# Bayesian Maximum Likelihood

- Bayesians describe the mapping from prior beliefs about $\theta$, summarized in $p(\theta)$, to new posterior beliefs in the light of observing the data, $Y^{data}$.

- General property of probabilities:

$$p\left(Y^{data}, \theta\right) = \left\{ \begin{array}{c} p\left(Y^{data}|\theta\right) \times p(\theta) \\ p\left(\theta|Y^{data}\right) \times p\left(Y^{data}\right) \end{array} \right. ,$$

which implies Bayes' rule:

$$p\left(\theta|Y^{data}\right) = \frac{p\left(Y^{data}|\theta\right) p(\theta)}{p\left(Y^{data}\right)},$$

mapping from prior to posterior induced by $Y^{data}$.

**Bayesian Maximum Likelihood ...**

● Properties of the posterior distribution, $p\left(\theta|Y^{data}\right)$.

  – The value of $\theta$ that maximizes $p\left(\theta|Y^{data}\right)$ ('mode' of posterior distribution).

  – Graphs that compare the marginal posterior distribution of individual elements of $\theta$ with the corresponding prior.

  – Probability intervals about the mode of $\theta$ ('Bayesian confidence intervals')

  – Other properties of $p\left(\theta|Y^{data}\right)$ helpful for assessing model 'fit'.

**Bayesian Maximum Likelihood ...**

- Computation of mode sometimes referred to as 'Basyesian maximum likelihood':

$$\theta^{\mathrm{mod}\,e} = \arg \max_{\theta} \left\{ \log \left[ p \left( Y^{data} | \theta \right) \right] + \sum_{i=1}^{N} \log \left[ p_i \left( \theta_i \right) \right] \right\}$$

maximum likelihood with a penalty function.

- Shape of posterior distribution, $p \left( \theta | Y^{data} \right)$, obtained by Metropolis-Hastings algorithm.
  - Algorithm computes
  $$\theta \left( 1 \right), ..., \theta \left( N \right),$$

  which, as $N \to \infty$, has a density that approximates $p \left( \theta | Y^{data} \right)$ well.

  - Marginal posterior distribution of any element of $\theta$ displayed as the histogram of the corresponding element $\{ \theta \left( i \right), i = 1, .., N \}$

# Metropolis-Hastings Algorithm (MCMC)

- We have (except for a constant):

$$f\left(\underbrace{\theta}_{N\times 1}|Y\right) = \frac{f\left(Y|\theta\right)f\left(\theta\right)}{f\left(Y\right)}.$$

- We want the marginal posterior distribution of $\theta_i$ :

$$h\left(\theta_i|Y\right) = \int_{\theta_{j\neq i}} f\left(\theta|Y\right)d\theta_{j\neq i}, \ i = 1, ..., N.$$

- MCMC algorithm can approximate $h\left(\theta_i|Y\right)$.

- Obtain ($V$ produced automatically by gradient-based maximization methods):

$$\theta^{\mod e} \equiv \theta^* = \arg\max_\theta f\left(Y|\theta\right)f\left(\theta\right), \ V \equiv \left[-\frac{\partial^2 f\left(Y|\theta\right)f\left(\theta\right)}{\partial\theta\partial\theta'}\right]^{-1}_{\theta=\theta^*}.$$

**Metropolis-Hastings Algorithm (MCMC) ...**

- Compute the sequence, $\theta^{(1)}, \theta^{(2)}, ..., \theta^{(M)}$ ($M$ large) whose distribution turns out to have pdf $f(\theta|Y)$.

  - $\theta^{(1)} = \theta^*$

  - to compute $\theta^{(r)}$, for $r > 1$

    * step 1: select candidate $\theta^{(r)}$, $x$,

    $$\text{draw } \underbrace{x}_{N \times 1} \text{ from } \theta^{(r-1)} + kN \left( \overbrace{\underbrace{0}_{N \times 1}, V}^{\text{'jump' distribution'}} \right), \; k \text{ is a scalar}$$

    * step 2: compute scalar, $\lambda$ :
    $$\lambda = \frac{f(Y|x) f(x)}{f\left(Y|\theta^{(r-1)}\right) f\left(\theta^{(r-1)}\right)}$$

    * step 3: compute $\theta^{(r)}$ :
    $$\theta^{(r)} = \begin{cases} \theta^{(r-1)} & \text{if } u > \lambda \\ x & \text{if } u < \lambda \end{cases}, \; u \text{ is a realization from uniform } [0, 1]$$

**Metropolis-Hastings Algorithm (MCMC) ...**

- Approximating marginal posterior distribution, $h\left(\theta_i | Y\right)$, of $\theta_i$

  – compute and display the histogram of $\theta_i^{(1)}, \theta_i^{(2)}, ..., \theta_i^{(M)}$, $i = 1, ..., N$.

