*Some background on choice e.g.*

Consider a general family of distributions:

*vary w/ covariates*

*subfamily of exponential family that includes Binomial/Bernoulli and Poisson*

$$\log f(y_i; \theta_i, \phi) = \frac{y_i \theta_i - b(\theta_i)}{a_i(\phi)} + c(y_i, \phi).$$

$$f(y_i; \theta_i, \phi) = \exp\left\{ \underbrace{\frac{y_i \theta_i}{a_i(\phi)} + c(y_i, \theta_i)}_{(*)} - \underbrace{\frac{b(\theta_i)}{a_i(\phi)}}_{(**)} \right\}$$

recall exponential family w/ parameter $\underline{\theta} = (\theta_1, \dots, \theta_s)^T$ is of the form

$$f(y; \underline{\theta}) = h(y) \exp\left\{ \underbrace{\sum_{j=1}^{S} g_j(\underline{\theta}) T_j(y)}_{(*)} - \underbrace{B(\underline{\theta})}_{(**)} \right\}$$

Assumes: $T_1(y_i) = y_i$ $\quad + \quad g_1(\underline{\theta}) = \frac{\theta_i}{a_i(\phi)}$ (subfamily of exponential family).

Similar to single parameter exponential family except for dispersion term $a_i(\phi)$.

**Example (Normal model):** $\quad E y_i = \mu_i$

$$f(y_i; \mu_i, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left( -\frac{(y_i - \mu_i)^2}{2\sigma^2} \right)$$

$$\log f(y_i; \mu_i, \sigma) = \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(y_i - \mu_i)^2}{2\sigma^2}$$

$$= -\log(\sqrt{2\pi}\sigma) - \frac{y_i^2 - 2\mu_i y_i + \mu_i^2}{2\sigma^2}$$

$$= \frac{y_i \mu_i - \frac{\mu_i^2}{2}}{\sigma^2} - \log(\sqrt{2\pi}\,\sigma) - \frac{y_i^2}{2\sigma^2}$$

So,

$$a_i(\phi) = \sigma^2$$

$$\theta_i = \mu_i$$

$$b(\theta_i) = \frac{\mu_i^2}{2} = \frac{\theta_i^2}{2}$$

$$c(y_i, \phi) = -\log(\sqrt{2\pi}\,\sigma) - \frac{y_i^2}{2\sigma^2} \quad \text{(only depends on } \sigma, \text{ not } \mu_i)$$

We can learn something about this distribution by considering it's mean and variance. Because we don't have an explicit form of the density, we rely on two facts:

1. $\mathrm{E}\left[\frac{\partial \log f(Y_i;\theta_i,\phi)}{\partial \theta_i}\right] = 0.$

   " derivative of log density wrt parameter of interest $\theta_i$ has expection $0$ "

See HW3.

2. $\mathrm{E}\left[\frac{\partial^2 \log f(Y_i;\theta_i,\phi)}{\partial \theta_i^2}\right] + \mathrm{E}\left[\left(\frac{\partial \log f(Y_i;\theta_i,\phi)}{\partial \theta_i}\right)^2\right] = 0.$

These will also come up later when we talk about information matrix.

For $\log f(y_i; \theta_i, \phi) = \frac{y_i\theta_i - b(\theta_i)}{a_i(\phi)} + c(y_i, \phi),$

Using ①:

$$\frac{\partial}{\partial \theta_i} \log f(y_i; \theta_i, \phi) = \frac{1}{a_i(\phi)}\left(y_i - b'(\theta_i)\right)$$

$$E\left[\frac{1}{a_i(\phi)}\left(Y_i - b'(\theta_i)\right)\right] = 0 \implies b'(\theta_i) = E[Y_i] \implies$$ information about the mean is contained in $b'(\theta_i)$

E.g. Normal model
$$b(\theta_i) = \frac{\theta_i^2}{2} \implies b'(\theta_i) = \frac{2\theta_i}{2} = \theta_i = \mu_i = E[Y_i] \checkmark$$
↑ desired      ↑ from the family form

Using ②:

$$\frac{\partial^2}{\partial \theta_i^2} \log f(y_i; \theta_i, \phi) = -\frac{b''(\theta_i)}{a_i(\phi)} \implies E\left[\frac{\partial^2}{\partial \theta_i^2} \log f(Y_i; \theta_i, \phi)\right] = -\frac{b''(\theta_i)}{a(\phi)}$$

$$E\left[\left(\frac{\partial \log f(Y_i; \theta_i, \phi)}{\partial \theta_i}\right)^2\right] = E\left[\left(\frac{1}{a_i(\phi)}(Y_i - b'(\theta_i))\right)^2\right] = \frac{1}{a_i(\phi)^2} E\left[(Y_i - \mu_i)^2\right]$$

$$\implies -\frac{b''(\theta_i)}{a(\phi)} + \frac{1}{a_i(\phi)^2}\text{Var } Y_i = 0 \implies \text{Var}[Y_i] = a_i(\phi)b''(\theta_i)$$

