# LECTURE III.  A NORMAL APPROXIMATION THEOREM

The basic Lemmas I.3 and I.4 are elaborated to obtain a normal approxima-
tion theorem.  This is applied to the special case of a sum of independent
random variables, the study of which was begun in the first lecture, and also
to the study of the sum of a random diagonal.  The substantial shortcomings of
these results, which will be described later in this lecture, will be overcome
to some extent in later lectures.

Lemma 1:  Let $(W,W')$ be an exchangeable pair of real random variables
such that

(1)
$$E^W W' = (1-\lambda)W$$

with

(2)
$$0 < \lambda < 1$$

and let $h: R \rightarrow R$ be a bounded continuous function with bounded piecewise
continuous derivative $h'$.  Then, with $U_N h$ defined by (II.4), that is

(3)
$$(U_N h)(w) = e^{\frac{1}{2}w^2} \int_{-\infty}^{w} [h(x)-Nh]e^{-\frac{1}{2}x^2} dx$$

$$= -e^{\frac{1}{2}w^2} \int_{w}^{\infty} [h(x)-Nh]e^{-\frac{1}{2}x^2} dx,$$

where $Nh$ is defined by (II.2), that is

(4)
$$Nh = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} h(x)e^{-\frac{1}{2}x^2} dx,$$

we have

(5)     $Eh(W) = Nh + E(U_N h)'(W)[1 - \frac{1}{2\lambda} E^W(W'-W)^2]$

$+ \frac{1}{2\lambda} \int E(W'-W)(z - \frac{W+W'}{2})[\mathcal{I}\{z \leq W'\} - \mathcal{I}\{z \leq W\}]d(U_N h)'(z).$

Proof: We have to express the remainder in (I.45), where f was written instead of $U_N h$, as the sum of the last two terms on the right hand side of (5). In fact

(6)     $E[(U_N h)'(W) - \frac{1}{2\lambda}(W'-W)((U_N h)(W') - (U_N h)(W))]$

$= E(U_N h)'(W)[1 - \frac{1}{2\lambda} E^W(W'-W)^2]$

$- \frac{1}{2\lambda} E(W'-W)[(U_N h)(W')-(U_N h)(W)-(W'-W)(U_N h)'(W)].$

But, for $W < W'$,

(7)     $(U_N h)(W') - (U_N h)(W)-(W'-W)(U_N h)'(W)$

$= \int_W^{W'} [(U_N h)'(z) - (U_N h)'(W)]dz$

$= \int_W^{W'} (\int_W^z (U_N h)''(y)dy)dz = \int_W^{W'} (W'-y)(U_N h)''(y)dy,$

and, for $W' < W$,

(8)     $(U_N h)(W') - (U_N h)(W) - (W'-w)(U_N h)'(W)$

$= - \int_{W'}^W [(U_N h)'(z) - (U_N h)'(W)]dz$

$= \int_{W'}^W (\int_z^W (U_N h)''(y)dy)dz = \int_{W'}^W (y-W')(U_N h)''(y)dy.$

Now (7) and (8) can be combined to obtain, for all $W, W'$,

(9)     $(U_N h)(W')-(U_N h)(W)-(W'-W)(U_N h)'(W)$

$= \int(W'-y)[\mathcal{I}\{y \leq W'\} - \mathcal{I}\{y \leq W\}](U_N h)''(y)dy.$

Taking expectation, applying the exchangeability of W and W' again, and using (6) and (I.45), we obtain (5).

