

## Saddlepoint Approximations for Function of Means from Several Populations\*

BING-YI JING

(Department of Mathematics, Hong Kong University of Science and Technology, Hong Kong)

ZE-HUI LI

(Department of Mathematics, Lanzhou University, Lanzhou, 730000)

### Abstract

Saddlepoint approximations for marginal tail probabilities for a real-valued function of vector means from several populations are developed. The approximations are then shown to give great numerical accuracy, which are demonstrated in some numerical examples. Application to the bootstrap is also considered in order to avoid intensive Monte Carlo simulations.

**Keywords:** Bootstrap, Edgeworth expansion, saddlepoint approximation.

**AMS Subject Classification:** 62E20, 62G09.

### §1. Introduction

Saddlepoint approximations of the Lugannani-Rice type for marginal (and conditional) tail probabilities for a real-valued function of a random vector have recently been considered by different authors, for examples, Skovgaard (1987), Daniels and Young (1991), DiCiccio and Martin (1991), Wang (1993), Jing and Robinson (1994) and others. The method used by these authors is a common one: first find a 1-1 transformation from the random vector to a new random vector, and then integrate out the unwanted variables in the joint density of the new random vector by Laplace approximation, finally apply Temme's method to obtain the desired tail area probabilities. In this paper, we shall employ the same techniques to obtain saddlepoint approximations to a function of vector means from several populations. The approach taken here closely follows that of Jing and Robinson (1994) for the case of a function of vector means from one population. Several examples will be considered to illustrate the accuracy of the saddlepoint approximations. We shall also investigate the performance of these approximations in the bootstrap context.

### §2. The saddlepoint approximations for several populations

For simplicity, we shall only consider saddlepoint approximations to a function of vector means from two populations; the generalization to the case of more than two populations is straightforward.

Suppose we have two sets of observations  $\{X_1, \dots, X_m\}$  and  $\{Y_1, \dots, Y_n\}$ , where  $X_i$ 's and  $Y_j$ 's are  $d_1$ -dimensional and  $d_2$ -dimensional random vectors  $R^{d_1}$  and  $R^{d_2}$  (both  $d_1$  and  $d_2$  are positive integers) from two populations with distribution functions  $F$  and  $G$ , respectively. Without loss of generality, we assume that  $m \geq n$ . In many case, we are interested in comparing some quantities from the two populations. For instance, we may be interested in the difference of the means  $EX_1 - EY_1$ , the ratio of the means  $EX_1/EY_1$ , or the ratio of the variances  $\text{Var}(X_1)/\text{Var}(Y_1)$ , etc. Note that all these quantities can be expressed as functions of means from the two population, and can be treated together. First let us formulate the problem.

\*Research supported by National Nature Science Foundation of China 19871035 and Nature Science Foundation of Gansu Province B4.

Received 1999.10.12. Revised 2001.1.3.

Let  $a_1 = g_1(x, y)$  be a real-valued function from  $R^d$  to  $R^1$ , where  $x \in R^{d_1}$  and  $y \in R^{d_2}$  and  $d = d_1 + d_2$ . Let  $\bar{X} = \sum X_i/m$  and  $\bar{Y} = \sum Y_i/n$ . We wish to approximate

$$P(A_1 \geq a_1) \equiv P(g_1(\bar{X}, \bar{Y}) \geq a_1).$$

Clearly, the difference of two means  $EX_1 - EY_1$ , the ratio of two means  $EX_1/EY_1$ , and the ratio of two variances  $\text{Var}(X_1)/\text{Var}(Y_1)$  can all be put into the forms of  $g_1(\mu_x, \mu_y)$  which can be estimated by  $g_1(\bar{X}, \bar{Y})$ , where  $\mu_x = EX_1$  and  $\mu_y = EY_1$ . The purpose of this paper is to study the distribution function of  $g_1(\bar{X}, \bar{Y})$ . Note that  $g_1(\bar{X}, \bar{Y})$  can not be put into the usual smooth function of vector means from one single population, as  $X_i$ 's and  $Y_j$ 's are from different distributions. However, as can be seen from Theorem 1 below, the saddlepoint approximations for both cases are strikingly similar. The development here will closely follow that of Jing and Robinson (1994), the outline of which is sketched below.