- Other objects of interest:

  – mean and variance of posterior distribution $\theta$ :

$$E\theta \simeq \bar{\theta} \equiv \frac{1}{M} \sum_{j=1}^{M} \theta^{(j)}, \ Var\left(\theta\right) \simeq \frac{1}{M} \sum_{j=1}^{M} \left[\theta^{(j)} - \bar{\theta}\right] \left[\theta^{(j)} - \bar{\theta}\right]'.$$

**Metropolis-Hastings Algorithm (MCMC) ...**

● Some intuition

 – Algorithm is more likely to select moves into high probability regions than into low probability regions.

 – Set, $\left\{ \theta^{(1)}, \theta^{(2)}, ..., \theta^{(M)} \right\}$, populated relatively more by elements near mode of $f(\theta|Y)$.

 – Set, $\left\{ \theta^{(1)}, \theta^{(2)}, ..., \theta^{(M)} \right\}$, also populated (though less so) by elements far from mode of $f(\theta|Y)$.

**Metropolis-Hastings Algorithm (MCMC) ...**

- Practical issues

    - what value should you set $k$ to?

        * set $k$ so that you accept (i.e., $\theta^{(r)} = x$) in step 3 of MCMC algorithm are roughly 27 percent of time

    - what value of $M$ should you set?

        * a value so that if $M$ is increased further, your results do not change

            · in practice, $M = 10,000$ (a small value) up to $M = 1,000,000$.

    - large $M$ is time-consuming. Could use Laplace approximation (after checking its accuracy) in initial phases of research project.

# Laplace Approximation to Posterior Distribution

- In practice, Metropolis-Hasting algorithm very time intensive. Do it last!

- In practice, Laplace approximation is quick, essentially free and very accurate.

- Let $\theta \in R^N$ denote the $N-$dimensional vector of parameters and

$$g\left(\theta\right) \equiv \log f\left(y|\theta\right) f\left(\theta\right),$$

$f\left(y|\theta\right)$ ~likelihood of data

$f\left(\theta\right)$ ~prior on parameters

$\theta^*$ ~maximum of $g\left(\theta\right)$ (i.e., mode)

**Laplace Approximation to Posterior Distribution ...**

- Second order Taylor series expansion about $\theta = \theta^*$ :

$$g\left(\theta\right) \approx g\left(\theta^*\right) + g_\theta\left(\theta^*\right)\left(\theta - \theta^*\right) - \frac{1}{2}\left(\theta - \theta^*\right)' g_{\theta\theta}\left(\theta^*\right)\left(\theta - \theta^*\right),$$

  where

$$g_{\theta\theta}\left(\theta^*\right) = -\frac{\partial^2 \log f\left(y|\theta\right) f\left(\theta\right)}{\partial\theta\partial\theta'}\Big|_{\theta=\theta^*}$$

- Interior optimality implies:

$$g_\theta\left(\theta^*\right) = 0, \ \ g_{\theta\theta}\left(\theta^*\right) \text{ positive definite}$$

- Then,

$$f\left(y|\theta\right) f\left(\theta\right) \simeq f\left(y|\theta^*\right) f\left(\theta^*\right) \exp\left\{-\frac{1}{2}\left(\theta - \theta^*\right)' g_{\theta\theta}\left(\theta^*\right)\left(\theta - \theta^*\right)\right\}.$$

**Laplace Approximation to Posterior Distribution ...**

● Note

$$\frac{1}{(2\pi)^{\frac{N}{2}}} \left| g_{\theta\theta}\left(\theta^*\right) \right|^{\frac{1}{2}} \exp\left\{ -\frac{1}{2}\left(\theta - \theta^*\right)' g_{\theta\theta}\left(\theta^*\right)\left(\theta - \theta^*\right) \right\}$$

$$= \text{ multinormal density for } N - \text{dimensional random variable } \theta$$

with mean $\theta^*$ and variance $g_{\theta\theta}\left(\theta^*\right)^{-1}$ .

● So, posterior of $\theta_i$ (i.e., $h\left(\theta_i | Y\right)$) is approximately
$$\theta_i \sim N\left(\theta_i^*, \left[g_{\theta\theta}\left(\theta^*\right)^{-1}\right]_{ii}\right).$$

● This formula for the posterior distribution is essentially free, because $g_{\theta\theta}$ is computed as part of gradient-based numerical optimization procedures.

