

Testing Uniformity Based on Regression and EDF*

Namhyun Kim¹⁾

Abstract

Some tests of the goodness of fit of the uniform distribution between 0 and 1 are presented. The powers of the tests under certain alternatives are examined. As a result, the statistic based on the difference between the order statistics and the modal value of them gives good powers. We also give modifications of the statistic without using the extensive tables of the critical points.

Keywords: Goodness of fit; order statistics; uniform distribution.

1. 서론

표본 X_1, \dots, X_n 에 대해서 귀무가설

H_0 : 확률표본이 분포함수 $F(x)$ 에서 추출되었음.'

을 검정하는 고전적인 일변량 적합도 검정 문제를 고려한다. 이를 위한 방법은 크게 χ^2 -검정법, 경험적 누적분포함수 (empirical cumulative distribution function, EDF)나 확률 플롯 (probability plot) 등의 그래프를 이용하는 방법, 경험적 분포함수와 모분포함수 (population distribution function)의 차이를 보는 방법, 확률 플롯의 직선의 정도를 살펴보는 회귀 또는 상관계수 검정법, 적률 (moment)을 이용하는 방법 등으로 나눌 수 있다 (D'Agostino와 Stephens, 1986).

그래프를 이용한 방법은 자료의 전체적인 형태를 파악하는데 도움을 주고 사용하기 편리하다는 장점이 있다. 그러나 이 방법은 적합도의 정도를 수치화하는 다른 구체적인 검정법의 보조적인 수단으로 이용되는 것이 적절하다. 위에서 언급한 그래프 중 EDF는 검정하고자 하는 분포의 누적분포함수와의 차이를 눈으로 판단한다. 그러나 일반적으로 이들은 곡선이므로 곡선과 곡선의 차이를 보아야 하는 어려움이 있다. 따라서 이보다는 직선으로부터 벗어난 정도를 보는 것이 더욱 편리할 것이다. 이를 위하여, 가정된 분포가 사실이라면 EDF가 거의 직선이 되도록 척도를 조정해 놓은 것이 확률플롯이라고 할 수 있다. 따라서 확률플롯에서는 그 직선의 정도를 봄으로써 분포에 대한 적합도

* This work was supported by 2006 Hongik University Research Fund. This work was done while the author was on sabbatical (Sep. 2006–Aug. 2007).

1) Professor, Department of Science, Hongik University, 72-1 Sangsu-Dong, Mapo-Gu, Seoul 121-791, Korea.

E-mail : nhkim@hongik.ac.kr

검정을 수행할 수 있고 이와 관련된 검정방법이 회귀 또는 상관계수 검정법이다. 따라서 회귀 또는 상관계수 검정법은 대부분 직관적으로 이해하기 쉬운 편이나 다른 검정법에 비해서 상대적으로 통계적 성질에 대한 규명이 부족한 편이다.

본 논문에서는 U_1, \dots, U_n 이 균일분포 $U(0, 1)$ 에서의 확률표본인지, 즉,

$$H_0 : F(x) := U(0, 1)$$

인지를 검정하는 문제를 다룬다. 여기서 $:=$ 는 왼쪽의 분포함수가 오른쪽 분포의 분포함수와 동일함을 의미한다.

여기에서는 앞에서 언급한 여러 가지 검정방법 중 특히 $U(0, 1)$ 에 대한 회귀 또는 상관계수 검정에 대해서 고려하고자 한다. 적합도 검정을 위한 확률플롯에서는 주로 순서통계량 $X_{(i)}$ 과 적절한 i 의 함수 T_i , 예를 들면 $F^{-1}(i/(n+1))$ 에 대한 그래프를 작성한다. 앞에서 언급한 바와 같이, 만일 가설이 사실이면 이 그래프는 거의 직선을 나타내게 되고, 이 직선의 정도를 보는 것이 회귀검정, 특히 $\mathbf{X} = (X_{(1)}, \dots, X_{(n)})$ 와 $\mathbf{T} = (T_1, \dots, T_n)$ 의 상관계수를 보는 것이 상관계수 검정이다. 회귀검정에 대해서 좀 더 살펴보자.

