

7. Point Estimation

Method of Moments Estimation

Definition 7.1

Let X_1, \dots, X_n be a random sample then the **sample moments** are defined as

$$m_r = \frac{1}{n} \sum_{i=1}^n X_i^r$$

- The method of moments estimators arise from equating the sample moments m_r with the population moments $\mu'_r = E[X^r]$.
- In most situations we can write the population moments in terms of the parameter vector $\theta = (\theta_1, \dots, \theta_k)$.

Definition 7.2

Suppose that X_1, \dots, X_n is a random sample from a distribution with parameter vector $\theta = (\theta_1, \dots, \theta_k)$. Let μ'_r be the moments of X and assume that μ'_k exists and is finite. Let m_r be the sample moments of X_1, \dots, X_n .

The **method of moments estimator** of θ is given by the solution to the k simultaneous equations

$$m_r = \mu'_r(\theta_1, \dots, \theta_k) \quad r = 1, \dots, k$$

Maximum Likelihood Estimation

Definition 7.3

Suppose $x = (x_1, \dots, x_n)$ are the observed values of n iid random variables from a family of distributions indexed by the unknown, possibly vector, parameter $\theta \in \Theta$ and let the likelihood be $L(\theta | x)$. Then the **maximum likelihood estimate** of θ is a value $\hat{\theta}(x)$ such that

$$L(\hat{\theta}(x) | x) \geq L(\theta | x) \quad \forall \theta \in \Theta$$

The **maximum likelihood estimator** of θ is the random variable $\hat{\theta}(X)$.

- If the likelihood $L(\boldsymbol{\theta}; \mathbf{x})$ is differentiable in $\boldsymbol{\theta}$, then an interior maximum of L can be given by solving the k equations

$$\frac{\partial}{\partial \theta_i} L(\theta_1, \dots, \theta_k; \mathbf{x}) = 0 \quad i = 1, \dots, k.$$

- The maximum value found this way is **not** necessarily the mle.
- Usually easier to maximize the *log-likelihood*

$$l(\boldsymbol{\theta} | \mathbf{x}) = \log L(\boldsymbol{\theta} | \mathbf{x})$$

- In the *iid* case

$$l(\boldsymbol{\theta} | \mathbf{x}) = \sum_{i=1}^n \log f_X(x_i | \boldsymbol{\theta})$$

Theorem 7.1

If $\hat{\theta}$ is the maximum likelihood estimate of a parameter θ and $\eta = g(\theta)$ is a transformation of the parameter then the maximum likelihood estimate of η is

$$\hat{\eta} = g(\hat{\theta}).$$

Theorem 7.2

Let X_1, \dots, X_n be a sample from a population with density $f(x | \theta)$ and let $T(\mathbf{X})$ be the minimal sufficient statistic for θ . Then the maximum likelihood estimator of θ depends on the sample only through the value of $T(\mathbf{X})$.

Newton–Raphson Method

- The mle is a solution to

$$U(\hat{\theta}) = \left. \frac{\partial l(\theta | \mathbf{x})}{\partial \theta} \right|_{\theta=\hat{\theta}} = 0$$

- Observed information matrix

$$J(\hat{\theta}) = - \left. \frac{\partial^2 l(\theta | \mathbf{x})}{\partial \theta \partial \theta^T} \right|_{\theta=\hat{\theta}}$$

- From Taylor's Theorem

$$U(\hat{\theta}) \approx U(\theta^*) - J(\theta^*)(\hat{\theta} - \theta^*)$$

- Hence

$$\hat{\theta} \approx \theta^* + J^{-1}(\theta^*)U(\theta^*)$$

1. Choose a reasonable starting value $\hat{\theta}^{(0)}$.

2. Update your estimate

$$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} + J^{-1}(\hat{\theta}^{(k)}) U(\hat{\theta}^{(k)})$$

3. Terminate the algorithm when

$$\|\hat{\theta}^{(k)} - \hat{\theta}^{(k+1)}\| < \varepsilon.$$

for some pre-determined tolerance $\varepsilon > 0$.

