# 1.3 Likelihoods for Regression Models

We will start with linear regression and then talk about more general models.

$\downarrow$
nonlinear
GLM

## 1.3.1 Linear Model

Consider the familiar linear model

$$Y_i = \boldsymbol{x}_i^\top \boldsymbol{\beta} + \epsilon_i, \qquad i = 1, \ldots, n,$$

where $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$ are known nonrandom vectors.

$$E\left[\varepsilon_i\right] = 0 \quad \text{and} \quad Var\left[\varepsilon_i\right] = \sigma^2$$

often estimate $\beta$ by $\hat{\beta}_{OLS}$, which does not require a distribution for $\varepsilon_i$.

For likelihood-based estimation, We need a distribution for $\varepsilon_i$! Start w/ $\varepsilon_i \sim N(0, \sigma^2)$.

$$\Rightarrow L(\boldsymbol{\beta}, \sigma | \{Y_i, \boldsymbol{x}_i\}_{i=1}^n) = \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}\sigma}\right) \exp\left(-\frac{(Y_i - x_i^\top \beta)^2}{2\sigma^2}\right)$$

$$= \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - x_i^\top \beta)^2\right)$$

take log,
derivatives, set = 0,
solve

$\hat{\beta}_{MLE} = (X^\top X)^{-1} X^\top Y$      same as $\hat{\beta}_{OLS}$!

$\hat{\sigma}^2_{MLE} = \frac{1}{n} \sum_{i=1}^n (Y_i - x_i^\top \beta)^2$    (only asymptotically unbiased).

What do you do when $\epsilon_i$ are not Gaussian?

- transform data so $\epsilon_i$ look Gaussian.

- use a different distribution for $\epsilon_i$!

**Example (Venice sea levels):** The annual maximum sea levels in Venice for $1931-1981$ are :

we know maxima are not Gaussian!



Approach 1: OLS   $E[\epsilon_i]=0$ , $Var[\epsilon_i]=\sigma^2$   No distributional assumption.

Approach 2: Assume $\epsilon_i \sim$ Gumbel (extreme value dsn), use ML

$$f(y; \sigma) = \frac{1}{\sigma} \exp\left(-\frac{y}{\sigma}\right) \exp\left(-\exp\left(-\frac{y}{\sigma}\right)\right).$$

$$\Rightarrow L\left(\beta, \sigma \mid \{y_i, x_i\}_{i=1}^n\right) = \prod_{i=1}^n f\left(y_i - x_i^T\beta\right) = \prod_{i=1}^n \frac{1}{\sigma} \exp\left(-\frac{y_i - x_i^T\beta}{\sigma}\right) \exp\left(-\exp\left(-\frac{y_i - x_i^T\beta}{\sigma}\right)\right)$$

YOUR TURN : Fit both approaches to the venice data.

OLS

$\hat{\beta_0} = 104.8$   $\hat{\beta_1} = .56$   $(SE .177)$

MLE, GUMBEL

$\hat{\beta_0} = 96.8$ , $\hat{\beta_1} = 0.56$   $\left(\begin{array}{c} SE \\ .136 \end{array}\right)$

OLS vs MLE? If EV Model is correct, more efficient (note: standard errors).

$\beta_0$ difference:   $\overset{\text{Gumbel}}{E[\epsilon_i]} = 0.577\sigma = 0.577\hat{\sigma}_{MLE} = 0.577(14.5) \neq 0$

$$96.8 + .577(14.5) = 105.1$$

## 1.3.2 Additive Errors Nonlinear Model

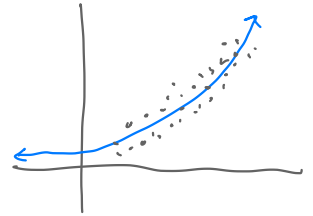Previous example had ① linear trend, ② Non-Gaussian errors.

Non-linear additive model:
$$Y_i = g(\underline{x}_i, \beta) + \varepsilon_i$$

Often interested in $\varepsilon_i \sim N(0, \sigma^2)$ but $g(\underline{x}_i, \beta) \neq \underline{x}_i^T \beta \Rightarrow$ ML required.

① non-linear trend, ② Gaussian errors.