α 는 위치모수, β 는 척도모수로 $F(x) = F_0((x - \alpha)/\beta)$ 이고 W_1, \dots, W_n 은 $F_0(x)$ 에서의 표본, $W_{(1)}, \dots, W_{(n)}$ 은 순서통계량이라고 하면,

$$X_i = \alpha + \beta W_i, \quad i = 1, \dots, n$$

이므로 $m_i = E(W_{(i)})$, E 는 기댓값을 나타낼 때

$$E(X_{(i)}) = \alpha + \beta m_i \quad (1.1)$$

이다. 따라서 $X_{(i)}$ vs $m_{(i)}$ 의 플롯은 \mathbf{X} 가 $F(x)$ 에서의 표본이면 균사적으로 절편 α , 기울기 β 인 직선을 따른다. 만일 $m_{(i)}$ 의 계산이 용이하지 않은 경우에는 다른 i 의 함수 T_i 를 사용하게 되고 식 (1.1)의 모형은

$$X_{(i)} = \alpha + \beta T_{(i)} + \epsilon_i$$

로 대체된다. 오차항 ϵ_i 는 $T_i = m_i$ 일 때는 $E(\epsilon_i) = 0$ 이다. 자주 사용되는 T_i 는 $T_i = F^{-1}(i/(n+1))$ 이다. 이에 대한 자세한 사항은 D'Agostino와 Stephens (1986, Ch. 5, 5.1–5.4)를 참고로 하였다.

귀무가설이 $H_0 : F(x) := U(0, 1)$ 인 경우에는 $\alpha = 0$, $\beta = 1$ 로 기지이므로 식 (1.1)은

$$E(U_{(i)}) = m_i \quad (1.2)$$

이 된다. 여기서 $m_i = i/(n+1)$ 이다. 이 경우 $\mathbf{U} = (U_{(1)}, \dots, U_{(n)})$ 과 $\mathbf{m} = (m_1, \dots, m_n)$ 의 상관계수

$$R(\mathbf{U}, \mathbf{m}) = \frac{\sum_{i=1}^n (U_{(i)} - \bar{U})(m_i - \bar{m})}{\left(\sum_{i=1}^n (U_{(i)} - \bar{U})^2 \sum (m_i - \bar{m})^2 \right)^{1/2}} \quad (1.3)$$

만으로는 귀무가설 H_0 를 검정하기에 적절하지 않다. 왜냐하면 \mathbf{U} 가 구간 $[0, 1]$ 이 아닌 임의의 구간 $[a, b]$ 에서 균일분포를 따르더라도 $R(\mathbf{U}, \mathbf{m})$ 은 거의 1에 가깝기 때문이다. 여기서 $\bar{U} = \sum_i U_i/n$, $\bar{m} = \sum_i m_i/n$ 이다. 일반적으로 상관계수는 두 변수에 대해서 정의한다. 여기서 \mathbf{m} 은 상수이고 변수는 아니라 편의상 유사한 정의를 사용하도록 하자. 회귀검정에서도 마찬가지 이유로 직선의 정도만을 보는 것은 충분하지 않다. 따라서 이를 적절히 보완할 수 있는 방법이 필요하다. 2절에서는 이에 대한 몇 가지 방법을 고려하고 이들의 검정력을 비교한다.

2. 통계량 및 검정력 비교

$U(0, 1)$ 의 검정에 상관계수를 이용하기 위하여 고려할 수 있는 가장 간단한 방법은 $(0, 0)$ 과 $(1, 1)$ 을 표본에 추가하여 상관계수를 살펴보는 것이다. 또한 회귀검정에서 많이 고려되는 방법은 모형 (1.2)의 잔차에 해당하는 $r_i = U_{(i)} - m_i$ 를 이용하는 것이다. r_i 를 이용한 통계량들은 여러 가지 제안되어왔다.

$$C^+ = \max_i r_i, \quad C^- = \max_i (-r_i), \quad C = \max(C^+, C^-), \quad K = C^+ + C^-$$

와 같은 통계량이 Durbin (1969), Brunk (1962), Stephens (1969, 1970) 등에 의해 연구되었다. 잔차 또는 잔차의 함수의 합에 기반을 둔 통계량도 당연히 고려되었고 이러한 통계량으로 $T_1 = \sum |r_i|/n$, $T_2 = \sum r_i^2/n$ 과 같은 통계량이 Hegazy와 Green (1975)에 의해 연구되었다. 이들은 r_i 대신

$$v_i = U_{(i)} - \frac{i-1}{n-1} \quad (2.1)$$

을 이용한 통계량 $T'_1 = \sum |v_i|/n$, $T'_2 = \sum v_i^2/n$ 도 고려하였고 주로 ' $H_0 : F(x) := U(a, b)$, a, b 는 미지'인 경우를 중심으로 통계량의 검정력을 비교하였다. 식 (2.1)의 v_i 는 $U_{(i)}$ 가 중앙값 (median)을 제외하고는 치우친 (skewed) 분포를 가지므로 $U_{(i)}$ 분포의 기대값인 $m_i = i/(n+1)$ 대신 최빈값 (mode) $\xi_i = (i-1)/(n-1)$ 을 고려한 것이다.