**Laplace Approximation to Posterior Distribution ...**

- Marginal likelihood of data, $y$, is useful for model comparisons. Easy to compute using the Laplace approximation.

- Property of Normal distribution:

$$\int \frac{1}{(2\pi)^{\frac{N}{2}}} \left| g_{\theta\theta}\left(\theta^*\right) \right|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(\theta - \theta^*\right)' g_{\theta\theta}\left(\theta^*\right) \left(\theta - \theta^*\right) \right\} d\theta = 1$$

- Then,

$$\int f\left(y|\theta\right) f\left(\theta\right) d\theta \simeq \int f\left(y|\theta^*\right) f\left(\theta^*\right) \exp \left\{ -\frac{1}{2} \left(\theta - \theta^*\right)' g_{\theta\theta}\left(\theta^*\right) \left(\theta - \theta^*\right) \right\} d\theta$$

$$= \frac{f\left(y|\theta^*\right) f\left(\theta^*\right)}{\frac{1}{(2\pi)^{\frac{N}{2}}} \left| g_{\theta\theta}\left(\theta^*\right) \right|^{\frac{1}{2}}} \int \frac{1}{(2\pi)^{\frac{N}{2}}} \left| g_{\theta\theta}\left(\theta^*\right) \right|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(\theta - \theta^*\right)' g_{\theta\theta}\left(\theta^*\right) \left(\theta - \theta^*\right) \right\} d\theta$$

$$= \frac{f\left(y|\theta^*\right) f\left(\theta^*\right)}{\frac{1}{(2\pi)^{\frac{N}{2}}} \left| g_{\theta\theta}\left(\theta^*\right) \right|^{\frac{1}{2}}}.$$

**Laplace Approximation to Posterior Distribution ...**

- Formula for marginal likelihood based on Laplace approximation:

$$f\left(y\right) = \int f\left(y|\theta\right) f\left(\theta\right) d\theta \simeq \left(2\pi\right)^{\frac{N}{2}} \frac{f\left(y|\theta^*\right) f\left(\theta^*\right)}{\left|g_{\theta\theta}\left(\theta^*\right)\right|^{\frac{1}{2}}}.$$

- Suppose $f(y|Model\ 1) > f(y|Model\ 2)$. Then, posterior odds on Model 1 higher than Model 2.

- 'Model 1 fits better than Model 2'

- Can use this to compare across two different models, or to evaluate contribution to fit of various model features: habit persistence, adjustment costs, etc.

# Generalized Method of Moments

- Express your econometric estimator into Hansen's GMM framework and you get standard errors

  - Essentially, *any* estimation strategy fits (see Hamilton)

- Works when parameters of interest, $\beta$, have the following property:

$$E \underbrace{u_t}_{N \times 1} \left( \underbrace{\beta}_{n \times 1} \right) = 0, \ \beta \text{ true value of some parameter(s) of interest}$$

$$u_t\left(\beta\right) \sim \text{stationary stochastic process (and other conditions)}$$

  - $n = N$ : 'exactly identified'

  - $n < N$ : 'over identified'

**Generalized Method of Moments ...**

– Example 1: mean

$$\beta = Ex_t,$$

$$u_t(\beta) = \beta - x_t.$$

– Example 2: mean and variance

$$\beta = \begin{bmatrix} \mu & \sigma \end{bmatrix},$$

$$Ex_t = \mu, E(x_t - \mu)^2 = \sigma^2.$$

then,

$$u_t(\beta) = \begin{bmatrix} \mu - x_t \\ (x_t - \mu)^2 - \sigma^2 \end{bmatrix}.$$

**Generalized Method of Moments  ...**

– Example 3: mean, variance, correlation, relative standard deviation

$$\beta = \begin{bmatrix} \mu_y & \sigma_y & \mu_x & \sigma_x & \rho_{xy} & \lambda \end{bmatrix}, \ \lambda \equiv \sigma_x/\sigma_y,$$

where

$$Ey_t = \mu_y, \ E\left(y_t - \mu_y\right)^2 = \sigma_y^2$$

$$Ex_t = \mu_x, \ E\left(x_t - \mu_x\right)^2 = \sigma_x^2$$

$$\rho_{xy} = \frac{E\left(y_t - \mu_y\right)\left(x_t - \mu_x\right)}{\sigma_y\sigma_x}.$$

then

$$u_t\left(\beta\right) = \begin{bmatrix} \mu_x - x_t \\ \left(x_t - \mu_x\right)^2 - \sigma_x^2 \\ \mu_y - y_t \\ \left(y_t - \mu_y\right)^2 - \sigma_y^2 \\ \sigma_y\sigma_x\rho_{xy} - \left(y_t - \mu_y\right)\left(x_t - \mu_x\right) \\ \sigma_y\lambda - \sigma_x \end{bmatrix}.$$