Theorem 1: Under the assumptions of Lemma 1,

(10) $\quad |Eh(W)-Nh|$

$$\leq 2 \sup|h-Nh|\sqrt{E[1-E^W \frac{1}{2\lambda} (W'-W)^2]^2} + \frac{1}{4\lambda} \sup|h'|E|W'-W|^3$$

and, for all real $w_0$,

(11) $\quad |P\{W \leq w_0\} - \Phi(w_0)|$

$$\leq 2\sqrt{E[1 - \frac{1}{2\lambda} E^W(W'-W)^2]^2} + (2\pi)^{-\frac{1}{4}} \sqrt{\frac{1}{\lambda} E|W'-W|^3}.$$

<u>Proof</u>: We recall that by Lemma II.3,

(12) $$\sup|(U_N h)'| \leq 2 \sup|h-Nh|$$

and

(13) $$\sup|(U_N h)''| \leq 2 \sup|h'|.$$

Thus from Lemma 1 we obtain

(14) $\quad |Eh(W)-Nh| = |E(U_N h)'(W)[1 - \frac{1}{2\lambda} E^W(W'-W)^2]$

$$+ \frac{1}{2\lambda} \int E(W'-W)(z - \frac{W+W'}{2})[\mathcal{I}\{z \leq W'\} - \mathcal{I}\{z \leq W\}](U_N h)''(z)dz|$$

$$\leq \sup|(U_N h)'|E|1 - \frac{1}{2\lambda} E^W(W'-W)^2|$$

$$+ \sup|(U_N h)''|\frac{1}{2\lambda} E\int_{W \wedge W'}^{W \vee W'} |W'-W||z - \frac{W+W'}{2}|dz$$

$$\leq 2 \sup|h-Nh|\sqrt{E[1 - \frac{1}{2\lambda} E^W(W'-W)^2]^2}$$

$$+ 2 \sup|h'| \cdot \frac{1}{2\lambda} E \frac{|W'-W|^3}{4},$$

which is (10). Here I have written

(15) $$a \wedge b = \min(a,b)$$

and

(16) $$a \vee b = \max(a,b).$$

In order the prove (11), we shall apply (10) with $h = h_{w_0,\Delta}$ defined for $\Delta > 0$ by

(17)
$$h_{w_0,\Delta}(w) = \begin{cases} 1 & \text{if } w \leq w_0 \\ 1 - \dfrac{w-w_0}{\Delta} & \text{if } w_0 \leq w \leq w_0 + \Delta \\ 0 & \text{if } w \geq w_0 + \Delta. \end{cases}$$

Then

(18)
$$Eh_{w_0-\Delta,\Delta}(W) \leq P\{W \leq w_0\} \leq Eh_{w_0,\Delta}(W),$$

and

(19)
$$\sup|h_{w_0,\Delta} - Nh_{w_0,\Delta}| \leq 1$$

and

(20)
$$\sup|h'_{w_0,\Delta}| \leq \frac{1}{\Delta}.$$

Thus, using (10) we obtain

(21)
$$P\{W \leq w_0\} \leq Eh_{w_0,\Delta}(W)$$

$$\leq Nh_{w_0,\Delta} + 2\sqrt{E[1 - \frac{1}{2\lambda} E^W(W'-W)^2]^2} + \frac{1}{4\lambda\Delta} E|W'-W|^3$$

$$\leq \Phi(w_0) + \frac{\Delta}{\sqrt{2\pi}} + 2\sqrt{E[1 - \frac{1}{2\lambda} E^W(W'-W)^2)]^2} + \frac{1}{4\lambda\Delta} E|W'-W|^3.$$

The upper bound on the right hand side is minimized by

(22)
$$\Delta = \frac{(2\pi)^{\frac{1}{4}}}{2} \sqrt{\frac{1}{\lambda} E|W'-W|^3}$$

yielding

(23)
$$P\{W \leq w_0\} \leq \Phi(w_0) + 2\sqrt{E[1 - \frac{1}{2\lambda} E^W(W'-W)^2]^2}$$

$$+ (2\pi)^{-\frac{1}{4}} \sqrt{\frac{1}{\lambda} E|W'-W|^3}.$$

The bound for the error in (11) follows from (23) and the corresponding lower bound.

