



## 2. Statistical Modelling (4)

Statistical Modelling & Machine Learning

Jaejik Kim

Department of Statistics  
Sungkyunkwan University

STA3036

# Models for Discrete or Non-normal $Y$ Variables

- ▶ Classical regression models assume continuous  $Y$  and normal error distribution.
- ▶ When  $Y$  is a discrete random variable or it does not have normal distribution, classical regression models do not work properly.
  - ▶ The range of  $\mu = E(Y) = \mathbf{X}\beta$ .
  - ▶ Statistical inference due to normality assumption. ] normality assumption을 만족하지 못하면 통계적 추론을 할 수가 없다.
- ▶ E.g., suppose that  $Y$  has Bernoulli response  $Y_i = 0$  or  $1$ .  
 $\mu_i = E(Y_i) = P(Y_i = 1) \in [0, 1]$   
 $Var(Y_i) = \mu_i(1 - \mu_i)$  (not constant). ) violates the normality assumption

# Generalized Linear Model

⇒ normality assumption이 충족되지 않았을 때 사용하는 Model

- ▶ Generalized linear model (GLM): Extension of classical linear model. <sup>\*</sup> GLM에서 MLE를 통해 parameters를 구하게 되면 closed form으로 나오지 않기 때문에 numerical method로 찾아야 한다.

※ 범주형 자료분석에서  
다용

- ▶ 3 components of GLM:

1. **Systematic component:**  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ . :  $i$ 번째 observation마다 분포의 parameter가 바뀌는 의미.
2. **Random component:**  $Y_i$ 's are independent random variables with  $E(Y_i) = \mu_i$  and pdf (pmf) in the exponential family as follows:

GLM에서 필요한 조건이나 상환들이 exponential family에 대해서만 적용가능하다.

$$p(y_i; \theta_i, \phi) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{a_i(\phi)} + c(y_i, \phi) \right\}, \quad (1)$$

- ▶  $\theta_i$ : Location parameter (usually our interest). :  $\mathcal{M}_i$ 에 대한 함수로 표현될 수 있음
- ▶  $\theta_i$  can be expressed as some function of  $\mu_i = E(Y_i)$ .
- ▶  $\phi$ : Scale parameter (nuisance parameter).

위 2개를 연결해주는 함수 3.

- 3. **Link function:** The link between the systematic and random components.

$$g(\mu_i) = \eta_i, \quad \mathcal{M}_i \text{가 } \theta_i \text{에 대한 함수이므로, } \eta_i \text{와 } \theta_i \text{를 연결시킨다.}$$

where  $g$  is one-to-one and differentiable. <sup>\*</sup> random component의 range를 systematic component에 맞추는 역할

# Exponential Family

- ▶ Exponential family: A set of distributions whose pdf (pmf) satisfies the format of (1).

- ▶ Distributions in Exponential family: Normal, exponential, Bernoulli, binomial, Poisson, gamma, geometric, etc.

GLM 사용 가능

Negative Binomial, ...

- ▶ E.g., Normal distribution:  $Y_i \sim^{indep.} N(\mu_i, \sigma^2)$ .

nuisance parameter

$$\begin{aligned} p(y_i; \mu_i, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_i - \mu_i)^2}{2\sigma^2} \right\} \\ &= \exp \left\{ \left[ y_i \mu_i - \frac{\mu_i^2}{2} \right] \frac{1}{\sigma^2} - \frac{y_i^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right\} \end{aligned}$$

- ▶  $\theta_i = \mu_i, \phi = \sigma^2,$
  - ▶  $a_i(\phi) = \phi, b(\theta_i) = \theta_i^2/2, c(y_i, \phi) = -[y_i^2/\phi + \log(2\pi\phi)]/2.$

# Exponential Family

- E.g., Binomial distribution:  $Y_i \sim^{indep.} \text{Binom}(n_i, p_i)$ .

$$\begin{aligned} p(y_i; p_i) &= \binom{n_i}{y_i} p_i^{y_i} (1 - p_i)^{n_i - y_i} \\ &= \binom{n_i}{y_i} p_i^{y_i} (1 - p_i)^{n_i} (1 - p_i)^{-y_i} \\ &= \binom{n_i}{y_i} \left( \frac{p_i}{1 - p_i} \right)^{y_i} (1 - p_i)^{n_i} \\ &= \exp \left[ y_i \log \left( \frac{p_i}{1 - p_i} \right) + n_i \log (1 - p_i) + \log \binom{n_i}{y_i} \right] \end{aligned}$$

$$\begin{aligned} \therefore \theta_i &= \log \left( \frac{p_i}{1 - p_i} \right) \\ \emptyset &= 1 \\ a_i(\emptyset) &= \emptyset \\ b(\theta_i) &= n_i \log (1 + e^{\theta_i}) \\ c(y_i, \emptyset) &= \log \binom{n_i}{y_i} \end{aligned}$$

