

Supplemental Material to Talk-5

1 The normal consistency factor for the range, d_2

1.1 The joint pdf of order statistics

Theorem 1. Let X_1, X_2, \dots, X_n be a random sample with continuous cdf $F(x)$ and pdf $f(x)$. Let $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ be the order statistics of a random sample. Then the joint pdf of $X_{(i)}$ and $X_{(j)}$ for $1 \leq i < j \leq n$ is given by

$$f_{(i,j)}(u, v) = \frac{n!}{(i-1)! \times 1! \times (j-i-1)! \times 1! \times (n-j)!} \times \\ \left[F(u) \right]^{i-1} f(u) \left[F(v) - F(u) \right]^{j-i-1} f(v) \left[1 - F(v) \right]^{n-j}$$

for $-\infty < u < v < \infty$.

Sketch Proof. For more details, refer to Theorem 5.4.6 in Casella and Berger¹ or Exercise 6.3-10 in Hogg *et al.*².

The value of $f_{(i,j)}(u, v)\Delta u\Delta v$ is approximated by $P(u \leq X_{(i)} < u + \Delta u, v \leq X_{(j)} < v + \Delta v)$ which is the probability that $i - 1$ random variables are less than u , one random variable is in $[u, u + \Delta u)$, $j - i - 1$ random variables are in $[u + \Delta u, v)$, random variable is in $[v, v + \Delta v)$, and the remaining $n - j$ random variables are greater than or equal to $v + \Delta v$.

Now we have the five subintervals, $(-\infty, u)$, $[u, u + \Delta u)$, $[u + \Delta u, v)$, $[v, v + \Delta v)$, and $[v + \Delta v, \infty)$, and the probabilities of falling to each of the subintervals are

¹CASELLA, G./BERGER, R. L. Statistical Inference. 2nd edition. Pacific Grove, CA: Duxbury, 2002.

²HOGG, R. V./TANIS, E. A./ZIMMERMAN, D. L. Probability and Statistical Inference. 9th edition. Pearson, 2015.

$F(u)$, $F(u + \Delta u) - F(u)$, $F(v) - F(u + \Delta u)$, $F(v + \Delta v) - F(v)$, and $1 - F(v + \Delta v)$, respectively.

The probability $P(u \leq X_{(i)} < u + \Delta u, v \leq X_{(j)} < v + \Delta v)$ is obtained by the multinomial distribution and we thus have

$$f_{(i,j)}(u, v) \Delta u \Delta v \approx \frac{n!}{(i-1)! \times 1! \times (j-i-1)! \times 1! \times (n-j)!} \times \\ [F(u)]^{i-1} [F(u + \Delta u) - F(u)] [F(v) - F(u + \Delta u)]^{j-i-1} \times \\ [F(v + \Delta v) - F(v)] [1 - F(v + \Delta v)]^{n-j}.$$

Then we have

$$f_{(i,j)}(u, v) \approx \frac{n!}{(i-1)! (j-i-1)! (n-j)!} \times [F(u)]^{i-1} \times \frac{F(u + \Delta u) - F(u)}{\Delta u} \times \\ [F(v) - F(u + \Delta u)]^{j-i-1} \times \frac{F(v + \Delta v) - F(v)}{\Delta v} \times [1 - F(v + \Delta v)]^{n-j}.$$

In the limit as $\Delta u \rightarrow 0^+$ and $\Delta v \rightarrow 0^+$, we have the result. \square

1.2 Transformations of variables

First, we briefly review the transformation of rectangular coordinates (x, y) to polar coordinates (r, θ) . If the function $f(x, y)$ is defined on the domain D , then we have

$$\iint_D f(x, y) dx dy = \iint_{D^*} f(r \cos \theta, r \sin \theta) r dr d\theta. \quad (1)$$

This popular transformation is used in statistics to prove that

$$\int_{-\infty}^{\infty} \phi(z) dz = 1, \quad (2)$$

where

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \cdot e^{-z^2/2}.$$

To prove (2), it suffices to show that

$$\int_{-\infty}^{\infty} \phi(x) dx \int_{-\infty}^{\infty} \phi(y) dy = 1,$$

which implies that we need to show

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy = 2\pi. \quad (3)$$

Using (1), the integral in (3) becomes

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)/2} dx dy = \int_0^{\infty} \int_0^{2\pi} e^{-r^2/2} r dr d\theta.$$

Since

$$\int_0^{\infty} e^{-r^2/2} r dr = \left[-e^{-r^2/2} \right]_0^{\infty} = 1,$$

we have

$$\int_0^{2\pi} \int_0^{\infty} e^{-r^2/2} r dr d\theta = \int_0^{2\pi} 1 d\theta = 2\pi.$$

Next, we consider the transformation of bivariate continuous-type random variables with a joint pdf $f(x_1, x_2)$ defined on S_X . Then it is well known that the probability $P[(X_1, X_2) \in D]$ is calculated using the double integral

$$P[(X_1, X_2) \in D] = \iint_D f(x_1, x_2) dx_1 dx_2.$$

Consider new random variables, $Y_1 = u_1(X_1, X_2)$ and $Y_2 = u_2(X_1, X_2)$ where $y_1 = u_1(x_1, x_2)$ and $y_2 = u_2(x_1, x_2)$ have only one-to-one transformations, that is, their inverse transforms exist and are given by $x_1 = v_1(y_1, y_2)$ and $x_2 = v_2(y_1, y_2)$. Then the probability $P[(X_1, X_2) \in D]$ is also calculated using the double y_1 and y_2

$$P[(Y_1, Y_2) \in D^*] = \iint_{D^*} f(v_1(y_1, y_2), v_2(y_1, y_2)) \cdot |J| dx_1 dx_2.$$

Thus, the joint pdf of Y_1 and Y_2 is given by

$$g(y_1, y_2) = f(v_1(y_1, y_2), v_2(y_1, y_2)) \cdot |J|,$$

where $(y_1, y_2) \in S_Y$ and J is the Jacobian determinant of order two

$$J = \det \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{bmatrix} = \det \begin{bmatrix} \frac{\partial v_1(y_1, y_2)}{\partial y_1} & \frac{\partial v_1(y_1, y_2)}{\partial y_2} \\ \frac{\partial v_2(y_1, y_2)}{\partial y_1} & \frac{\partial v_2(y_1, y_2)}{\partial y_2} \end{bmatrix}.$$