Now let us return to the special case of a sum of independent random variables. We need only evaluate the expectations occurring in the bounds on the right hand side of (10) and (11) in the situation discussed at the end of Lecture I starting before (I.51). I shall state the result as

Corollary 1:   Let $X_1,\ldots,X_n$ be independent real random variables with

(24) $$EX_i = 0, \quad EX_i^2 = \sigma_i^2$$

and

(25) $$\Sigma \sigma_i^2 = 1.$$

Then, for all real $w_0$,

(26) $$|P\{W \leq w_0\} - \Phi(w_0)| \leq \sqrt{\Sigma(EX_i^4 - \sigma_i^4)} + (2\pi)^{-\frac{1}{4}} \sqrt{2\Sigma(E|X_i|^3 + 3\sigma_i^3)}$$

and for any bounded continuous functions h:  R → R with bounded piecewise continuous derivative h',

(27) $$|Eh(W) - Nh| \leq \sup|h - Nh| \sqrt{\Sigma(EX_i^4 - \sigma_i^4)} + \frac{1}{2} \sup|h'| \Sigma(E|X_i|^3 + 3\sigma_i^3).$$

Proof:   We shall need the fact that, if $X_i^*$ is independent of $X_i$ and has the same distribution as $X_i$ then

(28) $$E^X(X_i^* - X_i)^2 = \sigma_i^2 + X_i^2$$

and

(29) $$E|X_i^* - X_i|^3 \leq 2(E|X_i|^3 + 3(EX_i^2)E|X_i|) \leq 2(E|X_i|^3 + 3\sigma_i^3).$$

Of course (29) could be improved substantially.  Then the expectations on the right hand side of (10) and (11) can be bounded by using

(30) $$E[1 - \frac{n}{2} E^W(W'-W)^2]^2 = E[1 - \frac{n}{2} E^W(X_I^* - X_I)^2]^2$$

$$\leq E[1 - \frac{n}{2} E^X(X_I^* - X_I)^2]^2 = E[1 - \frac{1}{2} \Sigma E^X(X_i^* - X_i)^2]^2$$

$$= E[\frac{1}{2} \Sigma(X_i^2 - \sigma_i^2)]^2 = \frac{1}{4} \Sigma(EX_i^4 - \sigma_i^4)$$

and

(31) $$E|W'-W|^3 = E|X_I^* - X_I|^3 = \frac{1}{n} \Sigma E|X_i^* - X_i|^3 \leq \frac{2}{n} \Sigma(E|X_i|^3 + 3\sigma_i^3).$$

Substituting in (10) and (11) we obtain (26) and (27).

Let us look at this corollary in the special case of a sum of independent identically distributed random variables.  Let $Y_1, Y_2,\ldots$ be independent

identically distributed random variables with mean 0 and variance 1 and let

(32)
$$W = \frac{1}{\sqrt{n}} \Sigma \, Y_i .$$

Then, by applying Corollary 1 to the

(33)
$$X_i = \frac{1}{\sqrt{n}} Y_i$$

we see that, with

(34)
$$\beta_3 = E|Y_i|^3$$

and

(35)
$$\beta_4 = EY_i^4 ,$$

we have

(36)
$$|P\{W \le w_0\} - \Phi(w_0)| \le \sqrt{\frac{\beta_4 - 1}{n}} + (2\pi n)^{-\frac{1}{4}} \sqrt{2(\beta_3 + 3)}$$

and

(37)
$$|Eh(W) - Nh| \le \sup|h - Nh| \sqrt{\frac{\beta_4 - 1}{n}} + \frac{1}{2\sqrt{n}} \sup|h'|(\beta_3 + 3).$$

We see that the bound in (37) is of the right order of magnitude $n^{-\frac{1}{2}}$ provided $\beta_4 < \infty$, but the bound in (36) is too large, being of the order of $n^{-\frac{1}{4}}$. These faults will be repaired in later lectures.

