

Quiz

- ▶ What are the study designs for case-control study and cohort study?
- ▶ State the invariance property.
- ▶ If we care about RR (of disease), what study should we use?

Over Dispersion

- ▶ Definition: the variance of the response Y exceeds the nominal variance. $Y \sim \text{Bin}(m, \pi) \rightarrow \text{var}_n(Y) = m\pi(1 - \pi)$.
- ▶ Dispersion parameter:

$$a(\phi) = \frac{\text{var}(Y)}{\text{var}_n(Y)}.$$

- ▶ Example: we hypothesize that the support rate of NE Patriots is constant across 5 midwestern states. That is, the proportion of people in the populations of those states who would root for the team is constant. We perform polls by randomly sampling $m = 200$ people in each of the 5 states: Wisconsin 57, Michigan 113, Illinois 56, Iowa 121, Minnesota 153.

$$Y_i \sim \text{Bin}(200, \pi) \quad i=1, \dots, 5 \quad \text{var}_n(Y_i) = 50, \quad \text{var}(Y_i) = 1801$$

$$\hat{\pi} = \frac{\sum Y_i}{200 \times 5} = 0.5$$

$$\text{var}(Y_i) = \frac{1}{5} \sum_{i=1}^5 (Y_i - \bar{Y})^2 = 1801$$

$$\text{var}_n(Y_i) = 200 \times 0.5 \times 0.5 = 50$$

- ▶ Over dispersion is very common in practice (under-dispersion is relatively rare)
- ▶ It means the data may not exactly follow the hypothetical distribution
- ▶ Grouped data may have over dispersion
- ▶ Binary data ($m = 1$) do NOT have over dispersion
- ▶ Count data may also have over dispersion

Source of Over-Dispersion

- Intra-class correlation:

For $Y = \sum_{i=1}^m Y_i$, $Y_i \sim \text{Bin}(1, \pi)$, we have $\text{corr}(Y_i, Y_j) = \rho$. Then Y may NOT follow $\text{Bin}(m, \pi)$.

$$\begin{aligned}\mathbb{E}(Y) &= m\pi \\ \text{var}(Y) &= m\pi(1 - \pi)[1 + (m - 1)\rho]\end{aligned}$$

if $m > 1$ & $\rho > 0$, $1 + (m - 1)\rho > 1$



► Hierarchical sampling:

Assume m individuals form m/k clusters, each with cluster size k .

Each cluster has its own event rate which follows a distribution.

For example, assume $Y_i | p_i \sim \text{Bin}(k, p_i)$ and $p_i \sim (\pi, \delta\pi(1 - \pi))$ for $i = 1, \dots, m/k$. Then $Y = \sum_{i=1}^{m/k} Y_i$ actually has

$$\mathbb{E}(Y) = m\pi$$

$$\text{var}(Y) = m\pi(1 - \pi)[1 + (k - 1)\delta]$$

$$\text{Var}(Y) = (\text{Var}(Y | p_i)) + \text{Var}(E(Y | p_i))$$

► In particular, if $k = m$ and $p_i \sim \text{Beta}(a, b)$, we have $Y \sim \text{Beta} - \text{Binomial}(m, \pi, \delta)$, where

$$\mathbb{E}(Y) = m\pi$$

$$\text{var}(Y) = m\pi(1 - \pi)[1 + (m - 1)\delta]$$

and $\pi = \frac{a}{a+b}$ and $\delta = \frac{1}{1+a+b}$.

Parameter Estimation

To model over-dispersed binomial data, we assume

$$\mathbb{E}(y_i) = m_i \pi_i, \quad \text{var}(y_i) = m_i \pi_i (1 - \pi_i) \phi.$$

We need to estimate β and ϕ .

If ϕ is a constant independent of m_i

- ▶ Quasi-likelihood estimator: $\hat{\beta}_Q = \hat{\beta}_{MLE}$
- ▶ $\text{var}(\hat{\beta}_Q) = \phi \text{var}(\hat{\beta}_{MLE})$

Estimating Constant Dispersion

Assume the dispersion parameter ϕ is a constant for independent m_i

- Under the dispersion model, the generalized Pearson χ^2 is

$$G = \sum_{i=1}^n \frac{(y_i - m_i \hat{\pi}_i)^2}{m_i \hat{\pi}_i (1 - \hat{\pi}_i) \phi} \sim \chi^2(n-p)$$

$$\frac{\hat{\beta}_R}{\hat{\tau}_i} = \frac{\hat{\beta}_{ML}}{\hat{\tau}_i}$$

dep orig

- ϕ can be estimated from

$$\hat{\phi} = G_0 / (n-p)$$

where G_0 is the generalized Pearson χ^2 from the original model fitting without over-dispersion.

