

□ Solution: Will formulate model using moments { formula page }:

$$E(y_i) = \mu_i$$

$$\text{var}(y_i) = \left( \frac{\phi_i}{w_i} \right) \cdot v(\mu_i) \text{ for some function } v$$

$$g(\mu_i) = x_i^T \beta \text{ for link function } g$$

and use quasi-likelihood to estimate  $\beta$  and  $\phi$

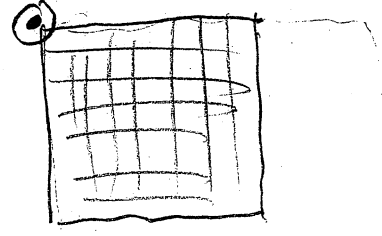
□ Later; now look at Meron data { meron.pdf }

$Y_i = \text{tree. } g_2 = \text{binary response}$

□  $(x, y)$ : geographical trend

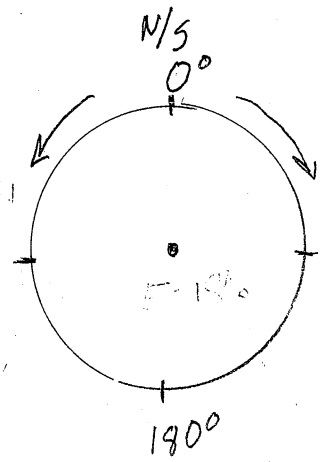
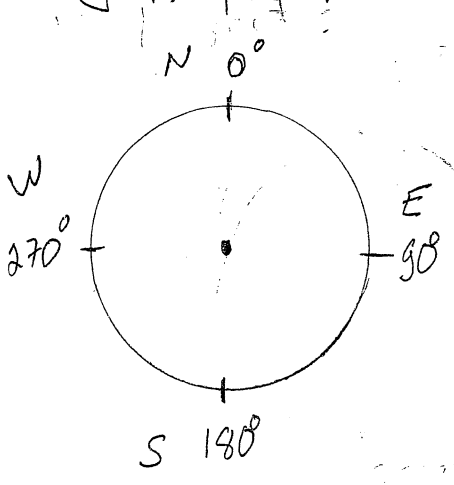
□ grazing variables: →

Kibbutz Sasa

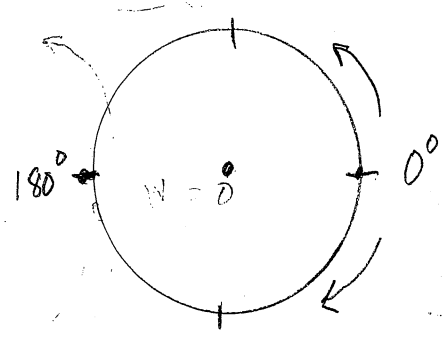


E/W

□ Aspect:



South: more sun  
(=> fewer trees) W

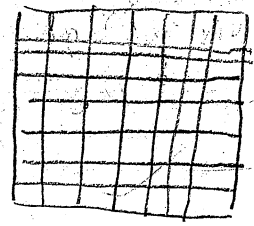


West: more moisture from Mediterranean

Remarks: ... response

Full data set:  $75 \times 80 = 6,000$  pixels

⇒ spatial dependence



Original response:

original  $y = \begin{cases} \% \text{ tree cover} \\ \% \text{ shrub cover} \\ \% \text{ grass} \end{cases}$

- multivariate response

constraint:

$\mu_1 + \mu_2 + \mu_3 = 1$

⇒ compositional data

→ vegetation index

Could define vegetation index

$y = 1$  (all grass), ...,  $7$  (all trees)

⇒  $y$  is ordinal variable, ordered polychotomous regression

Simplification

$y = \begin{cases} 1 & \text{if trees dominant } (\% \text{ tree} > \begin{cases} \% \text{ shrub, } \% \text{ grass} \\ \text{or} \\ \% \text{ shrub} + \% \text{ grass} \end{cases}) \\ 0 & \text{if not} \end{cases}$

□ Exploratory Analysis of Meron data (31)  
{meron.preliminary.pdf}

{p.1}

→ confounding of grazing type & intensity

{p.2} (like loc-host and south in med fly data)

- ew. aspect only up to 100

{p.3} (obs. in a region  $\approx$  sloping towards the east)

{p.4} Effect on trees:

- Fewer trees (Tree.92 = 0) when N/S aspect high  
(N/S = 0  $\Leftrightarrow$  north; so high N/S means  $\approx$  south and more sun)

- more trees (Tree.92 = 1) when E/W aspect high  
(E/W = 0  $\Leftrightarrow$  east, so high E/W aspect means  $\approx$  west  $\Rightarrow$  water from sea)

- higher altitude  $\approx$  lower road distance ( $r = -0.23$ )

↓  
fewer people

{p.5} □ linear regression results.