1.3 Distribution of the range

Let Z_1, Z_2, \dots, Z_n be a random sample from a standard normal distribution with pdf $\phi(z)$ and cdf $\Phi(z)$. For notational convenience, we denote $X_1 = Z_{(1)}$ and $X_2 = Z_{(n)}$. Using Theorem 1, we have the joint pdf of X_1 and X_2

$$f_{(1,n)}(x_1, x_2) = n(n-1) \phi(x_1) \phi(x_2) [\Phi(x_2) - \Phi(x_1)]^{n-2}.$$

The goal is to derive the distribution of the range of the sample, $X_2 - X_1 = Z_{(n)} - Z_{(1)}$. Next we consider the new random variables given by $Y_1 = X_1$ and $Y_2 = X_2 - X_1$. Notice that the random variable Y_2 is the *range*. The inverse transforms are easily obtained by $x_1 = y_1$ and $x_2 = y_1 + y_2$. Then, using the Jacobian transformation discussed earlier, the joint pdf of Y_1 and Y_2 , denoted by $g(y_1, y_2)$, is given by

$$g(y_1, y_2) = n(n-1) \phi(y_1) \phi(y_1 + y_2) [\Phi(y_1 + y_2) - \Phi(y_1)]^{n-2} |J|,$$

where $-\infty < y_1 < \infty$, $y_2 > 0$ and

$$J = \det \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{bmatrix} = \det \begin{bmatrix} \frac{\partial y_1}{\partial y_1} & \frac{\partial y_1}{\partial y_2} \\ \frac{\partial(y_1 + y_2)}{\partial y_1} & \frac{\partial(y_1 + y_2)}{\partial y_2} \end{bmatrix} = \det \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} = 1.$$

Thus, since the Jacobian determinant is equal to one, the pdf of Y_2 (the range) is just the marginal pdf of Y_2 which is given by

$$\begin{aligned} g_2(y_2) &= \int_{-\infty}^{\infty} g(y_1, y_2) dy_1 \\ &= n(n-1) \int_{-\infty}^{\infty} \phi(y_1) \phi(y_1 + y_2) [\Phi(y_1 + y_2) - \Phi(y_1)]^{n-2} dy_1. \end{aligned} \quad (4)$$

Note that the cdf of Y_2 can be easily obtained by

$$G_2(y_2) = n \int_{-\infty}^{\infty} \phi(y_1) [\Phi(y_1 + y_2) - \Phi(y_1)]^{n-1} dy_1.$$

1.4 The k -th moment of the range

Using the pdf of the range in (4), we can obtain the k -th moment of the range, $Y_2 = Z_{(n)} - Z_{(1)}$, by calculating the expectation as follows:

$$\begin{aligned}
 E(Y_2^k) &= \int_0^\infty y_2^k g_2(y_2) dy_2 \\
 &= n(n-1) \int_0^\infty y_2^k \int_{-\infty}^\infty \phi(y_1) \phi(y_1 + y_2) [\Phi(y_1 + y_2) - \Phi(y_1)]^{n-2} dy_1 dy_2 \\
 &= n(n-1) \int_{-\infty}^\infty \int_0^\infty y_2^k \phi(y_1) \phi(y_1 + y_2) [\Phi(y_1 + y_2) - \Phi(y_1)]^{n-2} dy_2 dy_1 \\
 &= n(n-1) \int_{-\infty}^\infty \left\{ \int_0^\infty y_2^k [\Phi(y_1 + y_2) - \Phi(y_1)]^{n-2} \phi(y_1 + y_2) dy_2 \right\} \phi(y_1) dy_1.
 \end{aligned}$$

It should be noted that if we replace (y_1, y_2) with (x, w) , then the expression for $E(Y_2^k)$ is identical to Equation (1) of Harter³ which is given by

$$E(W^k) = n(n-1) \int_{-\infty}^\infty \left\{ \int_0^\infty w^k [\Phi(x+w) - \Phi(x)]^{n-2} \phi(x+w) dw \right\} \phi(x) dx, \quad (5)$$

where $W = Z_{(n)} - Z_{(1)}$.

Clearly, the expression for the k -th moment of the range requires the evaluation of a complicated double integral. Fortunately, for the case where $k = 1$ which is the expectation, we can derive an alternative formula involving only a single integral. The derivation of this formula will require the application of three different lemmas which we state and prove below.

Lemma 1. *Let X be a continuous random variable with cdf $F(x)$. If $E(|X|^k)$ exists, then we have*

$$(i) \lim_{x \rightarrow \infty} x^k \{1 - F(x)\} = 0 \quad \text{and} \quad (ii) \lim_{x \rightarrow -\infty} |x|^k F(x) = 0.$$

Proof. (i) For $x > 0$, we have

$$0 \leq x^k \{1 - F(x)\} = x^k \int_x^\infty dF(t) = \int_x^\infty x^k dF(t) \leq \int_x^\infty t^k dF(t).$$

³HARTER, H. LEON Tables of Range and Studentized Range. The Annals of Mathematical Statistics, 31 1960, Nr. 4 <URL: <http://dx.doi.org/10.1214/aoms/1177705684>>.

Now, if we can show that the last term $\int_x^\infty t^k dF(t) \rightarrow 0$ in the limit as $x \rightarrow \infty$, then this will complete the proof because we just showed that $0 \leq x^k \{1 - F(x)\} \leq \int_x^\infty t^k dF(t)$ for $x > 0$. In order to prove that $\int_x^\infty t^k dF(t) \rightarrow 0$ in the limit as $x \rightarrow \infty$, note that we also have

$$\begin{aligned} \int_x^\infty t^k dF(t) &= \int_{-\infty}^\infty |t|^k dF(t) - \int_{-\infty}^x |t|^k dF(t) \\ &= E(|X|^k) - \int_{-\infty}^x |t|^k dF(t). \end{aligned}$$

Since $E(|X|^k)$ exists and $\lim_{x \rightarrow \infty} \int_{-\infty}^x |t|^k dF(t) = E(|X|^k)$, we have

$$\int_x^\infty t^k dF(t) = E(|X|^k) - \int_{-\infty}^x |t|^k dF(t) \rightarrow 0$$

in the limit as $x \rightarrow \infty$.