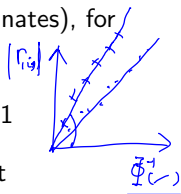
- Similarly, ϕ can be also estimated by $D_0 / (n-p)$ where D_0 is the deviance of the original model without over-dispersion.

Model Checking of Over Dispersion

Half Normal Plot:

- ▶ Order the absolute value of residuals (Pearson or deviance residuals) obtained from the Binomial model without dispersion
- ▶ Plot $|r_{(i)}|$ (y coordinates) against $\Phi^{-1}\left(\frac{n+i+0.5}{2n+1.125}\right)$ (x coordinates), for $i = 1, \dots, n$
- ▶ Reference line is a straight line through origin with slope 1
- ▶ Linear deviation from the reference line indicates constant over-dispersion
- ▶ Empirical slope is roughly $\sqrt{\phi}$

$$|r_{(i)}|$$



$$\frac{|r_{(i)}|}{\sqrt{\phi}}$$

Hypothesis Testing

$$H_0 : \beta = \beta_0 \text{ vs } H_1 : \beta \neq \beta_0$$

► Wald Test

$$TS = \underbrace{(\hat{\beta}_Q - \beta_0)}_{\parallel \hat{\beta}_{MLE}}^T \underbrace{\text{var}^{-1}(\hat{\beta}_Q)}_{\parallel \phi \cdot \text{var}^{-1}(\hat{\beta}_{MLE})} \underbrace{(\hat{\beta}_Q - \beta_0)}_{\parallel \hat{\beta}_{MLE}} = \underline{\underline{TS_W / \phi}}$$

where TS_W is the Wald Test statistic for the model ignoring dispersion.

► Score Test

$$TS = \underbrace{Q(\beta_0)}_{\parallel \frac{S(\beta_0)}{\phi}}^T \underbrace{\text{var}^{-1}(Q(\beta_0))}_{\parallel \phi \cdot \text{var}^{-1}(S(\beta_0))} \underbrace{Q(\beta_0)}_{\parallel \frac{S(\beta_0)}{\phi}} = \underline{\underline{TS_S / \phi}}$$

where TS_S is the Score Test statistic for the model ignoring dispersion.

Assume (β_1, β_2) is the coefficient vector.

$$H_0 : \beta_2 = \mathbf{0} \text{ vs } H_1 : \beta_2 \neq \mathbf{0}$$

Deviance Analysis

- ▶ Model 1: $\eta = \mathbf{Z}_1\beta_1$, Deviance D_1
- ▶ Model 2: $\eta = \mathbf{Z}_1\beta_1 + \mathbf{Z}_2\beta_2$, Deviance D_2
- ▶ Pearson χ^2 statistic G_0
- ▶ Calculate D_1, D_2, G_0 from binomial model without dispersion.
- ▶ $\frac{D_1 - D_2}{\phi} \sim \chi^2(p_2)$ and $\frac{G_0}{\phi} \sim \chi^2(n - p_1 - p_2)$ are approximately independent
- ▶ Therefore, the test statistic is

$$\frac{(D_1 - D_2)/p_2}{G_0/(n - p_1 - p_2)} \sim F(p_2, n - p_1 - p_2) \quad (\text{under } H_0)$$

Example

The table below listed the number of rats surviving the 21-day lactation period (Y), and the number of rats alive at four days in the same litter (m) in control and treated groups.

Group		1	2	3	4	5	6	7	8
Control (X=0)	Y	13	12	9	9	8	8	12	11
	m	13	12	9	9	8	8	13	12
Treated (X=1)	Y	12	11	10	9	10	9	9	8
	m	12	11	10	9	11	10	10	9

Group		9	10	11	12	13	14	15	16
Control (X=0)	Y	9	9	8	11	4	5	7	7
	m	10	10	9	13	5	7	10	10
Treated (X=1)	Y	8	4	7	4	5	3	3	0
	m	9	5	9	7	10	6	10	7