

2/03/16 ALSM 2
INFERENCE

Algorithm: To calculate β 's

Usually $\hat{\theta} = \arg \max_{\theta} \log(\text{like}(\theta))$ usually differentiate this function and equate to 0

Normal Distribution

link g : identity function [mapping \mathbb{R} onto \mathbb{R}]

$$p(y|\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y-\theta)^2}{2\sigma^2}\right]$$

likelihood wrt some parameter β : $\text{lik}(\beta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y_i - \theta_i)^2}{2\sigma^2}\right]$

$\theta_i = x_i^T \beta$ - combination of mean.

$$\log(\text{lik}(\beta)) = -N \log(\sqrt{2\pi}\sigma) - \sum (y_i - x_i^T \beta)^2 / 2\sigma^2$$

$$\text{SSE} = \sum (y_i - x_i^T \beta)^2$$

$$\hat{\beta} = \arg \min_{\beta} \text{SSE} \quad \text{(now we get our } \beta\text{'s)}$$

$$\text{SSE} = (\|y - X\beta\|)^2 \quad (\text{norm})^2$$

$$\vec{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1n} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nn} \end{bmatrix}$$

When we compute derivative of SSE: $\frac{d(\text{SSE})}{d\beta} = \frac{d \{ (y - X\beta)^T (y - X\beta) \}}{d\beta}$

$$= \frac{d}{d\beta} [y^T y - y^T X\beta - (X\beta)^T y + (X\beta)^T X\beta]$$

$$= 0 + (y^T X)^T - X^T y + 2X^T X \beta = 0$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

1. Differentiate by Method of Lagrange

2. Set $\nabla \text{Lagrange} = 0$

3. Solving system of equations defined by $\nabla \text{Lagrange} = 0$

Newton-Raphson

-locking at $\log(\text{like}(\theta))$ near $\hat{\theta}$

$$\log \text{like}(\theta) = \log \text{like}(\theta^i) + (\theta - \theta^i)' \nabla \log \text{like}(\theta^i) + \frac{1}{2} (\theta - \theta^i)' [\nabla^2 \log \text{like}(\theta^i)] (\theta - \theta^i) + \text{remainder}$$

$\nabla \log \text{like}(\theta^i)$: gradient at θ^i

$\nabla^2 \log \text{like}(\theta^i)$: Hessian matrix at θ^i : H_{θ^i}

After differentiation: $\nabla \log \text{like} \approx H_{\theta^i} (\theta - \theta^i) + \nabla \log \text{like}(\theta^i) = 0$

Newton-Raphson Method

Initialize $\theta^i = \theta^0$

$$\theta^{(i+1)} = \theta^{(i)} - [H_{\theta^i}]^{-1} \nabla \log \text{like}(\theta^i) \quad (\text{Inverse} \times \text{gradient})$$

Until convergence

- If loglike is quadratic, NR will find solution in one step (i.e. normal distribution)
- In GLM, only one extreme, no local extremes
- The quantity $[-H]$ (Hessian - observed information) determines the sharpness of the peak in the likelihood function and its maximum

Alternative to Newton-Raphson: Method of Scoring

$$I = E[-H] \quad (\text{Expected Fisher Information})$$

Replaces J by I in NR method

Multinomial Distribution and Poisson Random Variables

- Consider 2 independent random variables that follow a poisson distribution, i.e.

$$y_1 \sim P_{\theta_1}(\lambda_1) \quad \text{Show that } n = y_1 + y_2 \sim P_{\theta_1}(\lambda_1 + \lambda_2)$$

$$y_2 \sim P_{\theta_2}(\lambda_2)$$

$$p(y_1 | \lambda_1) = \frac{\lambda_1^{y_1}}{y_1!} \exp(-\lambda_1)$$

$$p(n) = \sum_{y_1=0}^{\infty} \sum_{y_2=0}^{\infty} p(n, y_1, y_2) \quad \text{Partitioning}$$

$$= \sum_{y_1=0}^{\infty} p(n | y_1, y_2) p(y_1, y_2) \quad (\text{Bayes}) \quad (\text{independent})$$

$$= \sum_{y_1=0}^{\infty} p(n - y_1, y_2) p(y_1) p(y_2) \quad \leftarrow$$

10/3/16

Atam

$$\text{Integrate: } = \sum_{y_2=0}^n P(y_1, n-y_1) P_{y_2}(y_2)$$

$$= \sum_{y_2=0}^n \frac{\lambda_1^{n-y_1}}{(n-y_1)!} \exp(-\lambda_1) \frac{\lambda_2^{y_2} \exp(-\lambda_2)}{y_2!}$$

$$= \sum_{y_2=0}^n \frac{\lambda_1^{n-y_1}}{(n-y_1)!} \frac{\lambda_2^{y_2}}{y_2!} \exp(-\lambda_1 - \lambda_2) \quad \text{expecting: } p(n) = \frac{(\lambda_1 + \lambda_2)^n}{n!} \exp(-(\lambda_1 + \lambda_2))$$

$$= \sum_{y_2=0}^n \frac{n!}{y_2! (n-y_2)!} \lambda_1^{n-y_2} \lambda_2^{y_2}$$

Show

$$\text{Joint probability } p(y_1, y_2 | n) = \frac{n!}{y_1! y_2!} \theta_1^{y_1} \theta_2^{y_2}$$

$$\text{Given } y_i \sim P(\lambda_i) \quad \forall i=1, 2 \quad n = y_1 + y_2 + y_3 \quad \text{with } \theta_i = \frac{\lambda_i}{\sum_{j=1}^3 \lambda_j}$$

$$p(y_1, y_2 | n) = \frac{p(n | y_1, y_2) p(y_1, y_2)}{p(n)} \quad \text{Bayes}$$

$$= \frac{\delta(n - (y_1 + y_2)) \pi_{j=1}^3 \frac{\lambda_j^{y_j}}{y_j!} \exp(-\lambda_j)}{(\sum_{j=1}^3 \lambda_j)^n \exp(-(\lambda_1 + \dots + \lambda_3))}$$

Limitation of likelihood function

- Assumes responses are independent
- Outliers in data can give a zero, all data except one group voting for parameter β
- May not be robust because of outliers
- Hard to remove outliers in large datasets
- Some problem with Bayesian approach

Linear Regression Scott Roper

$$p(y | x, \beta) = N(y, x^T \beta, \sigma^2) \quad \text{Hypothesis / Model}$$

$$\epsilon \in \epsilon_i = y_i - x_i^T \beta$$

$$\epsilon = y - x^T \beta \sim N(0, \sigma^2)$$

Take empirical distribution

$$p(\epsilon) = \frac{1}{n} \sum_{i=1}^n \delta(\epsilon, \epsilon_i)$$

↑ Dirac function

4

Get empirical distribution as close as possible to model distribution
i.e. get β for a minimizing distance between two functions!

Scott: Sum (inner product): $L.E \{ z_{i,j}^n p(y_i | x_i, \beta) \}$

allows for "or" to Scott β

can have multiple gaussian for creation of pdf