## Multinomial Distribution
- Generalization/Extension to Binomial Distribution
- Binomial is a joint probability distribution

$$P(y_1,\ldots,y_J \mid \theta_1,\ldots,\theta_J,n) = \frac{n!}{y_1!\ldots y_J!}\,\theta_1^{y_1}\theta_2^{y_2}\ldots\theta_J^{y_J}$$

$J$: # of categories
$\theta_i$ are the respective probabilities of the categories $\theta_1+\ldots+\theta_J=1$
$n = y_1+y_2+\ldots+y_J$

When $J=2$  $P(y_1,y_2\mid\theta_1,\theta_2) = \frac{n!}{y_1!y_2!}\theta_1^{y_1}\theta_2^{y_2}$   $n=y_1+y_2$  $\theta_1+\theta_2=1$

$P(y_i\mid\theta,n) = \frac{n!}{y!(n-y)!}\theta^y(1-\theta)^{n-y}$  Binomial Distribution

- Even if multinomial is not a member of the exponential family, we can control $\vec{\theta}\,\vec{y}$ collected over N groups via a set of parameters $\vec{\beta}$

## Nominal Logistic Regression
The outcomes of experiments are in $J$ categories and there is no natural order amongst the response categories. One category is arbitrarily chosen as a reference category.
e.g. $\theta_1$. Then the logits for the other categories are defined by:

$$\text{Logit}(\theta_j) = \left[\frac{\theta_j}{\theta_1}\right] = x^T\beta_j \quad \forall j=2,\ldots,J$$

having the constraint $\sum_{j=1}^J \theta_j =1$
When the estimates $\hat\beta_j$ are computed, then:

$$\hat\theta_j = \hat\theta_1\,\exp(x^T\hat\beta_j) \quad \forall j=2,\ldots,J$$

$$\hat\theta_1 = \frac{1}{1+\sum_{j=2}^J \exp(x^T\hat\beta_j)}$$

or $\hat\theta_j = \frac{\exp(x^T\hat\beta_j)}{1+\sum_{j=2}^J \exp(x^T\hat\beta_j)}$   $\forall j=2,\ldots,J$   Softmax function → converting output to probability

$\hat\theta_1 + \hat\theta_2\exp(x^T\hat\beta_2)+\ldots+\hat\theta_J\exp(x^T\hat\beta_J)=1$

## Alligators Example

Linear model makes proportional to explanatory variables:

$$\log\left[\frac{\theta_j}{\theta_1}\right] = \beta_j^{indoor} + \beta_j^{sex} + \beta_j \, size + \beta_j \angle 1 + \beta_j \angle 2 + \beta_j \angle 3 \qquad \forall j = 2,\ldots,5$$

$$\theta_j = \frac{\exp(x^T \beta_j)}{1 + \sum_{j=2}^{J} \exp(x^T \beta_j)} \qquad \forall j = 2,\ldots,5 \quad \text{with } x = [1, sex, size, \angle 1, \angle 2, \angle 3]$$

$$\beta_j = [\beta_j^{indoor}, \beta_j^{sex}, \beta_j^{size}, \beta_j^{\angle 1}, \beta_j^{\angle 2}, \beta_j^{\angle 3}]$$

$$\theta_1 = \frac{1}{1 + \sum_{j=2}^{J} \exp(x^T \beta_j)}$$

## Function to maximize

$$\mathcal{L} = \prod_{i=1}^{16} P(y_{1i}, y_{2i}, y_{3i}, y_{4i}, y_{5i} \mid \theta_{1i}, \theta_{2i}, \theta_{3i}, \theta_{4i}, \theta_{5i})$$

$$n_i = y_{1i} + y_{2i} + y_{3i} + y_{4i} + y_{5i} \qquad \rightarrow \text{group sizes can be different}$$

$$\theta_{1i} = \frac{1}{1 + \sum_{j=2}^{J} \exp(x_i^T \hat{\beta})}$$

$$\theta_{ji} = \frac{\exp(x_i^T \hat{\beta}_j)}{1 + \sum_{j=2}^{J} \exp(x_i^T \hat{\beta}_j)}$$

$$\mathcal{L} = \prod_{i=1}^{16} \frac{n_i}{y_{1i}! \, y_{2i}! \, y_{3i}! \, y_{4i}! \, y_{5i}!} \; \theta_{1i}^{y_{1i}} \, \theta_{2i}^{y_{2i}} \, \theta_{3i}^{y_{3i}} \, \theta_{4i}^{y_{4i}} \, \theta_{5i}^{y_{5i}} \qquad \substack{J=4 \\ \theta_{j,i}} \quad d.f = 4 \times 6$$

$\uparrow$ saturated model

$\ell(\theta_{1i}, \theta_{2i}, \theta_{3i}, \theta_{4i}, \theta_{5i}) \quad d.f = 6 \times 4$

$\ell(\{\beta_j\} \; j = 2, \ldots 5)$

## ODDS RATIO

$\qquad\qquad\qquad\qquad\qquad$ i = group $\qquad$ j = vote choice

$$OR \text{ (example)} = \frac{\beta_{j=2, i=1}}{\beta_{j=1, i=1}} \qquad \text{ratio for group 1 women (18-23) } sex = 0 \; Age1 = 0 \; Age2 = 0$$

$$\frac{\beta_{j=2, i=4}}{\beta_{j=1, i=4}} \qquad \text{ratio for group 2 men (18-23) } Sex = 1 \; Age1 = 0 \; Age2 = 0$$

Category 2 = $y_2$ $\quad$ j=1 reference category

How to use $\exp(\text{coitres})$ to get OR ratio

1.5 $\Rightarrow$ indicates the option is more important to women than men

- Assumed use of Softmax function