Hegazy와 Green (1975)은 주로 복합귀무가설에서 위의 통계량을 살펴보았으나, 본 논문에서는 이들의 통계량을 단순귀무가설 $H_0 : F(x) := U(0, 1)$ 을 위한 통계량으로서 좀 더 자세히 살펴보려 한다.

EDF 통계량은 경험적 분포함수 $F_n(x) = (1/n) \sum_{i=1}^n I(X_i \leq x)$ (I 는 표시함수 (indicator function))와 H_0 에서의 모분포함수 $F(x)$ 와의 차이를 보는 것으로 차이의 최대를 보는 Kolmogorov-Smirnov 통계량 $D = \sup_x |F_n(x) - F(x)|$, Cramér-von Mises 통계량, 이의 변형된 형태인 Anderson-Darling (1952) 통계량 등이 대표적이다.

위의 $T_2 = \sum r_i^2/n$, $T'_2 = \sum v_i^2/n$ 은 EDF 통계량인 Cramér-von Mises W^2 -통계량

$$W^2 = \sum \left(U_{(i)} - \frac{2i-1}{2n} \right)^2 + \frac{1}{12n}$$

과 형태가 유사하다. 균일분포의 경우 누적분포함수가 $F(x) = x$ 이므로 잔차제곱합으로 구성된 T_2 나 $F_n(x)$ 와 $F(x)$ 의 차이의 제곱의 적분에서 구한 W^2 통계량은 유사한 형태를 갖는다.

위의 세 통계량 T_2 , T'_2 , W^2 은 모두

$$\sum (U_{(i)} - p_i(c))^2, \quad p_i(c) = \frac{i - c}{n - 2c + 1}$$

의 형태로 각각 $c = 0$ (T_2), $c = 1$ (T'_2), $c = 1/2$ (W^2)에 해당한다. 결국 위 통계량들의 비교는 $U(0, 1)$ 의 경우 확률 플롯에서 어떤 플로팅 위치 (plotting position)가 합리적인지를 살펴보는 문제로 귀결된다. 플로팅 위치는 상당히 오랫동안 연구되어 온 주제 중 하나이다. 이에 대해서는 Barnett (1975), Harter (1984), Looney와 Gullledge (1985) 등을 참고한다.

본 논문에서는 $U(0, 1)$ 의 적합도 검정 통계량으로

$$G_p(c) = \sum |U_{(i)} - p_i(c)|^p, \quad p = 1, 2, \quad c = 0, 1/2, 1, \quad (p = 3, c = 1), \quad (p = 4, c = 1)$$

을 고려한다. $G_p(c)$ 는 $T_1, T'_1, T_2, T'_2, W^2$ 의 모든 통계량을 포함하므로 본 논문에서는 통계량의 형태와 플로팅 위치에 따른 검정력의 변화를 살펴볼 것이다.

또한 $(0, 0)$, $(1, 1)$ 을 자료에 포함한 상관계수를 고려한다. 비교를 위하여 식 (1.3)의 상관계수와 가장 널리 알려진 EDF 통계량인 Kolmogorov-Smirnov 통계량 D 도 포함하였다. 상관계수 또한 플로팅 위치에 의존한다. 따라서 $c = 0$, $c = 1$ 일 때의 두 가지 플로팅 위치 $p_i(c)$ 를 고려하였다. 즉, $\mathbf{U}^+ = (0, U_{(1)}, \dots, U_{(n)}, 1)$, $\mathbf{m}^+ = (0, m_1, \dots, m_n, 1)$, $\boldsymbol{\xi}^+ = (0, \xi_1, \dots, \xi_n, 1)$ 라고 할 때 $R(\mathbf{U}, \mathbf{m})$, $R(\mathbf{U}^+, \mathbf{m}^+)$, $R(\mathbf{U}, \boldsymbol{\xi})$, $R(\mathbf{U}, \boldsymbol{\xi}^+)$ 을 비교한다. 상관계수는 모두 식 (1.3)과 유사하게 정의된다.