# Link Function

- ▶ For  $Y_i$  with a certain distribution in the exponential family, various link functions exist.  
*이중에 알맞는걸 선택*
- ▶ E.g., for a binary random variable  $Y$ , we want to map  $\mathbb{R}$  ( $\eta = \mathbf{x}\beta$ ) to  $[0, 1]$  (the range of  $\mu = E(Y)$ ).
  - ▶ All cdf can map  $\mathbb{R}$  to  $[0, 1]$ . *cdf의 구간은  $(-\infty, \infty)$ 이고, 값의 범위는  $[0, 1]$ 이므로*

$$\mu = F(\eta) \Rightarrow F^{-1}(\mu) = \eta.$$

★ 가장 유명

- ▶ **Logit link:**  $\log \frac{\mu}{1-\mu}$  (inverse of unit logistic cdf).
- ▶ **Probit link:**  $\Phi^{-1}(\mu) = \eta$  (inverse of standard normal cdf).
- ▶ **log-log link:**  $\log(-\log(\mu)) = \eta$  (inverse of Gumbel cdf).

# Canonical Link Function

- ▶ For each distribution, there is one link function that is mathematically convenient  $\Rightarrow$  Canonical link function.  
주어진 분포에 대해서 많이 사용되는 link function
- ▶ Canonical Link:  $\theta = \eta$ .  $\theta_i(\mu_i) = \eta_i$

Distribution	Canonical Link
Normal distribution	$g(\mu) = \mu$
Bernoulli distribution	$g(\mu) = \log(\mu/(1 - \mu))$
Poisson distribution	$g(\mu) = \log \mu$
Gamma distribution	$g(\mu) = \mu^{-1}$

- ▶ Choice of link should be made on
  - ▶ model fit,
  - ▶ model interpretability,
  - ▶ mathematical convenience (canonical link or not).



# Maximum Likelihood Estimation in GLM

Model parameter를 Estimate 하는 단계

- ▶ For independent  $Y_1, \dots, Y_n$ , the log-likelihood function is

$$l(\theta; \mathbf{y}) = \sum_{i=1}^n \log p(y_i; \theta_i),$$

- ▶  $\theta_i = \theta_i(\mu_i)$ .
- ▶  $g(\mu_i) = \eta_i = \mathbf{x}_i^\top \beta$ . ↑ "MLE의 함수도 MLE다"
- ▶ By the invariance property of MLE, MLE of  $\mu_i$  and  $\theta_i$  can be obtained by the MLE  $\hat{\beta}$  of the model parameters.

$$\hat{\beta} = \max_{\beta} l(\beta; \mathbf{y}).$$

- ▶ No closed-form solution exists  $\Rightarrow$  Numerical method (Newton-Raphson or Fisher scoring).

# Logistic Regression

- ▶ Suppose that  $Y$  is a binary output variable.
- ▶ Canonical link function:  $g(\mu_i) = \frac{\log(\mu_i/(1 - \mu_i))}{= \theta_i}$ , where  $\mu_i = E(Y_i) = p_i$ .

$$\log \frac{p_i}{1 - p_i} = \frac{\mathbf{x}_i^\top \boldsymbol{\beta}}{= \eta_i},$$

where  $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip})^\top$  &  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^\top$ .

- ▶  $p_i = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \Rightarrow \frac{p_i}{1 - p_i} = \exp(\mathbf{x}_i^\top \boldsymbol{\beta}).$

# Estimation of $\beta$

- Likelihood function:

$$\begin{aligned}\max_{\beta} L(\beta; \mathbf{y}) &= \max_{\beta} \left[ \prod_{i:y_i=1} p_i \prod_{i':y_{i'}=0} (1 - p_{i'}) \right] \\ &= \max_{\beta} \prod_{i:y_i=1} \frac{e^{\mathbf{x}_i^{\top} \beta}}{1 + e^{\mathbf{x}_i^{\top} \beta}} \prod_{i':y_{i'}=0} \frac{1}{1 + e^{\mathbf{x}_{i'}^{\top} \beta}}.\end{aligned}$$

⇒ Numerical method (Iteratively Reweighted Least Squares)

⇒  $\hat{\beta}$  (MLE).

