# Lec 4: Generalized Linear Model

Yanjun Han

Sept. 26, 2023

# Generalized Linear model.

<u>Setting</u>. For $i = 1, 2, \cdots, n$, let $y_i \overset{ind}{\sim} p_{\theta_i}(y_i) = \exp\left(\langle \theta_i, T(y_i)\rangle - A(\theta_i)\right)h(y_i)$,

where $\theta_i = \left(\langle x_i, \beta_1\rangle, \langle x_i, \beta_2\rangle, \cdots, \langle x_i, \beta_d\rangle\right) \in \mathbb{R}^d$

- $x_i \in \mathbb{R}^p$ : feature / covariate
- $(\beta_1, \cdots, \beta_d) \in \mathbb{R}^{p \times d}$ : regression coefficients
- written in matrix form : $\theta_i = \beta^T x_i$

<u>MLE</u>.

$$\hat{\beta} = \underset{\beta}{\arg\max} \ \prod_{i=1}^{n} p_{\theta_i}(y_i)$$

$$= \underset{\beta}{\arg\max} \ \sum_{i=1}^{n} \left(\langle \beta^T x_i, T(y_i)\rangle - A(\beta^T x_i)\right)$$

$$= \underset{\beta}{\arg\max} \ \underbrace{\text{Tr}\left(\sum_{i=1}^{n} T(y_i) x_i^T \cdot \beta\right)}_{\text{linear in } \beta} - \underbrace{\sum_{i=1}^{n} A(\beta^T x_i)}_{\text{convex in } \beta}$$

Estimating equation $(d=1)$ : $\sum_{i=1}^{n} T(y_i) x_i = \sum_{i=1}^{n} A'(\hat{\beta}^T x_i) x_i$.

The computation of MLE is a convex problem, thus efficient.

$$\boxed{\text{In R: } \quad \text{model} \leftarrow \text{glm}( y \sim X, \text{family}).}$$

<u>Examples</u>. 1. Linear regression.

$$y_i \sim N(\theta_i, 1) = N(\beta^T x_i, 1) \qquad \textcolor{red}{\mathbb{R}^{n \times p}}$$
$$\Rightarrow \hat{\beta} = \underset{\beta}{\arg\min} \ \sum_{i=1}^{n} (y_i - \beta^T x_i)^2 = \underset{\beta}{\arg\min} \ \|y - X\beta\|_2^2$$

2. Logistic regression.

$$y_i \sim \text{Bern}\left(\frac{1}{1 + e^{-\theta_i}}\right) = \text{Bern}\left(\frac{1}{1 + e^{-\beta^T x_i}}\right)$$
$$\Rightarrow \hat{\beta} = \underset{\beta}{\arg\max} \ \sum_{i=1}^{n} \left(y_i \log \frac{1}{1 + e^{-\beta^T x_i}} + (1 - y_i)\log \frac{e^{-\beta^T x_i}}{1 + e^{-\beta^T x_i}}\right)$$
$$= \underset{\beta}{\arg\max} \ \sum_{i=1}^{n} \left(y_i \beta^T x_i - \log(1 + e^{\beta^T x_i})\right)$$

## 2'. Probit model.

$$Y_i \sim \text{Bern}(\Phi(\theta_i)) = \text{Bern}(\Phi(\beta^T x_i)),$$

where $\Phi$ is the standard normal CDF:

$$\Phi(t) = \int_{-\infty}^{t} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

MLE:
$$\hat{\beta} = \underset{\beta}{\text{argmax}} \sum_{i=1}^{\hat{n}} \left( Y_i \log \Phi(\beta^T x_i) + (1-Y_i) \log (1-\Phi(\beta^T x_i)) \right)$$

**Lemma**. The above objective is concave in $\beta$.

**Pf**. For $f(x) = \log \Phi(x)$:
$$f'(x) = \frac{\varphi(x)}{\Phi(x)}, \quad f''(x) = \frac{\varphi' \Phi - \varphi^2}{\Phi^2} = -\frac{(x\Phi + \varphi)\varphi}{\Phi^2}.$$

Gaussian Mills ratio:
$$1 - \Phi(x) < \frac{\varphi(x)}{x}, \quad x > 0$$
$$\implies x\Phi(x) + \varphi(x) > 0, \ x < 0 \implies f''(x) < 0.$$

> In an exponential family, there could be more than
> one parametrizations such that the MLE computation
> in the corresponding GLM is a convex problem.