이들 통계량의 검정력을 몇 가지 대립가설에서 비교한다. 대립가설은 Stephens (1974)에서 균일분포의 검정을 위하여 고려한 $(0, 1)$ 에서의 몇 가지 분포를 택한다. 이들의 분포함수는 다음과 같다.

$$\begin{aligned} A : \quad F(z) &= 1 - (1 - z)^k, & 0 \leq z \leq 1, \\ B : \quad F(z) &= 2^{k-1}z^k, & 0 \leq z \leq 0.5; \\ &F(z) = 1 - 2^{k-1}(1 - z)^k, & 0.5 \leq z \leq 1, \\ C : \quad F(z) &= 0.5 - 2^{k-1}(0.5 - z)^k, & 0 \leq z \leq 0.5; \\ &F(z) = 0.5 + 2^{k-1}(z - 0.5)^k, & 0.5 \leq z \leq 1. \end{aligned}$$

분포 A 는 $U(0, 1)$ 보다 0쪽의 확률이 높은 분포이고, B, C 는 대칭이며 각각 0.5 또는 0과 1쪽의 확률이 높은 분포이다. 그 정도는 물론 모두 k 에 따라 결정된다. 본 논문에서는 Stephens (1974)과 마찬가지로 분포 A , $k = 1.5, 2$ 분포 B , $k = 1.5, 2, 3$, 분포 C , $k = 1.5, 2$ 의 대립가설과 $n = 10, 20, 40$ 을 고려하였고 여기에 균일분포 $U(0.25, 0.75)$, $U(0, 0.75)$, $U(0.25, 1)$ 또한 추가하였다. 이는 균일분포이나 범위가 적절하지 않을 경우 상관계수의 이용이 불합리함을 확인하기 위함이다.

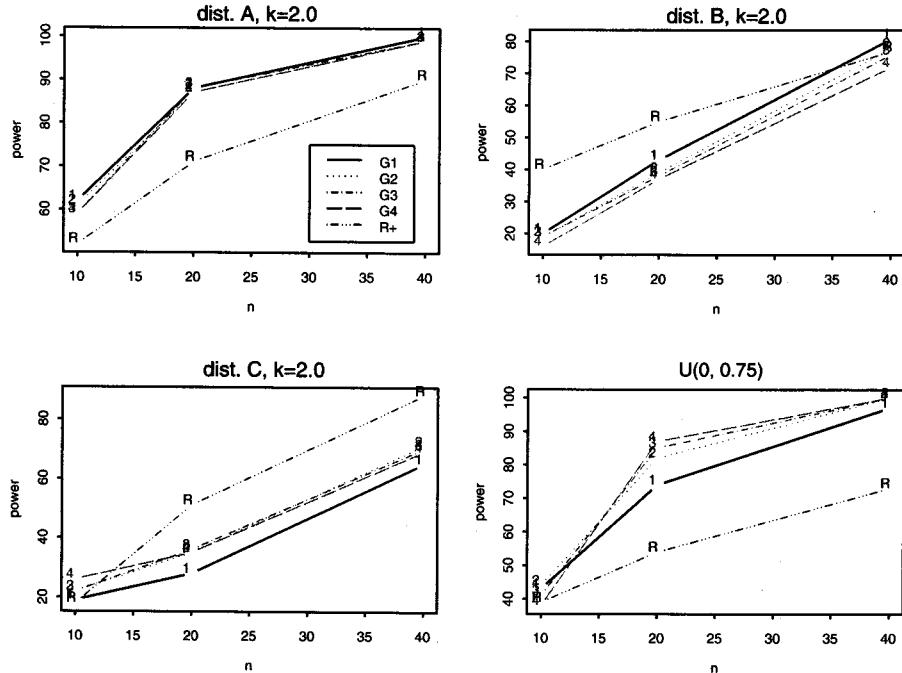
표 2.1: 검정력 비교 ($\alpha = 0.1$)