(ii) For $x < 0$, we have

$$0 \leq |x|^k F(x) = |x|^k \int_{-\infty}^x dF(t) = \int_{-\infty}^x |x|^k dF(t) \leq \int_{-\infty}^x |t|^k dF(t).$$

Now, if we can show that the last term $\int_{-\infty}^x |t|^k dF(t) \rightarrow 0$ in the limit as $x \rightarrow -\infty$, then this will complete the proof because we just showed that $0 \leq |x|^k F(x) \leq \int_{-\infty}^x |t|^k dF(t)$ for $x < 0$. In order to prove that $\int_{-\infty}^x |t|^k dF(t) \rightarrow 0$ in the limit as $x \rightarrow -\infty$, note that we also have

$$\begin{aligned} \int_{-\infty}^x |t|^k dF(t) &= \int_{-\infty}^\infty |t|^k dF(t) - \int_x^\infty |t|^k dF(t) \\ &= E(|X|^k) - \int_x^\infty |t|^k dF(t). \end{aligned}$$

Since $E(|X|^k)$ exists and $\lim_{x \rightarrow -\infty} \int_x^\infty |t|^k dF(t) = E(|X|^k)$, we have

$$\int_{-\infty}^x |t|^k dF(t) = E(|X|^k) - \int_x^\infty |t|^k dF(t) \rightarrow 0$$

in the limit as $x \rightarrow -\infty$. □

Lemma 2. Let X be a continuous random variable with cdf $F(x)$. Then we have

$$E(X) = \int_0^\infty [1 - F(x)] dx - \int_{-\infty}^0 F(x) dx \quad (6)$$

$$= \int_0^\infty [1 - F(x) - F(-x)] dx. \quad (7)$$

Proof. We have

$$\begin{aligned}
 E(X) &= \int_{-\infty}^{\infty} x dF(x) \\
 &= \int_{-\infty}^0 x dF(x) + \int_0^{\infty} x dF(x) \\
 &= \int_{-\infty}^0 x dF(x) - \int_0^{\infty} x d[1 - F(x)].
 \end{aligned} \tag{8}$$

Using integration by parts (*i.e.*, $\int u dv = uv - \int v du$), we have

$$\int_{-\infty}^0 x dF(x) = \left[xF(x) \right]_{-\infty}^0 - \int_{-\infty}^0 F(x) dx \tag{9}$$

and

$$\int_0^{\infty} x d[1 - F(x)] = \left[x\{1 - F(x)\} \right]_0^{\infty} - \int_0^{\infty} [1 - F(x)] dx. \tag{10}$$

Applying Lemma 1 to both (9) and (10), we have

$$\int_{-\infty}^0 x dF(x) = - \int_{-\infty}^0 F(x) dx \tag{11}$$

and

$$\int_0^{\infty} x d[1 - F(x)] = - \int_0^{\infty} [1 - F(x)] dx. \tag{12}$$

Next, we can substitute Equations (11) and (12) into (8), in order to obtain (6). Finally, using a change of integration variable technique, we know that $\int_{-\infty}^0 F(x) dx = \int_0^{\infty} F(-x) dx$. Therefore, (6) is equivalent to (7) which completes the proof. \square

It should be noted that Lemma 2 is also valid for discrete random variables.

Lemma 3. Let X_1, X_2, \dots, X_n be a random sample with cdf $F(x)$. Let $F_{(j)}(x)$ denote the cdf of the j -th order statistic $X_{(j)}$. Then we have

$$(i) F_{(n)}(x) = [F(x)]^n \quad \text{and} \quad (ii) F_{(1)}(x) = 1 - [1 - F(x)]^n. \tag{13}$$

Proof. (i) Since $X_{(n)} = \max_{1 \leq i \leq n} X_i$ and X_1, X_2, \dots, X_n are independent, we have

$$\begin{aligned} F_{(n)}(x) &= P\left[\max_{1 \leq i \leq n} X_i \leq x\right] \\ &= P[X_1 \leq x, X_2 \leq x, \dots, X_n \leq x] \\ &= P[X_1 \leq x] P[X_2 \leq x] \cdots P[X_n \leq x] \\ &= [F(x)]^n. \end{aligned}$$

(ii) Similarly, since $X_{(1)} = \min_{1 \leq i \leq n} X_i$, we have

$$\begin{aligned} 1 - F_{(1)}(x) &= 1 - P\left[\min_{1 \leq i \leq n} X_i \leq x\right] \\ &= P\left[\min_{1 \leq i \leq n} X_i > x\right] \\ &= P[X_1 > x, X_2 > x, \dots, X_n > x] \\ &= P[X_1 > x] P[X_2 > x] \cdots P[X_n > x] \\ &= [1 - F(x)]^n. \end{aligned}$$

Thus, we have

$$F_{(1)}(x) = 1 - [1 - F(x)]^n.$$

This completes the proof. □

Theorem 2. Let X_1, X_2, \dots, X_n be a random sample with cdf $F(x)$. Then the expectation of the range is given by

$$E[X_{(n)} - X_{(1)}] = \int_{-\infty}^{\infty} \left\{ 1 - [F(x)]^n - [1 - F(x)]^n \right\} dx.$$

Proof. Using Lemma 2, we have

$$E[X_{(n)}] = \int_0^{\infty} [1 - F_{(n)}(x) - F_{(n)}(-x)] dx$$

and

$$E[X_{(1)}] = \int_0^{\infty} [1 - F_{(1)}(x) - F_{(1)}(-x)] dx.$$

Applying (13) in Lemma 3 to the integral above, we obtain

$$E[X_{(n)}] = \int_0^\infty \left\{ 1 - [F(x)]^n - [F(-x)]^n \right\} dx$$

and

$$E[X_{(1)}] = \int_0^\infty \left\{ [1 - F(x)]^n dx - 1 + [1 - F(-x)]^n \right\} dx.$$

Thus, we have

$$\begin{aligned} E[X_{(n)} - X_{(1)}] &= \int_0^\infty \left\{ 1 - [F(x)]^n - [F(-x)]^n - [1 - F(x)]^n + 1 - [1 - F(-x)]^n \right\} dx \\ &= \int_0^\infty \left\{ 1 - [F(x)]^n - [1 - F(x)]^n + 1 - [F(-x)]^n - [1 - F(-x)]^n \right\} dx \\ &= \int_0^\infty \left\{ 1 - [F(x)]^n - [1 - F(x)]^n \right\} dx + \int_0^\infty \left\{ 1 - [F(-x)]^n - [1 - F(-x)]^n \right\} dx. \end{aligned}$$

