نکند)، نسبت احتمال به اندازه ضریب  $e^\beta$  تغییر می‌کند (مثلا اگر  $\beta_1 = -0.2$  و  $e^{-0.2} = 0.818$  بدین معناست که احتمال نسبی انتخاب خودرو شخصی نسبت به تاکسی، ۱۸.۲٪ کمتر از قبل می‌شود).



## ❖ نرخ نهایی جانشینی (Marginal Rate of Substitution (MRS))

یکی دیگر از تحلیل‌های حاصل از مدل لوجیت شرطی، محاسبه و تحلیل **نرخ نهایی جانشینی** است.

**نرخ نهایی (حاشیه‌ای) جانشینی** بیانگر تغییر در گزینه‌ها به ازای ثابت ماندن مطلوبیت است.

در مثال حمل‌ونقل مدل زیر را در نظر بگیرید،

$$U_{ij} = \alpha_j + \beta_1 C_{ij} + \beta_2 T_{ij} + \varepsilon_{ij}$$

تغییرات مطلوبیت را به ازای تغییر در هزینه حمل‌ونقل و تغییر در زمان صرف شده، حساب کرده و برابر با صفر قرار

می‌دهیم؛

$$\Delta U_{ij} = \beta_1 \Delta C_{ij} + \beta_2 \Delta T_{ij} = 0$$

نرخ نهایی جانشینی برابر است با؛

$$MRS = -\frac{\Delta C_{ij}}{\Delta T_{ij}} = \frac{\beta_2}{\beta_1}$$



## ❖ نرخ نهایی جانشینی

نرخ نهایی جانشینی برابر است با؛

$$MRS = -\frac{\Delta C_{ij}}{\Delta T_{ij}} = \frac{\beta_2}{\beta_1}$$

$\beta_2$  بیانگر اثر یک گزینه غیر پولی (زمان صرف شده) است، و  $\beta_1$  بیانگر اثر یک گزینه پولی (هزینه حمل و نقل) نسبت  $MRS$  بیانگر آن است که اگر یک واحد در زمان صرف شده کاسته شود، چقدر به هزینه (تمایل به

پرداخت) اضافه خواهد شد؛ که برابر است با  $\frac{\beta_2}{\beta_1}$

رابطه  $MRS$  را به صورت روبرو نیز می‌توان نوشت؛  
$$\Delta C_{ij} = -\frac{\beta_2}{\beta_1} \Delta T_{ij}$$

این رابطه بیانگر **میزان تمایل به پرداخت هزینه بابت تغییر در زمان تلف شده** است.

مثلاً اگر  $\frac{\beta_2}{\beta_1} = 1.5$  آنگاه می‌توان گفت یک واحد کاهش در زمان تلف شده ( $\Delta T_{ij} = -1$ )، امکان افزایش

هزینه حمل و نقل به میزان  $1/5$  برابر را بدون تغییر در میزان مطلوبیت، فراهم می‌کند.

یعنی اگر قیمت  $1000$  تومان باشد؛ به ازای یک واحد کاهش در زمان، کاربر حاضر است  $1500$  تومان بپردازد.



## ❖ ناهمگونی در ترجیحات (Preference heterogeneity)

- تا اینجا فرض بر این بود که همه افراد با ویژگی‌های مختلف، دارای ترجیحات و تمایل به پرداخت یکسان هستند.
- به همین دلیل ضرایب  $\beta$  برای همه افراد یکسان در نظر گرفته شد.
- این ضرایب می‌توانند متغیر باشند (یعنی میزان اهمیت یک متغیر برای افراد مختلف، متفاوت باشد).
- مثلاً،  $\beta_1$  که ضریب هزینه حمل‌ونقل است می‌تواند با افزایش درآمد، افزایش یابد.
- یعنی با افزایش درآمد، افراد شیوه حمل‌ونقل گران‌تر را انتخاب می‌کنند.
- یا با افزایش درآمد، تمایل دارند شیوه حمل‌ونقلی که زمان سفر کمتری دارد را انتخاب کنند.
- اگر فرضیه این باشد که افراد با افزایش درآمد به دنبال انتخاب شیوه حمل‌ونقلی با زمان سفر کمتر هستند، می‌توان این موضوع را با وارد کردن حاصلضرب زمانِ صرف شده و درآمد، لحاظ نمود.