대립가설	n	D	$G_1(1)$	$G_1(\frac{1}{2})$	$G_1(0)$	$G_2(1)$	W^2	$G_2(0)$	$G_3(1)$	$G_4(1)$	$R(1)$	$R(0)$	$R^+(1)$	$R^+(0)$
$A, k = 1.5$	10	26	28	29	26	28	26	25	25	28	13	13	27	20
	20	38	47	44	41	45	44	42	42	45	22	17	35	35
	40	64	70	72	72	70	70	71	73	72	40	33	47	45
$A, k = 2.0$	10	52	62	61	60	61	59	58	59	59	21	20	52	43
	20	81	88	87	86	87	86	85	88	87	43	39	71	69
	40	98	100	99	99	99	99	100	99	99	71	68	90	90
$B, k = 1.5$	10	7	10	8	5	11	7	4	13	10	11	11	20	12
	20	14	15	12	7	15	11	9	15	16	13	10	26	20
	40	22	29	22	16	27	23	18	27	27	23	19	35	28
$B, k = 2.0$	10	12	20	10	3	19	9	4	19	16	13	12	40	23
	20	24	43	32	15	39	29	19	38	37	25	19	55	47
	40	58	81	73	61	78	72	61	76	72	45	44	77	74
$B, k = 3.0$	10	26	56	29	8	49	23	9	44	40	21	18	80	57
	20	66	92	85	65	88	80	66	86	85	41	41	95	93
	40	97	100	100	99	100	100	100	100	100	74	74	100	100
$C, k = 1.5$	10	18	14	15	20	17	15	19	16	17	20	21	13	20
	20	26	17	21	25	18	23	23	17	20	29	22	22	31
	40	34	23	29	31	28	32	40	28	31	46	42	38	46
$C, k = 2.0$	10	32	19	26	31	22	27	33	22	26	33	34	18	35
	20	48	28	41	50	35	44	53	36	35	61	51	51	61
	40	73	65	77	81	71	76	83	70	69	92	89	88	91
$U(0.25, 0.75)$	10	17	59	32	7	52	20	7	41	38	11	9	92	66
	20	97	99	95	78	99	95	83	100	100	11	10	100	100
	40	100	100	100	100	100	100	100	100	100	9	10	100	100
$U(0, 0.75)$	10	34	43	40	34	44	33	30	41	38	8	10	39	25
	20	91	74	68	64	82	76	71	85	87	12	10	54	45
	40	100	97	97	95	100	100	100	100	100	10	10	73	60
$U(0.25, 1)$	10	32	40	43	35	43	35	31	43	41	11	10	40	27
	20	92	74	70	64	80	74	69	82	84	10	10	58	46
	40	100	98	97	97	100	100	99	100	100	11	10	70	61

표 2.1과 그림 2.1에는 몇 가지 대립가설에서의 검정력을 비교결과를 제시하였다. 유의수준은 $\alpha = 0.1$ 이다. 표의 숫자는 여러가지 통계량에서 유의한 결과를 보인 표본의 비율을 퍼센트로 표시한 것이다. 각 표본은 S-plus 6.1을 이용하여 추출하였고 표본 수 $N = 1,000$ 을 이용하였다. 각 통계량의 기각값은 모의실험 (simulation)을 통해 구하였으며 이때 표본수는 $N = 10,000$ 을 이용하였다.

검정력을 비교한 결과 다음과 같은 사실을 볼 수 있다. 표 2.1에서는 편의상 $R(\mathbf{X}, \mathbf{m})$, $R(\mathbf{X}, \boldsymbol{\xi})$ 을 $R(0), R(1)$ 으로 $R(\mathbf{X}^+, \mathbf{m}^+)$, $R(\mathbf{X}^+, \boldsymbol{\xi}^+)$ 을 $R^+(0), R^+(1)$ 으로 표시하였다.

첫째, 통계량 $G_1(c), G_2(c)$ 의 결과에서 플로팅 위치에 따른 검정력을 살펴보자. 그러면 최빈값 ξ_i 에 대한 플롯의 경우 ($c = 1$) 전반적으로 분포 B 에서 높은 검정력을 보이고 기대값 m_i 에 대한 플롯 ($c = 0$)은 분포 C 에서 높은 검정력을 보인다. 분포 A 에서는 플로

그림 2.1: 검정력 비교 ($\alpha = 0.1$)

팅 위치에 따른 검정력의 변화가 크게 눈에 띄지 않는다. 전반적으로는 최빈값 ξ_1 에 대한 플로팅($c = 1$)이 대부분의 대립가설에서 높은 검정력을 보여준다.

둘째, $G_p(1)$, $p = 1, 2, 3, 4$ 를 보면 분포 B에서 p 가 커짐에 따라 검정력이 감소하는 경향이 약간 보이기는 하나 p 에 따른 검정력의 변화는 크게 눈에 띄는 바가 없다. 따라서 $p = 1$ 또는 $p = 2$ 를 사용함이 적절해 보인다 (그림 2.1 참조).

셋째, Kolmogorov-Smirnov 통계량 D는 분포 A와 B에서는 $G_p(c)$ 통계량보다는 전반적으로 검정력이 낮으나 분포 C에서는 유사한 검정력을 나타낸다. 범위가 $(0, 1)$ 이 아닌 균일분포에서는 표본크기가 작을때 ($n = 10$)를 제외하고는 검정력이 우수하다.