Using the change of the integration variable technique for the last term in the above, we have

$$\int_0^\infty \left\{ 1 - [F(-x)]^n - [1 - F(-x)]^n \right\} dx = \int_{-\infty}^0 \left\{ 1 - [F(x)]^n - [1 - F(x)]^n \right\} dx.$$

It is immediate from this result that we have

$$\begin{aligned} E[X_{(n)} - X_{(1)}] &= \int_0^\infty \left\{ 1 - [F(x)]^n - [1 - F(x)]^n \right\} dx + \int_{-\infty}^0 \left\{ 1 - [F(x)]^n - [1 - F(x)]^n \right\} dx \\ &= \int_{-\infty}^\infty \left\{ 1 - [F(x)]^n - [1 - F(x)]^n \right\} dx, \end{aligned}$$

which completes the proof. \square

It should be noted that the above lemmas and theorems are also valid for non-normal distributions. But, we use the results specifically in the case of the normal distribution. Now suppose that we have a random sample from a *standard* normal distribution, Z_1, Z_2, \dots, Z_n , and we want to calculate the expectation of the sample range. Then we have

$$E[Z_{(n)} - Z_{(1)}] = \int_{-\infty}^\infty \left\{ 1 - [\Phi(z)]^n - [1 - \Phi(z)]^n \right\} dz.$$

Note that the integrand, $1 - [\Phi(z)]^n - [1 - \Phi(z)]^n$, is an even function due to the fact that $\Phi(-z) = 1 - \Phi(z)$ which allows for the simplification of the expectation:

$$d_2 = E[Z_{(n)} - Z_{(1)}] = 2 \int_0^\infty \left\{ 1 - [\Phi(z)]^n - [1 - \Phi(z)]^n \right\} dz.$$

Thus, the estimator $(X_{(n)} - X_{(1)})/d_2$ is unbiased for σ if a random sample, X_1, X_2, \dots, X_n , is from a normal distribution.

If X_i are from $N(\mu, \sigma^2)$, then we have $R = X_{(n)} - X_{(1)} = \sigma(Z_{(n)} - Z_{(1)})$, where Z_i are from $N(\mu, 1)$. Note that d_3 is defined by $\text{Var}(R/d_3) = \sigma^2$. Thus, using $\text{Var}(R) = \sigma^2 \text{Var}(Z_{(n)} - Z_{(1)})$, we have

$$d_3 = \sqrt{\text{Var}(Z_{(n)} - Z_{(1)})} = \sqrt{\text{Var}(W)} = \sqrt{E[W^2] - \{E(W)\}^2},$$

where $W = Z_{(n)} - Z_{(1)}$. Then, using (5), we can calculate d_3 .

The values of d_2 and d_3 can be easily calculated using the R programming language and we illustrate these below.

R code implementation for $d_2 = E[Z_{(n)} - Z_{(1)}]$

```

1  # =====
2  # Calculating d2 factors
3  # -----
4  # =====
5  # Using a single integral
6  # -----
7  d2A = function(n) {
8      integrand = function(x) {
9          Phi = pnorm(x)
10         return(1 - (1-Phi)^n - Phi^n)
11     }
12     tmp = integrate(integrand, lower=0, upper=Inf)
13     return(2*tmp$value)
14 }
15
16 # =====
17 # Using a double integral
18 # -----
19 d2B = function(n,k=1) {
20     f= function(x,y) {
21         (y^k) * exp(-(x^2+(x+y)^2)/2) * (pnorm(x+y)-pnorm(x))^(n-2)
22     }
23     tmp = integrate(function(y) {
24         sapply(y,function(y) {
25             integrate(function(x){f(x,y)},-Inf,Inf)$value})
26         }, 0, Inf)
27     return( n*(n-1)/2/pi*tmp$value )
28 }
29
30 # =====
31 # Using a simulation
32 # -----
33 d2C = function(n, iter=1000) {
34     m1 = 0
35     for ( i in 1:iter ) {
36         r = diff( range(rnorm(n)) )
37         m1 = m1 + r / iter
38     }
39     return(m1)
40 }
41 # -----
42
43 # =====
44 # Example 1
45 # -----
46 d2A(25)
47 d2B(25)
48 d2C(25) # Not accurate enough with iter=1000
49 d2C(25, iter=100000)
50
51 # =====
52 # Example 2
53 # -----

```

```

54 system.time(d2A(25))
55 system.time(d2B(25))
56 system.time(d2C(25))
57
58 # =====
59 # Example 3
60 # -----
61 system.time(for(i in 1:100) d2A(25))
62 system.time(for(i in 1:100) d2B(25))
63 system.time(for(i in 1:100) d2C(25))
64 system.time(for(i in 1:100) d2C(25, iter=10000))
65
66 # =====
67 # Calculating d3 factors
68 # -----
69 # =====
70 # Using a single integral
71 # (Maybe impossible)
72 # -----
73
74 # =====
75 # Using a double integral
76 # -----
77 d3B = function(n) {
78   s2 = d2B(n,k=2) - (d2A(n))^2
79   return( sqrt(s2) )
80 }
81
82 # =====
83 # Using a simulation
84 # -----
85 d3C = function(n, iter=1000) {
86   m1 = 0 ; m2 = 0
87   for ( i in 1:iter ) {
88     r = diff( range(rnorm(n)) )
89     m1 = m1 + r / iter
90     m2 = m2 + r^2 / iter
91   }
92   s2 = iter/(iter-1) * (m2-m1^2)
93   return( sqrt(s2) )
94 }
95 # -----
96
97 # =====
98 # Example 4
99 # -----
100 d3B(25)
101 d3C(25) # Not accurate enough
102 d3C(25,iter=100000)

```

2 The bias correction factor for the sample standard deviation, c_4

It is well known that the sample variance is unbiased under the normal distribution assumption. That is

$$E[S^2] = \sigma^2,$$

where $S^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n-1)$. Given that the variance estimate is unbiased, a natural question arises. *Is the sample standard deviation also unbiased?* We already know that, *if* the sample standard deviation is unbiased, then the relation below holds

$$E[S] = \sigma,$$

where $S = [\sum_{i=1}^n (X_i - \bar{X})^2 / (n-1)]^{1/2}$. As it turns out, the sample standard deviation is *not* unbiased under the normal distribution assumption but, by calculating the expectation of S explicitly, we can derive a bias-corrected expression which is unbiased. The derivation follows below.