For any two vectors  $a$  and  $b$ , we shall use  $\langle a, b \rangle$  to denote their inner product. Let  $K_x(\theta) = \log(Ee^{\langle \theta, X_1 \rangle})$  and  $K_y(\tau) = \log(Ee^{\langle \tau, Y_1 \rangle})$ , where  $\theta \in R^{d_1}$  and  $\tau \in R^{d_2}$ , be the cumulant generating functions of  $X_1$  and  $Y_1$ , respectively. By saddlepoint approximations for the density of  $\bar{X}$  and  $\bar{Y}$  (see Daniels (1987), for instance), we have

$$\begin{aligned} f_{\bar{X}, \bar{Y}}(x, y) &= \frac{e^{-m[\langle \theta, x \rangle - K_x(\theta)]}}{(2\pi/m)^{d_1/2} [\det(K_x''(\theta))]^{1/2}} \frac{e^{-n[\langle \tau, y \rangle - K_y(\tau)]}}{(2\pi/n)^{d_2/2} [\det(K_y''(\tau))]^{1/2}} (1 + O(n^{-1})) \\ &= \frac{e^{-N\Lambda(x, y)}}{(2\pi/N)^{d/2} \Delta^{1/2}} (1 + O(n^{-1})), \end{aligned}$$

where the saddlepoints  $\theta$  and  $\tau$  are solutions of  $K_x'(\theta) = x$  and  $K_y'(\tau) = y$ , respectively, and

$$\begin{aligned} \Lambda(x, y) &= \frac{m}{N} \langle \theta, x \rangle - K_x(\theta) + \frac{n}{N} \langle \tau, y \rangle - K_y(\tau), \\ \Delta &= \det \left[ \frac{N}{m} K_x''(\theta) \right] \det \left[ \frac{N}{n} K_y''(\tau) \right]. \end{aligned}$$

Now let us construct a 1-1 transformation

$$\begin{cases} a_1 = g_1(x_1, \dots, x_{d_1}, y_1, \dots, y_{d_2}) \\ a_2 = g_2(x_1, \dots, x_{d_1}, y_1, \dots, y_{d_2}) \\ \dots \dots \dots \\ a_d = g_d(x_1, \dots, x_{d_1}, y_1, \dots, y_{d_2}), \end{cases}$$

that is,  $a = g(x, y) \equiv g(z)$ . Write the inverse transformation as  $z = g^{-1}(a)$ . Denote the absolute value of the jacobian of the transformation as  $J(a) = |\det(\partial z / \partial a)|$ . So the joint density of  $A \equiv g(\bar{X}, \bar{Y}) \equiv (g_1(\bar{X}, \bar{Y}), \dots, g_d(\bar{X}, \bar{Y}))$  is  $f_A(a) = f_{\bar{X}, \bar{Y}}(x(a), y(a))J(a)$ . Hence the marginal density for  $A_1 = g_1(\bar{X}, \bar{Y})$  is

$$f_{A_1}(a) = \int_{R^{d-1}} f_A(a) da_2 \dots da_d = \frac{e^{-NH(a_1)} \tilde{J}}{\tilde{\Delta}^{1/2} [\det(\tilde{L}_{22})]^{1/2}} (1 + O(n^{-1})), \quad (1)$$

where

$$\begin{aligned} L(a) &= \Lambda(z(a)), \\ L_{22} &= \partial L(a) / \partial (a_2, \dots, a_d), \\ H(a_1) &= \inf_{a_2, \dots, a_d} L(a) \equiv L(a_1, \tilde{a}_2, \dots, \tilde{a}_d), \end{aligned}$$

and write  $\tilde{a} = (a_1, \tilde{a}_2, \dots, \tilde{a}_d)$ , (i.e.  $\tilde{a}$  minimizes  $L(a)$  for fixed  $a_1$ ) and also  $\tilde{J}$ ,  $\tilde{L}_{22}$  and  $\tilde{\Delta}$  are evaluated at  $\tilde{a}$ .

To get a saddlepoint approximation for the tail probability  $P(A_1 \geq a_1)$ , we can integrate  $f_{A_1}(a)$  in (1). By applying the Temme's method as in Jing and Robinson (1994), we get the following theorem.