Finally let us look briefly at the case of a sum of a random diagonal. Let $\{a_{ij}\}_{i,j \in \{1...n\}}$ be a square array of numbers such that, for every $i \in \{1,...,n\}$,

(38)
$$\sum_j a_{ij} = 0$$

and, for every $j \in \{1,...,n\}$,

(39)
$$\sum_i a_{ij} = 0,$$

and also

(40)
$$\sum_i \sum_j a_{ij}^2 = n-1.$$

Let $\Pi$ be a random permutation of $\{1,...,n\}$ such that, for every one to one

$\pi:\ \{1,\ldots,n\} \to \{1,\ldots,n\},$

(41) $$P\{\Pi=\pi\} = \frac{1}{n!}.$$

We shall be interested in the distribution of

(42) $$W = \Sigma\ a_{i\Pi(i)}.$$

In order to apply Theorem 1, we introduce an ordered pair (I,J) of random variables independent of $\Pi$ with

(43) $$P\{I=i\ \&\ J=j\} = \frac{1}{n(n-1)}$$

for every $i,j\ \epsilon\ \{1,\ldots,n\}$ with $i \neq j$, and define the random permutation $\Pi'$ by

(44) $$\Pi'(i) = \begin{cases} \Pi(i) & \text{if}\ \ i \notin \{I,J\} \\ \Pi(J) & \text{if}\ \ i = I \\ \Pi(I) & \text{if}\ \ i = J. \end{cases}$$

Finally we define

(45) $$W' = \Sigma\ a_{i\Pi'(i)} = W + \alpha_{I,J,\Pi(I),\Pi(J)}$$

where

(46) $$\alpha_{i,j,k,\ell} = a_{i\ell} + a_{jk} - a_{ik} - a_{j\ell}.$$

Clearly $(\Pi,\Pi')$ is an exchangeable pair since we can first choose the pair (I,J), then the $\pi(i)$ for $i \notin \{I,J\}$ and then, of the two remaining possibilities, choose one to be $\Pi$ and the other to be $\Pi'$ with equal probability. It follows that $(W,W')$ is also an exchangeable pair.

We must verify the condition (1) which is part of the hypotheses of Theorem 1 and then bound the quantities occurring on the right hand side of (10) and (11). By (45) and (46),

(47) $E^{\Pi}W' = W + E^{\Pi}[a_{I\Pi(J)} + a_{J\Pi(I)} - a_{I\Pi(I)} - a_{J\Pi(J)}]$

$\qquad = W + \dfrac{1}{n(n-1)} \sum\limits_{i} \sum\limits_{j \neq i} [a_{i\Pi(j)} + a_{j\Pi(i)} - a_{i\Pi(i)} - a_{j\Pi(j)}]$

$\qquad = W + \dfrac{1}{n(n-1)} [-\sum\limits_{i} a_{i\Pi(i)} - \sum\limits_{j} a_{j\Pi(j)} - (n-1)\sum\limits_{i} a_{i\Pi(i)} - (n-1)\sum\limits_{j} a_{j\Pi(j)}]$

$\qquad = (1 - \dfrac{2}{n-1})W.$

Thus (1) holds with

(48)
$$\lambda = \frac{2}{n-1} \ .$$

We shall bound the first term on the right-hand side of (10) or (11) by using the inequality

(49)
$$E[1 - \frac{1}{2\lambda} E^W(W'-W)^2]^2 \le E[1 - \frac{1}{2\lambda} E^\Pi(W'-W)^2]^2$$

which is true because $\Pi$ determines W.

In order to bound the right-hand side of (49) we first compute

(50)
$$E^\Pi(W'-W)^2 = \frac{1}{n(n-1)} \sum_{i,j} (a_{i\Pi(j)} + a_{j\Pi(i)} - a_{i\Pi(i)} - a_{j\Pi(j)})^2$$

$$= \frac{2}{n(n-1)} [\sum_{i,j} a^2_{i\Pi(j)} + \sum_{i,j} a^2_{i\Pi(i)} + \sum_{i,j} (a_{i\Pi(j)} a_{j\Pi(i)} + a_{i\Pi(i)} a_{j\Pi(j)})]$$

$$- \frac{4}{n(n-1)} [\sum_{i,j} a_{i\Pi(j)} a_{i\Pi(i)} + \sum_{i,j} a_{i\Pi(j)} a_{j\Pi(j)}]$$

$$= \frac{2}{n(n-1)} [(n-1) + (n+2)\sum_i a^2_{i\Pi(i)} + \sum_{i,j} (a_{i\Pi(j)} a_{j\Pi(i)} + a_{i\Pi(i)} a_{j\Pi(j)})]$$

where the sums over the pair of indices i,j are subject to the restriction $i \ne j$. The final equality uses (38), (39), and (40). From (50) we can compute