It is also known that $(n-1)S^2/\sigma^2$ has the chi-square distribution with $n-1$ degrees of freedom which is equivalent to the gamma distribution with $\alpha = (n-1)/2$ and $\theta = 2$. Now, it is well known that

$$E[Y^c] = \frac{\Gamma(\alpha + c)\theta^c}{\Gamma(\alpha)}$$

when Y has the gamma distribution with parameters α and θ . Clearly, for $c = 1/2$, we have

$$E[\sqrt{Y}] = \frac{\Gamma(\alpha + 1/2)\sqrt{\theta}}{\Gamma(\alpha)}.$$

Now let $Y = (n-1)S^2/\sigma^2$ so that Y has the gamma distribution with parameters $\alpha = (n-1)/2$ and $\theta = 2$. Then using the relation above, we obtain

$$E[\sqrt{(n-1)S^2/\sigma^2}] = \frac{\Gamma(n/2)\sqrt{2}}{\Gamma(n/2 - 1/2)}.$$

We then take the previous expectation and simplify under the square-root on the left hand side of the above equation and obtain

$$\frac{\sqrt{n-1}}{\sigma} \cdot E[S] = \frac{\Gamma(n/2)\sqrt{2}}{\Gamma(n/2 - 1/2)}.$$

But this implies that

$$E[S] = c_4 \sigma$$

where

$$c_4 = \sqrt{\frac{2}{n-1}} \cdot \frac{\Gamma(n/2)}{\Gamma(n/2 - 1/2)}.$$

Thus, the estimator S/c_4 is unbiased for σ .

The bias correction factor c_4 can be easily calculated using the R programming language and we illustrate this below.

R code implementation for $c_4 = E[S]/\sigma$

```

1  # =====
2  # Calculating c4 factors
3  # -----
4
5  # =====
6  # Using an original gamma
7  # -----
8  c4A = function(n) sqrt(2/(n-1))*gamma(n/2)/gamma(n/2-1/2)
9
10 # =====
11 # Using a log-gamma
12 # -----
13 c4B = function(n) {
14     tmp = lgamma(n/2) - lgamma(n/2-1/2)
15     sqrt(2/(n-1)) * exp(tmp)
16 }
17
18 # =====
19 # Using a simulation
20 # -----
21 c4C = function(n, iter=1000) {
22     m1 = 0
23     for ( i in 1:iter ) {
24         m1 = m1 + sd(rnorm(n)) / iter
25     }
26     return(m1)
27 }
28 # -----
29
30 # =====
31 # Example 1
32 # -----
33 c4A(10)
34
35 c4B(10)
36
37 c4C(10) # Not accurate enough with iter=1000
38 c4C(10, iter=100000)
39
40 # =====
41 # Example 2
42 # -----
43 c4A(350) # value out of the calculation of gamma function
44
45 c4B(350)
46
47 c4C(350) # Not accurate enough with iter=1000
48 c4C(350, iter=100000)

```

3 The power function of statistical hypothesis testing

In this section, we perform the two-sided hypothesis testing for $H_0 : \mu = \mu_0$ versus $H_1 : \mu \neq \mu_0$. We assume that X_1, X_2, \dots, X_n are from the normal with mean μ and variance σ^2 . This test is well known as z -test (when σ is known) or t -test (when σ is unknown) in the statistics literature. For example, see Section 8.3 of Casella and Berger⁴.

3.1 When the variance is known

When the variance σ^2 is known, the power function of the z -test is then given by

$$\begin{aligned} K_z(\mu) &= P\left[\frac{|\bar{X} - \mu_0|}{\sigma/\sqrt{n}} > z_{\alpha/2}\right] \\ &= 1 - P\left[\frac{|\bar{X} - \mu_0|}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right] \\ &= 1 - P\left[-z_{\alpha/2} \leq \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right]. \end{aligned} \quad (14)$$

Notice that X_i are from $N(\mu, \sigma^2)$, not from $N(\mu_0, \sigma^2)$. Thus, $(|\bar{X} - \mu|)/(\sigma/\sqrt{n})$ is distributed as the standard normal distribution, $N(0, 1)$, but $(|\bar{X} - \mu_0|)/(\sigma/\sqrt{n})$ is not.

We rewrite the power function in (14) by

$$\begin{aligned} K_z(\mu) &= 1 - P\left[-z_{\alpha/2} + \frac{\mu_0}{\sigma/\sqrt{n}} \leq \frac{\bar{X}}{\sigma/\sqrt{n}} \leq z_{\alpha/2} + \frac{\mu_0}{\sigma/\sqrt{n}}\right] \\ &= 1 - P\left[-z_{\alpha/2} + \frac{\mu_0 - \mu}{\sigma/\sqrt{n}} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2} + \frac{\mu_0 - \mu}{\sigma/\sqrt{n}}\right] \\ &= 1 - P\left[-z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right] \\ &= 1 - P\left[-z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}} \leq Z \leq z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right], \end{aligned}$$

where $Z \sim N(0, 1)$. Thus, we have

$$K_z(\mu) = 1 - \left\{ \Phi\left(z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) - \Phi\left(-z_{\alpha/2} - \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) \right\},$$

⁴CASELLA, G./BERGER, R. L. Statistical Inference. 2nd edition. Pacific Grove, CA: Duxbury, 2002.

where $\Phi(\cdot)$ is the cdf of the standard normal distribution. Using the relation $\Phi(-z) = 1 - \Phi(z)$, we can rewrite the above by

$$K_z(\mu) = 1 - \left\{ 1 - \Phi\left(-z_{\alpha/2} + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) - 1 + \Phi\left(z_{\alpha/2} + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) \right\}.$$

Therefore, the power function of the z -test is given by

$$\begin{aligned} K_z(\mu) &= 1 + \Phi\left(-z_{\alpha/2} + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) - \Phi\left(z_{\alpha/2} + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) \\ &= 1 - \Phi\left(z_{\alpha/2} + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right) + \Phi\left(-z_{\alpha/2} + \frac{\mu - \mu_0}{\sigma/\sqrt{n}}\right). \end{aligned}$$

3.2 When the variance is not known

Theorem 3. *Let X_1, X_2, \dots, X_n be a random sample from a normal distribution with μ and variance σ^2 . Let \bar{X} and S^2 denote the sample mean and variance, respectively. Then the following t -test statistic under the local alternative $H_1 : \mu = \mu_0 + \delta\sigma/\sqrt{n}$*

$$\frac{\bar{X} - \mu_0}{S/\sqrt{n}}$$

has a non-central t -distribution with $n - 1$ degrees of freedom and non-centrality δ .