Example: exponential growth model

$$g(x, \underline{\beta}) = \beta_0 \exp(\beta_1 x)$$
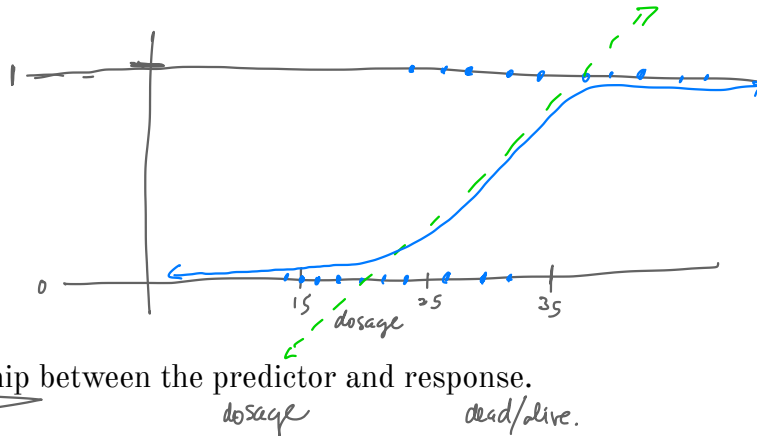
## 1.3.3 Generalized Linear Models

↗ of response.

Regression: build a relationship between a parameter (mean) & covariates.

LM's: stochastic element is additive w/ mean.

GLM's: stochastic element is different.

Imagine an experiment where individual mosquitos are given some dosage of pesticide. The response is whether the mosquito lives or dies. The data might look something like:

$=0$    $=1$

| $x$ (dosage) | $y$ (0=lives, 1=dies) |
|---|---|
| 15 | 0 |
| 17 | 0 |
| 18 | 1 |
| 20 | 0 |
| 21 | 1 |
| ⋮ | ⋮ |

**Goal:** Model the relationship between the predictor and response.

dosage          dead/alive.

Sounds like regression!

Big difference: $Y_i$'s are not continuous. They only take values of 0 or 1.

**Question:** What would a curve of best fit look like? Would we want a function that only takes values in $\{0,1\}$.

It seems sensible to have a curve which takes values near 0 for low doses & near 1 for high doses and intermediate values for middle doses.

what does this curve represent? Probability.

**Refined Goal:** Model relationship between predictor (dosage) + probability of success in response

Let's build a <u>sensible</u> model. <u>Note</u>: We don't observe the <u>probability</u>. (mosquito dies).

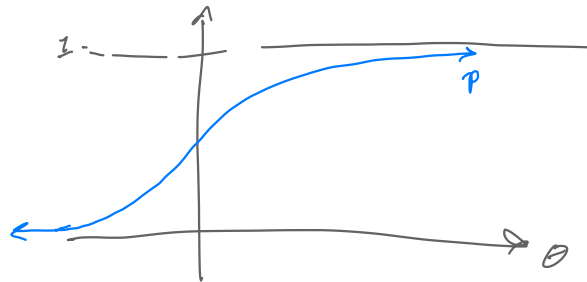**Step 1:** Find a function that behaves the way we want.

like the blue curve.

Consider the logistic function,

$$p = \frac{\exp(\theta)}{1 + \exp(\theta)},$$

As $\theta \to \infty,\ p \to 1$

$\quad\ \ \theta \to -\infty,\ p \to 0$

$\quad\ \ \theta = 0,\ p = \frac{1}{2}$



By changing $\theta$, we can change location, slope, direction of this function.