**Theorem 1** Under some regularity conditions, we have

$$P(A_1 \geq a_1) = 1 - \Phi(\hat{w}\sqrt{N}) - \frac{\phi(\hat{w}\sqrt{N})}{\sqrt{N}} \left( \frac{1}{\hat{w}} - \frac{1}{\hat{u}} + O(n^{-3/2}) \right), \quad (2)$$

where

$$\begin{aligned}\hat{w} &= \sqrt{2H(a_1)} \operatorname{sign}(a_1 - \alpha_1), \\ \hat{u} &= \sqrt{\tilde{\Delta}} \det(\tilde{L}_{22}) \tilde{L}_1 / \tilde{J},\end{aligned}$$

where  $\tilde{L}_1 = \partial L(\tilde{a}) / \partial a_1$  and  $\alpha_1 = \{a_1 : \inf_{a_1} H(a_1)\}$ .

To apply Theorem 1, one can take the following easy steps.

Step 1: For fixed  $a_1$ , one can solve  $(\tilde{\theta}, \tilde{\tau}, \tilde{a}_2, \dots, \tilde{a}_d)$  from

$$\begin{cases} K'_x(\theta) = x(a) \\ K'_y(\tau) = y(a) \\ L_2(a) = \frac{m}{N} \left( \frac{\partial x(a)}{\partial(a_2, \dots, a_d)} \right)^T \theta + \frac{n}{N} \left( \frac{\partial y(a)}{\partial(a_2, \dots, a_d)} \right)^T \tau = 0. \end{cases}$$

Note here we have  $(2d - 1)$  unknowns and  $(2d - 1)$  equations.

Step 2: Write  $\tilde{a} = (a_1, \tilde{a}_2, \dots, \tilde{a}_d)$ , then calculate

$$\begin{aligned}H(a_1) &= L(\tilde{a}) = \frac{m}{N} [\langle \tilde{\theta}, x(\tilde{a}) \rangle - K_x(\tilde{\theta})] + \frac{n}{N} [\langle \tilde{\tau}, y(\tilde{a}) \rangle - K_y(\tilde{\tau})], \\ \tilde{J}(a_1) &= \left| \det \left( \frac{\partial z(\tilde{a})}{\partial a} \right) \right|, \\ \tilde{L}_1 &= \frac{m}{N} \left( \frac{\partial x(\tilde{a})}{\partial a_1} \right)^T \tilde{\theta} + \frac{n}{N} \left( \frac{\partial y(\tilde{a})}{\partial a_1} \right)^T \tilde{\tau}, \\ \tilde{L}_{22} &= \frac{m}{N} \tilde{L}_{x22} + \frac{n}{N} \tilde{L}_{y22},\end{aligned}$$

where

$$\begin{aligned}\tilde{L}_{x22} &= \left( \frac{\partial x(\tilde{a})}{\partial(a_2, \dots, a_d)} \right)^T (K''_x(\tilde{\theta}))^{-1} \left( \frac{\partial x(\tilde{a})}{\partial(a_2, \dots, a_d)} \right) + M_x(\tilde{a}, \tilde{\theta}), \\ M_x(a, \theta) &= \begin{pmatrix} \langle \theta, \frac{\partial^2 x}{\partial a_2 \partial a_2} \rangle & \dots & \langle \theta, \frac{\partial^2 x}{\partial a_2 \partial a_d} \rangle \\ \dots & \dots & \dots \\ \langle \theta, \frac{\partial^2 x}{\partial a_d \partial a_2} \rangle & \dots & \langle \theta, \frac{\partial^2 x}{\partial a_d \partial a_d} \rangle \end{pmatrix}_{(d-1) \times (d-1)},\end{aligned}$$

and  $\tilde{L}_{y22}$  and  $M_y(\tilde{a}, \tilde{\tau})$  are similarly defined.

Step 3: Calculate  $\hat{w}$  and  $\hat{u}$ , and apply (2) in Theorem 1.

An alternative saddlepoint approximation can also be obtained, which does not depend on the transformation  $a = g(z)$ . see Jing and Robinson (1994), DiCiccio and Martin (1991) for more details in related cases. We shall not pursue this here. But it should be pointed out that it is often easier to construct simple 1-1 transformation  $a = g(z)$  (which is often easy to find) and then apply Theorem 1 than using the alternative formula.