- interpretation of slope, aspect, slope  $\times$  aspect coeff's - later

- cattle at low intensity??

- fertilizer?

- {p.3: cattle, low  $\approx$   $\Leftrightarrow$

{ alt. high  
E/W aspect  $\approx$  high  
(w: more moisture)

□ Before fitting GLM to Meron, another link function (32)

□ So far formula page

Canonical links  $\eta_i = \beta_0 + \beta_1 x_i = g(\mu_i) = \mu_i$  (normal)

$\{g(\mu_i) = \theta_i\}$

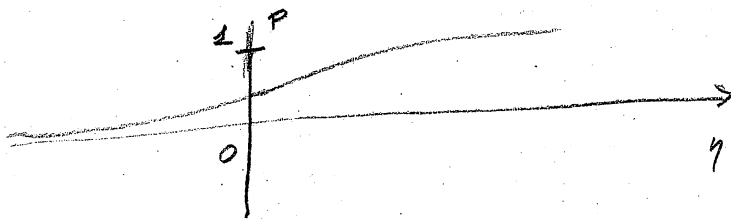
$= g(\lambda_i) = \log(\lambda_i)$  (Poisson)

$= g(p_i) = \log \frac{p_i}{1-p_i}$  (Bernoulli)

□ Another link function for Bernoulli:

Want  $g: (0, 1) \rightarrow (-\infty, \infty)$

$g^{-1}: (-\infty, \infty) \rightarrow (0, 1)$



$g^{-1}(\eta) = \Phi(\eta)$ ,  $\Phi = N(0, 1)$  CDF

$g(p) = \Phi^{-1}(p)$ , probit link function  
("probability unit")

□ Interpretation of logit:  $\log\left(\frac{p}{1-p}\right) = \log$  "odds"

□ Interpretation of probit?

$$y_i \sim \text{Ber}(p_i),$$

$$p_i = \Phi(\eta_i), \quad \eta_i = x_i^T \beta$$

$$= P(Z_i \leq \eta_i), \quad Z_i \sim N(0, 1).$$

□ interpretation / motivation:

$y = 1$  when underlying latent variable

$W$  exceeds an unknown

threshold  $c$ .

$$(Z = -W :$$

$$W > c$$

$$\Leftrightarrow$$

$$Z < -c$$

EG  $W =$  "heart risk index"

$$W \sim N(\beta_0 + \beta_1 x_1 + \beta_2 x_2, \sigma^2)$$

where

$$x_1 = \text{BMI}$$

( $\beta_i$ : probably  $> 0$ )

$$x_2 = \text{cholesterol level}$$

$Y = 1$  if has heart attacks.

Model:  $Y = 1_{\{W > c\}}$  where  $c$  is unknown threshold

Then

$$P(Y=1) = P(W > c)$$

$$= P(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon > c), \quad \epsilon \sim N(0, \sigma^2)$$

$$= P(-\epsilon < \frac{\beta_0 + \beta_1 x_1 + \beta_2 x_2 - c}{\sigma})$$

$$= P(Z < \underbrace{\left(\frac{\beta_0 - c}{\sigma}\right)}_{\beta_0} + \underbrace{\left(\frac{\beta_1}{\sigma}\right)}_{\beta_1} x_1 + \underbrace{\left(\frac{\beta_2}{\sigma}\right)}_{\beta_2} x_2), \quad Z = -\epsilon \sim N(0, 1)$$

$\Rightarrow$  Probit model

Note:  $\beta_j > 0 \iff \beta_j > 0$

□ logit or probit? Saw (p. 21), notes

- that for any  $\underline{x} = (x_1, \underline{x}_{(2)}, x_p)$ ,

and  $\frac{p}{1-p}$  = "odds for  $y=1$ ",

$$\frac{\text{odds}(x_1 + \Delta, \underline{x}_{(2)})}{\text{odds}(x_1, \underline{x}_{(2)})} = e^{\beta_1 \Delta}$$

Q: How interpret  $\beta_1$ ?

$$\Rightarrow \frac{\text{odds}(x_1 + \Delta, \underline{x}_{(2)}) - \text{odds}(x_1, \underline{x}_{(2)})}{\Delta} = \text{odds}(x_1, \underline{x}_{(2)}) \frac{[e^{\beta_1 \Delta} - 1]}{\Delta}$$

$$\begin{aligned} \Delta \rightarrow 0 \\ \Rightarrow \frac{\partial}{\partial x_1} \text{odds}(x_1, \underline{x}_{(2)}) &= \text{odds}(x_1, \underline{x}_{(2)}) \cdot \underbrace{\frac{d}{dx} e^{\beta_1 x}}_{\beta_1} \Big|_{x=0} \end{aligned}$$

□ So  $\beta_j$  gives change in odds as fcn of  $x_j$  on a multiplicative scale,

for any  $\underline{x}$ .

