

Maximum-likelihood fitting of the logistic regression models

STAT 206
19 Mar 21
Extra Lecture

$(y_i | p_i) \sim \text{Bernoulli}(p_i) \quad (i = 1, \dots, n)$ ①

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \sum_{j=1}^k x_{ij} \beta_j = \beta_0 + (\beta) x_i$$

[LR] = regression structure

$y = (y_1, \dots, y_n)$ observed $X = (x_1, \dots, x_n)$

in which

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}$$

is the matrix of predictor values and

$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$ is the vector of logistic regression coefficients

$$\theta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

The sampling distribution is

$$p(y | \beta_0, \beta \text{ [LR] } \mathcal{B}) =$$

$$\prod_{i=1}^n p(y_i | \beta_0, \beta \text{ [LR] } \mathcal{B})$$

is the vector of unknown parameters

$\{\theta\} = k+1$
 $\dim(\theta) =$ unknowns
 $(n > k+1)$

and this is

$$p(z|\beta, \beta [LR] \mathcal{B}) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i} \quad \text{expressed} \quad (2)$$

in terms of the probabilities $p_i = P(\Sigma_i = 1 | \mathcal{B})$

so the likelihood function is $\mathcal{P} = (p_1, \dots, p_n)$

$$l(\mathcal{P} | z [LR] \mathcal{B}) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

$$= \prod_{i=1}^n p_i^{y_i} (1-p_i)^{-y_i} (1-p_i)^{1-y_i} = \prod_{i=1}^n (1-p_i) \left(\frac{p_i}{1-p_i} \right)^{y_i}$$

and the log likelihood function is then

$$\begin{aligned}
 \ell(\mathcal{P} | z [LR] \mathcal{B}) &= \left[\sum_{i=1}^n \log(1-p_i) \right] + \left[\sum_{i=1}^n y_i \log \left(\frac{p_i}{1-p_i} \right) \right] \\
 &= \left[\sum_{i=1}^n \log(1-p_i) \right] + \left\{ \sum_{i=1}^n y_i \left[\beta_0 + \sum_{j=1}^k x_{ij} \beta_j \right] \right\} \\
 &= \left[\sum_{i=1}^n \log(1-p_i) \right] + \left(\beta_0 \sum_{i=1}^n y_i \right) + \left(\sum_{i=1}^n y_i \sum_{j=1}^k x_{ij} \beta_j \right)
 \end{aligned}$$

interchanging the order of summation in the ⁽³⁾
 last term gives

$$\sum_{i=1}^n \gamma_i \sum_{j=1}^k x_{ij} \beta_j = \sum_{i=1}^n \sum_{j=1}^k x_{ij} \gamma_i \beta_j$$

$$= \sum_{j=1}^k \sum_{i=1}^n x_{ij} \gamma_i \beta_j = \sum_{j=1}^k \beta_j \left(\sum_{i=1}^n x_{ij} \gamma_i \right)$$

so we "see" that there is a $(k+1)$ dimensional
 set of minimal sufficient statistics (MSS),
 namely $\left\{ \underbrace{\left(\sum_{i=1}^n \gamma_i \right)}_{S_y} \text{ and } \underbrace{\left(\sum_{i=1}^n x_{ij} \gamma_i \right)}_{(S_{xy})_j} \text{ for each } j=1, \dots, k \right\} \textcircled{*}$

but wait: what about the annoying
 $\sum_{i=1}^n \log(1-p_i)$ term? Not so fast with
 the MSS claim, although $\textcircled{*}$ was a
 highly useful observation: those $(k+1)$

quantities can be computed once, before ⁽⁴⁾ the iterative search for the MLEs begins