Kolmogorov-Smirnov 통계량은 Stephens (1974)에서도 다른 EDF 통계량과 검정력을 비교하였다. 이를 토대로 Stephens (1974)은 Kolmogorov-Smirnov 통계량 D보다는 Cramér-von Mises 통계량 W^2 이나 Anderson-Darling 통계량 A^2 이 유용하다고 결론내리고 있다. 여기서도 이와 유사한 검정력 결과를 보여 주고 있다.

넷째, 자료에 $(0, 0)$, $(1, 1)$ 을 추가한 R^+ 는 예상과 같이 대부분의 대립가설에서 검정력을 높이는 효과가 있다. R^+ 의 플로팅 위치에 따른 경향은 G_p -통계량의 경우와 유사하고 마찬가지로 $R^+(1)$ 의 경우 $R^+(0)$ 보다 좀 더 우수한 검정력을 보여준다. 또한 R^+ 의

검정력 증가는 일반적으로 표본크기가 증가하면서 그 정도가 약해지고 있다. 표본크기에 따라 $(0, 0)$, $(1, 1)$ 의 가중값을 변화함으로써 검정력의 향상을 기대해 볼 수도 있을 것이다. 이에 대해서는 좀 더 세심한 연구가 필요하다. R 의 경우 위치, 척도에 무관한 통계량이므로 범위가 다른 균일분포의 대립가설에서 검정력은 유의수준과 유사하게 나타난다.

다섯째, $G_p(1)$ 과 $R^+(1)$ 을 비교하면 $R^+(1)$ 은 대칭인 대립가설에서는 검정력이 매우 우수하나 비대칭 분포에서는 검정력이 매우 떨어진다 (그림 2.1 참조). 따라서 총괄검정(omnibus test)의 경우는 이용하기에 무리가 있으나 특수한 대립가설에서는 좋은 검정통계량이 될 수 있다.

위의 결과를 토대로 볼 때 $U(0, 1)$ 의 검정을 위해서는 $G_1(1)$ 또는 $G_2(1)$ 이 적절해 보인다. 앞에서 언급한 바와 같이 $G_2(1)$ 은 플로팅 위치를 제외하면 W^2 과 유사한 통계량으로 $G_1(1)$ 보다는 이론의 전개가 좀 더 용이할 것이다. 즉, $U(0, 1)$ 의 검정에서는 순서통계량의 분포의 최빈값 ξ_i 를 플로팅 위치로 이용하는 것이 좀 더 합리적이며 따라서 W^2 의 변형된 형태인 $G_2(1)$ 을 제안한다. 3절에서는 $G_2(1)$ 의 구체적인 이용방법에 대해 살펴본다.

3. 수정된 통계량 및 구체적인 검정의 적용방법

2절에서 제안한 $G_2(1)$ -통계량은 앞에서 언급한 바와 같이 W^2 -통계량과 매우 밀접한 관계에 있다. W^2 -통계량의 귀무가설에서의 극한분포는, Z_1, Z_2, \dots 이 $N(0, 1)$ 을 따르는 독립인 확률변수일 때

$$\sum_{j=1}^{\infty} \frac{Z_j^2}{(j\pi)^2} \quad (3.1)$$

과 같이 독립인 확률변수의 무한합으로 표현된다 (Shorack과 Wellner, 1986). 이와 관련된 내용은 de Wet과 Venter (1973)에 의해서도 연구되었다. 식 (3.1)의 분포의 분위수는 Anderson과 Darling (1952) 또는 Shorack과 Wellner (1986, 3장 표4 또는 5장 표1) 등에서 찾을 수 있다. 그러나 유한한 표본크기 n 에 대해서는 정확한 값을 구하기 힘들기 때문에 이론적인 결과보다는 몬테칼로 방법 (Monte Carlo method) 등을 통하여 기각값을 구하는 것이 일반적이다. Hegazy와 Green (1975)은 2절의 T_2, T'_2 의 적률을 구하고 이를 이용하여 T_2 와 T'_2 의 귀무가설에서의 분포를 피어슨 분포로 근사하여 기각값을 구하였으나 이도 역시 주어진 표본크기 n 에 대하여 각각 구해야 하는 단점이 있다. 이러한 단점을 보완하기 위하여 Stephens (1970)은 W^2 을 포함한 EDF 통계량들의 수정된 형태를 제안하였다. 이 절에서는 $G_2(1)$ -통계량에 대해서 이와 유사한 형태를 찾아보려 한다.