U So:

(36)

- Motivation for probit:

possible latent variable

- Motivation for logit:

- "log odds" interpretation, of  $\beta$ , for any  $x$

- also useful in retrospective sampling.



$D$   $Y = \text{rare disease} = \begin{cases} 1 \\ 0 \end{cases}$ ,  $X = \text{risk factor} = \begin{cases} 1 \\ 0 \end{cases}$

Want  $P(Y=1|X=1)$  vs  $P(Y=1|X=0)$

Prospective sample:

100 individuals w  $X=0$   
100 ———  $X=1$

observe  $Y$  & compare

Problem: might get all  $Y_i = 0$

Retrospective sample

100 with  $Y=0$   
100 ———  $Y=1$

compute  $X$  & compare

Problem Can only estimate  $P(X|Y)$ ; want  $P(Y|X)$

Sol'n: Use data to estimate odds ratio

$$\frac{\text{odds}(X=1|Y=1)}{\text{odds}(X=1|Y=0)} = \frac{P(X=1|Y=1) / [1 - P(X=1|Y=1)]}{P(X=1|Y=0) / [1 - P(X=1|Y=0)]}$$

$$= \frac{P(X=1|Y=1) \cdot P(X=0|Y=0)}{P(X=1|Y=0) \cdot P(X=0|Y=1)} \cdot \frac{P(Y=1)}{P(Y=1)} \cdot \frac{P(Y=0)}{P(Y=0)}$$

$$= \frac{P(Y=1, X=1) \cdot P(Y=0, X=0)}{P(Y=0, X=1) \cdot P(Y=1, X=0)}$$

$$= \frac{P(Y=1, X=1) / P(Y=0, X=1)}{P(Y=1, X=0) / P(Y=0, X=0)} \cdot \frac{P(X=1) / P(X=1)}{P(X=0) / P(X=0)}$$

$$= \frac{P(Y=1|X=1) / P(Y=0|X=1)}{P(Y=1|X=0) / P(Y=0|X=0)}$$

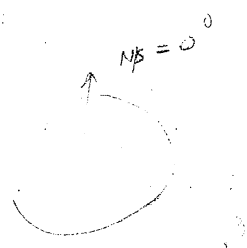
$$= \frac{\text{odds}(Y=1|X=1)}{\text{odds}(Y=1|X=0)}$$

Example: Linear-logistic & probit regression

- (1) Regressors are logistic functions
- (2) Probit is logistic, but with more variables are significant
- (3) compare logit & probit

Effect of slope?

I. change in predicted prob of 10% ( $\hat{\beta} > 0$ )



1.1  
 1.2  
 1.3  
 1.4  
 1.5  
 1.6  
 1.7  
 1.8  
 1.9  
 2.0

□ Now: GLIM for Meron { meron.glim.fit.pdf } (39)

Linear regression (from before) significant  $\hat{\beta}_i$ :

alt, water.dist, cattle.low  $\hat{\beta}_i > 0$   
 stream.dist, slope: ns.aspect  $\hat{\beta}_i < 0$

Logistic regression

alt, water.dist  $\hat{\beta}_i > 0$   
 stream.dist, slope: ns.aspect  $\hat{\beta}_i < 0$

Probit regression

alt, water.dist, ns.aspect  $\hat{\beta}_i > 0$   
 stream.dist, slope: ns.aspect  $\hat{\beta}_i < 0$

- □ → Deviance Residuals
  - Null Deviance Residual Deviance
  - Dispersion parameter (=  $\phi$ )
  - →  $\chi^2$  test
  - Null, Residual deviance
  - $\chi^2$  test
  - AIC
- } later  
} later  
} trees

□  $\chi^2$  test

□  $\chi^2$  test

□  $\chi^2$  test

Question (probit results):

A = N/S aspect {A=0: North}

S = slope

$$\hat{\beta}_A = 0.2$$

$$\hat{\beta}_S = 4.8$$

$$\hat{\beta}_{A \times S} = -0.18$$

Is facing north (A=0) good or bad for trees?

Answer depends on slope

$$\Phi^{-1}(p) \approx c + 4.8 S + 0.02 A - 0.18 A \times S$$

$$= c + 0.02 A + [4.8 - 0.18 A] \cdot S$$

- Need a graph {meron.interaction.graphs.pdf}

- For n = 6,000, assuming  $y_i$  independent

solid line is top curve

Top graph: facing North is good for trees and

Bottom graph: facing West " " " "

IF the slope is large enough