Thoughts:
- Variance depends on $i$
- Var$[Y_i]$ positive $\implies$ $b''(\theta_i)$ positive ∀ values of $\theta_i$
  $\implies b(\theta_i)$ strictly convex
  $b'(\theta_i)$ monotone increasing so $b'^{-1}$ exists.

**Example (Bernoulli model):**

$$f(y_i; p_i) = p_i^{y_i}(1-p_i)^{1-y_i}$$

$$\log f(y_i; p_i) = y_i \log p_i + (1-y_i)\log(1-p_i)$$

$$= y_i \boxed{\log \frac{p_i}{1-p_i}} + \boxed{\log(1-p_i)}$$

comparing to general form:

$$\frac{y_i \boxed{\theta_i} - \boxed{b(\theta_i)}}{a_i(\phi)} + c(y_i, \phi)$$

$$\Rightarrow \theta_i = \log\left(\frac{p_i}{1-p_i}\right) \implies p_i = \frac{\exp(\theta_i)}{1+\exp(\theta_i)}$$

$$b(\theta_i) = -\log(1-p_i)$$

$$= -\log\left(1 - \frac{\exp(\theta_i)}{1+\exp(\theta_i)}\right)$$

$$= -\log\left(\frac{1}{1+\exp(\theta_i)}\right)$$

$$= \log(1+\exp(\theta_i))$$

Notice: $b'(\theta_i) = \frac{1}{1+\exp(\theta_i)}\exp(\theta_i) = p_i$ ✓

$$a_i(\phi) = 1.$$

$$a_i(\phi)\, b''(\theta_i) = -\left(-\left(1+\exp(\theta_i)\right)^{-2}\exp(\theta_i)\exp(\theta_i) + \left(1+\exp(\theta_i)\right)^{-1}\exp(\theta_i)\right)$$

$$= \frac{(1+\exp(\theta_i))\exp(\theta_i) - \exp(\theta_i)\exp(\theta_i)}{(1+\exp(\theta_i))^2} = \frac{\exp(\theta_i)}{(1+\exp(\theta_i))^2}$$

$$= \frac{\exp(\theta_i)}{(1+\exp(\theta_i))} \cdot \frac{1}{1+\exp(\theta_i)}$$

$$= p_i(1-p_i) = \text{Var } Y_i$$

Finally, back to modelling. Our **goal** is to build a relationship between the mean of $Y_i$ and covariates $x_i$.

$$\text{choose} \quad x_i^T \beta = g(\mu_i) = g(E[Y_i]).$$

what to choose for $g$?

$$\mu_i = g^{-1}(x_i^T \beta)$$

we know $\mu_i = b'(\theta_i)$

$\qquad\qquad$ we know this exists!

$$\implies b'(\theta_i) = g^{-1}(x_i^T \beta) \quad \text{OR} \quad \theta_i = b'^{-1}(g^{-1}(x_i^T \beta))$$

If we choose $g^{-1} = b'$, then this will clean up nicely! i.e. $\theta_i = x_i^T \beta$.

"canonical" or "natural" link function.

$$\implies \text{log-likelihood is}$$

$$\ell\left(\beta, \phi \mid \{Y_i, x_i\}_{i=1}^n\right) = \sum_{i=1}^n \left\{ \frac{Y_i \, x_i^T \beta - b(x_i^T \beta)}{a_i(\phi)} + c(Y_i, \phi) \right\}$$

**Example (Bernoulli model, cont'd):**

$$b(\theta_i) = \log\left(1 + \exp(\theta_i)\right)$$

$$b'(\theta_i) = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)} = \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \quad \text{canonical link}$$

$\downarrow$ mean

$$p_i = \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)}$$

(same form as before)

# 1.4 Marginal and Conditional Likelihoods

Consider a model which has $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$, where $\boldsymbol{\theta}_1$ are the parameters of interest and $\boldsymbol{\theta}_2$ are nuisance parameters.