$$U_{ij} = \alpha_j + \beta_1 C_{ij} + \beta_2 T_{ij} + \gamma_2 T_{ij} I_{ij} + \dots$$



## ❖ ناهمگونی در ترجیحات

$U_{ij} = \alpha_j + \beta_1 C_{ij} + \beta_2 T_{ij} + \gamma_2 T_{ij} I_{ij} + \dots$   $I_{ij}$  درآمد فرد  $i$  را نشان می‌دهد که گزینه  $j$  را انتخاب می‌کند.

این معادله را به صورت روبرو بازنویسی می‌کنیم.

$$U_{ij} = \alpha_j + \beta_1 C_{ij} + (\beta_2 + \gamma_2 I_{ij}) T_{ij} + \dots$$

ضریب  $T_{ij}$  برابر با  $(\beta_2 + \gamma_2 I_{ij})$  است که متغیر می‌باشد، و بستگی به سطح درآمد فرد دارد.

اگر  $\beta_1$  نیز تابع درآمد باشد،  $U_{ij}$  به زیر خواهد بود؛

$$U_{ij} = \alpha_j + (\beta_1 + \gamma_1 I_{ij}) C_{ij} + (\beta_2 + \gamma_2 I_{ij}) T_{ij} + \dots$$

در این حالت، تمایل به پرداخت برای کاهش زمان صرف شده (نرخ نهایی جانشینی) برابر است با؛

$$MRS = \frac{\beta_2 + \gamma_2 I_{ij}}{\beta_1 + \gamma_1 I_{ij}}$$

بنابراین، در این حالت **نرخ نهایی جانشینی** یا تمایل به پرداخت برای کاهش زمان تلف شده، بستگی به سطح درآمد فرد  $i$  خواهد داشت.



## ❖ استقلال گزینه‌های نامربوط (Independence of Irrelevant Alternatives (IIA))

یکی از الزامات مدل لوجیت شرطی این است که انتخاب از بین مجموعه گزینه‌ها باید دارای شرط **استقلال گزینه‌های نامربوط** باشد.

○ فرض **استقلال گزینه‌های نامربوط** بیان می‌کند که؛  
✓ انتخاب بین دو گزینه A و B نباید وابسته به ویژگی‌های گزینه نامربوط C باشد.

مثلا فرض کنید افراد می‌توانند بین گزینه **اتوبوس** و **خودرو شخصی** انتخاب کنند.

با معرفی یک شیوه جدید، مثلا **مترو**، اگر فرض استقلال گزینه‌های نامربوط برقرار باشد، نباید تغییری در ترجیح نسبی افراد برای انتخاب بین دو گزینه **اتوبوس** و **خودرو شخصی** ایجاد شود.

✓ یعنی اگر یک نفر بین دو گزینه **اتوبوس** و **خودرو شخصی**، گزینه **اتوبوس** را انتخاب کرده بود؛ نباید بعد از معرفی گزینه **مترو**، گزینه **خودرو شخصی** را انتخاب کند!



## ❖ استقلال گزینه‌های نامربوط (Independence of Irrelevant Alternatives (IIA))

توجه شود که در مدل **لوجیت شرطی** فرض کرده‌ایم که  $\varepsilon_{ij}$  ها iid هستند، یعنی،

- تصمیمات افراد مستقل از هم است (**استقلال افراد**) و
- گزینه‌ها نیز مستقل از هم هستند (**استقلال گزینه‌ها**)

استقلال تصمیمات افراد، منطقی است، اما اگر گزینه‌ها شباهت داشته باشند، فرض استقلال گزینه‌ها نقض می‌شود.

$$\ln PR = \ln \frac{P(Y_i = j)}{P(Y_i = k)} = (\mathbf{x}_{ij} - \mathbf{x}_{ik})' \boldsymbol{\beta}$$

برای دو گزینه  $j$  و  $k$ ، لگاریتم نسبت احتمال برابر است با؛

مشاهده می‌شود که در صورت برقراری فرض استقلال گزینه‌ها،  $PR$  تابعی از گزینه‌های  $j$  و  $k$  است (و تابع سایر گزینه‌ها نیست).