### §3. Some numerical examples

#### 3.1 An exact case

Let  $F = N(0, 1)$  and  $G = N(0, 1)$ . To find out the distribution for  $\mathbf{P}(\bar{X} + \bar{Y} \geq a_1)$ , we choose  $d = 2$ .  $a_1 = x + y$  and  $a_2 = y$ . Applying Theorem 1, we get

$$\mathbf{P}(\bar{X} + \bar{Y} \geq a_1) = 1 - \Phi(a_1 / \sqrt{1/m + 1/n}).$$

In this case, the saddlepoint approximation is exact.

#### 3.2 Two-sample studentized- $t$ statistic

Let  $X_1, \dots, X_m \sim N(0, 1)$  and  $Y_1, \dots, Y_n \sim N(0, 1)$ . Define the two-sample studentized  $t$  statistic as

$$T = \frac{\bar{X} - \bar{Y}}{\sqrt{S_x^2/m + S_y^2/n}},$$

where  $S_x^2 = \sum(X_i - \bar{X})^2/(m-1)$ ,  $S_y^2 = \sum(Y_j - \bar{Y})^2/(n-1)$ . It is known that when  $m = n$ ,  $T$  follows Student- $t$  distribution with  $(2n-2)$  degrees of freedom. So in the following, we shall apply the saddlepoint approximation (2) for  $P(T \geq a_1)$  and compare the results with the exact probabilities.

To apply Theorem 1, we need to find out a 1-1 transformation and also the cumulant generating functions. For the first part, we choose  $d = 4$  and the transformation to be

$$\begin{cases} a_1 = (x_1 - y_1)/\sqrt{(x_2 - x_1^2)/(m-1) + (y_2 - y_1^2)/(n-1)} \\ a_2 = x_2 - x_1^2 \\ a_3 = y_1 \\ a_4 = y_2 - y_1^2. \end{cases}$$

Its inverse transformation is

$$\begin{cases} x_1 = a_3 + a_1\sqrt{a_2/(m-1) + a_4/(n-1)} \\ x_2 = a_2 + x_1^2 \\ y_1 = a_3 \\ y_2 = a_4 + a_3^2. \end{cases}$$

It is easy to see that the Jacobian of this transformation is  $\sqrt{a_2/(m-1) + a_4/(n-1)}$ . The cumulant generating functions of both distributions are

$$\begin{aligned} K_x(\theta_1, \theta_2) &= \frac{\theta_1^2}{2(1-2\theta_2)} - \log\left(\frac{1-2\theta_2}{2}\right), \\ K_y(\tau_1, \tau_2) &= \frac{\tau_1^2}{2(1-2\tau_2)} - \log\left(\frac{1-2\tau_2}{2}\right). \end{aligned}$$

Some numerical results are presented in Table 1 to illustrate the performance of the saddlepoint approximations. For the case  $m = n$ , we know that  $T$  follows Student- $t$  distribution with  $(2n-2)$  degrees of freedom, so the exact percentiles can be found from the  $t$  table. To calculate approximate percentiles by saddlepoint approximations, we can invert  $P(A_1 \geq a_1) = \alpha$  for various values of  $\alpha$ , using saddlepoint approximations (2) in Theorem 1. For comparison purposes, we also include percentiles obtained from the Edgeworth expansions for  $T$ , which have been derived in Hall and Martin (1988). It takes the following form,

$$P(T \leq x) = \Phi(x) + p_1(x)\phi(x) + p_2(x)\phi(x) + O(n^{-3/2}), \quad (3)$$

where

$$\begin{aligned} p_1(x) &= \frac{1}{6}\gamma(2x^2 + 1), \\ p_2(x) &= x\left[\frac{1}{12}\kappa(x^2 - 3) - \frac{1}{18}\gamma^2(x^4 + 2x^2 - 3) - \frac{1}{4}(ax^2 + 3b)\right], \end{aligned}$$