## 3. Poisson regression.

$$Y_i \sim \text{Poi}(e^{\theta_i}) = \text{Poi}(e^{\beta^T x_i})$$
$$\implies \hat{\beta} = \underset{\beta}{\text{argmax}} \sum_{i=1}^{\hat{n}} (T(Y_i)\beta^T x_i - A(\beta^T x_i))$$
$$= \underset{\beta}{\text{argmax}} \sum_{i=1}^{\hat{n}} (Y_i \beta^T x_i - e^{\beta^T x_i}).$$

## 4. Multinomial logit regression.

Recall that
$$\theta = (\theta_1, \cdots, \theta_k)$$
$$T(y) = (\mathbb{1}(y=1), \mathbb{1}(y=2), \cdots, \mathbb{1}(y=k))$$
$$A(\theta) = \log(e^{\theta_1} + \cdots + e^{\theta_k})$$

Model : $\quad \mathbb{P}(y_i = j \mid x_i) = \dfrac{e^{\beta_j^T x_i}}{e^{\beta_1^T x_i} + e^{\beta_2^T x_i} + \cdots + e^{\beta_k^T x_i}}$ .

MLE :

$$\hat{\beta} = \underset{\beta}{argmax} \; \sum_{i=1}^{\hat{n}} \left( \mathbb{1}(y_i = 1)\beta_1^T x_i + \mathbb{1}(y_i = 2)\beta_2^T x_i + \cdots \right.$$
$$\left. + \mathbb{1}(y_i = k)\beta_k^T x_i - \log\left( \sum_{j=1}^{k} e^{\beta_j^T x_i} \right) \right)$$
$$= \underset{\beta}{argmax} \; \sum_{j=1}^{k} \beta_j^T \sum_{i : y_i = j} x_i - n \log\left( \sum_{j=1}^{k} e^{\beta_j^T x_i} \right).$$

Note : the MLE is not unique, as $(\beta_1, \cdots, \beta_k)$ and
$(\beta_1 + c, \cdots, \beta_k + c)$ give the same objective.
So we can assume that $\beta_1 = 0$.


4' <u>Ordered logit model (ordinal regression)</u>.

Suppose $y_i$ could take $k$ values with <u>ordered</u> relationship.

Model :
$$\log \frac{\mathbb{P}(y_i \le j)}{\mathbb{P}(y_i > j)} = \alpha_j + \beta^T x \quad (j = 1, 2, \cdots, k-1)$$

or equivalently,
$$\mathbb{P}(y_i \le j) = \frac{1}{1 + e^{-(\alpha_j + \beta^T x)}}.$$

Proportional odds assumption : the difference in the log-odds
$$\log \frac{\mathbb{P}(y_i \le j+1)}{\mathbb{P}(y_i > j+1)} - \log \frac{\mathbb{P}(y_i \le j)}{\mathbb{P}(y_i > j)}$$
is independent of $x$. More on this in Lecture 5.