To deal with the $\sum_{i=1}^n \log(1-p_i)$ term, we need to solve the equation $\left\{ \begin{array}{l} \ell \text{ for logit} \\ \text{not likelihood} \end{array} \right.$

$$\log\left(\frac{p_i}{1-p_i}\right) = \ell_i \leftarrow \text{backwards for } p_i :$$

$$\frac{p_i}{1-p_i} = e^{\ell_i}$$

$$p_i = (1-p_i)e^{\ell_i} = e^{\ell_i} - p_i e^{\ell_i}$$

$$p_i(1 + e^{\ell_i}) = e^{\ell_i} \rightarrow p_i = \frac{e^{\ell_i}}{1 + e^{\ell_i}}$$

$$\text{and } 1-p_i = \frac{1}{1 + e^{\ell_i}}$$

so the pesky term is

$$\sum_{i=1}^n \log(1-p_i) = \sum_{i=1}^n \log\left(\frac{1}{1 + e^{\ell_i}}\right) = - \sum_{i=1}^n \log(1 + e^{\ell_i})$$

and the full loglikelihood function $l(\beta)$ in terms of the β_j 's

$$l(\beta_0, \beta | \mathbf{y}, (\mathbf{X}, \mathbf{R}) \beta) = - \sum_{i=1}^n \log \left\{ 1 + \exp(\beta_0 + \sum_{j=1}^k \beta_j x_{ij}) \right\} + \beta_0 \sum_{i=1}^n y_i + \sum_{j=1}^k \beta_j \left(\sum_{i=1}^n x_{ij} y_i \right)$$

Covariance information between β and \mathbf{y}

This clearly cannot be maximized in closed form; iterative numerical methods are needed to find the MLEs. The standard

numerical method, Fisher scoring, is just an application of Newton-Raphson, and usually converges quickly to a (local) maximum of the log likelihood function;

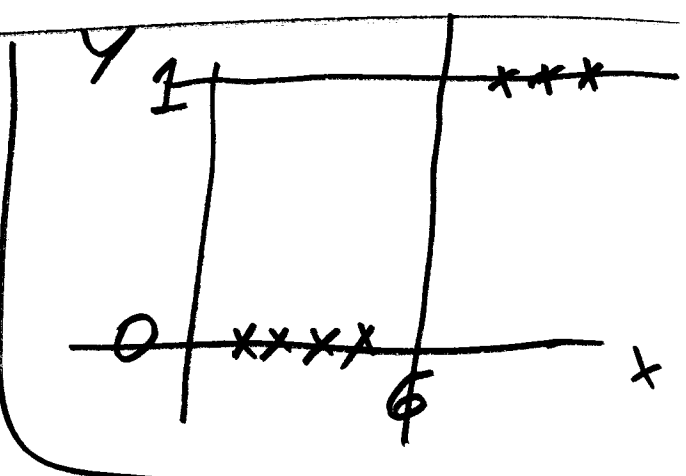
This method is also an instance of an approach called Iteratively Reweighted Least Squares (IRLS).

Convergence from last page The MLEs are issues fail to exist in this model in 2

situations: ① If $\begin{pmatrix} X \\ n \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ is not of full rank

ie., if there are linear dependencies among the columns of X ; R solves this by booting predictors out of the model that cause the

collinearity ② Here for this dataset there's a separating hyperplane



(in this case, the vertical line $x=6$) such that $(x < 6)$ perfectly predicts $y=0$ and $(x > 6)$ perfectly predicts $y=1$.

The log likelihood function may be maximized in this situation at $\beta = \pm \infty$ or ludicrous standard errors may emerge (to be continued) someday

$(Y | p \in \mathcal{B}) \sim \text{Bernoulli}(p)$
 (PMF)