Let $\theta = \beta_0 + \beta_1 x \implies p = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}.$

```
# understanding the logistic function
# first, theta just equals x
x <- seq(-7, 7, .1)
theta <- x
y <- exp(theta)/(1 + exp(theta))
ggplot() + geom_line(aes(x, y))

# now, let theta be a linear function of x
theta <- 1 + 3*x
y <- exp(theta)/(1 + exp(theta))
ggplot() + geom_line(aes(x, y))
```
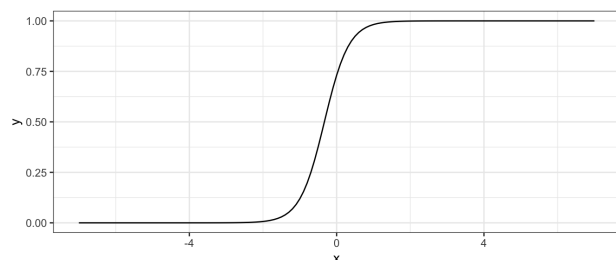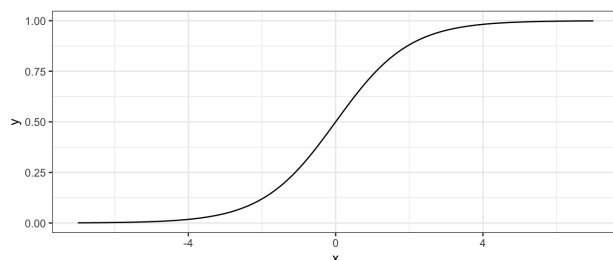


Now we can connect probabilities to covariate $x$!

We'd be done if we observed probabilities, but our response only takes values of 0 and 1.

**Step 2:** Build a stochastic mechanism to relate to a binary response.

Recall the Bernoulli distribution

$$Y = \begin{cases} 0 & \text{w. p.} \quad 1-p \\ 1 & \text{w. p.} \quad p. \end{cases}$$

biased
Coin flip example w/ $p = 0.75$. Flip coin, you will observe $0$ (tails) or $1$ (heads).

Aside: We could instead think about binomial dsn, which counts # of successes for $n$ trials.

$$X = \sum_{i=1}^{n} Y_i \quad , \quad Y_i \overset{iid}{\sim} Bern(p). \quad X \text{ takes values in } \{0, 1, \cdots, n\}.$$

$$P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}.$$

**Step 3:** Put Step 1 and Step 2 together.

$$Y_i \overset{ind}{\sim} Bernoulli\,(p_i) \qquad \qquad + \qquad \qquad p_i = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)} \quad , \quad \theta_i = \beta_0 + \beta_1 x_i$$

↑
outcome of $i^{th}$ observation (observed).

↑
prob of $i^{th}$ observation having success (unobserved).

<u>Goal</u>: estimate $\beta_0$ and $\beta_1$. Find the "best" estimates.

Fitting our model: Does OLS make sense? No.

What else can we do? Maximum likelihood!
↳ Find the parameters ($\beta$s) which make the density agree best w/ data we observed!

pmf of Bernoulli: $f(y_i ; p_i) = p_i^{y_i} (1-p_i)^{1-y_i}$

⟶ take $y_i$'s to estimate $p_i$'s.

Consider the likelihood contribution.

$$L_i(p_i|Y_i) = p_i^{Y_i}(1-p_i)^{(1-Y_i)} \qquad (Y_i\text{'s are } 0 \text{ or } 1).$$

So the log-likelihood contribution is

$$\ell_i(p_i) = Y_i \log p_i + (1-Y_i)\log(1-p_i) = \underbrace{\log(1-p_i) + Y_i \log \frac{p_i}{1-p_i}}_{(*)}$$

Recall, we said $\underbrace{p_i = \frac{\exp(\theta_i)}{1+\exp(\theta_i)}}$ was sensible.

Manipulating! $\quad p_i + p_i \exp(\theta_i) = e^{(\theta_i)}$

$$p_i = (1-p_i)\exp(\theta_i).$$

OR

$$\frac{p_i}{1-p_i} = \exp(\theta_i)$$

$$\log\left(\frac{p_i}{1-p_i}\right) = \theta_i \quad (1)$$

$$p_i \exp(-\theta_i) = 1-p_i$$

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

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

$$-\log(1+\exp(\theta_i)) = \log(1-p_i) \quad (2).$$

Plug in (1) and (2) into $(*)$.
Which gives us,

$$\ell_i(\theta_i) = \underbrace{-\log(1+\exp(\theta_i))} + Y_i\theta_i \qquad (\text{now in terms of } \theta_i \text{ not } p_i)$$

notice now the term w/ the data is "nice" for MLE things.
Why? Because log-likelihood + "sensible" function $p_i = \frac{\exp(\theta_i)}{1+\exp(\theta_i)}$ work well together.