*not what we are interested in performing inference for*

*when dimension of $\underline{\theta}_2$ is large MLE's of $\underline{\theta}_1$ can be biased for small samples and inconsistent in large samples.*

One way to improve estimation for $\boldsymbol{\theta}_1$ is to find a one-to-one transformation of the data $\boldsymbol{Y}$ to $(\boldsymbol{V}, \boldsymbol{W})$ such that either

$$f_Y(\underline{y}; \theta_1, \theta_2) = f_{W,V}(\underline{w}, \underline{v}; \theta_1, \theta_2)$$

$$= f_{W,V}(\underline{w} \mid \underline{v}; \theta_1, \theta_2) \, f_V(\underline{v}; \theta_1) \qquad \text{``marginal''}$$

*alternative likelihood*

OR $\qquad f_Y(\underline{y}; \underline{\theta}_1, \theta_2) = f(\underline{w} \mid \underline{v}; \theta_1) \, f_V(\underline{v}; \theta_1, \theta_2) \qquad \text{``conditional''}$

*either way looking to split density into a piece that doesn't depend on $\underline{\theta}_2$ (nuisance parameters)*

The key feature is that one component of each contains only the parameter of interest.

$\underline{\theta}_1$

**Example (Neyman-Scott problem):** Let $Y_{ij}, i = 1, \ldots, n, j = 1, 2$ be intependent normal random variables with possible different means $\mu_i$ but the same variance $\sigma^2$.

$$Y_{ij}, \quad \underbrace{i = 1, \ldots, n}_{n \text{ groups}}, \quad \underbrace{j = 1, 2}_{\substack{\text{ONLY 2} \\ \text{obs per group}}} \quad \overset{iid}{\sim} \quad N\left(\underset{\underset{\text{mean}}{\text{group}}}{\mu_i}, \sigma^2\right) \underset{\underset{\text{variance.}}{\text{common}}}{}$$

$$\theta = \left(\overbrace{\mu_1, \ldots, \mu_n}^{\text{nuisance params}}, \sigma^2\right)^T$$
$$\underset{n+1 \text{ parameters.}}{}$$

Our goal is to estimate $\sigma^2$. Should we be able to?

Yes: lots of groups!

No: only 2 obs per group.

Usual asymptotic assumptions as $n$ grows,
Here as $n$ grows, # groups grows $\Rightarrow$ # parameters grows.

Following the usual arguments,

$$\hat{\mu}_{i,\text{MLE}} = \frac{Y_{i1} + Y_{i2}}{2}$$

$$\hat{\sigma}^2_{\text{MLE}} = \frac{1}{2n} \sum_{i=1}^{n} \sum_{j=1}^{2} (Y_{ij} - \hat{\mu}_{i,\text{MLE}})^2$$

$$= \frac{1}{2n} \sum_{i=1}^{n} \left\{ \left( Y_{i1} - \frac{Y_{i1} + Y_{i2}}{2} \right)^2 + \left( Y_{i2} - \frac{Y_{i1} + Y_{i2}}{2} \right)^2 \right\}$$

$$= \frac{1}{2n} \sum_{i=1}^{n} \left\{ \underline{\left( \frac{Y_{i1} - Y_{i2}}{2} \right)^2} + \underline{\left( \frac{Y_{i2} - Y_{i1}}{2} \right)^2} \right\}$$
$$\text{terms are equivalent}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{1}{4} \left( Y_{i1} - Y_{i2} \right)^2$$

$$\mathrm{E}[\hat{\sigma}^2_{\mathrm{MLE}}] = E\left[\frac{1}{n}\sum_{i=1}^{n}\frac{1}{4}\left(Y_{i_1}-Y_{i_2}\right)^2\right]$$

$$(iid) \quad = \frac{1}{4}E\left[\left(Y_{i_1}-Y_{i_2}\right)^2\right]$$

$$= \frac{1}{4}E\left[\left((Y_{i_1}-\mu_i)-(Y_{i_2}-\mu_i)\right)^2\right]$$

$$= \frac{1}{4}E\left[(Y_{i_1}-\mu_i)^2 - 2(Y_{i_1}-\mu_i)(Y_{i_2}-\mu_i)+(Y_{i_2}-\mu_i)^2\right]$$

$$= \frac{1}{4}\left[\sigma^2 - 0 + \sigma^2\right]$$

$$= \frac{\sigma^2}{2} \neq \sigma^2 \ !$$

So as $n \to \infty$, $\hat{\sigma}^2_{MLE} \xrightarrow{p} \frac{\sigma^2}{2}$ by WLLN

This seems bad.

Happens because # of nuisance parameters grows w/ $n$ (# groups).

A reworking of the data seems more promising. Let,

$$V_i = \frac{Y_{i1}-Y_{i2}}{\sqrt{2}} \qquad \text{and} \qquad W_i = \frac{Y_{i1}+Y_{i2}}{\sqrt{2}}$$