



ETC3580: Advanced Statistical Modelling

Week 4: Binomial and proportion
responses

Outline

1 Binomial responses

2 Proportion responses

Binomial responses

Binomial distribution

Y is binomially distributed $B(m, p)$ if

$$P(Y = y) = \binom{m}{y} p^y (1 - p)^{m-y}$$

Y = number of "successes" in m independent trials, each with probability p of success.

- Likelihood

$$L = \prod_{i=1}^n \binom{m_i}{y_i} p_i^{y_i} (1 - p_i)^{m_i - y_i}$$

- Binomial regression also uses a logit link

$$p_i = e^{\eta_i} / (1 + e^{\eta_i})$$

$$\eta_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_q x_{i,q}$$

Binomial likelihood

Likelihood

$$L = \prod_{i=1}^n \binom{m_i}{y_i} p_i^{y_i} (1 - p_i)^{m_i - y_i}$$

$$\begin{aligned}\log L(\beta) &= \sum_{i=1}^n \left[\log \binom{m_i}{y_i} + y_i \log(p_i) + (m_i - y_i) \log(1 - p_i) \right] \\&= \sum_{i=1}^n \left[\log \binom{m_i}{y_i} + y_i \eta_i - y_i \log(1 + e^{\eta_i}) - (m_i - y_i) \log(1 + e^{\eta_i}) \right] \\&= \sum_{i=1}^n \left[\log \binom{m_i}{y_i} + y_i \eta_i - m_i \log(1 + e^{\eta_i}) \right]\end{aligned}$$

Binomial responses

In R:

- `glm` needs a two-column matrix of success and failures. (So rows sum to m).

```
fit <- glm(cbind(successes, failures) ~  
  x1 + x2,  
  family=binomial, data=df)
```

- Everything else works the same as for binary regression.

Overdispersion

- If mean correctly modelled, but observed variance larger than model, we called the data “overdispersed”. [Same for underdispersion.]
- Concept of overdispersion irrelevant for OLS and logistic regression because there cannot be any more variance than what is modelled.
- For binomial regression: $y_i \sim B(m_i, p_i)$, $E(y_i) = m_i p_i$, $\text{Var}(y_i) = m_i p_i (1 - p_i)$.
- If model correct, $D = -2 \log L \sim \chi^2_{n-q}$.
So $D > n - q$ indicates overdispersion.

Overdispersion

$D > n - q$ can also be the result of:

- missing covariates or interaction terms
- negligence of non-linear effects
- large outliers
- sampling from clusters
- non-independence
- m small (χ^2 approximation fails)

Overdispersion

Solution 1: Drop strict binomial assumption and let $E(y_i) = m_i p_i$, $\text{Var}(y_i) = \phi m_i p_i (1 - p_i)$.

Pearson residuals

$$r_i = \frac{y_i - m_i \hat{p}_i}{\sqrt{m_i \hat{p}_i (1 - \hat{p}_i)}}$$

Simple estimate of dispersion parameter

Estimate
$$\hat{\phi} = \frac{1}{n} \sum_{i=1}^n r_i^2.$$

- OK for estimation and standard errors on coefficients.
- But no proper inference via deviance.

Overdispersion

Solution 2: Define a “quasi-likelihood” that behaves like the log-likelihood but allows for

$$\text{Var}(y_i) = \phi m_i p_i (1 - p_i).$$

```
fit <- glm(cbind(successes, failures) ~  
  x1 + x2,  
  family=quasibinomial, data=df)
```

- Must use F -tests rather than χ^2 tests for comparing models. (Approximation only.)

Inference for GLMs

Model	<i>F</i> -test	χ^2 test
Normal OLS	Exact	–
Binary Logistic	Approx	Better approx
Binary Probit	Approx	Better approx
Binomial Logistic	Approx	Better approx
Quasibinomial logistic	Approx	–

Outline

1 Binomial responses

2 Proportion responses

Proportion responses: three approaches

Suppose $y \in [0, 1]$.

logitNormal model

$$\log(y/(1 - y)) \mid \mathbf{x} \sim N(\mathbf{x}'\boldsymbol{\beta}, \sigma^2)$$

quasiBinomial model

$y \mid \mathbf{x}$ has mean p and variance $\phi p(1 - p)$
where $\log(p/(1 - p)) = \mathbf{x}'\boldsymbol{\beta}$

