

Binary Regression

Hoff Chapter 10 -11, ISL Chapter 4

October 23, 2014

Seedling Survival

Tropical rain forests have up to 300 species of trees per hectare, which leads to difficulties when studying processes which occur at the community level. To gain insight into species responses, a sample of seeds were selected from a suite of eight species selected to represent the range of regeneration types which occur in this community.

Name	Size	Cotyledon type
Ardisia	3	H
C. biflora	7	H
Gouania	1	E
Hirtella	8	H
Inga	4	H
Maclura	2	E
C. racemosa	6	H
Strychnos	5	E

Size = 1 smallest to 8 largest
E = Epigeal - cotyledons
H = Hypogeal - seed food reserves

Experimental Design

This representative community was then placed in experimental plots manipulated to mimic natural conditions

- ▶ 8 PLOTS: 4 in forest gaps, 4 in understory conditions
- ▶ Each plot split in half: mammals were excluded from one half with a CAGE
- ▶ 4 subplots within each CAGE/NO CAGE
- ▶ 6 seeds of each SPECIES plotted in each SUBPLT
- ▶ 4 LITTER levels applied to each SUBPLT
- ▶ LIGHT levels at forest floor recorded
- ▶ SURV an indicator of whether they germinated and survived was recorded

Which variables are important in determining whether a seedling will survive? Are there interactions that influence survival probabilities?

Modeling Survival

Distribution for Survival of a single Seedling is Bernoulli

$$E[\text{SURV}_i \mid \text{covariates}] = \pi_i$$

How should we relate covariates to probability of survival?

For example, probability of survival may depend on whether there was a CAGE to prevent animals from eating the seedling or LIGHT levels.

- ▶ Naive approach: Regress SURV on CAGE and LIGHT

$$\hat{\pi}_i = \hat{\beta}_0 + \hat{\beta}_1 \text{CAGE}_i + \hat{\beta}_2 \text{LIGHT}_i$$

- ▶ Problems:
 - ▶ Fitted values of probabilities are not constrained to (0, 1)
 - ▶ Variances are not constant $\pi_i(1 - \pi_i)$ under Bernoulli model
- ▶ Unbiased?

Logistic Regression

To build in the necessary constraints that the probabilities are between 0 and 1 convert to log-odds or “logits”

- ▶ Odds of survival: $\pi_i / (1 - \pi_i)$

$$\text{logit}(\pi_i) \stackrel{\text{def}}{=} \log \left(\frac{\pi_i}{1 - \pi_i} \right) = \beta_0 + \beta_1 \text{CAGE}_i + \beta_2 \text{LIGHT}_i = \eta_i$$

- ▶ η_i is the linear predictor
- ▶ logit is the *link* function that relates the mean π_i to the linear predictor η_i
- ▶ Generalized Linear Models (GLMs)
- ▶ Find Maximum Likelihood Estimates (optimization problem)

Logits

To convert from the linear predictor η to the mean π , use the inverse transformation:

- ▶ $\log \text{ odds (SURV} = 1) = \eta$
- ▶ $\text{odds (SURV} = 1) = \exp(\eta) = \omega$
- ▶ $\pi = \text{odds} / (1 + \text{odds}) = \omega / (1 + \omega)$
- ▶ $\omega = \pi / (1 - \pi)$

Can go in either direction

Interpretation of Coefficients

$$\omega_i = \exp(\beta_0 + \beta_1 \text{CAGE}_i + \beta_2 \text{LIGHT}_i)$$

- ▶ When all explanatory variables are 0 (CAGE= 0, LIGHT= 0), the odds of survival are $\exp(\beta_0)$
- ▶ The ratio of odds (or odds ratio) at $X_j = A$ to odds at $X_j = B$, for fixed values of the other explanatory variables is

$$\text{Odds ratio} = \frac{\omega_A}{\omega_b} = \exp(\beta_j(A - B))$$

$$\text{Odds ratio} = \frac{\omega_A}{\omega_b} = \exp(\beta_j) \text{ if } A - B = 1$$

$$\text{Odds}(X_j = A) = \exp(\beta_j) \cdot \text{Odds}(X_j = B)$$

- ▶ Coefficients are log odds ratios

- ▶ use `glm()` rather than `lm()`
- ▶ model formula as before
- ▶ need to specify family (and link if not default)

```
seeds.glm0 = glm(SURV ~ 1, data=seeds, family=binomial)
seeds.glm1 = glm(SURV ~ CAGE + LIGHT,
                  family = binomial,
                  data=seeds)
plot(seeds.glm1)
summary(seeds.glm1)
anova(seeds.glm1, seeds.glm0, test="Chi")
```


Estimates

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-1.437e+00	7.095e-02	-20.258	< 2e-16	***
CAGE	7.858e-01	8.753e-02	8.977	< 2e-16	***
LIGHT	-4.683e-07	8.762e-08	-5.344	9.09e-08	***

- ▶ Coefficient for the dummy variable CAGE= 0.79
- ▶ If CAGE increases by 1 unit (No CAGE to CAGE) the odds of survival change by $\exp(.79) = 2.2$
- ▶ The odds of survival in a CAGE are 2.2 times higher than odds of survival in the open.