$P(Y=y | p \in \mathcal{B}) =$

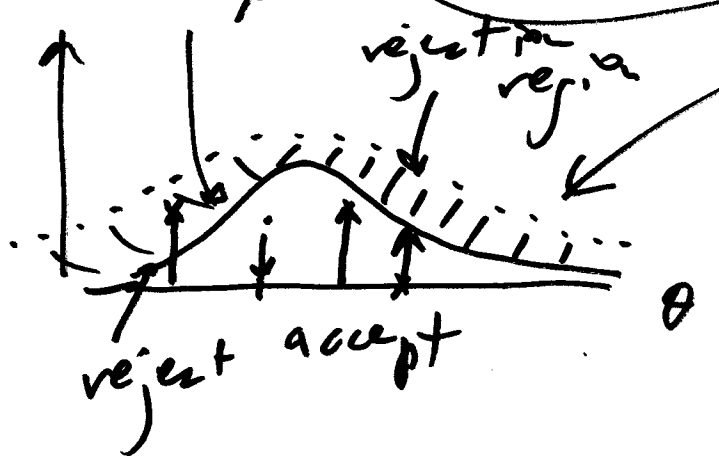
$E(Y)$

$p^y (1-p)^{1-y}$

target PDF
 $p(\theta)$

$I_{\{0,1\}}(y)$

$$= \begin{cases} p & \text{if } y=1 \\ 1-p & 0 \\ 0 & \text{else} \end{cases}$$



envelope function

function $G(\theta)$

① $G(\theta) \geq p(\theta)$

② it's integrable

$\int_{\mathbb{R}} G(\theta) d\theta < \infty$

normalizable

$g(\theta) = \frac{G(\theta)}{\int_{\mathbb{R}} G(\theta) d\theta}$

③ g should be easy to sample from

to sample from

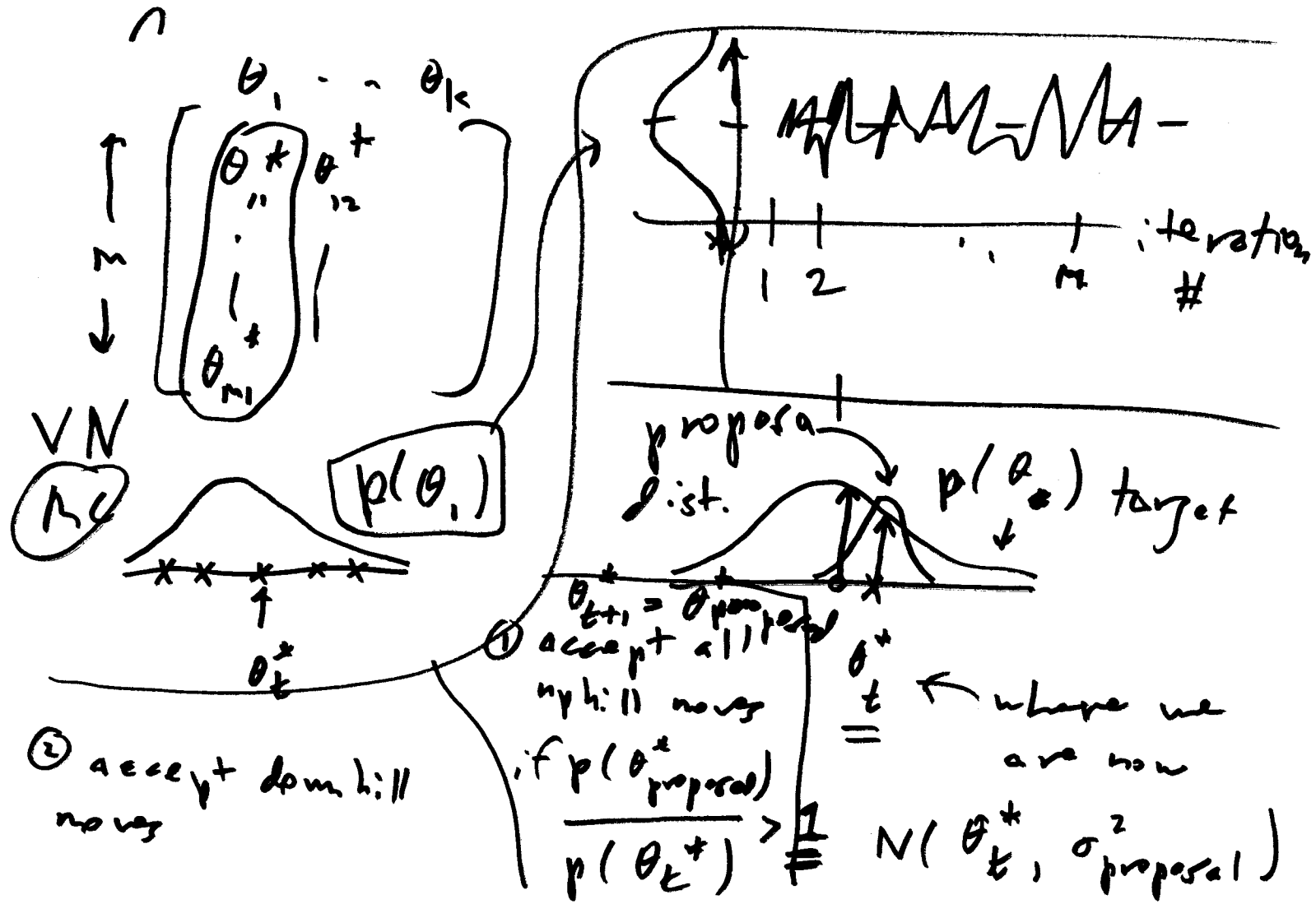
v. N: $\boxed{I \& D}$ but low efficiency for large k ⑧