<u>NOT</u> a coincidence.

So the log-likelihood is

$$\ell(\theta_1,\dots,\theta_n) = \sum_{i=1}^{n} \ell_i(\theta_i)$$

$$= \sum_{i=1}^{n} \left\{-\log(1+\exp(\theta_i)) + Y_i\theta_i\right\}$$

$$\Rightarrow \ell(\beta_0,\beta_1) = \sum_{i=1}^{n}\left\{-\log(1+\exp(\beta_0+\beta_1 x_i)) + Y_i(\beta_0+\beta_1 x_i)\right\}$$

To optimize? *Must be done numerically.*

```
## data on credit default
data("Default", package = "ISLR")
head(Default)
```

```
##    default student   balance      income
## 1       No      No  729.5265 44361.625
## 2       No     Yes  817.1804 12106.135
## 3       No      No 1073.5492 31767.139
## 4       No      No  529.2506 35704.494
## 5       No      No  785.6559 38463.496
## 6       No     Yes  919.5885  7491.559
```

*optimizing likelihood numerically.*

```
## fit model with ML
m0 <- glm(default ~ balance, data = Default, family = binomial)
tidy(m0) |> kable()
```

*broom package.*

$\hat{\beta_0}, \hat{\beta_1}$     $se(\hat{\beta_0}), se(\hat{\beta_1})$     $\frac{\hat{\beta_i}}{se(\hat{\beta_i})}$

data's are 0,1

$H_0: \beta_i = 0$  $i=1,2.$
$H_a: \beta_i \neq 0$

| term        | estimate    | std.error | statistic | p.value |
|-------------|-------------|-----------|-----------|---------|
| (Intercept) | -10.6513306 | 0.3611574 | -29.49221 | 0       |
| balance     | 0.0054989   | 0.0002204 | 24.95309  | 0       |

```
glance(m0) |> kable()
```

| null.deviance | df.null | logLik    | AIC      | BIC      | deviance | df.residual | nobs  |
|---------------|---------|-----------|----------|----------|----------|-------------|-------|
| 2920.65       | 9999    | -798.2258 | 1600.452 | 1614.872 | 1596.452 | 9998        | 10000 |

$\ell(\hat{\beta_0}, \hat{\beta_1}).$

```
## plot the curve
x_new <- seq(0, 2800, length.out = 200)
theta <- m0$coefficients[1] + m0$coefficients[2]*x_new
p_hat <- exp(theta)/(1 + exp(theta))

ggplot() +
  geom_point(aes(balance, as.numeric(default) - 1), alpha = 0.5, data
        = Default) +
  geom_line(aes(x_new, p_hat), colour = "blue") +
  ylab("Probability of Defaulting")
```

never outside of [0,1] → valid probabilites!

In general, a GLM is three pieces:

1. The random component

probability dsn from
exponential family.

link
function.

2. The systemic component

A function relating the parameter
of interest (mean!) to $\theta$

$$E[Y] = \vec{g}^{-1}(\eta)$$

3. A linear predictor

$$\theta = X\beta.$$

Ex: Logistic Regression

$$Y_i \sim \text{Binom}(p_i)$$

$$p_i = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)} = \vec{g}^{-1}(\theta_i).$$

Note $Y_i \sim \text{Bern}(p_i)$
$$E[Y_i] = p_i$$

$$\theta_i = \underline{x}_i \beta$$

Explanation:

describe the generating
mechanism of <u>observed data</u>

transforming linear relationship
to be on a scale that makes
sense for the parameter of interest

"linking" linear relationship to mean.

describing how $\theta$ is a linear
function of predictor variables.

Remarks:

① Standard formulation denotes function by $\vec{g}^{-1}$: $\quad p = \vec{g}^{-1}(\theta) = \dfrac{\exp(\theta)}{1 + \exp(\theta)}.$

$$\Rightarrow \theta = g(p) = \log\left(\frac{p}{1-p}\right).$$

② Parameter of interest is still the mean, just like linear regression.

③ Theoretical reasons for exponential family ... relationship btw/ param of interest & Variance.

Example (Poisson regression): $\quad \nearrow$ for count data.