Confidence Intervals

- ▶ MLEs are approximately normally distributed (large samples)
 - ▶ mean β_j
 - ▶ estimated variance $SE(\beta_j)^2$
- ▶ Asymptotic posterior distribution for β_j is $N(\hat{\beta}_j, SE(\beta_j)^2)$
- ▶ $(1 - \alpha)100\%$ CI based on normal theory:

$$\hat{\beta}_j \pm Z_{\alpha/2} SE(\beta_j)$$

- ▶ 95% CI for coefficient for CAGE:

$$0.79 \pm 1.960.088 = (0.62, 0.96)$$

- ▶ Exponentiate to obtain interval for odds ratio:
 $\exp(0.62), \exp(0.96) = (1.85, 2.62)$

The odds of survival in a CAGE are 1.85 to 2.62 times higher than odds of survival in the open (with probability 0.95).

Deviance

The concept of Deviance replaces Sum-of-Squares in GLMs

- ▶ residual deviance = $-2 \log$ likelihood at MLEs
- ▶ null deviance = residual deviance under model with constant mean (Total Sum of Squares)
- ▶ analysis of deviance
- ▶ change in (residual) deviance has an asymptotic χ^2 distribution with degrees of freedom based on the change in parameters

```
> anova(seeds.glm1, test="Chi")
```

Analysis of Deviance Table

Terms added sequentially (first to last)

	Df	Deviance	Resid. Df	Resid. Dev	P(> Chi)
NULL			3071	3426.1	
CAGE	1	91.6	3070	3334.5	1.065e-21
LIGHT	1	35.5	3069	3299.0	2.548e-09

Other Variables

```
> seeds.glm3 = glm(SURV ~ SPECIES + CAGE + log(LIGHT) +  
  factor(LITTER), data=seeds, family=binomial)  
> summary(seeds.glm3)
```

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-0.83171	0.25280	-3.290	0.00100	**
SPECIESC. biflora	0.18742	0.17000	1.103	0.27024	
SPECIESC. racemosa	0.84686	0.16320	5.189	2.11e-07	***
SPECIESGouania	-2.63505	0.36003	-7.319	2.50e-13	***
SPECIESHirtella	1.14131	0.16246	7.025	2.14e-12	***
SPECIESInga	0.90420	0.16295	5.549	2.87e-08	***
SPECIESMaclura	-2.63505	0.36003	-7.319	2.50e-13	***
SPECIESStrychnos	-1.18652	0.21736	-5.459	4.80e-08	***
CAGE	0.92548	0.09654	9.586	< 2e-16	***
log(LIGHT)	-0.08997	0.02001	-4.496	6.92e-06	***
factor(LITTER)1	0.09099	0.13494	0.674	0.50013	
factor(LITTER)2	0.24257	0.13425	1.807	0.07079	.
factor(LITTER)4	-0.13625	0.13531	-1.007	0.31396	

Analysis of Deviance

```
> anova(seeds.glm3, test="Chi")
```

Analysis of Deviance Table

Model: binomial, link: logit

Response: SURV

Terms added sequentially (first to last)

	Df	Deviance	Resid. Df	Resid. Dev	P(> Chi)
NULL			3071	3426.1	
SPECIES	7	572.2	3064	2853.9	2.353e-119
CAGE	1	109.3	3063	2744.6	1.399e-25
log(LIGHT)	1	15.3	3062	2729.3	9.191e-05
factor(LITTER)	3	7.9	3059	2721.4	4.769e-02

Interactions?

The presence of a CAGE may be more important for survival for some species than others - implies an interaction

The odds of survival | Cage compared to odds of survival | no Cage depend on SPECIES

Fit model with upto 4 way interactions:

```
seeds.glm4 = glm(SURV~SPECIES*CAGE*log(LIGHT)*LITTER,  
                 data=seeds, family=binomial)
```

The analysis of deviance test suggests that there are three way interactions

Hierarchical Model

So far we have not taken into account all the sources of variation or information about the experimental design.

- ▶ SPECIES (size & cotyledon type) and LITTER are randomized to sub-plots. Expect that survival of seedlings in the same sub-plot may be related, which suggests a sub-plot random effect.
- ▶ sub-plots are nested within CAGE within plots (so expect that sub-plots in the same CAGE are correlated, as well as sub-plots within the same plot may have a similar survival.
- ▶ Plots characteristics may affect survival (light levels)

How to model?