Proof. Recall the definition of a non-central t -distribution. Let $Z \sim N(0, 1)$ and V has a chi-square distribution with r degrees of freedom. Suppose that Z and V are independent. Then the quotient below has a non-central t -distribution with r degrees of freedom and non-centrality μ :

$$\frac{Z + \mu}{\sqrt{V/r}}.$$

Let $V = (n - 1)S^2/\sigma^2$ for convenience. Then V has a chi-square distribution with $n - 1$ degrees of freedom. See Theorem 5.3.1 of CASELLA, G./BERGER, R. L.

Statistical Inference. 2nd edition. Pacific Grove, CA: Duxbury, 2002. We have

$$\begin{aligned}\frac{\sqrt{n}(\bar{X} - \mu_0)}{S} &= \frac{\sqrt{n}(\bar{X} - \mu_0)/\sigma}{\sqrt{V/(n-1)}} \\ &= \frac{\sqrt{n}(\bar{X} - \mu)/\sigma + \sqrt{n}(\mu - \mu_0)/\sigma}{\sqrt{V/(n-1)}} \\ &= \frac{Z + \sqrt{n}(\mu - \mu_0)/\sigma}{\sqrt{V/(n-1)}},\end{aligned}$$

where $Z = (\bar{X} - \mu)/(\sigma/\sqrt{n})$ and $Z \sim N(0, 1)$. Thus, under the local $H_1 : \mu = \mu_0 + \delta\sigma/\sqrt{n}$, we have

$$\frac{\bar{X} - \mu_0}{S/\sqrt{n}} = \frac{Z + \delta}{\sqrt{V/(n-1)}}.$$

Since S^2 and \bar{X} are independent, V and Z are also independent. This completes the proof. \square

Next, we want to obtain the power function for testing $H_0 : \mu = \mu_0$ versus $H_1 : \mu \neq \mu_0$. Since $\mu = \mu_0 + \{\sqrt{n}(\mu - \mu_0)/\sigma\} \cdot \{\sigma/\sqrt{n}\}$, it is immediate upon using Theorem 3 that $(\bar{X} - \mu_0)/(S/\sqrt{n})$ under H_1 has the non-central t -distribution with $\nu = n - 1$ degrees of freedom and non-centrality $\delta = (\mu - \mu_0)/(\sigma/\sqrt{n})$.

The critical region for testing $H_0 : \mu = \mu_0$ versus $H_1 : \mu \neq \mu_0$ is given by

$$\frac{|\bar{X} - \mu_0|}{S/\sqrt{n}} > t_{\alpha/2}.$$

For convenience, we let $T_{n-1}(\delta) = (\bar{X} - \mu_0)/(S/\sqrt{n})$. Then the critical region can be rewritten as $|T_{n-1}(\delta)| > t_{\alpha/2}$. Then the power function is given by

$$\begin{aligned}K_t(\mu) &= P(|T_{n-1}(\delta)| > t_{\alpha/2}) \\ &= P(T_{n-1}(\delta) > t_{\alpha/2}) + P(T_{n-1}(\delta) < -t_{\alpha/2}) \\ &= 1 - \Phi_{\nu, \delta}(t_{\alpha/2}) + \Phi_{\nu, \delta}(-t_{\alpha/2}),\end{aligned}$$

where $\Phi_{\nu, \delta}(\cdot)$ is the cdf of the non-central t -distribution with $\nu = n - 1$ degrees of freedom and non-centrality $\delta = (\mu - \mu_0)/(\sigma/\sqrt{n})$.

3.3 Simulation

As a pilot study, we carried out the Monte Carlo simulation to obtain the empirical power curve of the t -test for testing $H_0 : \mu = 1/2$ versus $H_1 : \mu \neq 1/2$ with the

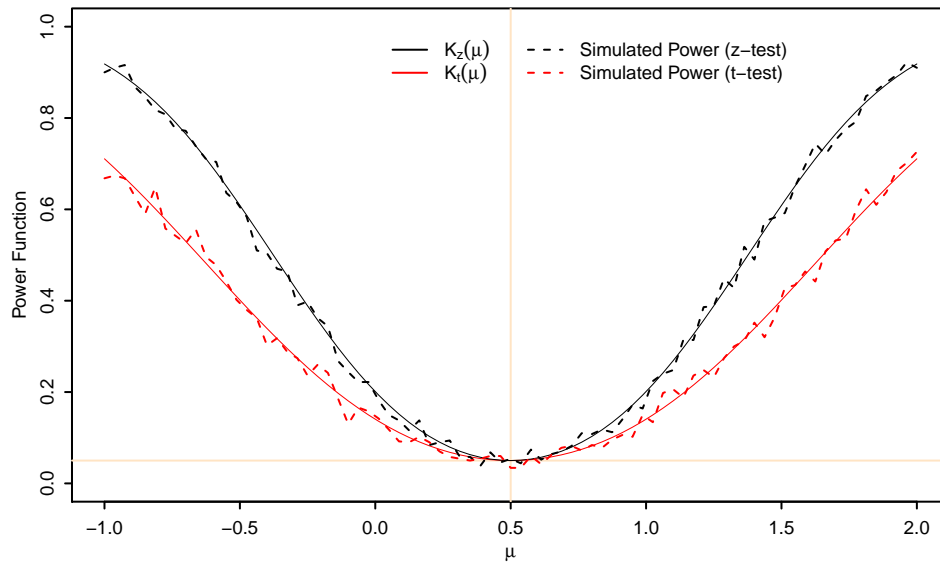


Figure 1: Empirical power curve along with theoretical power curves.

significance level $\alpha = 0.05$. A random sample of size $n = 5$ was generated from the normal distribution with mean μ and $\sigma = 1$, where μ varies from -1 to 2 . This sampling was replicated $I = 1,000$. In Figure 1, we plotted the empirical power curve (blue solid line) along with the theoretical power curves of the z -test (dotted red line) and the t -test (dotted blue line).