بنابراین، در انتخاب میان دو گزینه  $j$  و  $k$ ، سایر گزینه‌ها نامربوط هستند.

یعنی وقتی می‌خواهیم ترجیحات فرد  $i$  را در خصوص این دو گزینه بررسی کنیم، حضور یا عدم حضور سایر گزینه‌ها اهمیتی ندارد.

• به عبارت دیگر، احتمال نسبی دو گزینه، ربطی به ویژگی‌های سایر گزینه‌ها ندارد.

• یا ریسک نسبی دو گزینه، مستقل از سایر گزینه‌هاست.



## ❖ خلاصه مفروضات مهم مدل‌های MNL و Conditional Logit

1. استقلال گزینه‌های نامربوط ((Independence of Irrelevant Alternatives (IIA)
2. جملات خطای تصادفی در میان افراد، مستقل هستند.
3. جملات خطای تصادفی در میان گزینه‌ها، مستقل هستند.
4. همه جملات خطای تصادفی ( $\varepsilon_{ij}$  ها) دارای توزیع مشابه هستند.
5. جملات خطا از توزیع گامبل پیروی می‌کنند.
6. ناهمگنی پنهان وجود ندارد. همه افراد با ویژگی‌های مختلف، دارای ترجیحات یکسان (همگن) هستند (تفاوت‌ها از طریق متغیرها لحاظ شده است).

# مدل رگرسیون داده‌های گسسته

فصل  
یازدهم

## Discrete data Regression Model

ریسک و شانس	مقدمه
لوجیت چندگانه	پروبیت دوگانه
لوجیت شرطی	لوجیت دوگانه (لوجستیک)
لوجیت آشیانه‌ای و ترکیبی	معیارهای نیکویی برآزش

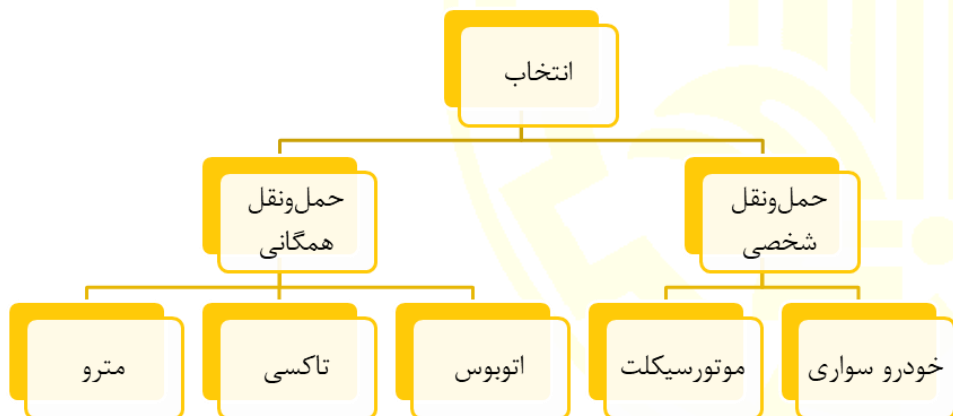


# مدل‌های لوجیت – لوجیت آشیانه‌ای

مدل **لوجیت آشیانه‌ای (Nested Logit)** برای مواردی طراحی شده است که فرض **استقلال گزینه‌های نامربوط (IIA)** بین گروهی از گزینه‌ها نقض می‌شود.