**Generalized Method of Moments ...**

– Example 4: New Keynesian Phillips curve

$$\pi_t = 0.99 E_t \pi_{t+1} + \gamma s_t,$$

or,

$$\pi_t - 0.99\pi_{t+1} - \gamma s_t = \eta_{t+1}$$

where,

$$\eta_{t+1} = 0.99\left(E_t \pi_{t+1} - \pi_{t+1}\right) \Longrightarrow E_t \eta_{t+1} = 0$$

Under Rational Expectations : $\eta_{t+1} \perp$ time $t$ information, $z_t$

$$u_t\left(\gamma\right) = \left[\pi_t - 0.99\pi_{t+1} - \gamma s_t\right] z_t$$

**Generalized Method of Moments ...**

• Inference about $\beta$

  – Estimator of $\beta$ in exactly identified case $(n = N)$

    ∗ Choose $\hat{\beta}$ to mimick population property of true $\beta$,

$$Eu_t(\beta) = 0.$$

    ∗ Define:

$$g_T(\beta) = \frac{1}{T} \sum_{t=1}^{T} u_t(\beta).$$

    ∗ Solve

$$\hat{\beta} : \ g_T \left( \underbrace{\hat{\beta}}_{N \times 1} \right) = \underbrace{0}_{N \times 1}.$$

**Generalized Method of Moments ...**

– Example 1: mean

$$\beta \;=\; Ex_t,$$

$$u_t\left(\beta\right) \;=\; \beta - x_t.$$

Choose $\hat{\beta}$ so that

$$g_T\left(\hat{\beta}\right) = \frac{1}{T}\sum_{t=1}^{T} u_t\left(\hat{\beta}\right) = \hat{\beta} - \frac{1}{T}\sum_{t=1}^{T} x_t = 0$$

and $\hat{\beta}$ is simply sample mean.

**Generalized Method of Moments ...**

– Example 4 in exactly identified case

$$Eu_t\left(\gamma\right) = E\left[\pi_t - 0.99\pi_{t+1} - \gamma s_t\right]z_t, \ z_t \sim \text{scalar}$$

choose $\hat{\gamma}$ so that

$$g_T\left(\hat{\beta}\right) = \frac{1}{T}\sum_{t=1}^{T}\left[\pi_t - 0.99\pi_{t+1} - \hat{\gamma}s_t\right]z_t = 0,$$

or. standard instrumental variables estimator:

$$\hat{\gamma} = \frac{\frac{1}{T}\sum_{t=1}^{T}\left[\pi_t - 0.99\pi_{t+1}\right]z_t}{\frac{1}{T}\sum_{t=1}^{T}s_t z_t}$$

**Generalized Method of Moments ...**

– Key message:

* In exactly identified case, GMM does not deliver a new estimator you would not have thought of on your own

· means, correlations, regression coefficients, exactly identified IV estimation, maximum likelihood.

* GMM provides framework for deriving asymptotically valid formulas for estimating sampling uncertainty.

**Generalized Method of Moments ...**

– Estimating $\beta$ in overidentified case $(N > n)$

   ∗ Cannot exactly implement sample analog of $Eu_t(\beta) = 0$ :

$$g_T\left(\underbrace{\hat{\beta}}_{n\times 1}\right) = \underbrace{0}_{N\times 1}$$

   ∗ Instead, 'do the best you can':

$$\hat{\beta} = \arg\min_{\beta} g_T(\beta)' W_T g_T(\beta),$$

   where

   $W_T$ ~ is a positive definite weighting matrix.

   ∗ GMM works for any positive definite $W_T$, but is most efficient if $W_T$ is inverse of estimator of variance-covariance matrix of $g_T\left(\hat{\beta}\right)$ :

$$(W_T)^{-1} = E g_T\left(\hat{\beta}\right) g_T\left(\hat{\beta}\right)'.$$

**Generalized Method of Moments ...**

– This choice of weighting matrix very sensible:

∗ weight heavily those moment conditions (i.e., elements of $g_T\left(\hat{\beta}\right)$) that are precisely estimated

∗ pay less attention to the others.

