

Logistic regression
(Bernoulli outcome)

$$(Y_i | p_i) \sim \text{Bernoulli}(p_i)$$

$(\text{if diabetes}) (y_i = 1), \dots, n$ ← 768
 $y = (y_1, \dots, y_n)$

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

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

$$p(\beta | \mathcal{X}, \mathcal{Y}) \leftarrow \text{prior for } \beta$$

likelihood function

$$l(\beta | \mathcal{Y}, \mathcal{X}, \mathcal{Z}, \mathcal{B}) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

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

$n \times k$

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

log-likelihood function

$$ll(\beta | \mathcal{Y}, \mathcal{X}, \mathcal{Z}, \mathcal{B}) = \sum_{i=1}^n \left[y_i \log p_i + (1-y_i) \log (1-p_i) \right]$$

some algebra reveals that

$$\log \left(\frac{p_i}{1-p_i} \right) = L_i \Leftrightarrow p_i = \frac{1}{1 + \exp(-L_i)}$$

Bayes's Theorem

$$p(\beta | \mathcal{Z} \cap \mathcal{B}) =$$

$$c p(\beta | \mathcal{Z} \cap \mathcal{B}) \cdot \mathcal{L}(\beta | \mathcal{Y} \cap \mathcal{Z} \cap \mathcal{B})$$

log rule

log posterior

log likelihood

$$\log p(\beta | \mathcal{Z} \cap \mathcal{B}) = c + \mathcal{L}(\beta | \mathcal{Y} \cap \mathcal{Z} \cap \mathcal{B})$$

$$+ \log p(\beta | \mathcal{Z} \cap \mathcal{B})$$

log prior

Consider,

eg., a Normal prior for β

$$p(\beta | \mathcal{Z}_n \cap \mathcal{B}) \sim N_k(\mu_{\beta}, \Sigma_{\beta})$$

$$p(\beta | \mathcal{Z}_n \cap \mathcal{B}) = c \exp \left[-\frac{1}{2} (\beta - \mu_{\beta})^T \Sigma_{\beta}^{-1} (\beta - \mu_{\beta}) \right]$$

take $\Sigma_{\beta} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_k^2 \end{bmatrix}$ and $\mu_{\beta} = \underline{0}$

then

$$\log p(\beta | \mathcal{Z} | N, \mathcal{B}) = -\frac{1}{\gamma} \sum_{j=1}^k \left(\frac{\beta_j}{\sigma_j} \right)^2$$

for some $\gamma > 0$

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

$$\|\beta\|^2 = \beta^T \beta = \sum_{j=1}^k \beta_j^2$$

log posterior

$$\log p(\beta | \mathcal{Z} | \Sigma | \mathcal{B}) =$$

log likelihood

$$L(\beta | \mathcal{Z} | \Sigma | \mathcal{B})$$

$$- \frac{1}{\gamma} \sum_{j=1}^k \left(\frac{\beta_j}{\sigma_j} \right)^2$$

← log prior

Consider estimating β not by maximizing the log likelihood but by maximizing

the log posterior: result is $\hat{\beta}_{MAP}$ called maximum a posteriori

but look what the prior is doing: it's

shrinking the $\hat{\beta}_{MLE}$ vector toward $\underline{0}$

Why might this be a good idea?

imagine k (# predictors) growing with
 n (# observations) fixed turns out
(makes sense)

that $\hat{\beta}_{MLE}$ becomes more & more

unstable, in the sense that

$SE\left[\left(\hat{\beta}_{MLE}\right)_j\right] \uparrow$ sharply

$$\begin{array}{c|c} \mathbf{X} & \hat{\mathbf{y}} \\ \hline y_1 & \hat{y}_1 \\ \vdots & \vdots \\ y_n & \hat{y}_n \end{array}$$

general
prediction
problem

a natural way to measure
how good $\hat{\mathbf{y}}$ is in predicting \mathbf{y}

is with $RMS\hat{E} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$

In terms of random variables, the square of
of MSE is $MSE(\underline{Y}, \hat{\underline{Y}}) = E[(\underline{Y} - \hat{\underline{Y}})^2]$

fact: $MSE = [\text{bias}(\hat{\underline{Y}})]^2 + V(\hat{\underline{Y}})$

excellent
idea: if we deliberately bias $\hat{\underline{\beta}}_{MLE}$
by shrinking it toward $\underline{0}$, the
variance of the resulting $\hat{\underline{\beta}}_{MAP}$ may
go down more than the $[\text{bias}(\hat{\underline{\beta}}_{MAP})]^2$
goes up, yielding a reduction in (R)MSE

This turns out to work better & better
as $k \uparrow$