به عبارت دیگر، برخی گزینه‌ها به دلایل غیرقابل مشاهده (مانند راحتی، قابلیت دسترسی) با هم همبستگی دارند. لوجیت آشیانه‌ای این گزینه‌ها را در یک «آشیانه (nest)» قرار می‌دهد و اجازه می‌دهد خطاهای آنها با هم همبسته باشند. فرآیند تصمیم‌گیری به صورت درختی (دو سطحی یا چندسطحی) مدل می‌شود:

- **سطح بالا:** انتخاب بین آشیانه‌ها (مثلاً «حمل‌ونقل عمومی» در برابر «حمل‌ونقل شخصی»)
- **سطح پایین:** انتخاب بین گزینه‌های درون یک آشیانه (مثلاً اتوبوس در برابر مترو یا تاکسی)



به عبارت دیگر، در مدل لوجیت آشیانه‌ای (Nested Logit)، فرآیند انتخاب به صورت سلسله‌مراتبی (درختی) مدل می‌شود: ابتدا فرد تصمیم می‌گیرد کدام آشیانه (دسته از گزینه‌های مشابه) را انتخاب کند. سپس درون آن آشیانه، یک گزینه خاص را برمی‌گزیند.



روابط برای آشیانه  $m$  (که شامل  $l$  گزینه باشد):

احتمال بدون قید انتخاب گزینه  $j$  که در آشیانه  $m$  قرار دارد:  $P_{ij} = P(m) \times P(j|m)$

$$P(m) = \frac{e^{\lambda_m I_m}}{\sum_{k=1}^M e^{\lambda_k I_k}}$$

۱- احتمال انتخاب آشیانه  $m$

$M$ : تعداد آشیانه‌ها

$$I_m = \ln \sum_{l \in \text{Nest } m} e^{v_{il}/\lambda_m}$$

۲- مقدار شامل (Inclusive value) برای آشیانه  $m$

$$P(j|m) = \frac{e^{v_{ij}/\lambda_m}}{\sum_{l \in \text{Nest } m} e^{v_{il}/\lambda_m}}$$

۳- احتمال شرطی انتخاب گزینه  $j$  در آشیانه  $m$

در اسلایدهای بعد، توضیح مختصری درباره هر یک از این روابط ارائه شده است.



## روابط برای آشیانه $m$ (که شامل $l$ گزینه باشد):

احتمال بدون قید انتخاب گزینه  $z$  که در آشیانه  $m$  قرار دارد:  $P_{ij} = P(m) \times P(j|m)$

$$P(m) = \frac{e^{\lambda_m I_m}}{\sum_{k=1}^M e^{\lambda_k I_k}} \quad \text{۱- احتمال انتخاب آشیانه } m$$

- پارامتر  $\lambda_m$  پارامتر ناهمگونی (dissimilarity) در آشیانه  $m$  نام دارد.
- محدوده مجاز آن  $0 < \lambda_m \leq 1$  است. اگر  $\lambda_m = 1$ ، مدل به لوجیت ساده (بدون آشیانه) تبدیل می‌شود (فرض IIA برقرار است).
- هرچه  $\lambda_m$  کوچک‌تر باشد، همبستگی درون آشیانه بیشتر است (گزینه‌های درون آشیانه جانشین‌های نزدیکی برای یکدیگر هستند).
- در واقع به کمک جمله  $\lambda_m I_m$  سعی می‌کنیم **مطلوبیت سطح بالا (بین آشیانه‌ها)** برای آشیانه  $m$  را محاسبه کنیم.
- مقدار  $I_m$  لگاریتم مجموع مطلوبیت‌های گزینه‌های موجود در آشیانه  $m$  (**مطلوبیت سطح پایین**) را محاسبه می‌کند.



# مدل‌های لوجیت – لوجیت آشیانه‌ای

## روابط برای آشیانه $m$ (که شامل $l$ گزینه باشد):

احتمال بدون قید انتخاب گزینه  $l$  که در آشیانه  $m$  قرار دارد:  $P_{ij} = P(m) \times P(j|m)$

۲- مقدار شامل (Inclusive value) برای آشیانه  $m$

$$I_m = \ln \sum_{l \in Nest\ m} e^{V_{il}/\lambda_m}$$

**مقدار شامل:** معیاری است که مطلوبیت مورد انتظار فرد از انتخاب بهترین گزینه درون آشیانه  $m$  را خلاصه می‌کند. به عبارت دیگر، این مقدار نشان می‌دهد که آشیانه  $m$  به طور کلی چقدر برای فرد جذاب است، با در نظر گرفتن تمام گزینه‌های درون آن و ویژگی‌هایشان.