**Generalized Method of Moments  ...**

    – Estimator of $W_T^{-1}$

        * Note:

$$Eg_T\left(\hat{\beta}\right)g_T\left(\hat{\beta}\right)'$$

$$= \frac{1}{T^2}E\left[u_1\left(\hat{\beta}\right)+u_2\left(\hat{\beta}\right)+...+u_T\left(\hat{\beta}\right)\right]\left[u_1\left(\hat{\beta}\right)+u_2\left(\hat{\beta}\right)+...+u_T\left(\hat{\beta}\right)\right]'$$

$$= \frac{1}{T}[\frac{T}{T}Eu_t\left(\hat{\beta}\right)u_t\left(\hat{\beta}\right)'+\frac{T-1}{T}Eu_t\left(\hat{\beta}\right)u_{t+1}\left(\hat{\beta}\right)'+...+\frac{1}{T}Eu_t\left(\hat{\beta}\right)u_{t+T-1}\left(\hat{\beta}\right)'$$

$$+\frac{T-1}{T}Eu_t\left(\hat{\beta}\right)u_{t-1}\left(\hat{\beta}\right)'+\frac{T-2}{T}Eu_t\left(\hat{\beta}\right)u_{t-2}\left(\hat{\beta}\right)'+..+\frac{1}{T}Eu_t\left(\hat{\beta}\right)u_{t-T+1}\left(\hat{\beta}\right)']$$

$$= \frac{1}{T}\left[C\left(0\right)+\sum_{r=1}^{T-1}\frac{T-r}{T}\left(C\left(r\right)+C\left(r\right)'\right)\right],$$

      where

$$C\left(r\right)=Eu_t\left(\hat{\beta}\right)u_{t-r}\left(\hat{\beta}\right)'$$

      * $W_T^{-1}$ is '$\frac{1}{T}\times$spectral density matrix at frequency zero, $S_0$, of $u_t\left(\hat{\beta}\right)$'

**Generalized Method of Moments ...**

– Conclude:

$$W_T^{-1} = Eg_T\left(\hat{\beta}\right)g_T\left(\hat{\beta}\right) = \frac{1}{T}\left[C\left(0\right) + \sum_{r=1}^{T-1}\frac{T-r}{T}\left(C\left(r\right) + C\left(r\right)'\right)\right] = \frac{S_0}{T}.$$

– $W_T^{-1}$ estimated by

$$\widehat{W_T^{-1}} = \frac{1}{T}\left[\hat{C}\left(0\right) + \sum_{r=1}^{T-1}\frac{T-r}{T}\left(\hat{C}\left(r\right) + \hat{C}\left(r\right)'\right)\right] = \frac{1}{T}\hat{S}_0,$$

imposing whatever restrictions are implied by the null hypothesis, i.e., (as in ex. 4)

$$C\left(r\right) = 0,\ r > R \text{ some } R.$$

– which is 'Newey-West estimator of spectral density at frequency zero'
  * Problem: need $\hat{\beta}$ to compute $W_T^{-1}$ and need $W_T^{-1}$ to compute $\hat{\beta}$!!

  · Solution - first compute $\hat{\beta}$ using $W_T = I$, then iterate...

**Generalized Method of Moments ...**

● Sampling Uncertainty in $\hat{\beta}$.

    – The exactly identified case

    – By the Mean Value Theorem, $g_T\left(\hat{\beta}\right)$ can be expressed as follows:

$$g_T\left(\hat{\beta}\right) = g_T\left(\beta_0\right) + D\left(\hat{\beta} - \beta_0\right),$$

    where $\beta_0$ is the true value of the parameters and

$$D = \frac{\partial g_T\left(\beta\right)}{\partial \beta'}|_{\beta=\beta^*}, \text{ some } \beta^* \text{ between } \beta_0 \text{ and } \hat{\beta}.$$

    – Since $g_T\left(\hat{\beta}\right) = 0$ and $g_T\left(\beta_0\right) \overset{a}{\sim} N\left(0, S_0/T\right)$, it follows:

$$\hat{\beta} - \beta_0 = -D^{-1}g_T\left(\beta_0\right),$$

    so

$$\hat{\beta} - \beta_0 \overset{a}{\sim} N\left(0, \frac{\left(D'S_0^{-1}D\right)^{-1}}{T}\right)$$

52

**Generalized Method of Moments ...**

– The overidentified case.

  ∗ An extension of the ideas we have already discussed.

  ∗ Can derive the results for yourself, using the 'delta function method' for deriving the sampling distribution of statistics.

  ∗ Hamilton's text book has a great review of GMM.