Stephens (1970)은 경험적인(empirical) 방법으로 W^2 의 수정된 형태인

$$W^{2*} = (W^2 - 0.4/n + 0.6/n^2)(1 + 1/n) \quad (3.2)$$

을 얻고 이 값이 주어진 유의수준에서 근사분포의 기각값보다 크면 H_0 를 기각하는 방법을 제안하였다. 2절의 $G_2(1)$ 에 대해서 식 (3.2)과 유사한 표현을 얻기위해서 $G_2(1)$ 과

W^2 사이의 관계를 살펴보자.

$$G_2(1) = \sum_{i=1}^n \left(U_{(i)} - \frac{i-1}{n-1} \right) = W^2 + D$$

이고 여기서 D 는

$$D = \sum_{i=1}^n \left(\frac{i-0.5}{n} - \frac{i-1}{n-1} \right) \left(U_{(i)} - \frac{i-1}{n-1} \right) - \frac{1}{12n}$$

이다. D 는 확률변수에 의존하므로 이의 기대값을 계산하면

$$\begin{aligned} E(D) &= \sum_{i=1}^n \left(\frac{i-0.5}{n} - \frac{i-1}{n-1} \right) \left(\frac{i}{n+1} - \frac{i-1}{n-1} \right) - \frac{1}{12n} \\ &= \frac{2n+1}{3(n-1)^2} - \frac{n+1}{2(n-1)^2} - \frac{1}{12n} \end{aligned}$$

이고 이를 반영하면 수정된 $G_2(1)$ 으로

$$G_2^M(1) = \left(G_2(1) - \frac{19}{60} \frac{1}{n} + \frac{n+1}{2(n-1)^2} - \frac{2n+1}{3(n-1)^2} + \frac{0.6}{n^2} \right) \left(1 + \frac{1}{n} \right) \quad (3.3)$$

을 얻는다.

두번째 방법으로 $E(D)$ 의 각 항이 $O(1/n)$ 이므로 식 (3.2)에서 $1/n$ 항의 계수만을 바꾸어 주어진 유의수준 α 에 좀 더 가까운 경험적인 유의수준 α' 을 주는 계수를 모의실험 (simulation)을 통하여 찾은 결과

$$G_2^*(1) = \left(G_2(1) - \frac{0.75}{n} + \frac{0.6}{n^2} \right) \left(1 + \frac{1}{n} \right) \quad (3.4)$$

을 얻었다. 식 (3.3)과 식 (3.4)를 이용하여 모의실험을 통하여 얻은 유의수준 α' 을 표 3.1에 나타내었다. 그 결과 $G_2^M(1)$ 보다는 $G_2^*(1)$ 가 좀 더 주어진 유의수준 α 에 가까운 값을 갖는다. 이 경우 각 표본은 S-plus 6.1을 이용하여 추출하였고 표본 수 $N = 10,000$ 을 이용하였다.

이 결과를 바탕으로 균일분포 $U(0,1)$ 의 검정을 위하여, 주어진 자료에 대해서 통계량 $G_2(1)$ 그리고 $G_2^*(1)$ 을 계산한 후 $G_2^*(1)$ 이 표 3.2의 기각값을 초과하면 귀무가설을 기각하는 검정방법을 제안한다.

위의 검정법을 D'Agostino와 Stephens (1986, 부록 A)의 Leghorn chick data에 적용해 보자. 이 자료는 생후 21일된 흰 leghorn 병아리 20마리의 무게 X 를 그램으로 측정한 것으로 원자료는 Bliss (1967)에서 발췌한 것이다. 자료가 평균 $\mu = 200$, 표준편차 $\sigma = 35$ 인 정규분포를 따르는지를 검정하고자 한다. 우선 확률적분변환을 적용하여 얻은 Z 를 표 3.3에 표시하였다. 이 값을 이용하여 $G_2(1)$ 과 $G_2^*(1)$ 을 계산하면 $G_2(1) = 0.2$, $G_2^*(1) = 0.172$ 이다. 따라서 표 3.2에 따르면 $p > 0.15$ 이고 귀무가설 H_0 는 기각하지 못한다. 실제로 p 값은 $p \approx 0.33$ 이다.