- $\hat{\beta} > 0, p(x) \uparrow$ ,  
 $\hat{\beta} < 0, p(x) \downarrow$ . (Not linear relationship).

# Multinomial Logistic Regression

Multinomial Distribution  $Y \rightarrow$  nominal categorical  
Multinomial Logistic Regression을 사용하자

- ▶ Logistic regression with  $K$  classes ( $K > 2$ )  $\Rightarrow$  Multinomial logistic regression:

$$\log \frac{\mu_k}{\mu_K} = \ln \frac{P(Y = k | X = x)}{P(Y = K | X = x)} = x^\top \beta_k$$

small "k" above  $\mu_k$ , capital "K" below  $\mu_K$

for  $k = 1, \dots, K - 1$  with  $\sum_{k=1}^K P(Y = k | X = x) = 1$ .  
 $k$ 가 변할 때마다  $\beta_k$ , 식의 계수가 변한다.

- ▶  $x = (1, x_1, \dots, x_p)^\top$  &  $\beta_k = (\beta_{k0}, \beta_{k1}, \dots, \beta_{kp})^\top$ .
- ▶ The choice of denominator, the  $K$ th class, is arbitrary.

# Multinomial Logistic Regression

- By solving for  $P(Y = k|X = x)$ , we have

$$P(Y = k|X = x) = \frac{\exp(x^\top \beta_k)}{1 + \sum_{j=1}^{K-1} \exp(x^\top \beta_j)}$$

for  $k = 1, \dots, K - 1$  and

$$P(Y = K|X = x) = \frac{1}{1 + \sum_{j=1}^{K-1} \exp(x^\top \beta_j)}.$$

$\Rightarrow$  ML estimation (numerical method)

$\Rightarrow \hat{\beta}_k, k = 1, \dots, K - 1$  (MLE).

# Ordinal Response: Cumulative Logit Model

- ▶ Ordinal data: Categories are ordered (e.g., good, medium, bad).
- ▶ Suppose that response  $Y$  takes ordered category values  $k = 1, \dots, K$ , let  $p_k = P(Y = k | \mathbf{X})$ .  
*Y가 k번째 category에 속할 확률*
- ▶ Cumulative probability:

$$\gamma_k = \sum_{j=1}^k p_j = P(Y \leq k | \mathbf{X}), \quad k = 1, \dots, K.$$

*small case k*      *Capital case K*       $\therefore \gamma_K = 1$  Since it cumulates all  $K$  possibilities

- ▶ Cumulative logit: 일반적인 Logistic Regression에서는  $p_i$ 를 활용해서 Log-odds를 만들었는데, Cumulative Logit에서는 Log-odds를 만드는 대신에  $\gamma_k$ 이다.  
*부*  $p_i + (1 - p_i) = \gamma_k + (1 - \gamma_k) = 1$

$$\log \frac{\gamma_k}{1 - \gamma_k} = \log \frac{p_1 + \dots + p_k}{p_{k+1} + \dots + p_K}, \quad k = 1, \dots, K - 1.$$

*k<sup>th</sup> order 보다 큰 observations들의 합*

# Ordinal Response: Cumulative Logit Model

## ► Cumulative Logit Model (Proportional odds model):

$$\log \frac{\gamma_{ik}}{1 - \gamma_{ik}} = \alpha_k + \mathbf{x}_i^\top \boldsymbol{\beta}, \quad i = 1, \dots, n, \quad k = 1, \dots, K - 1.$$

$\alpha_k$  is dependent on  $k$ .  
 $\boldsymbol{\beta}$  coefficients are the same for all  $k$ .  
 $\gamma_{ik}$  is the probability of being in category  $k$  or smaller.

독립된  $\beta$ 를 가정하지 않으면 order가 역전되지 않게 하기 위한 constraints들이 많이 생기고, 이러한 constraints들이 많아서 optimize하기 상당히 까다롭다.  
 동일한  $\beta$ 를 갖는 이유: 경해진 order가 역전되지 않게 하기 위해서 기울기를 고정하고  $\alpha$ 만 움직이게 해야 한다.

- $\alpha_k$  is increasing in  $k$  because  $\gamma_k$  is increasing in  $k$  for fixed  $\mathbf{x}$ .
- This model has the same effects  $\boldsymbol{\beta}$  for each logit model (The  $K - 1$  logistic curves have the same shape).

형태는 모두 동일하되 위치만 다르다.