R code implementation for Figure 1

```

1  #=====
2  pdf(file="Two-side-t-test.pdf", width=5.0, height=3.0, encoding="TeXtext.enc")
3  par(mfrow=c(1,1), mar=c(5, 5, 1, 1), omi=c(0,0,0,0), cex=0.6, mex=0.6)
4  #-----
5  iter = 500
6  mu0 = 0.5; sigma=1; alpha=0.05 ; n = 5
7  MU = seq(-1,2, l=81)
8
9  # Power function : Assuming sigma=known
10 Kz = function(mu, alpha, mu0, sigma, n) {
11     z.cut = qnorm(1-alpha/2)
12     tmp = (mu-mu0)/(sigma/sqrt(n))
13     pnorm(z.cut + tmp, lower.tail=FALSE) + pnorm(-z.cut + tmp)
14 }
15
16 # Power function : Assuming sigma = unknown
17 Kt = function(mu, alpha, mu0, sigma, n) {
18     t.cut = qt(1-alpha/2, df=n-1)
19     ncp = (mu-mu0)/(sigma/sqrt(n))
20     pt(t.cut, df=n-1, ncp=ncp, lower.tail=FALSE) + pt(-t.cut, df=n-1, ncp=ncp)
21 }
22 powerz = Kz(MU, alpha=alpha, mu0=mu0, sigma=sigma, n=n)
23 powert = Kt(MU, alpha=alpha, mu0=mu0, sigma=sigma, n=n)
24
25
26 #=====
27 # Simulation Approach
28 #-----
29 nMU = length(MU)
30
31 # sigma known
32 sim.powerz = numeric(nMU)
33 z.cut = qnorm(1-alpha/2)
34 for ( j in 1:nMU ) {
35     for ( i in 1:iter ) {
36         x = rnorm(n, mean=MU[j], sd=sigma)
37         s = sigma      # sigma is known
38         test.stat = abs((mean(x)-mu0)/(s/sqrt(n)))
39         if (test.stat>z.cut) sim.powerz[j] = sim.powerz[j] + 1/iter
40     }
41 }
42
43 # sigma unknown
44 sim.powert = numeric(nMU)
45 t.cut = qt(1-alpha/2, df=n-1)
46 for ( j in 1:nMU ) {
47     for ( i in 1:iter ) {
48         x = rnorm(n, mean=MU[j], sd=sigma)
49         s = sd(x)      # sigma is UNKNOWN
50         test.stat = abs((mean(x)-mu0)/(s/sqrt(n)))
51         if (test.stat>t.cut) sim.powert[j] = sim.powert[j] + 1/iter
52     }

```

```
53 }  
54  
55 # Plot  
56 plot(NA,NA, xlim=range(MU), ylim=c(0,1), type="n",  
57       xlab=expression(mu), ylab="Power Function")  
58 abline(h=alpha, v=mu0, col="bisque")  
59 lines(MU, powert, type="l", lty=1, lwd=0.5, col="red")  
60 lines(MU, powerz, type="l", lty=1, lwd=0.5, col="black" )  
61 lines(MU, sim.powerz, lty=2, lwd=1.0, col="black")  
62 lines(MU, sim.powert, lty=2, lwd=1.0, col="red")  
63 legend (0.0,1, legend=c(expression(K[z](mu)), expression(K[t](mu)) ),  
64        horiz=FALSE, bty="n", lty=c(1,1), col=c("black", "red") )  
65 legend (0.5,1,legend=c("Simulated Power (z-test)", "Simulated Power (t-test)"),  
66        horiz=FALSE, bty="n", lty=c(2,2), lwd=1.0, col=c("black","red") )
```