**تفسیر شهودی:** فرض کنید می‌خواهید تصمیم بگیرید که با شیوه حمل‌ونقل عمومی جابجا شوید یا حمل‌ونقل شخصی (دو آشیانه). هر آشیانه شامل گزینه‌های مختلفی است:

- آشیانه «حمل‌ونقل عمومی»: تاکسی، اتوبوس، مترو.
- آشیانه «حمل‌ونقل شخصی»: پیاده، موتورسیکلت، خودرو.

شما ابتدا نمی‌دانید کدام آشیانه را انتخاب کنید. اما اگر درون آشیانه حمل‌ونقل عمومی، یک گزینه (مثلاً مترو) خیلی عالی باشد، آن آشیانه برای شما جذاب‌تر می‌شود، حتی اگر بقیه گزینه‌های آن متوسط باشند.

مقدار شامل دقیقاً همین «جذابیت کلی» را اندازه می‌گیرد، که ترکیبی از مطلوبیت همه گزینه‌های درون آشیانه (با وزن نمایی) است. ✓ هر چه بهترین گزینه درون آشیانه مطلوب‌تر باشد، یا هر چه تعداد گزینه‌های خوب در آن آشیانه بیشتر باشد،  $I_m$  بزرگتر خواهد بود.



روابط برای آشیانه  $m$  (که شامل  $l$  گزینه باشد):

احتمال بدون قید انتخاب گزینه  $j$  که در آشیانه  $m$  قرار دارد:  $P_{ij} = P(m) \times P(j|m)$

$$P(j|m) = \frac{e^{V_{ij}/\lambda_m}}{\sum_{l \in \text{Nest } m} e^{V_{il}/\lambda_m}}$$

۳- احتمال شرطی انتخاب گزینه  $j$  در آشیانه  $m$



## مثال کاربردی

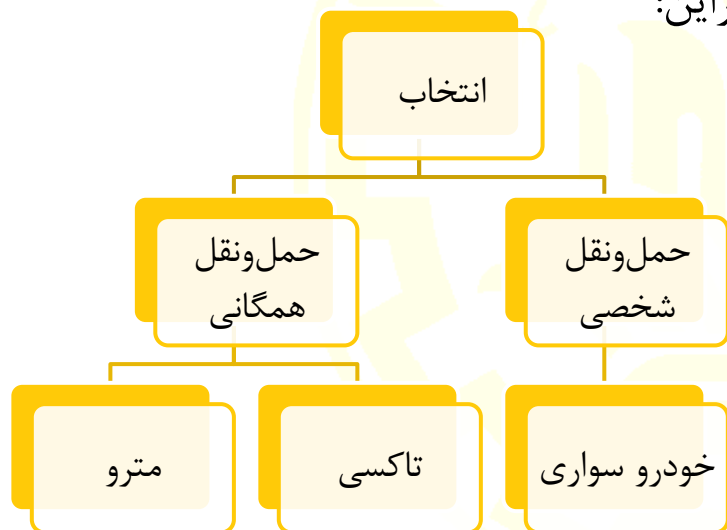
فرض کنید افراد بین سه گزینه برای رفتن به فرودگاه انتخاب می‌کنند:

(۱) خودرو شخصی، (۲) تاکسی، (۳) مترو.

تاکسی و مترو هر دو جزو «حمل‌ونقل عمومی» محسوب می‌شوند و ممکن است خطاهای آن‌ها به دلیل عواملی مانند «پرهیز از رانندگی» یا «دسترسی به ایستگاه» با هم همبسته باشد. بنابراین:

آشیانه ۱: خودرو شخصی (تک گزینه)

آشیانه ۲: حمل‌ونقل عمومی (شامل تاکسی و مترو)



با برآورد مدل، اگر  $\lambda_2 = 0.5$  به دست آید، یعنی همبستگی نسبتاً بالایی بین تاکسی و مترو وجود دارد و لوجیت ساده (که این همبستگی را نادیده می‌گیرد) نتایج اریب‌دار خواهد داد.



مدل لوجیت ترکیبی (Mixed Logit) یا لوجیت با پارامترهای تصادفی (Random Parameters Logit)

منعطف‌ترین و پیشرفته‌ترین نسخه مدل لوجیت است که برای مدلسازی رویدادهای انتخاب گسسته استفاده می‌شود.