MLE : $(\hat{\alpha}, \hat{\beta}) = \underset{(\alpha, \beta)}{argmax} \; \sum_{i=1}^{n} \left( \sum_{j=1}^{k} \mathbb{1}(y_i = j) \log \mathbb{P}(y_i = j) \right)$
$$= \underset{(\alpha, \beta)}{argmax} \; \sum_{i=1}^{n} \left( \sum_{j=1}^{k} \mathbb{1}(y_i = j) \times \right.$$
$$\left. \log\left( \frac{1}{1 + e^{-(\alpha_j + \beta^T x)}} - \frac{1}{1 + e^{-(\alpha_{j-1} + \beta^T x)}} \right) \right)$$

where $\quad \alpha_0 \overset{\Delta}{=} 0, \quad \alpha_k \overset{\Delta}{=} +\infty$ .

$$\boxed{\text{Exercise (HW): show that the log-likelihood is concave in } (\alpha, \beta).}$$

## Variance of MLE.

In the sequel we assume that $d = 1$ for simplicity, i.e. $\beta \in \mathbb{R}^{?}$.

F.O.C. for MLE:
$$0 = \sum_{i=1}^{n} \left( T(y_i) - A'(x_i^T \hat{\beta}^{MLE}) \right) x_i$$
$$= \sum_{i=1}^{n} \left( A'(x_i^T \beta) - A'(x_i^T \hat{\beta}^{MLE}) \right) x_i$$
$$+ \underbrace{\sum_{i=1}^{n} \left( T(y_i) - A'(x_i^T \beta) \right) x_i}_{\color{blue}{Cov(\cdot) = \sum_{i=1}^{n} A''(x_i^T \beta) x_i x_i^T}}$$

Delta method (Taylor expansion):
$$\text{first term} \approx \left( \sum_{i=1}^{n} A''(x_i^T \beta) x_i x_i^T \right) (\beta - \hat{\beta}^{MLE})$$

$$\boxed{Cov_\beta(\hat{\beta}^{MLE}) \approx \left( \sum_{i=1}^{n} A''(x_i^T \beta) x_i x_i^T \right)^{-1}}$$

## Fisher information.

<u>Def</u>. For a (regular) class of probability distributions $(p_\theta)_{\theta \in \mathbb{R}^d}$, the Fisher information at $\theta = \theta_0$ is defined as
$$I(\theta_0) = \mathbb{E}_{\theta_0} \left[ - \frac{\partial^2 \log p_\theta(y)}{\partial \theta^2} \Big|_{\theta = \theta_0} \right]$$

<div style="border:1px solid blue; color:blue;">

Side note: $\dot{\ell}_{\theta_0}(y) = \frac{\partial \log p_\theta(y)}{\partial \theta} \Big|_{\theta = \theta_0}$ <span style="color:red">(score)</span>

$\mathbb{E}_{\theta_0} [\dot{\ell}_{\theta_0}(y)] = 0$

$Cov_{\theta_0}(\dot{\ell}_{\theta_0}(y)) = I(\theta_0)$

</div>

In GLM: $\ell_\beta(x, y) = \sum_{i=1}^{n} \log p_{\theta_i}(y_i) = \sum_{i=1}^{n} (T(y_i)\beta^T x_i - A(\beta^T x_i))$
$$+ \text{Const}(x, y)$$

$\dot{\ell}_\beta(x, y) = \frac{\partial}{\partial \beta} \ell_\beta(x, y) = \sum_{i=1}^{n} \underbrace{(T(y_i) - A'(\beta^T x_i))}_{\text{has mean zero}} x_i$

$\ddot{\ell}_\beta(x, y) = \frac{\partial}{\partial \beta} \dot{\ell}_\beta(x, y) = -\sum_{i=1}^{n} A''(\beta^T x_i) x_i x_i^T$

$\implies I(\beta) = \mathbb{E}[-\ddot{\ell}_\beta(x, y)] = \sum_{i=1}^{n} A''(\beta^T x_i) x_i x_i^T$ .

---

(Asymptotic) Cramér-Rao bound: $I(\theta)^{-1}$ is the "best" covariance of any asymptotically unbiased estimator $\hat{\theta}$ for $\theta$ as $n \to \infty$.

---

Asymptotic efficiency of MLE: $\hat{\theta}^{MLE}$ asymptotically achieves the Cramér-Rao bound.