① Poisson $(\lambda)$.

③ $\theta = X\beta$

② $\lambda = \vec{g}^{-1}(\theta) = \exp(\theta) \quad (\lambda > 0).$

$$\Downarrow$$

$$\theta = \log(\lambda).$$

Some background on choice of $g$.

Consider a general <u>family</u> of distributions:

subfamily of exponential family dsns that includes Binomial, Poisson, etc.

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

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

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

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
$$T_1(y_i) = y_i$$
$$g_1(\underline{\theta}) = \frac{\theta_i}{a_i(\phi_i)}$$

subfamily of exponential family.
similar to single param exp family except dispersion term $a_i(\phi)$.

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

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

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

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

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

$$\theta_i = \mu_i$$
$$a_i(\phi) = \sigma^2$$
$$b(\theta_i) = \frac{\mu_i^2}{2} = \frac{\theta_i^2}{2}$$
$$c(y_i, \phi) = -\log\left( \sqrt{2\pi}\,\sigma \right) - \frac{y_i^2}{2\sigma^2} \quad \left( \text{depends on } \sigma^2, \text{ not } \mu_i \right).$$

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:

HW 2

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

These will come up again later
(Thursday?).

when we talk about information matrix.

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

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:

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

$$\Rightarrow E\left[\frac{1}{a_i(\phi)}\left(y_i - b'(\theta_i)\right)\right] \overset{\text{fact①}}{=} 0 \qquad \Rightarrow b'(\theta_i) = E[Y_i]. \Rightarrow \begin{matrix} \text{information about the} \\ \text{mean is contained in } b'(\theta_i). \end{matrix}$$

E.g. Normal model                    , from firm

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

Using ②:

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

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

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

Thoughts:
- Variance can depend on $i$
- Variance $Var[Y_i]$ positive $\Rightarrow b''(\theta_i)$ positive $\forall$ values of $\theta_i$
$$\Rightarrow b(\theta_i) \text{ strictly convex}$$
$$b'(\theta_i) \text{ monotone increasing} \Rightarrow b'^{-1} \text{ exists.}$$

**Example (Bernoulli model):**

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

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

$$= y_i \left[ \log \frac{p_i}{1 - p_i} - \left[ -\log(1 - p_i), \right] \right] + 0$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxx}}_{\textcircled{1}}$$

comparing to general form:

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

$$\theta_i = \log\left(\frac{p_i}{1 - p_i}\right) \Rightarrow p_i = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)}. \qquad a_i(\theta) = 1$$

$$c(y_i, \phi) = 0$$

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

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

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

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

$$b'(\theta_i) = \frac{1}{1 + \exp(\theta_i)} \cdot \exp(\theta_i) = p_i \overset{\checkmark}{=} E[Y_i]$$

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

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

$$= \frac{\exp(\theta_i)}{\left(1 + \exp(\theta_i)\right)^2} = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)} \cdot \frac{1}{1 + \exp(\theta}$$

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

Finally, back to modelling. Our **goal** is to build a relationship between the mean of $Y_i$ and covariates $\boldsymbol{x}_i$.

choose $\quad \boldsymbol{x}_i^T \boldsymbol{\beta} = g\left(E[Y_i]\right)$

What to choose for $g$?

$$E[Y_i] = g^{-1}\left(\boldsymbol{x}_i^T \boldsymbol{\beta}\right)$$

We know $\quad E[Y_i] = b'(\theta_i)$

we know this exists!

$$\implies b'(\theta_i) = g^{-1}\left(\boldsymbol{x}_i^T \boldsymbol{\beta}\right) \qquad \text{OR} \qquad \theta_i = b'^{-1}\left(g^{-1}\left(\boldsymbol{x}_i^T \boldsymbol{\beta}\right)\right)$$

If we choose $\underbrace{g^{-1} = b'}$, then this will clean up nicely! $\quad \theta_i = \boldsymbol{x}_i^T \boldsymbol{\beta}$.

"canonical" or "natural" link function.