این مدل دو محدودیت اصلی لوجیت شرطی را برطرف می‌کند:

1. **نقض فرض IIA**: اجازه می‌دهد همبستگی بین گزینه‌های مشابه وجود داشته باشد.

2. **لحاظ نمودن ناهمگونی مشاهده‌نشده در ترجیحات**: اجازه می‌دهد ضرایب مدل (میزان اهمیت هر

ویژگی) در بین افراد جامعه متفاوت باشد (تا از این طریق، تفاوت‌های پنهان میان افراد (مانند توانایی

بدنی) نیز در مدلسازی ترجیحات آن‌ها لحاظ گردد). زیرا این تفاوت‌ها می‌تواند موجب شود که

ارزش‌گذاری یک گزینه بین افراد، متفاوت باشد.



$U_{ij}$  را مطلوبیت انتخاب گزینه  $j$  برای فرد  $i$  و به صورت روبرو در نظر می‌گیریم.

$$U_{ij} = x'_{ij}\beta_i + \varepsilon_{ij}$$

که در آن، بردار ویژگی‌های گزینه‌ها و افراد، و  $\varepsilon_{ij}$  جمله خطاست (با توزیع مقدار حدی گامبل).

✓ تفاوت مهم این است که در اینجا،  $\beta_i$  ها بردار ضرایب تصادفی مربوط به فرد  $i$  بوده و از یک توزیع مشخص  $f(\beta|\theta)$  مانند توزیع نرمال (یا لگ‌نرمال، یکنواخت، و غیره) پیروی می‌کند:

$$\beta_i \sim f(\beta|\theta)$$

✓  $\theta$  پارامترهای توزیع (مانند میانگین و واریانس است).

$$\beta_i = \mu + \eta_i ; \eta_i \sim N(0, \sigma^2)$$

احتمال انتخاب گزینه  $j$  توسط فرد  $i$  در مدل لوجیت ترکیبی، برابر است با انتگرال احتمال‌های لوجیت استاندارد بر روی تمامی حالت‌های ممکن برای مقدار  $\beta$  که براساس چگالی احتمال مقادیر  $\beta$  وزن دار شده است:

$$P_{ij} = \int \frac{e^{x'_{ij}\beta_i}}{\sum_{k=1}^m e^{x'_{ij}\beta_i}} f(\beta|\theta) d\beta$$

که در آن،  $f(\beta|\theta)$  تابع چگالی مقادیر  $\beta$  است.

این انتگرال یک پاسخ به صورت فرم بسته (closed form) ندارد و باید با استفاده از روش‌های شبیه‌سازی مانند مونت کارلو آن را تقریب زد.



## مثال کاربردی

یک مدل انتخاب شیوه سفر را به صورت زیر در نظر بگیرید. در اینجا، از یک مدل لوجیت ترکیبی (با پارامتر تصادفی) استفاده شده تا ضرایب مربوط به دو متغیر زمان سفر و هزینه سفر را به صورت پارامتر تصادفی (و نه پارامتر ثابت) مدل کند:

$$U_{ij} = \beta_i^{time} T_{ij} + \beta_i^{cost} C_{ij} + \varepsilon_{ij}$$
$$\beta_i^{time} \sim N(\mu_{time}, \sigma_{time}^2)$$
$$\beta_i^{cost} \sim N(\mu_{cost}, \sigma_{cost}^2)$$

در اینجا فرض شده است که مقدار پارامترهای دو متغیر، به صورت تصادفی (براساس یک توزیع نرمال) در میان افراد جامعه تغییر می‌کند.

یعنی، مثلاً همه افراد به طور متوسط، ترجیح می‌دهند هزینه سفر کمتر باشد.

اما میزان اهمیت متغیر هزینه، می‌تواند در میان افراد، متفاوت باشد.

امید ریاضی توزیع ( $\mu$ ) میانگین میزان حساسیت افراد به هر متغیر را نشان می‌دهد.

واریانس توزیع ( $\sigma^2$ ) میزان پراکندگی این ترجیحات در میان جامعه را نشان می‌دهد.