(51)
$$E(W'-W)^2 = EE^\Pi(W'-W)^2 = \frac{4}{n-1} = 2\lambda$$

since

(52)
$$E \sum_i a^2_{i\Pi(i)} = \frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{n} a^2_{ik} = \frac{n-1}{n}$$

and

(53)
$$E \sum_{i,j} a_{i\Pi(j)} a_{j\Pi(i)} = E \sum_{i,j} a_{i\Pi(i)} a_{j\Pi(j)}$$

$$= \frac{1}{n(n-1)} \sum_{i,j} \sum_{k,\ell} a_{ik} a_{j\ell} = - \frac{1}{n(n-1)} \sum_{i,j} \sum_k a_{ik} a_{jk}$$

$$= \frac{1}{n(n-1)} \sum_{i,k} a^2_{ik} = \frac{1}{n},$$

where the sum over the pair of indices $k,\ell$ is also subject to the restriction $k \ne \ell$. Substituting $f(w) = w$ in (I.38) and using (51) we obtain

(54)
$$EW^2 = 1,$$

which indicates that the scale has been chosen in a natural way.

From (50) it is not hard to see that for some absolute constants C and C'

(55)  $\text{Var } E^{\Pi}(W'-W)^2$

$$\leq C\left[\frac{1}{n^2} \sum_i \text{Var } a_{i\Pi(i)}^2 + \frac{1}{n^4} \sum_{i,j} \text{Var}(a_{i\Pi(j)}a_{j\Pi(i)} + a_{i\Pi(i)}a_{j\Pi(j)})\right]$$

$$\leq \frac{C'}{n^3} \sum_{i,j} a_{ij}^4$$

and consequently, by (48) and (49),

(56)                      $$E[1 - \frac{1}{2\lambda} E^W(W'-W)^2]^2 \leq \frac{C'}{n} \sum_{i,j} a_{ij}^4.$$

Similarly, for use in the final remainder term of (10) or (11) we observe that

(57)                      $$E|W'-W|^3$$

$$= E|a_{I\Pi(J)} + a_{J\Pi(I)} - a_{I\Pi(I)} - a_{J\Pi(J)}|^3$$

$$\leq \frac{C}{n^2} \sum_{i,j} |a_{ij}|^3.$$

Thus in the present case (10) and (11) are specialized to

(58)                      $$|Eh(W) - Nh|$$

$$\leq \frac{C'}{\sqrt{n}} \sup|h| \sqrt{\sum_{i,j} a_{ij}^4} + \frac{C'}{n} \sup|h'| \sum_{i,j} |a_{ij}|^3$$

and

(59)                      $$|P\{W \leq w_0\} - \Phi(w_0)|$$

$$\leq \frac{C'}{\sqrt{n}} \left[\sqrt{\sum_{i,j} a_{ij}^4} + \sqrt{\sum_{i,j} |a_{ij}|^3}\right].$$

In Lemma 1 an identity was derived that makes the conditions for approximate normality of the random variable W fairly clear intuitively, at least in a rough way. Provided that the scale has been chosen in such a way that

(60)                      $$2\lambda = E(W'-W)^2,$$

the first remainder term is small if $\frac{1}{2\lambda} E^W(W'-W)^2$ is nearly constant, and the second remainder term is small if $|W'-W|$ is small, or more precisely if $E|W'-W|^3/E[W'-W]^2$ is small. These rough statements were made precise in Theorem 1. This was specialized in Corollary 1 to sums of independent random

variables and then specialized further in (36) and (37) to the identically
distributed case.  Finally the case of a sum of a random diagonal was consider-
ed briefly.  A special case of this is the permutation distribution of a sample
covariance and, more particularly, the distribution of the mean of a sample
from a finite population.