# Likelihood Theory

First: review likelihood theory for conditional and full maximum likelihood estimation.

Suppose data is $X = (Y, Z)$; write density of $X$ as

$$f_X(x|\theta) = f_{Y|Z}(y|z, \theta) f_Z(z|\theta)$$

Notation. Full log-likelihood is

$$\ell_X(\theta) = \log f_X(X|\theta).$$

Conditional log-likelihood is

$$\ell_{Y|X}(\theta) = \log f_{Y|Z}(Y|Z, \theta).$$

Marginal log-likelihood is

$$\ell_Z(\theta) = \log f_Z(Z|\theta).$$

Example: fitting AR($p$). Have

$$X = (X_0, \ldots, X_{T-1})$$
$$Y = (X_p \ldots, X_{T-1})$$
$$Z = (X_0, \ldots, X_{p-1})$$

To find full, marginal or conditional mle study roots of corresponding score function:

Example: conditional score is

$$U_{Y|Z}(\theta) = \nabla \ell_{Y|Z}(\theta)$$

where $\nabla$ denotes gradient vector with $j$th component

$$U_{Y|Z;j}(\theta) = \frac{\partial \log f(y|z, \theta)}{\partial \theta_j}$$

Usually, conditional mle $\widehat{\theta}_c$ solves

$$U_{Y|Z}(\widehat{\theta}_c) = 0$$

Hessian: matrix of second derivatives, entries

$$\frac{\partial^2 \log f(y|z, \theta)}{\partial \theta_j \partial \theta_k}$$

Should be negative definite at maximum.

General considerations of likelihood theory:

- Expected value of score function is 0 at true parameter value.

- Expected second derivative of log-likelihood is negative definite and grows with sample size (usually linearly)

- Typically expected second derivative large compared to standard deviation of score.

Evidence for above from general theory:

Drop subscripts on densities for convenience.

Differentiate the identity

$$1 = \int f(y|z, \theta) dy$$

with respect to $\theta_j$ (the $j$th component of $\theta$) and pull the derivative under the integral sign to get

$$
\begin{aligned}
0 &= \int \frac{\partial f(y|z, \theta)}{\partial \theta_j} dy \\
&= \int \frac{\partial \log f(y|z, \theta)}{\partial \theta_j} f(y|z, \theta) dy \\
&= \mathsf{E}_\theta(U_{Y|Z;j}(\theta)|Z)
\end{aligned}
$$

Conclude

$$\mathsf{E}_\theta \left[ U_{Y|Z}(\theta)|Z \right] = 0.$$

Conclusion (1) follows since

$$\mathsf{E}_\theta(U_{Y|Z;j}(\theta)|Z) = 0$$

take expected values to see that

$$\mathsf{E}_\theta(U_{Y|Z;j}(\theta)) = 0$$

Also: other two scores $U_X(\theta)$ and $U_Z(\theta)$ have mean 0 (when $\theta$ is the true value of $\theta$).

Differentiate identity again wrt $\theta_k$ to get

$$0 = \int \frac{\partial^2 \log f(y|z,\theta)}{\partial\theta_j \partial\theta_k} f(y|z,\theta)dy$$
$$+ \int \frac{\partial \log f(y|z,\theta)}{\partial\theta_j} \frac{\partial \log f(y|z,\theta)}{\partial\theta_k} f(y|z,\theta)dy$$

We define the conditional Fisher information matrix $I_{Y|Z}(\theta)$ to have $jk$th entry

$$\mathsf{E}\left[-\frac{\partial^2 \ell}{\partial\theta_j \partial\theta_k}|Z\right]$$

and get

$$I_{Y|Z}(\theta|Z) = \mathsf{Var}_\theta(U_{Y|Z}(\theta)|Z)$$

Corresponding identities based on $f_X$ and $f_Z$:

$$I_X(\theta) = \mathsf{Var}_\theta(U_X(\theta))$$

and

$$I_Z(\theta) = \mathsf{Var}_\theta(U_Z(\theta))$$

For first bit of (2) note Variances are non-negative definite.