## ► Two observations with input vector $\mathbf{x}_1$ and $\mathbf{x}_2$ , respectively.

$$\log \frac{\gamma_{1k}}{1 - \gamma_{1k}} - \log \frac{\gamma_{2k}}{1 - \gamma_{2k}} = \log \frac{\gamma_{1k}/(1 - \gamma_{1k})}{\gamma_{2k}/(1 - \gamma_{2k})} = (\mathbf{x}_1 - \mathbf{x}_2)^\top \boldsymbol{\beta}.$$

$\mathbf{x}_1, \mathbf{x}_2$  모두  $k$ 번째 category에 속하기 때문에 동일한  $\alpha_k$ 값을 가지므로 소거된다.

- Log cumulative odds ratio does not depend on  $k$ , only on  $(\mathbf{x}_1 - \mathbf{x}_2)$ .
- The ratio of odds of being in the  $k$ th or smaller category under two different inputs is the same for all categories.  $\Rightarrow$  Proportional odds model.

# Maximum Likelihood Estimation

~ one-hot encoding

- ▶ Let  $\mathbf{y}_i = (y_{i1}, \dots, y_{iK})^\top$ ,  $i = 1, \dots, n$ , where  $y_{ik} = 1$  if the  $i$ th obs. is in the  $k$ th ordered category. Otherwise,  $y_{ik} = 0$ .
- ▶ Likelihood function:

$$\begin{aligned} L(\boldsymbol{\alpha}, \boldsymbol{\beta}; \mathbf{y}) &= \prod_{i=1}^n \left[ \prod_{k=1}^K p_k^{y_{ik}} \right] = \prod_{i=1}^n \left[ \prod_{k=1}^K (\gamma_k - \gamma_{k-1})^{y_{ik}} \right] \\ &= \prod_{i=1}^n \left[ \prod_{k=1}^K \left\{ \frac{\exp(\alpha_k + \mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\alpha_k + \mathbf{x}_i^\top \boldsymbol{\beta})} - \frac{\exp(\alpha_{k-1} + \mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\alpha_{k-1} + \mathbf{x}_i^\top \boldsymbol{\beta})} \right\}^{y_{ik}} \right]. \end{aligned}$$

$$\begin{aligned} \log \frac{\eta_k}{1 - \eta_k} &= \alpha_k + \mathbf{x}_i^\top \boldsymbol{\beta} \\ \Rightarrow \eta_k &= \frac{\exp(\alpha_k + \mathbf{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\alpha_k + \mathbf{x}_i^\top \boldsymbol{\beta})} \end{aligned}$$



# Poisson Regression

- ▶  $Y$ : Count data  $\Rightarrow Y_1, \dots, Y_n \sim^{indep.} \text{Poisson}(\mu_i)$ .
- ▶ Canonical link function:  $g(\mu_i) = \log(\mu_i) = \eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ .
- ▶ Model:  $\log \mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ ,  $i = 1, \dots, n$ .
- ▶ Log-likelihood function:

$$\begin{aligned} l(\boldsymbol{\beta}; \mathbf{y}) &= \sum_{i=1}^n \overbrace{[y_i \log(\mu_i) - \mu_i - \log(y_i!)]}^{\text{exponential family form of poisson distribution}} \\ &= \sum_{i=1}^n \left[ y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \exp(\mathbf{x}_i^\top \boldsymbol{\beta}) - \log(y_i!) \right]. \end{aligned}$$

# Poisson Regression

포아송의 단점!

- ▶ Overdispersion: Data variation is higher than model's expectation.
- ▶ Overdispersion in Poisson regression is typical because  $E(Y_i) = Var(Y_i) = \mu_i$ .  
각의 모든 상황에서 분산은 평균보다 큰데 포아송은 평균과 분산이 동일해서
- ▶ To solve the overdispersion problem, the negative binomial model can be considered.  $Y \sim NB(\mu, \alpha)$

$$p(y) = \frac{\Gamma(y + \alpha^{-1})}{y! \Gamma(\alpha^{-1})} \left( \frac{\alpha \mu}{1 + \alpha \mu} \right)^y \left( \frac{1}{1 + \alpha \mu} \right)^{\alpha^{-1}}, \quad y = 0, 1, 2, \dots$$

- ▶  $E(Y) = \mu$  and  $Var(Y) = \mu + \alpha \mu^2$ .
- ▶ Negative binomial model:  $\log \mu_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ ,  $i = 1, \dots, n$ .

