

Recap: overdispersion etc.

- ▶ saturated model: $y_i \sim \text{Bin}(n_i, p_i)$, $\tilde{p}_i = y_i/n_i$,

$$\ell(\tilde{p}) = \sum \{y_i \log(y_i/n_i) + (n_i - y_i) \log(1 - y_i/n_i)\}$$

$$(y_{ij} - \bar{y}_i)$$
$$\hat{\mu}_i = \bar{y}_i$$

- ▶ what's the saturated model for linear regression? what is the maximized log-likelihood for this model?

$$\bar{y}_i = \mu_i + \varepsilon_i$$

$$i = 1, \dots, n$$

- ▶ with binomial data, large-ish n_i , residual deviance compares regression model to saturated model
- ▶ if it's too large, we have the wrong model
- ▶ lack of independence among individual Bernoullis; a few outliers; wrong predictors
- ▶ estimate $\tilde{\phi} = X^2/(n - p)$
- ▶ inflate variance $\hat{\beta} \sim N(\beta, \tilde{\phi}(X^T W X))$

ELM p. 43,4

ELM p. 45

instead of $N(\beta, X^T W X)$

... overdispersion

```
> summary(bmod)
```

```
Call:
```

```
glm(formula = cbind(survive, total - survive) ~ location + period,  
     family = binomial, data = troutegg)
```

```
...
```

```
period8      -2.3256      0.2429  -9.573 < 2e-16 ***  
period11     -2.4500      0.2341 -10.466 < 2e-16 ***
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(Dispersion parameter for binomial family taken to be 1)
```

```
Null deviance: 1021.469  on 19  degrees of freedom  
Residual deviance:   64.495  on 12  degrees of freedom  
AIC: 157.03
```

```
> summary(bmod2)
```

```
Call:
```

```
glm(formula = cbind(survive, total - survive) ~ location + period,  
     family = quasibinomial, data = troutegg)
```

```
period8      -2.3256      0.5609  -4.146 0.001356 **  
period11     -2.4500      0.5405  -4.533 0.000686 ***
```

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
(Dispersion parameter for quasibinomial family taken to be 5.330358)
```

... overdispersion

SM §10.6, p.512

- ▶ $Y \mid \epsilon \sim \text{Bin}(m, \epsilon p)$
- ▶ $E(\epsilon) = 1, \quad \text{var}(\epsilon) = \xi$
- ▶ $E(Y) = E\{E(Y \mid \epsilon)\} = E(mp\epsilon) = mp$
- ▶ $\text{var}(Y) = \text{var}\{E(Y \mid \epsilon)\} + E\{\text{var}(Y \mid \epsilon)\}$

$$\begin{aligned} & \text{var}(mp\epsilon) + E[m p \epsilon (1 - p \epsilon)] \\ & = m^2 p^2 \xi + mp - mp^2 E\epsilon^2 \end{aligned}$$

- ▶ $\text{var}(Y) = m\{p(1 - p) + \xi p^2(m - 1)\}$ ntbc
- ▶ variance is larger than $mp(1 - p)$ see also ELM p.44
- ▶ can't be detected if $m = 1$ m plays the role of n

... overdispersion

SM §10.6, p.512

- ▶ $Y \mid \epsilon \sim \text{Bin}(m, \epsilon p)$
- ▶ $E(\epsilon) = 1, \quad \text{var}(\epsilon) = \xi$
- ▶ $E(Y) = E\{E(Y \mid \epsilon)\} = E(mp\epsilon) = mp$
- ▶ $\text{var}(Y) = \text{var}\{E(Y \mid \epsilon)\} + E\{\text{var}(Y \mid \epsilon)\}$

$$= mp(1-p) \cdot \phi \quad \phi = ?$$

- ▶ $\text{var}(Y) = m\{p(1-p) + \xi p^2(m-1)\}$ ntbc
- ▶ variance is larger than $mp(1-p)$ see also ELM p.44
- ▶ can't be detected if $m = 1$ m plays the role of n_i

... matched case-control studies

- ▶ suppose we have 1 : M matching one case, M matched controls

- ▶ for person i in matched set j , we have

$$y_{ij}, x_{ij}, \quad i = 0, 1, \dots, M$$

$$j = 1, \dots, n \quad y_{ij} \\ i = 0, \dots, M$$

- ▶ model:

$$\log \frac{p_j(x_{ij})}{1 - p_j(x_{ij})} = \alpha_j + x_{ij}^T \beta$$

$n+p$
 pars

Mn people

- ▶ different intercept for each matched set confounding variables
- ▶ same effect of covariates across patients and sets β
- ▶ data: in matched set j , we have 1 case (person 0) and M controls (persons 1, ..., M)
- ▶

$$\begin{aligned} \Pr(y_{0j} = 1 \mid \sum_{i=1}^M y_{ij} = 1) &= \frac{\Pr(y_{0j} = 1, y_{1j} = 0, \dots, y_{Mj} = 0)}{\Pr(y_{1j} = 0, \dots, y_{Mj} = 0)} \\ &= \frac{\exp(x_{0j}^T \beta)}{\sum_{i=0}^M \exp(x_{ij}^T \beta)} \end{aligned}$$

... matched case-control studies

- ▶ suppose we have 1 : M matching one case, M matched controls
- ▶ for person i in matched set j , we have

$$y_{ij}, x_{ij}, \quad i = 0, 1, \dots, M$$

- ▶ model:

$$\log \frac{p_j(x_{ij})}{1 - p_j(x_{ij})} = \alpha_j + x_{ij}^T \beta$$

- ▶ different intercept for each matched set confounding variables
- ▶ same effect of covariates across patients and sets β
- ▶ data: in matched set j , we have 1 case (person 0) and M controls (persons 1, ..., M)

Fix j !

$$\Pr(y_{0j} = 1 \mid \sum_{i=1}^M y_{ij} = 1) = \frac{\Pr(y_{0j} = 1, y_{1j} = 0, \dots, y_{Mj} = 0)}{\Pr(y_{1j} = 0, \dots, y_{Mj} = 0)}$$

$$P_n(\text{case is case} \mid \text{ex. 1 case}) = \frac{\exp(x_{0j}^T \beta)}{\sum_{i=0}^M \exp(x_{ij}^T \beta)}$$

... matched case-control studies

$$P_n(Y_0 = 1 \mid Y_0 + Y_1 = 1)$$

$$\begin{matrix} 1 & 0 & | & 1 \\ \hline (1 & 0) & \in & \Sigma_1 \\ (0 & 1) & \in & \end{matrix}$$

$$= P(Y_0 = 1, Y_1 = 0)$$

$$= \frac{e^{\alpha + \beta x_0}}{e^{\alpha + \beta x_0} + e^{\alpha + \beta x_1}}$$

$$P_n(Y_0 = 1, Y_1 = 0) + P_n(Y_0 = 0, Y_1 = 1)$$

$$= \frac{e^{\alpha + \beta x_0}}{1 + e^{\alpha + \beta x_0}} \cdot \frac{1}{1 + e^{\alpha + \beta x_1}}$$

$$\frac{e^{\alpha + \beta x_0} / (1 + e^{\alpha + \beta x_0}) \cdot 1 / (1 + e^{\alpha + \beta x_1})}{\dots}$$