Evidence for other assertions in special case:

Model $X_t = \rho X_{t-1} + \epsilon_t$. For $Y = (X_1, \ldots, X_{T-1})$ and $Z = X_0$ we find

$$U_{Y|Z}(\rho, \sigma) = \begin{bmatrix} \dfrac{\sum_1^{T-1}(X_t - \rho X_{t-1})X_{t-1}}{\sigma^2} \\ \dfrac{\sum_1^{T-1}(X_t - \rho X_{t-1})^2}{\sigma^3} - \dfrac{T-1}{\sigma} \end{bmatrix}$$

Differentiating again gives the matrix of second derivatives

$$\begin{bmatrix} -\dfrac{\sum_1^{T-1} X_{t-1}^2}{\sigma^2} & -2\dfrac{\sum_1^{T-1}(X_t - \rho X_{t-1})X_{t-1}}{\sigma^3} \\ -2\dfrac{\sum_1^{T-1}(X_t - \rho X_{t-1})X_{t-1}}{\sigma^3} & -3\dfrac{\sum_1^{T-1}(X_t - \rho X_{t-1})^2}{\sigma^4} + \dfrac{T-1}{\sigma^2} \end{bmatrix}$$

Taking conditional expectations given $X_0$ gives

$$I_{Y|Z}(\rho, \sigma) = \begin{bmatrix} \dfrac{\sum_1^{T-1} \mathsf{E}[X_{t-1}^2|X_0]}{\sigma^2} & 0 \\ 0 & \dfrac{2(T-1)}{\sigma^2} \end{bmatrix}$$

To compute $W_k \equiv \mathsf{E}[X_k^2|X_0]$ write

$$X_k = \rho X_{k-1} + \epsilon_k$$

and get

$$W_k = \rho^2 W_{k-1} + \sigma^2$$

with $W_0 = X_0^2$.

Can check that

$$W_k = \rho^{2k} W_0 + \frac{\sigma^2(1 - \rho^{2k})}{1 - \rho^2}$$

so

$$W_k \to W_\infty = \frac{\sigma^2}{1 - \rho^2}$$

as $k \to \infty$.

It follows that

$$\frac{1}{T} I_{Y|Z}(\rho, \sigma) \rightarrow \begin{bmatrix} \frac{1}{1-\rho^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{bmatrix}$$

Notice although conditional Fisher information might have been expected to depend on $X_0$ it does not, at least for long series.

This is evidence for (2) and (3). Notice standard deviation of score is like square root of Fisher information, so small relative to FI.

# Large Sample Theory: Conditional Likelihood

Data $X = (Y, Z)$. Cond'l likelihood, score, Fisher information and mle: $\ell_{Y|Z}(\theta)$, $U_{Y|Z}(\theta)$, $\mathcal{I}_{Y|Z}(\theta)$ and $\widehat{\theta}$. In general standard maximum likelihood theory expected to apply to these conditional objects:

1) $P(\ell_{Y|Z}(\theta_0) > \ell_{Y|Z}(\theta)) \to 1$ as the "sample size" (often measured by the Fisher information) tends to infinity.

2) Can strengthen 1) to get $\widehat{\theta}$ is consistent (converges to true value as Fisher information converges to infinity).