# Survival Model

- ▶ Survival data:  $T$  is the survival time until death or failure.  
*살아있는 시간, 기간*
- ▶ Censoring: Property of survival data
  - ▶ It occurs when the outcome of a particular patient or component is unknown at the end of the study.  
*=> 피험자의 중도 포기 또는 관측된 관측기간 때문에 실험의 끝을 알 수 없는 경우.*
- ▶ Let the survival time  $T$  have a pdf  $f(t)$  and the cdf  $F(t)$ .
  - ▶  $F(t)$ : The fraction of the population dying by time  $t$ .
  - ▶  $1 - F(t)$ : Survival function (fraction still surviving at time  $t$ ).
  - ▶  $h(t)$ : Hazard function (instantaneous risk). *이에서의 순간사망률 (사력...?)*
  - ▶  $h(t)dt$ : Prob. of dying in the next small time interval  $dt$  given survival to time  $t$ . *이까지 살다가 이가 되자마자 사망할 확률*

$$h(t)dt = P(T \in [t, t + dt] | T > t) = \frac{f(t)dt}{1 - F(t)}$$

*이 작은 틈새 사망*      *생존시간 گذشته 한바 보다*

$$\Rightarrow h(t) = \frac{f(t)}{1 - F(t)}$$

# Proportional Hazard Model

- ▶ Proportional hazard model:

★  $X^T \beta = 0$  이면  $h(t|x) = \lambda(t)$  가 되므로,  
 $\lambda(t)$ 를 Baseline Hazard라고 칭한다.

$$h(t|x) = \lambda(t) \exp(\mathbf{x}^T \beta).$$

- ▶ Under this model, consider two observations with  $\mathbf{x}_1$  and  $\mathbf{x}_2$ , respectively.

$$\frac{h(t|\mathbf{x}_1)}{h(t|\mathbf{x}_2)} = \exp[(\mathbf{x}_1 - \mathbf{x}_2)^T \beta].$$

=> Independent from  $t$ . 위험률의 비율은 시간에 상관없이 일정하다.

- ▶ Proportional hazard: This ratio does not depend on  $t$ .
- ▶ From the proportional hazard model,

$$h(t) = f(t)/[1 - F(t)] = \lambda(t) \exp(\mathbf{x}^T \beta),$$

by taking integral on both sides,

$$-\log[1 - F(t)] = \Lambda(t) \exp(\mathbf{x}^T \beta),$$

우가 없으므로 상수항에 되면서 그대로 나옴

where  $\Lambda(t) = \int_{-\infty}^t \lambda(u) du$  (**Cumulative hazard**).

# ML Estimation of PH Model

- ▶ Survival function:

$$S(t) = 1 - F(t) = \exp\{-\Lambda(t) \exp(\mathbf{x}^\top \boldsymbol{\beta})\}.$$

- ▶ By minus derivative w.r.t.  $t$ ,

$$f(t) = \lambda(t) \exp\{\mathbf{x}^\top \boldsymbol{\beta} - \Lambda(t) \exp(\mathbf{x}^\top \boldsymbol{\beta})\}.$$

- ▶ Likelihood function:

- ▶ An object who died at time  $t$  contributes a factor  $f(t)$  to the likelihood. ] 각 t에서 관측된 죽음들은 f(t)만큼 기여했다.
- ▶ An object who censored at time  $t$  contributes  $S(t)$ . ] 각 t에서 censored된 사람은 S(t)만큼 기여했다.
- ▶ If the  $i$ th observation is died at time  $t$ ,  $w_i = 1$ . Otherwise,  $w_i = 0$ .

## ► Log-likelihood function:

$$\begin{aligned}
 \underline{l(\beta; \mathbf{t}, \mathbf{w})} &= \sum_{i=1}^n [w_i \log \underline{f(t_i)} + (1 - w_i) \log \underline{S(t_i)}] \\
 &= \sum_{i=1}^n \left[ w_i \{ \log \lambda(t_i) + \mathbf{x}_i^\top \beta \} - \Lambda(t_i) \exp(\mathbf{x}_i^\top \beta) \right].
 \end{aligned}$$

$\beta$ 에 대한 추정론 parametric method인데,  
 $\lambda(t)$ 에 대한 추정론 non-parametric 해서  
 $l(\beta; \mathbf{t}, \mathbf{w})$ 를 semi-parametric 이라 한다.

$\lambda(t_i)$ 를 추정해야 하는데,  $\lambda(t_i)$ 는 non-parametric 회귀 추정된다.

⇒ 모델의 해석 :  $X$ 가 치료여부에 대한 binary data인 경우, 만약  $\frac{h(t|X_1)}{h(t|X_2)} = \exp[(X_1 - X_2)^\top \beta]$ 에서  $X_1 = 1, X_2 = 0, \beta > 0$  라면  
 치료로 인한 발생률의  $e^\beta$  배로 증가된다고 할 수 있다.