Beta model

$y \mid \mathbf{x} \sim \text{Beta}(a, b)$
where $E(y \mid \mathbf{x}) = a/(a + b) = e^{\mathbf{x}'\boldsymbol{\beta}}/(1 + e^{\mathbf{x}'\boldsymbol{\beta}})$

logitNormal model

logitNormal model

$$\log(y/(1 - y)) \mid \mathbf{x} \sim N(\mathbf{x}'\beta, \sigma^2)$$

- Provided no empirical proportions are at either 0 or 1, we can compute a logit transformation of the observed proportions. $\log(y/(1 - y))$
- Then just fit a Gaussian linear regression using OLS.
- Back-transform the predictions using the inverse logit. $e^y/(1 + e^y)$

```
lm(log(y/(1-y)) ~ x1 + x2, data=df)
```

Quasi-binomial model

quasiBinomial model

$y \mid \mathbf{x}$ has mean p and variance $\phi p(1 - p)$
where $\log(p/(1 - p)) = \mathbf{x}'\beta$

- logit link keeps predicted proportions in $(0, 1)$
- Variance function $\phi p(1 - p)$ makes sense for proportions as greatest variation around $p = 0.5$ and least around $p = 0$ and $p = 1$.

```
glm(y ~ x1 + x2,  
    family=quasibinomial, data=df)
```

Beta regression

Beta model

$$\begin{aligned} & y \mid \mathbf{x} \sim \text{Beta}(a, b) \\ \text{where } & E(y \mid \mathbf{x}) = a/(a + b) = e^{\mathbf{x}'\beta} / (1 + e^{\mathbf{x}'\beta}) \end{aligned}$$

Beta density

$$f(y) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1}$$

where $y \in [0, 1]$ and $\Gamma(u) = \int_0^\infty x^{u-1} e^{-x} dx$.

$$\blacksquare \quad E(y) = \frac{a}{a+b} \quad \text{Var}(y) = \frac{ab}{(a+b)^2(a+b+1)}$$

Beta regression

Beta model

$$\begin{aligned} & y \mid \mathbf{x} \sim \text{Beta}(a, b) \\ \text{where} \quad & E(y \mid \mathbf{x}) = a/(a + b) = e^{\mathbf{x}'\beta} / (1 + e^{\mathbf{x}'\beta}) \end{aligned}$$

Beta density

$$f(y) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1}$$

where $y \in [0, 1]$ and $\Gamma(u) = \int_0^\infty x^{u-1} e^{-x} dx$.

$$\blacksquare E(y) = \frac{a}{a+b} \quad \text{Var}(y) = \frac{ab}{(a+b)^2(a+b+1)}$$

■ Reparameterize so $\mu = a/(a+b)$ and $\phi = a+b$.

Beta regression

Beta model

$$\begin{aligned} & y \mid \mathbf{x} \sim \text{Beta}(a, b) \\ \text{where} \quad & E(y \mid \mathbf{x}) = a/(a + b) = e^{\mathbf{x}'\beta} / (1 + e^{\mathbf{x}'\beta}) \end{aligned}$$

Beta density

$$f(y) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1}$$

where $y \in [0, 1]$ and $\Gamma(u) = \int_0^\infty x^{u-1} e^{-x} dx$.

$$\blacksquare E(y) = \frac{a}{a+b} \quad \text{Var}(y) = \frac{ab}{(a+b)^2(a+b+1)}$$

- Reparameterize so $\mu = a/(a + b)$ and $\phi = a + b$.
- Then $E(Y) = \mu$ and $\text{Var}(Y) = \mu(1 - \mu)/(1 + \phi)$.

Beta regression

Beta model

$$\begin{aligned} & y \mid \mathbf{x} \sim \text{Beta}(a, b) \\ \text{where } & E(y \mid \mathbf{x}) = a/(a + b) = e^{\mathbf{x}'\beta} / (1 + e^{\mathbf{x}'\beta}) \end{aligned}$$

Beta density

$$f(y) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} y^{a-1} (1-y)^{b-1}$$

where $y \in [0, 1]$ and $\Gamma(u) = \int_0^\infty x^{u-1} e^{-x} dx$.

$$\blacksquare E(y) = \frac{a}{a+b} \quad \text{Var}(y) = \frac{ab}{(a+b)^2(a+b+1)}$$

- Reparameterize so $\mu = a/(a + b)$ and $\phi = a + b$.
- Then $E(Y) = \mu$ and $\text{Var}(Y) = \mu(1 - \mu)/(1 + \phi)$.

```
mgcv::gam(y ~ x1 + x2, family=betar())
```