Factors for Constructing Variables Control Charts

Observations in	Chart for Averages				Chart for Standard Deviations						Chart for Ranges						
	Factors for Control Limits		Factors for Center Line		Factors for Control Limits						Factors for Center Line		Factors for Control Limits				
	Sample, n	A	A_2	A_3	c_4	$1/c_4$	Factors for Control Limits				Factors for Center Line		d_3	D_1	D_2	D_3	D_4
							B_3	B_4	B_5	B_6	d_2	$1/d_2$					
	2	2.121	1.880	2.659	0.7979	1.2533	0	3.267	0	2.606	1.128	0.8865	0.853	0	3.686	0	3.267
	3	1.732	1.023	1.954	0.8862	1.1284	0	2.568	0	2.276	1.693	0.5907	0.888	0	4.358	0	2.574
	4	1.500	0.729	1.628	0.9213	1.0854	0	2.266	0	2.088	2.059	0.4857	0.880	0	4.698	0	2.282
	5	1.342	0.577	1.427	0.9400	1.0638	0	2.089	0	1.964	2.326	0.4299	0.864	0	4.918	0	2.114
	6	1.225	0.483	1.287	0.9515	1.0510	0.030	1.970	0.029	1.874	2.534	0.3946	0.848	0	5.078	0	2.004
	7	1.134	0.419	1.182	0.9594	1.0423	0.118	1.882	0.113	1.806	2.704	0.3698	0.833	0.204	5.204	0.076	1.924
	8	1.061	0.373	1.099	0.9650	1.0363	0.185	1.815	0.179	1.751	2.847	0.3512	0.820	0.388	5.306	0.136	1.864
	9	1.000	0.337	1.032	0.9693	1.0317	0.239	1.761	0.232	1.707	2.970	0.3367	0.808	0.547	5.393	0.184	1.816
	10	0.949	0.308	0.975	0.9727	1.0281	0.284	1.716	0.276	1.669	3.078	0.3249	0.797	0.687	5.469	0.223	1.777
	11	0.905	0.285	0.927	0.9754	1.0252	0.321	1.679	0.313	1.637	3.173	0.3152	0.787	0.811	5.535	0.256	1.744
	12	0.866	0.266	0.886	0.9776	1.0229	0.354	1.646	0.346	1.610	3.258	0.3069	0.778	0.922	5.594	0.283	1.717
	13	0.832	0.249	0.850	0.9794	1.0210	0.382	1.618	0.374	1.585	3.336	0.2998	0.770	1.025	5.647	0.307	1.693
	14	0.802	0.235	0.817	0.9810	1.0194	0.406	1.594	0.399	1.563	3.407	0.2935	0.763	1.118	5.696	0.328	1.672
	15	0.775	0.223	0.789	0.9823	1.0180	0.428	1.572	0.421	1.544	3.472	0.2880	0.756	1.203	5.741	0.347	1.653
	16	0.750	0.212	0.763	0.9835	1.0168	0.448	1.552	0.440	1.526	3.532	0.2831	0.750	1.282	5.782	0.363	1.637
	17	0.728	0.203	0.739	0.9845	1.0157	0.466	1.534	0.458	1.511	3.588	0.2787	0.744	1.356	5.820	0.378	1.622
	18	0.707	0.194	0.718	0.9854	1.0148	0.482	1.518	0.475	1.496	3.640	0.2747	0.739	1.424	5.856	0.391	1.608
	19	0.688	0.187	0.698	0.9862	1.0140	0.497	1.503	0.490	1.483	3.689	0.2711	0.734	1.487	5.891	0.403	1.597
	20	0.671	0.180	0.680	0.9869	1.0133	0.510	1.490	0.504	1.470	3.735	0.2677	0.729	1.549	5.921	0.415	1.585
	21	0.655	0.173	0.663	0.9876	1.0126	0.523	1.477	0.516	1.459	3.778	0.2647	0.724	1.605	5.951	0.425	1.575
	22	0.640	0.167	0.647	0.9882	1.0119	0.534	1.466	0.528	1.448	3.819	0.2618	0.720	1.659	5.979	0.434	1.566
	23	0.626	0.162	0.633	0.9887	1.0114	0.545	1.455	0.539	1.438	3.858	0.2592	0.716	1.710	6.006	0.443	1.557
	24	0.612	0.157	0.619	0.9892	1.0109	0.555	1.445	0.549	1.429	3.895	0.2567	0.712	1.759	6.031	0.451	1.548
	25	0.600	0.153	0.606	0.9896	1.0105	0.565	1.435	0.559	1.420	3.931	0.2544	0.708	1.806	6.056	0.459	1.541

Critical values of t -distribution.

$$t_{\alpha}(r) = F^{-1}(1 - \alpha) \text{ or } P[T \leq t_{\alpha}(r)] = 1 - \alpha,$$

where $F^{-1}(\cdot)$ is the inverse CDF of the t distribution with r degrees of freedom.

r	$t_{0.40}(r)$	$t_{0.30}(r)$	$t_{0.25}(r)$	$t_{0.20}(r)$	$t_{0.10}(r)$	$t_{0.05}(r)$	$t_{0.025}(r)$	$t_{0.010}(r)$	$t_{0.005}(r)$	$t_{0.001}(r)$
1	0.325	0.727	1.000	1.376	3.078	6.314	12.706	31.821	63.657	318.309
2	0.289	0.617	0.816	1.061	1.886	2.920	4.303	6.965	9.925	22.327
3	0.277	0.584	0.765	0.978	1.638	2.353	3.182	4.541	5.841	10.215
4	0.271	0.569	0.741	0.941	1.533	2.132	2.776	3.747	4.604	7.173
5	0.267	0.559	0.727	0.920	1.476	2.015	2.571	3.365	4.032	5.893
6	0.265	0.553	0.718	0.906	1.440	1.943	2.447	3.143	3.707	5.208
7	0.263	0.549	0.711	0.896	1.415	1.895	2.365	2.998	3.499	4.785
8	0.262	0.546	0.706	0.889	1.397	1.860	2.306	2.896	3.355	4.501
9	0.261	0.543	0.703	0.883	1.383	1.833	2.262	2.821	3.250	4.297
10	0.260	0.542	0.700	0.879	1.372	1.812	2.228	2.764	3.169	4.144
11	0.260	0.540	0.697	0.876	1.363	1.796	2.201	2.718	3.106	4.025
12	0.259	0.539	0.695	0.873	1.356	1.782	2.179	2.681	3.055	3.930
13	0.259	0.538	0.694	0.870	1.350	1.771	2.160	2.650	3.012	3.852
14	0.258	0.537	0.692	0.868	1.345	1.761	2.145	2.624	2.977	3.787
15	0.258	0.536	0.691	0.866	1.341	1.753	2.131	2.602	2.947	3.733
16	0.258	0.535	0.690	0.865	1.337	1.746	2.120	2.583	2.921	3.686
17	0.257	0.534	0.689	0.863	1.333	1.740	2.110	2.567	2.898	3.646
18	0.257	0.534	0.688	0.862	1.330	1.734	2.101	2.552	2.878	3.610
19	0.257	0.533	0.688	0.861	1.328	1.729	2.093	2.539	2.861	3.579
20	0.257	0.533	0.687	0.860	1.325	1.725	2.086	2.528	2.845	3.552
21	0.257	0.532	0.686	0.859	1.323	1.721	2.080	2.518	2.831	3.527
22	0.256	0.532	0.686	0.858	1.321	1.717	2.074	2.508	2.819	3.505
23	0.256	0.532	0.685	0.858	1.319	1.714	2.069	2.500	2.807	3.485
24	0.256	0.531	0.685	0.857	1.318	1.711	2.064	2.492	2.797	3.467
25	0.256	0.531	0.684	0.856	1.316	1.708	2.060	2.485	2.787	3.450
26	0.256	0.531	0.684	0.856	1.315	1.706	2.056	2.479	2.779	3.435
27	0.256	0.531	0.684	0.855	1.314	1.703	2.052	2.473	2.771	3.421
28	0.256	0.530	0.683	0.855	1.313	1.701	2.048	2.467	2.763	3.408
29	0.256	0.530	0.683	0.854	1.311	1.699	2.045	2.462	2.756	3.396
30	0.256	0.530	0.683	0.854	1.310	1.697	2.042	2.457	2.750	3.385
∞	0.253	0.524	0.674	0.842	1.282	1.645	1.960	2.326	2.576	3.090

Cumulative Standard Normal Distribution, $\Phi(z) = P(Z \leq z)$

[illegible]