3) Usual Bartlett identities hold:

$$\mathsf{E}_{\theta} \left[ U_{Y|Z}(\theta) | Z \right] = 0$$

$$\mathcal{I}_{Y|Z}(\theta) \equiv \mathsf{Var} \left[ U_{Y|Z}(\theta) | Z \right]$$

$$= -\mathsf{E}_{\theta} \left[ \frac{\partial}{\partial \theta} U_{Y|Z}(\theta) | Z \right]$$

4) Score function is asymptotically normal:

$$\left\{\mathcal{I}_{Y|Z}(\theta)\right\}^{-1/2} U_{Y|Z}(\theta) \approx MVN(0, I)$$

5) Error in mle has approximate form

$$\widehat{\theta} - \theta \approx \left(\mathcal{I}_{Y|Z}(\theta)\right)^{-1} U_{Y|Z}(\theta)$$

6) The mle is approximately normal:

$$\left(\mathcal{I}_{Y|Z}(\theta)\right)^{1/2} \left(\widehat{\theta} - \theta\right) \approx MVN(0, I)$$

(where $I$ is the identity matrix).

7) The conditional Fisher information can be estimated by the observed information:

$$\left(\mathcal{I}_{Y|Z}(\theta)\right)^{-1} \left(-\frac{\partial}{\partial \theta} U_{Y|Z}(\widehat{\theta})\right) \to I$$

8) The log-likelihood ratio is approximately $\chi^2$:

$$2(\ell_{Y|Z}(\widehat{\theta}) - \ell_{Y|Z}(\theta_0)) \Rightarrow \chi_p^2$$

What if likelihood wrong?

# Non Gaussian series.

The fitting methods we have studied are based on the likelihood for a normal fit. However, the estimates work reasonably well even if the errors are not normal.

Example: AR(1) fit. We fit $X_t - \mu = \rho(X_{t-1} - \mu) + \epsilon_t$ using $\hat{\mu} = \bar{X}$ which is consistent for non-Gaussian errors. (In fact

$$(1 - \rho) \sum_0^{T-1} X_t + \rho X_{T-1} - X_0$$

$$= (T - 1)(1 - \rho)\mu + \sum_0^{T-1} \epsilon_t - \epsilon_0;$$

divide by $T$ and apply the law of large numbers to $\bar{\epsilon}$ to see that $\bar{X}$ is consistent.)

Outline of logic which follows. Assume: errors are iid mean 0, variance $\sigma^2$ and finite fourth moment $\mu_4 = E(\epsilon_t^4)$.

**Do not** assume errors have normal distribution.

1) The estimates of $\rho$ and $\sigma$ are consistent.

2) The score function satisfies

$$T^{-1/2}U(\theta_0) \Rightarrow MVN(0, B)$$

where

$$B = \begin{bmatrix} \frac{1}{1-\rho^2} & 0 \\ 0 & \frac{\mu_4 - \sigma^4}{\sigma^6} \end{bmatrix}$$

3) The matrix of second derivatives satisfies

$$\lim_{T \to \infty} -\frac{1}{T}\frac{\partial U}{\partial \theta} = \lim_{T \to \infty} -\frac{1}{T}\mathsf{E}\left(\frac{\partial U}{\partial \theta}\right) = A$$

where

$$A = \begin{bmatrix} \frac{1}{1-\rho^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{bmatrix}$$

4) If $\mathcal{I}$ is the (conditional) Fisher information then

$$\lim \lim_{T \to \infty} \frac{1}{T}\mathcal{I} = A$$

5) We can expand $U(\hat{\theta})$ about $\theta_0$ and get

$$T^{1/2}(\hat{\theta} - \theta_0) = \left[\frac{1}{T}I(\theta_0)\right]^{-1}\left[T^{-1/2}U(\theta_0)\right]$$
$$+ \text{negligible remainder}$$

6) So

$$T^{1/2}(\hat{\theta} - \theta_0)$$
$$\approx MVN(0, A^{-1}BA^{-1}) = MVN(0, \Sigma)$$

where

$$\Sigma = A^{-1}BA^{-1} = \begin{bmatrix} 1 - \rho^2 & 0 \\ 0 & \frac{\mu_4 - \sigma^4}{4\sigma^2} \end{bmatrix}$$

7) So $T^{1/2}(\hat{\rho} - \rho) \Rightarrow N(0, 1 - \rho^2)$ even for non-normal errors.