$\implies$ log-likelihood is

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

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

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

$$b'(\theta_i) = \frac{\exp(\theta_i)}{1 + \exp(\theta_i)}.$$

let $g^{-1} = b' \implies b'(\theta_i) = \dfrac{\exp\left(\boldsymbol{x}_i^T \boldsymbol{\beta}\right)}{1 + \exp(\boldsymbol{x}^T \boldsymbol{\beta})}$

$\qquad\qquad\qquad\qquad \parallel$

$\qquad\qquad\qquad\quad E[Y_i]$

$\qquad\qquad\qquad\qquad \parallel$

$\qquad\qquad\qquad\quad p_i$

$\qquad\qquad\quad$ (same as before).

# 1.4 Marginal and Conditional Likelihoods

Consider a model which has $\theta = (\theta_1, \theta_2)$, where $\theta_1$ are the parameters of interest and $\theta_2$ are nuisance parameters.

*not what we are interested in performing inference for.*

*When dimension of $\underline{\theta}_2$ is large, MLEs of $\underline{\theta}_1$ can be biased for small samples and inconsistent in large samples.*

One way to improve estimation for $\theta_1$ is to find a one-to-one transformation of the data $Y$ to $(V, W)$ such that either

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

*alternative likelihoods*

OR

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

*nuisance*
*parameters.*

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

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 $\sigma^2$.

$$Y_{ij}, \quad i=1,\ldots,n, \; j=1,2 \quad \overset{iid}{\sim} N(\mu_i, \sigma^2).$$

$\underbrace{\phantom{Y_{ij}}}_{\text{n groups}}$ $\underbrace{\phantom{i=1,n}}_{\substack{\text{ONLY two}\\\text{observations}\\\text{per group}}}$ $\underset{\substack{\text{group}\\\text{mean}}}{\uparrow}$ $\quad \underset{\text{common variance.}}{\nwarrow}$

$$\underline{\theta} = \left( \mu_{1}, \ldots, \mu_n, \; \sigma^2 \right)^T$$

$\underbrace{\phantom{(\mu_1, \ldots, \mu_n, \sigma^2)}}_{\text{n+1 parameters}}$

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, # of groups grows $\Rightarrow$ # parameters grows.

Following the usual arguments,

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

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

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

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

$$\underset{\text{equivalent}}{}$$

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

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

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

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

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

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

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

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

This seems bad.

Happens because the # of nuisance parameters grows as $n$ grows.

A reworking of the data seems more promising. Let,

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

Because $Y_{ij}$ are Gaussian,

$$V_i \sim N(0, \sigma^2) \text{ and } W_i \sim N(\sqrt{2}\mu_i, \sigma^2). \text{ Also, } V_i \perp W_i!$$

$$\begin{bmatrix} V_i \\ W_i \end{bmatrix} = \frac{1}{\sqrt{2}}\begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}\begin{bmatrix} Y_{i_1} \\ Y_{i_2} \end{bmatrix} \Rightarrow \text{Var}\left(\begin{bmatrix} V_i \\ W_i \end{bmatrix}\right) = \frac{1}{2}\begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}\sigma^2\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} = \sigma^2\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

consider the density of $V \;\&\; W$:

$$f_{VW}\left(\underline{v}, \underline{w}; \sigma^2, \mu_{(1-)}, \mu_n\right) \overset{\text{ind}}{=} \underbrace{f_V\left(\underline{v}; \sigma^2\right)}_{\text{no nuisance parameters!}} f_W\left(\underline{w}; \mu_{1}, ..., \mu_n, \sigma^2\right)$$

$$\Rightarrow \ell(\sigma \mid \underline{V}) = -n\log\sqrt{2\pi} - n\log\sigma - \frac{1}{2\sigma^2}\sum_{i=1}^{n}V_i^2$$

$$\frac{\partial\ell(\sigma\mid\underline{V})}{\partial\sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3}\sum_{i=1}^{n}V_i^2 \Rightarrow \hat{\sigma}^2_{\text{MLE}} = \frac{1}{n}\sum_{i=1}^{n}V_i^2$$

A marginal likelihood approach is simple provided you can find a statistic $V$ whose dsn is free of the nuisance parameter!

For conditional likelihoods, we can often exploit the existence of sufficient statistics for the nuisance parameters under the assumption that the parameter of interest is known.