---

<u>Bootstrap estimate for $\text{Cov}(\hat{\beta}^{MLE})$</u>: same as Lecture 3.

<u>Inference in GLM</u>.

<u>Deviance</u>. Deviance for data point $i$ is
$$D_i(\hat{\beta}; \beta) = D(x_i^T \hat{\beta}; x_i^T \beta) = 2(A(x_i^T \beta) - A(x_i^T \hat{\beta}) - A'(x_i^T \hat{\beta}) x_i^T(\beta - \hat{\beta})).$$
(a generalization of "training error" $(x_i^T \hat{\beta} - x_i^T \beta)^2$ in linear regression)

<u>Total deviance</u>. $D_+(\hat{\beta}; \beta) = \sum_{i=1}^{n} D_i(\hat{\beta}; \beta)$.

<u>Hoeffding's formula</u>.
$$D_+(\hat{\beta}^{MLE}; \beta) = 2\log \frac{p_{\hat{\beta}^{MLE}}(y^n | x^n)}{p_\beta(y^n | x^n)} \quad (\text{Pf: see HW})$$

<u>Deviance table</u>.

Setting: $\beta = ( \underset{p^{(1)} \times 1}{\beta^{(1)}} , \underset{p^{(2)} \times 1}{\beta^{(2)}} , \cdots , \underset{p^{(J)} \times 1}{\beta^{(J)}} )$
  $\underset{p \times 1}{\phantom{\beta}}$     with $\sum_{j=1}^{J} p^{(j)} = P$.

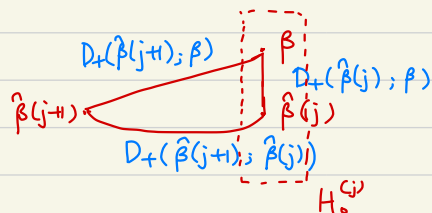Target: for each $j = 0, 1, \cdots, J$, test if $\beta^{(j+1)} = \cdots = \beta^{(J)} = 0$. ($H_o^{(j)}$)

Notation: let $\hat{\beta}(j)$ be the MLE assuming that $\hat{\beta}^{(j+1)} = \cdots = \hat{\beta}^{(J)} = 0$.
  ($\hat{\beta}(0) = 0$ , $\hat{\beta}(J) = \hat{\beta}^{MLE}$)

$\boxed{\begin{array}{l} \text{Deviance additivity theorem: if } \beta \in H_o^{(j)}, \\ \qquad D_+ ( \hat{\beta}(j+1); \hat{\beta}(j)) = D_+ ( \hat{\beta}(j+1); \beta) - D_+ ( \hat{\beta}(j) ; \beta) \\ \text{(Pf: see HW)} \end{array}}$

A pictorial illustration:



Generalized likelihood-ratio test:
  $D_+ ( \hat{\beta}(j+1) ; \hat{\beta}(j)) \sim \chi^2_{p^{(j+1)}}$ under $H_j : \beta^{(j+1)} = \cdots = \beta^{(J)} = 0$.

<u>Deviance table</u>.

| MLE | $2 \times$ log-likelihood | difference | compare with |
|---|---|---|---|
| $\hat{\beta}(0) = 0$ | $2\ell_{\hat{\beta}(0)}$ | | |
| $\hat{\beta}(1)$ | $2\ell_{\hat{\beta}(1)}$ | $D_+ ( \hat{\beta}(1); \hat{\beta}(0))$ | $\chi^2_{p^{(1)}}$ |
| $\hat{\beta}(2)$ | $2\ell_{\hat{\beta}(2)}$ | $D_+ ( \hat{\beta}(2); \hat{\beta}(1))$ | $\chi^2_{p^{(2)}}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $\hat{\beta}(J) = \hat{\beta}^{MLE}$ | $2\ell_{\hat{\beta}(J)}$ | $D_+ ( \hat{\beta}(J); \hat{\beta}(J-1))$ | $\chi^2_{p^{(J)}}$ |