8) On the other hand the estimate of $\sigma$ has a limiting distribution which will be different for non-normal errors (because it depends on $\mu_4$ which is $3\sigma^4$ for normal errors and something else in general for non-normal errors).

Here are details.

**Consistency**: One of our many nearly equivalent estimates of $\rho$ is

$$\widehat{\rho} = \frac{\sum (X_t - \bar{X})(X_{t-1} - \bar{X})}{\sum (X_t - \bar{X})^2}$$

Divide both top and bottom by $T$. You need essentially to prove

$$T^{-1} \sum (X_t - \mu)(X_{t-1} - \mu) \to C(1)$$

and

$$T^{-1} \sum (X_t - \mu)^2 \to C(0)$$

Each of these is correct and hinges on the fact that these linear processes are ergodic — long time averages converge to expected values. For these particular averages it is possible to compute means and variances and prove that the mean squared error converges to 0.

# Score function: asymptotic normality

The score function is

$$U(\rho, \sigma) = \begin{bmatrix} \frac{\sum X_{t-1}(X_t - \rho X_{t-1})}{\sigma^2} \\ \frac{\sum(X_t - \rho X_{t-1})^2}{\sigma^3} - \frac{T-1}{\sigma} \end{bmatrix}$$

If $\rho$ and $\sigma$ are the true values of the parameters then

$$U(\rho, \sigma) = \begin{bmatrix} \frac{\sum X_{t-1}\epsilon_t}{\sigma^2} \\ \frac{\sum \epsilon_t^2}{\sigma^3} - \frac{T-1}{\sigma} \end{bmatrix}$$

Claim that $T^{-1/2}U(\rho, \sigma) \Rightarrow MVN(0, B)$.

Proof: martingale central limit theorem.

Technically fix an $a \in R^2$, study $T^{-1/2}a^t U(\rho, \sigma)$.

Prove that limit is $N(0, a^t B a)$.

Here do only special cases $a = (1, 0)^t$ and $a = (0, 1)^t$.

The second of these is simply

$$T^{-1/2} \sum (\epsilon_i^2 - \sigma^2)/\sigma^3$$

which converges by usual CLT to $N(0, (\mu_4 - \sigma^4)/\sigma^6)$. For $a = (1,0)^t$ the claim is that

$$T^{-1/2} \sum X_{t-1} \epsilon_t \Rightarrow N(0, C(0)\sigma^2)$$

because $C(0) = \sigma^2/(1 - \rho^2)$.

To prove this assertion we define for each $T$ a martingale $M_{T,k}$ for $k = 1, \ldots, T$ where

$$M_{T,k} = \sum_1^k D_{T,i}$$

with

$$D_{T,i} = T^{-1/2} X_{i-1} \epsilon_i$$

The martingale property is that

$$\mathsf{E}(M_{T,k+1} | \epsilon_k, \epsilon_{k-1}, \ldots) = M_{T,k}$$

Martingale central limit theorem (Hall, P. and Heyde, C. C. (1980). *Martingale limit theory and its application.* New York: Academic Press.):

$$M_{T,T} \Rightarrow N(0, b)$$

provided that

$$\sum_k D_{T,k}^2 \to b$$

and provided that an analogue of Lindeberg's condition holds.

Here I check only the former condition:

$$\sum_k D_{T,k}^2 = \frac{1}{T} \sum X_{t-1}^2 \epsilon_t^2 \to \mathsf{E}(X_0^2 \epsilon_1^2) = C(0)\sigma^2$$

(by the ergodic theorem or you could compute means and variances).