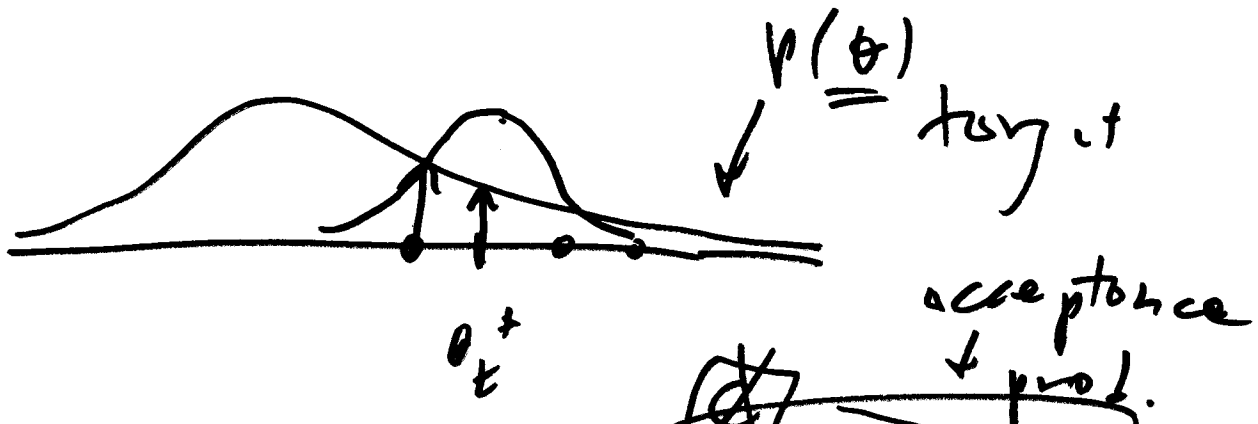
$\theta \sim (\theta_1, \dots, \theta_k)$

Metr. et al. / can't get rid of \boxed{D} , but how about relaxing \boxed{I} ?

from the target $p(\theta)$

Constructing a try 1st order Markov chain

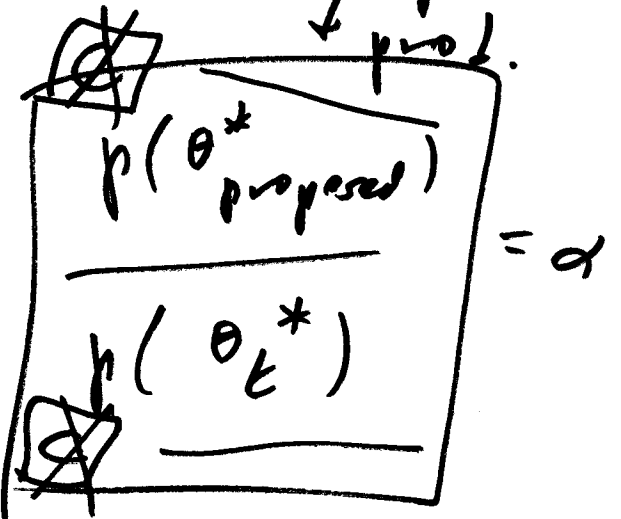




① accept all uphill moves:

$$\theta_{t+1}^* = \theta_{\text{proposed}}^*$$

if $\alpha \geq 1$ accept



\rightarrow accept with prob. $\min(\alpha, 1)$

② accept downhill moves with prob α

if reject, $\theta_{t+1}^* = \theta_t^*$

stay where you are (!)

$$p(\theta | -) = c \cdot p(\theta) \ell(\theta | -)$$