

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.
- ▶ 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).

Generalized Linear Model

- ▶ Generalized linear model (GLM): Extension of classical linear model.
- ▶ 3 components of GLM:
 1. **Systematic component**: $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$.
 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:

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

- ▶ θ_i : Location parameter (usually our interest).
 - ▶ θ_i can be expressed as some function of $\mu_i = E(Y_i)$.
 - ▶ ϕ : Scale parameter (nuisance parameter).
- 3. **Link function**: The link between the systematic and random components.

$$g(\mu_i) = \eta_i,$$

where g is one-to-one and differentiable.

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

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

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

Exponential Family

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

$$\begin{aligned} p(y_i; p_i) &= \binom{n_i}{y_i} p_i^{y_i} (1 - p_i)^{n_i - y_i} \\ &= \end{aligned}$$

- ▶ For Y_i with a certain distribution in the exponential family, various link functions exist.
- ▶ 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 $[0, 1]$.

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

- ▶ 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.
- ▶ Canonical Link: $\theta = \eta$.

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

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

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

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

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

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

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

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

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

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

- 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

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

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

for $k = 1, \dots, K - 1$ with $\sum_{k=1}^K P(Y = k|X = x) = 1$.

- ▶ $x = (1, x_1, \dots, x_p)^\top$ & $\beta_k = (\beta_{k0}, \beta_{k1}, \dots, \beta_{kp})^\top$.
- ▶ 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 $p_k = P(Y = k|\mathbf{X})$.
- ▶ Cumulative probability:

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

- ▶ Cumulative logit:

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

Ordinal Response: Cumulative Logit Model

- ▶ Cumulative Logit Model (Proportional odds model):

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

- ▶ α_k is increasing in k because γ_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).
- ▶ Two observations with input vector \mathbf{x}_1 and \mathbf{x}_2 , respectively.

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

- ▶ 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

- ▶ 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:

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

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:

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

Poisson Regression

- ▶ Overdispersion: Data variation is higher than model's expectation.
- ▶ Overdispersion in Poisson regression is typical because $E(Y_i) = \text{Var}(Y_i) = \mu_i$.
- ▶ To solve the overdispersion problem, the negative binomial model can be considered. $Y \sim \text{NB}(\mu, \alpha)$

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

- ▶ $E(Y) = \mu$ and $\text{Var}(Y) = \mu + \alpha \mu^2$.
- ▶ 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.
- ▶ Censoring: Property of survival data
 - ▶ It occurs when the outcome of a particular patient or component is unknown at the end of the study.
- ▶ 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).
 - ▶ $h(t)d(t)$: Prob. of dying in the next small time interval $d(t)$ given survival to time t .

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

Proportional Hazard Model

- ▶ Proportional hazard model:

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

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

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

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

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

by taking integral on both sides,

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

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

- ▶ 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.
- ▶ An object who censored at time t contributes $S(t)$.
- ▶ If the i th observation is died at time t , $w_i = 1$. Otherwise, $w_i = 0$.

- ▶ Log-likelihood function:

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