where  $\sigma_x^2$ ,  $\gamma_x$ ,  $\kappa_x$  are the variance, skewness and kurtosis for  $X$  with similar quantities defined for  $Y$ , and  $\sigma_{pool}^2 = \sigma_x^2/m + \sigma_y^2/n$  and

$$\begin{aligned} \gamma &= \left(\frac{\gamma_x}{m^2} - \frac{\gamma_y}{n^2}\right) / \sigma_{pool}^3, & \kappa &= \left(\frac{\kappa_x}{m^3} + \frac{\kappa_y}{n^3}\right) / \sigma_{pool}^4, \\ a &= \left(\frac{\sigma_x^4}{m^3} + \frac{\sigma_y^4}{n^3}\right) / \sigma_{pool}^3, & b &= 2\sigma_x^2\sigma_y^2\left(\frac{m+n}{m^2n^2}\right) / \sigma_{pool}^3. \end{aligned}$$

As can be seen from Table 1, the percentiles corresponding to saddlepoint approximations are extremely close to the "exact" ones, for both  $m = n = 5$  and  $m = n = 20$  cases. On the other hand, the percentiles

obtained from the Edgeworth expansions perform very poorly by comparison, particularly for the case small sample size  $m = n = 5$ .

Table 1 Comparisons of the percentiles of  $T$  under normal distributions

Probability	$m = n = 5$			$m = n = 20$		
	"Exact"	Saddlepoint	Edgeworth	"Exact"	Saddlepoint	Edgeworth
0.60	0.262	0.262	0.403	0.255	0.255	0.291
0.70	0.546	0.541	0.678	0.528	0.529	0.563
0.80	0.889	0.881	1.004	0.851	0.852	0.883
0.90	1.397	1.377	1.460	1.304	1.303	1.329
0.95	1.860	1.843	1.841	1.686	1.685	1.698
0.99	2.897	2.869	2.558	2.429	2.428	2.394
0.995	3.355	3.323	2.821	2.712	2.711	2.650
0.999	4.501	4.457	3.360	3.319	3.319	3.178
0.9995	5.041	4.993	3.569	3.566	3.565	3.383
0.9999	6.442	6.379	4.013	4.116	4.114	4.385

#### §4. Application to the bootstrap

Suppose we have two independent data sets  $\mathcal{X} = \{X_1, \dots, X_m\}$  and  $\mathcal{Y} = \{Y_1, \dots, Y_n\}$  from two unknown distributions  $F$  and  $G$ , respectively. Define  $T$  as in Example 2 except that  $(m-1)$  and  $(n-1)$  in the definitions of  $S_x^2$  and  $S_y^2$  will now be replaced by  $m$  and  $n$ . Then the bootstrap estimate of  $\mathbf{P}(T \geq a_1)$  is  $\mathbf{P}^*(T^* \geq a_1)$ , where

$$T^* = \frac{(\bar{X}^* - \bar{Y}^*) - (\bar{X} - \bar{Y})}{\sqrt{S_x^{*2}/m + S_y^{*2}/n}}$$

and  $\mathcal{X}^* = \{X_1^*, \dots, X_m^*\}$  and  $\mathcal{Y}^* = \{Y_1^*, \dots, Y_n^*\}$  are drawn with replacement from  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, and  $\bar{X}^*$ ,  $\bar{Y}^*$  are the bootstrap means,  $S_x^{*2}$  and  $S_y^{*2}$  are the bootstrap variances of  $\mathcal{X}^*$  and  $\mathcal{Y}^*$ , respectively. Also  $\mathbf{P}^*$  indicates the probability conditional on  $\mathcal{X}$  and  $\mathcal{Y}$ .

Generally,  $\mathbf{P}^*(T^* \geq a_1)$  can be estimated by Monte Carlo simulation. But now we apply the saddlepoint approximation (2) to avoid the intensive simulation, in a similar way to Davison and Hinkley (1988) and Daniels and Young (1991). Other related work includes DiCiccio and Martin (1991), Wang (1993), Jing and Robinson (1994) and so on. To apply the saddlepoint approximation (2), we choose the same transformation as in Example 2, except that we now replace  $m-1$  and  $n-1$  there to  $m$  and  $n$ , respectively. The empirical cumulant generating functions for  $\mathcal{X}$  and  $\mathcal{Y}$  are

$$\begin{aligned} \hat{K}_x(\theta_1, \theta_2) &= \log \left( \frac{1}{m} \sum_{i=1}^m \exp(\theta_1(x_i - \bar{x}) + \theta_2(x_i - \bar{x})^2) \right), \\ \hat{K}_y(\tau_1, \tau_2) &= \log \left( \frac{1}{n} \sum_{j=1}^n \exp(\tau_1(y_j - \bar{y}) + \tau_2(y_j - \bar{y})^2) \right). \end{aligned}$$