**Second derivative matrix and Fisher information**: matrix of negative second derivatives is

$$-\frac{\partial U}{\partial \theta} = \begin{bmatrix} \dfrac{\sum X_{t-1}^2}{\sigma^2} & 2\dfrac{\sum X_{t-1}(X_t - \rho X_{t-1})}{\sigma^3} \\ 2\dfrac{\sum X_{t-1}(X_t - \rho X_{t-1})}{\sigma^3} & 3\dfrac{\sum(X_t - \rho X_{t-1})^2}{\sigma^4} - \dfrac{T-1}{\sigma^2} \end{bmatrix}.$$

If you evaluate at the true parameter value and divide by $T$ the matrix and the expected value of the matrix converge to

$$A = \begin{bmatrix} \dfrac{C(0)}{\sigma^2} & 0 \\ 0 & \dfrac{2}{\sigma^2} \end{bmatrix}$$

(Again this uses the ergodic theorem or a variance calculation.)

**Taylor expansion**: next step — supposed to prove that a random vector has a MVN limit.

Usual tactic uses *Cramér-Wold device*: prove that each linear combination of entries in vector has a univariate normal limit.

Then $U(\widehat{\rho}, \widehat{\sigma}) = 0$ and Taylor's theorem is that

$$0 = U(\widehat{\rho}, \widehat{\sigma}) = U(\rho, \sigma) + \left[\frac{\partial U(\theta)}{\partial \theta}\right](\widehat{\theta} - \theta) + R$$

(Using $\theta^t = (\rho, \sigma)$; $R$ is remainder term — random variable with property that

$$P(\|R\|/\|U(\theta)\| > \eta) \to 0$$

for each $\eta > 0$.) Multiply through by

$$\left[\frac{\partial U(\theta)}{\partial \theta}\right]^{-1}$$

and get

$$T^{1/2}(\widehat{\theta} - \theta) = \left[-T^{-1}\frac{\partial U(\theta)}{\partial \theta}\right]^{-1} \left\{T^{-1/2}U(\rho, \sigma) + T^{-1/2}R\right\}.$$

It is possible with care to prove that

$$\left[-T^{-1}\frac{\partial U(\theta)}{\partial \theta}\right]^{-1}(T^{-1/2}R) \to 0$$

**Asymptotic normality**: consequence of Slutsky's theorem applied to Taylor expansion and results above for $U$ and $I$.

Slutsky's theorem: asymptotic distribution of $T^{1/2}(\hat{\theta} - \theta)$ same as that of

$$A^{-1}(T^{-1/2}U(\rho, \sigma))$$

which converges in distribution to

$$MVN(0, A^{-1}B(A^{-1})^t).$$

Now since $C(0) = \sigma^2/(1 - \rho^2)$

$$A^{-1}B(A^{-1})^t = \begin{bmatrix} 1 - \rho^2 & 0 \\ 0 & \frac{\mu_4 - \sigma^4}{4\sigma^4} \end{bmatrix}$$

**Behaviour of $\widehat{\rho}$:** pick off first component:

$$T^{1/2}(\widehat{\rho} - \rho) \Rightarrow N(0, 1 - \rho^2)$$

Notice answer same for normal and non-normal errors.

**Behaviour of $\widehat{\sigma}$:** on the other hand

$$T^{1/2}(\widehat{\sigma} - \sigma) \Rightarrow N(0, (\mu_4 - \sigma^4)/(4\sigma^2))$$

which has $\mu_4$ in it and will match the normal theory limit if and only if $\mu_4 = 3\sigma^4$.

**More general models:** For an ARMA$(p, q)$ model the parameter vector is

$$\theta = (a_1, \ldots, a_p, b_1, \ldots, b_q, \sigma)^t .$$

In general the matrices $B$ and $A$ are of the form

$$B = \left[ \begin{array}{cc} B_1 & 0 \\ 0 & \frac{\mu_4 - \sigma^4}{\sigma^6} \end{array} \right]$$

and

$$A = \left[ \begin{array}{cc} A_1 & 0 \\ 0 & \frac{2}{\sigma^2} \end{array} \right]$$

where $A_1 = B_1$ and $A_1$ is a function of the parameters $a_1, \ldots, a_p, b_1, \ldots, b_q$ only and is the same for both normal and non-normal data.