표 3.1: 수정된 통계량의 유의수준 α 에서 얻은 α' 비교

		α				
통계량	n	0.15	0.1	0.05	0.025	0.01
$G_2^M(1)$	10	0.158	0.110	0.061	0.031	0.013
	20	0.159	0.109	0.055	0.028	0.010
	30	0.159	0.107	0.056	0.029	0.011
	40	0.158	0.108	0.059	0.031	0.012
	50	0.156	0.103	0.053	0.024	0.011
$G_2^*(1)$	10	0.146	0.099	0.050	0.027	0.011
	20	0.146	0.100	0.052	0.029	0.012
	30	0.148	0.100	0.053	0.025	0.010
	40	0.150	0.099	0.052	0.027	0.011
	50	0.149	0.103	0.055	0.027	0.011

표 3.2: 수정된 통계량 $G_2^*(1)$ 과 기각값

$G_2^*(1)$ 의 백분위수	15%	10%	5%	2.5%	1.0%
$G_2^*(1) = \left(G_2(1) - \frac{0.75}{n} + \frac{0.6}{n^2} \right) \left(1 + \frac{1}{n} \right)$	0.284	0.347	0.461	0.581	0.743

표 3.3: Leghorn Chick Data

X	156 214	162 220	168 226	182 230	186 230	190 236	190 236	196 242	202 246	210 270
Z	0.104 0.655	0.139 0.716	0.180 0.771	0.304 0.804	0.345 0.804	0.388 0.848	0.388 0.848	0.455 0.885	0.523 0.906	0.612 0.997

4. 결론 및 토의

본 논문에서는 균일분포 $U(0,1)$ 에 대한 적합도 검정에 대해서 고려하였다. 이 경우는 일반적으로 사용되는 회귀나 상관계수 검정은 유용하지 않다. 따라서 회귀모형에서의 잔차에 기반을 둔 검정통계량을 고려하였고 이는 EDF 통계량인 Cramer-von Mises 통계량과 유사한 형태이다. 모의실험 결과 $U(0,1)$ 에서 이러한 형태의 통계량을 사용할 경우 순서통계량과 순서통계량의 분포의 최빈값과의 차이에 기반을 둔 통계량이 좀 더 우수한 검정력을 보여주었다.

이러한 형태의 통계량은 타분포나 모수가 미지인 경우뿐만 아니라 중도절단자료 (censored data)에 대한 적합도 검정에 대해서도 확장될 수 있다.

참고문헌

- Anderson, T. W. and Darling, D. A. (1952). Asymptotic theory of certain “goodness of fit” criteria based on stochastic processes. *The Annals of Mathematical Statistics*, **23**, 193–212.

- Barnett, V. (1975). Probability plotting methods and order statistics. *Applied Statistics*, **24**, 95–108.
- Bliss, C. I. (1967). *Statistics in Biology: Statistical Methods for Research in the Natural Sciences, Vol 1*. McGraw-Hill, New York.
- Brunk, H. D. (1962). On the range of the difference between hypothetical distribution function and Pyke's modified empirical distribution function. *The Annals of Mathematical Statistics*, **33**, 525–532.
- D'Agostino, R. B. and Stephens, M. A. (1986). *Goodness-of-fit Techniques*. Marcel Dekker, New York.
- de Wet, T. and Venter, J. H. (1973). Asymptotic distributions for quadratic forms with applications to tests of fit. *The Annals of Statistics*, **1**, 380–387.
- Durbin, J. (1969). Tests for serial correlation in regression analysis based on the periodogram of least-squares residuals. *Biometrika*, **56**, 1–15.
- Harter, H. L. (1984). Another look at plotting positions. *Communications in Statistics-Theory and Methods*, **13**, 1613–1633.
- Hegazy, Y. A. S. and Green, J. R. (1975). Some new goodness-of-fit tests using order statistics. *Applied Statistics*, **24**, 299–308.
- Looney, S. W. and Gulledge, T. R. Jr. (1985). Use of the correlation coefficient with normal probability plots. *The American Statistician*, **39**, 75–79.
- Shorack, G. R. and Wellner, J. A. (1986). *Empirical Processes with Applications to Statistics*. John Wiley & Sons , New York.
- Stephens, M. A. (1969). Results from the relation between two statistics of the Kolmogorov-Smirnov type. *The Annals of Mathematical Statistics*, **40**, 1833–1837.
- Stephens, M. A. (1970). Use of the Kolmogorov-Smirnov, Cramér-von Mises and related statistics without extensive tables. *Journal of the Royal Statistical Society, Ser. B*, **32**, 115–122.
- Stephens, M. A. (1974). EDF statistics for goodness of fit and some comparisons. *Journal of the American Statistical Association*, **69**, 730–737.

[Received June 2007, Accepted November 2007]