Fisher Scoring Method

Replace $J(\hat{\theta}^{(k)})$ with the **expected Fisher information** matrix evaluated at the current step

$$I(\hat{\theta}^{(k)}) = E_{\theta} \left[-\frac{\partial^2 l(\theta | \mathbf{x})}{\partial \theta \partial \theta^T} \right] \Bigg|_{\theta = \hat{\theta}^{(k)}}$$

- Choice of starting value is very important.
- It is possible that there are multiple solutions to the likelihood equation including minima and saddlepoints.
- It is often useful to start from multiple starting points to ensure that a global maximum is found.

The EM Algorithm

- Useful when there is missing data.
- \mathbf{y} is the observed (incomplete) data with log likelihood $l(\theta | \mathbf{y})$.
- \mathbf{x} is the complete data with log likelihood $l_c(\theta | \mathbf{x})$.
- Define the quantity

$$Q(\theta, \theta^*) = \mathbb{E}_{\theta^*} [l_c(\theta | \mathbf{X}) | \mathbf{y}]$$

- Usually easier to maximize $Q(\theta, \theta^*)$.

1. Define the complete data \mathbf{X} and its log likelihood $l_c(\theta | \mathbf{X})$
2. Choose an initial estimate $\hat{\theta}^{(0)}$.

3. E Step: Calculate

$$Q(\theta, \hat{\theta}^{(k)}) = \mathbb{E}_{\hat{\theta}^{(k)}} [l_c(\theta | \mathbf{X}) | \mathbf{y}]$$

4. M Step: Choose $\hat{\theta}^{(k+1)}$ to maximize $Q(\theta, \hat{\theta}^{(k)})$.

5. Iterate the E and M steps until

$$L(\hat{\theta}^{(k+1)} | \mathbf{y}) - L(\hat{\theta}^{(k)} | \mathbf{y}) < \varepsilon.$$

Theorem 7.3

If $\hat{\theta}^{(k)}$ is a sequence of iterates such that

$$Q(\hat{\theta}^{(k+1)}, \hat{\theta}^{(k)}) \geq Q(\hat{\theta}^{(k)}, \hat{\theta}^{(k)})$$

Then the incomplete data likelihood is monotone increasing

$$L(\theta^{(k)} | \mathbf{y}) \geq L(\hat{\theta}^{(k)} | \mathbf{y}).$$

Proof relies on [Jensen's Inequality](#)

If X is a random variable with finite mean and g is a concave function such that $E[|g(X)|] < \infty$ then

$$E[g(X)] \leq g(E[X]).$$

EM Algorithm in the Full Exponential Family

- Suppose that the density of the complete data \mathbf{X} is

$$f_{\mathbf{X}}(\mathbf{X} \mid \theta) = h(\mathbf{x})c(\theta) \exp \left\{ \sum_{i=1}^d \theta_i t_i(\mathbf{x}) \right\}.$$

- Then we can write

$$Q(\theta, \hat{\theta}^{(k)}) = \log c(\theta) + \sum_{i=1}^d \theta_i \mathbb{E}_{\hat{\theta}^{(k)}} [t_i(\mathbf{X}) \mid \mathbf{y}].$$

- Hence $\hat{\theta}^{(k+1)}$ satisfies

$$\begin{aligned} \mathbb{E}_{\hat{\theta}^{(k)}} [\mathbf{t}(\mathbf{X}) \mid \mathbf{y}] &= - \left. \frac{\partial \log(c(\theta))}{\partial \theta} \right|_{\theta=\hat{\theta}^{(k+1)}} \\ &= \mathbb{E}_{\hat{\theta}^{(k+1)}} [\mathbf{t}(\mathbf{X})]. \end{aligned}$$

EM Gradient Algorithm

- Do not maximize Q but take one step of a Newton–Raphson algorithm.

$$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} - \left[\frac{\partial^2 Q(\theta, \hat{\theta}^{(k)})}{\partial \theta \partial \theta^T} \right]_{\theta = \hat{\theta}^{(k)}}^{-1} \left[\frac{\partial Q(\theta, \hat{\theta}^{(k)})}{\partial \theta} \right]_{\theta = \hat{\theta}^{(k)}}.$$