<u>Model selection</u>. (Assuming $p^{(1)} = \cdots = p^{(J)} = 1$)

## 1. AIC (Akaike information criterion)

$$
\begin{aligned}
j^{AIC} &= \underset{j \in \{0,1,\cdots,J\}}{\arg\min} \; -D_+(j) + 2j \\
&= \underset{j \in \{0,1,\cdots,J\}}{\arg\min} \; -2\ell_{\hat\beta(j)} + 2j
\end{aligned}
$$

## 2. BIC (Bayesian information criterion).

$$
\begin{aligned}
j^{BIC} &= \underset{j \in \{0,1,\cdots,J\}}{\arg\min} \; -D_+(j) + j\ln n \\
&= \underset{j \in \{0,1,\cdots,J\}}{\arg\min} \; -2\ell_{\hat\beta(j)} + j\ln n
\end{aligned}
$$

## 3. Lasso.

$$
\hat\beta^{Lasso} = \underset{\beta}{\arg\min} \; -\frac{1}{n}\sum_{i=1}^{n}\log P_{x_i^T\beta}(y_i) + \lambda\|\beta\|_1
$$

- $\lambda$ is typically chosen by cross validation.

<u>Application:</u> Density estimation via <u>Lindsey's method</u>

Given i.i.d. $z_1, \cdots, z_n \sim p$, aim to fit
$$
p \approx p_\theta = \exp\left(\langle \theta, T(z)\rangle - A(\theta)\right) h(z)
$$
- known: $T(\cdot), h(\cdot)$     • unknown: $\theta \in \mathbb{R}^d$.

<u>Problem with MLE</u>: log-partition function $A(\theta)$ untractable (more in Lec 6)

<u>Lindsey's method</u>

- Suppose $Z \subseteq \mathbb{R}$, and $Z = Z_1 \cup Z_2 \cup \cdots \cup Z_K$, with
$$Z_K = [z_k - \tfrac{\Delta_k}{2}, z_k + \tfrac{\Delta_k}{2}].$$

- For small $\Delta_k$,
$$\mathbb{P}(z \in Z_k) = \int_{Z_k} p_\theta(z)\,dz$$
$$\approx \exp(\langle \theta, T(z_k) \rangle - A(\theta))\, h(z_k)\, \Delta_k =: p_k.$$

- For $y_k = \#\{z_i \in Z_k\}$, then
$$(y_1, \cdots, y_K) \sim \text{Multi}(n; (p_1, \cdots, p_K))$$

- Poisson trick: fit
$$y_k \overset{\text{ind.}}{\sim} \text{Poi}\left(e^{\langle \theta, T(z_k) \rangle + \log(h(z_k)\Delta_k) + \theta_0}\right)$$

  This is a Poisson GLM!

- <u>Poisson conditioning property</u>:

  if $y_i \overset{\text{ind.}}{\sim} \text{Poi}(\lambda_i)$, then
  $$(y_1, \cdots, y_K) \mid \sum_{k=1}^{K} y_k = n \sim \text{Multi}\left(n; \left(\tfrac{\lambda_1}{\sum_k \lambda_k}, \cdots, \tfrac{\lambda_K}{\sum_k \lambda_k}\right)\right)$$

  Therefore, $(y_1, \cdots, y_K) \mid \sum_{k=1}^{K} y_k = n \sim \text{Multi}(n; (q_1, \cdots, q_K))$, with
  $$q_k = \frac{\exp(\langle \theta, T(z_k) \rangle + \log(h(z_k)\Delta_k) + \theta_0)}{\sum_j \exp(\langle \theta, T(z_j) \rangle + \log(h(z_j)\Delta_j) + \theta_0)}$$
  $$\propto \exp(\langle \theta, T(z_k) \rangle)\, h(z_k)\, \Delta_k = p_k.$$

- Think: what does $\theta_0$ represent?