For comparison purposes, we generated two random samples of sizes  $m = n = 10$  from two independent normal distributions. The data are

$$\begin{aligned} X_i &: -1.35, -0.81, -0.24, 0.28, -1.55, -0.59, 0.50, 0.82, 0.96, -1.82; \\ Y_j &: -0.61, -0.82, 1.32, -0.27, 0.93, -0.99, -0.11, 1.27, 0.87, 0.04. \end{aligned}$$

We calculate the tail probability  $\mathbf{P}^*(T^* \geq a_1)$  for various values of  $a_1$  by the following three methods: "exact" probability obtained by 100,000 simulations; saddlepoint approximations; and Edgeworth expansions, which is the bootstrap version of (3) obtained by replacing all the population quantities by their sample counterparts. The results are presented in Table 2.

As can be seen from the table, the saddlepoint approximations give very accurate results, in particular, at the tails of the distributions, as one might have expected. On the other hand, Edgeworth expansions perform rather poorly at the tails. It is interesting to note, though, that Edgeworth expansions give better approximations at the center of the distribution here.

Table 2 Comparisons of the bootstrap probabilities  $P^*(T^* \geq a_1)$   
( $m = n = 10$ )

$a_1$	$P^*(T^* \geq a_1)$		
	"Exact"	Saddlepoint	Edgeworth
3.267	0.01	0.012	0.006
2.713	0.025	0.020	0.016
2.294	0.05	0.051	0.044
1.835	0.10	0.114	0.112
1.313	0.20	0.185	0.188
0.964	0.30	0.265	0.274
0.670	0.40	0.353	0.368
0.400	0.50	0.451	0.471
0.132	0.60	0.559	0.584
-0.1529	0.70	0.679	0.707
-0.4842	0.80	0.819	0.853
-0.9655	0.90	0.901	0.933
-1.3813	0.95	0.946	0.972
-1.7493	0.975	0.976	0.993

### References

- [1] Daniels, H.E., Uniform approximations for tail probabilities, *International Statistical Review*, **55**(1987), 37-48.
- [2] Daniels, H.E. and Young, G.A., Saddlepoint approximation for the studentized mean, with an application to the bootstrap, *Biometrika*, **78**(1991), 169-179.
- [3] Davison, A.C. and Hinkley, D.V., Saddlepoint approximations in resampling methods, *Biometrika*, **75**(1988), 417-431.
- [4] DiCiccio, T.J. and Martin, M.A., Approximations of marginal tail probabilities for a class of smooth functions with applications to Bayesian and conditional inference, *Biometrika*, **78**(1991), 891-902.
- [5] Hall and Martin, On the bootstrap and two-sample problems, *Austral. J. Statist.*, **30A**(1988), 179-192.
- [6] Jing, B.Y. and Robinson, J., Saddlepoint approximations for marginal and conditional probabilities of transformed variables, *Annals of Statistics*, **22**(3)(1994), 1115-1132.
- [7] Skovgaard, Ib.M., Saddlepoint expansions for conditional distributions, *J. Appl. Probab.*, **24**(1987), 875-887.
- [8] Wang, S., Saddlepoint approximations in conditional inference, *J. Appl. Probab.*, **30**(1993), 397-404.

## 几个总体均值函数的鞍点逼近

邢炳义

(香港科学技术大学, 香港)

李泽慧

(兰州大学, 兰州, 730000)

本文发展了几个总体的均值向量的函数的尾部概率的鞍点逼近方法, 并应用它于 Bootstrap 估计中, 以代替 Monte Carlo 模拟, 一些数值例子说明了其近似精度.

关键词: Bootstrap, Edgeworth 展开, 鞍点逼近.

学科分类号: 0212.