- Not guaranteed to be monotone
- Alternative is to let $0 < a^{(k)} \leq 1$ and let

$$\hat{\theta}^{(k+1)} = \hat{\theta}^{(k)} - a^{(k)} \left[\frac{\partial^2 Q(\theta, \hat{\theta}^{(k)})}{\partial \theta \partial \theta^T} \right]_{\theta = \hat{\theta}^{(k)}}^{-1} \left[\frac{\partial Q(\theta, \hat{\theta}^{(k)})}{\partial \theta} \right]_{\theta = \hat{\theta}^{(k)}}.$$

where $a^{(k)}$ is chosen to ensure that

$$Q(\hat{\theta}^{(k+1)}, \hat{\theta}^{(k)}) \geq Q(\hat{\theta}^{(k)}, \hat{\theta}^{(k)})$$

Monte Carlo EM Algorithm

- Cannot calculate $Q(\theta, \hat{\theta}^{(k)})$ in closed form.
- Suppose that we can sample $\mathbf{x}_1^{(k)}, \dots, \mathbf{x}_R^{(k)}$ from the conditional distribution $f_{\mathbf{X}|\mathbf{Y}}(\mathbf{x} | \mathbf{y}, \hat{\theta}^{(k)})$.

- Define

$$\hat{Q}(\theta, \hat{\theta}^{(k)}) = \frac{1}{R} \sum_{j=1}^R l_c(\theta | \mathbf{x}_j^{(k)})$$

- Now choose $\hat{\theta}^{(k+1)}$ to maximize $\hat{Q}(\theta, \hat{\theta}^{(k)})$.

- MCEM Algorithm will follow the EM “path” with random noise.
- Assessment of convergence is harder since convergence behaves like a monotone function with random noise.
- Variability of the random noise will depend on the value of R .
- It is common to take R small for the first steps and then increase R for later iterations.

Bayes Estimation

- In Bayesian statistics, the parameter θ is considered a random variable.
- The marginal distribution for θ is the probability distribution before any data is collected and so is called the **prior distribution** for θ and is commonly denoted $\pi(\theta)$.
- The model is now considered to be a conditional distribution of the data for a given value of the parameter.
- Hence the joint distribution of the data and parameter is the product of the likelihood and the prior.

- The conditional distribution of the parameter given the observed data is then

$$\pi(\theta | \mathbf{x}) = \frac{f(\mathbf{x} | \theta)\pi(\theta)}{\int f(\mathbf{x} | \theta)\pi(\theta) d\theta}$$

- This is called the **posterior distribution** for θ .
- The posterior distribution is then used to make inference about the random quantity θ .
- One obvious point estimator for θ is the mean of the posterior distribution

$$\hat{\theta} = E[\theta | \mathbf{x}] = \frac{\int \theta f(\mathbf{x} | \theta)\pi(\theta) d\theta}{\int f(\mathbf{x} | \theta)\pi(\theta) d\theta}$$

Conjugate Priors

Definition 7.4

Given a family \mathcal{F} of pdf's (or pmf's) $f(x | \theta)$ indexed by a parameter θ , then a family, Π of prior distributions is said to be **conjugate** for the family \mathcal{F} if the posterior distribution of θ is in the family Π for all $f \in \mathcal{F}$, all priors $\pi(\theta) \in \Pi$ and all possible data sets x .

- Conjugate families often make the mathematics of Bayesian statistics easier.
- They may not adequately describe the prior knowledge about the parameter.

Non-informative Priors

- How do we express prior ignorance about a parameter?
- If the set of possible values of θ is a finite interval then we may use a uniform distribution.
- Most parameter spaces, however are infinite.

Definition 7.5

Suppose that θ is a parameter with prior distribution $\pi(\theta)$. The prior distribution is called *improper* if

$$\int \pi(\theta) d\theta = \infty.$$

Definition 7.6

Suppose that $\mathbf{X} \sim f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})$. The *Fisher Information Matrix* is

$$I(\boldsymbol{\theta}) = -\mathbb{E} \left[\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \log f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) \right]$$

The *Jeffrey's Prior* for $\boldsymbol{\theta}$ is

$$\pi(\boldsymbol{\theta}) \propto \sqrt{|I(\boldsymbol{\theta})|}.$$

- Attempt to provide a general noninformative prior.
- The Jeffrey's Prior is often improper
- Usually, however,

$$\int f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} < \infty$$