Let $T_i$ be sufficient for a nuisance parameter $(\mu_i)$

Then the conditional distribution of the data given $\underline{T} = (T_1,..,T_n)$ doesn't depend on the nuisance parameters.
$(\underline{Y})$

$\Rightarrow$ We can look for the conditional dsn of data $|\underline{T}$

**Example (Exponential Families):** The structure of exponential families is such that it is often possible to exploit their properties to eliminated nuisance parameters. Let $Y$ have a density of the form

$$f(y; \boldsymbol{\eta}) = h(y) \exp\left\{ \sum_{i=1}^{s} \eta_i T_i(y) - A(\boldsymbol{\eta}) \right\},$$

then

If $\eta = (\underline{\theta_1}^T, \underline{\theta_2}^T)$,

(Thm 2.6
pg 104)

sufficient statistics

$$f(y; \underline{\theta_1}, \underline{\theta_2}) = h(y) \exp\left\{ \sum \theta_{1i} W_i + \sum \theta_{2j} V_j - A(\underline{\theta_1}, \underline{\theta_2}) \right\}$$

Then the conditional dsn of $\underline{W}$ given $\underline{V}$ is exponential family of the form:

$$f(\underline{w}|\underline{v}; \underline{\theta_1}) = q(\underline{w}) \exp\left\{ \sum_{i=1}^{r} \theta_{1i} w_i - A_{\underline{v}}(\underline{\theta_1}) \right\}$$

a conditional density
that doesn't depend on $\underline{\theta_2}$.

Thus, exponential families often provide an automatic procedure for finding $\boldsymbol{W}$ and $\boldsymbol{V}$.

**Example (Logistic Regression):** For binary $Y_i$, the standard logistic regression model is

$$P(Y_i = 1) = p_i(\boldsymbol{x}_i, \boldsymbol{\beta}) = \frac{\exp(\boldsymbol{x}_i^\top \boldsymbol{\beta})}{1 + \exp(\boldsymbol{x}_i^\top \boldsymbol{\beta})}$$

and the likelihood is

$$L(\boldsymbol{\beta}|\boldsymbol{Y}, \boldsymbol{X}) = \prod_{i=1}^{n} p_i(x_i, \beta)^{Y_i} \{1 - p_i(x_i, \beta)\}^{1-Y_i}$$

$$= \prod_{i=1}^{n} \left\{ \frac{\exp(x_i^\top \beta)}{1 + \exp(x_i^\top \beta)} \right\}^{Y_i} \left\{ \frac{1}{1 + \exp(x_i^\top \beta)} \right\}^{1-Y_i}$$

$$= \frac{\exp\left( \sum_{i=1}^{n} Y_i (x_i^\top \beta) \right\}}{\prod_{i=1}^{n} (1 + \exp(x_i^\top \beta))}$$

$$= c(X, \beta) \cdot \exp\left( \sum_{j=1}^{p} \beta_j \sum_{i=1}^{n} x_{ij} Y_i \right).$$

$$\Rightarrow T_j = \sum_{i=1}^{n} x_{ij} Y_i \qquad j = 1, \dots, p \quad \text{are sufficient for this exponential family model.}$$

Suppose $\theta_1 = \beta_k$ is the parameter of interest and the others are nuisance parameters.

$$\Rightarrow W_1 = T_k = \sum_{i=1}^{n} x_{ik} Y_i \quad \text{and} \quad V = (T_1, \dots, T_{k-1}, T_{k+1}, \dots, T_p)^T \qquad \overset{\text{joint}}{\underset{\text{marginal}}{\nearrow}}$$

and the conditional density $\quad P(T_k = t_k | T_1 = t_1, \dots, T_{k-1} = t_{k-1}, T_{k+1} = t_{k+1}, \dots, T_p = t_p)$

$$\vdots$$

$$= \frac{c(t_1, \dots, t_p) \exp(\beta_k t_k)}{\sum_u c(t_1, \dots, t_{k-1}, u, t_{k+1}, \dots, t_p) \exp(\beta_k u)} \quad \overset{\text{only depends}}{\underset{\text{on } \beta_k}{\longleftarrow}} \Rightarrow \text{we can maximize this to get } \hat{\beta}_k$$

there exists fast computational ways to compute.