

다수의 고장 원인을 갖는 기기의 신뢰성 모형화 및 분석*

Reliability Modeling and Analysis for a Unit with Multiple Causes of Failure

백상엽**, 임태진***, 이창훈**

Sang-Yeop Baek**, Tae-Jin Lim***, Chang-Hoon Lie**

Abstract

This paper presents a reliability model and a data-analytic procedure for a repairable unit subject to failures due to multiple non-identifiable causes. We regard a failure cause as a state and assume the life distribution for each cause to be exponential. Then we represent the dependency among the causes by a Markov switching model(MSM) and estimate the transition probabilities and failure rates by maximum likelihood(ML) method. The failure data are incomplete due to masked causes of failures. We propose a specific version of EM(expectation and maximization) algorithm for finding maximum likelihood estimator(MLE) under this situation. We also develop statistical procedures for determining the number of significant states and for testing independency between state transitions. Our model requires only the successive failure times of a unit to perform the statistical analysis. It works well even when the causes of failures are fully masked, which overcomes the major deficiency of competing risk models. It does not require the assumption of stationarity or independency which is essential in mixture models. The stationary probabilities of states can be easily calculated from the transition probabilities estimated in our model, so it covers mixture models in general. The results of simulations show the consistency of estimation and accuracy gradually increasing according to the difference of failure rates and the frequency of transitions among the states.

* 본 논문은 과학기술처에서 주관한 한국 원자력 연구소의 원자력 안전성 향상 연구 지원에 의해 수행되었음.

** 서울대학교 산업공학과 (Dept. of Industrial Eng. Seoul National Univ.)

*** 송실대학교 산업공학과 (Dept. of Industrial Eng. Soong-sil Univ.)

1. 서론

일반적으로 여러 부품으로 구성된 기기에는 다수의 고장 원인이 존재한다. 그리고 이러한 고장 원인에 따라 기기의 고장 시간을 분류해 보면 분포적 특성에 상당한 차이가 있음을 알 수 있다[11]. 따라서 여러 부품으로 구성된 기기에 대한 정확한 신뢰성 모형화를 위해서는 고장 시간만이 아니라 고장 원인에 따른 고장 시간의 분포 특성도 반영할 수 있는 신뢰성 모형이 요구된다고 할 수 있다[13].

기기의 신뢰성 분석에 있어 기기의 고장 시간 및 고장 원인을 모두 고려할 수 있는 기존의 방법론으로는 경쟁적 위험 모형(competitive risk model) 및 분포 혼합 모형(mixture model)을 들 수 있다[13]. 경쟁적 위험 모형에서는 계통 고장에 대해 서로 다른 두가지 이상의 원인이 존재하고 이러한 원인에 의해 특정한 고장 유형(failure mode)이 발생하는 상황에서, 가능한 고장 원인이 직렬 구조를 이루고 있고 그 중에서 가장 먼저 작용하는 고장 원인이 고장을 발생시킨 원인이라는 가정하에 모형화가 수행된다[16]. 분포 혼합 모형은 기기의 고장은 가능한 고장 원인 중에 일정한 확률로 하나의 원인이 작용하여 고장을 발생시킨다는 가정을 사용하고 있다[13]. 이러한 기존의 모형들은 연속되는 고장들의 발생 원인간에는 서로 독립이라는 가정을 사용하고 있다. 그러나 수리 가능한 기기에 있어 고장이 발생하면 기기가 완전히 대체되지 않는 한, 고장을 발생시키지 않은 고장 원인에 대해서는 아무런 조치도 수행하지 않았다고 볼 수 있으므로 이는 고장 원인간의 종

속성이 존재하는 경우라 할 수 있다. 또한 어떤 환경적 요인에 의해 다수의 고장 원인이 영향을 받는 경우에도 고장 원인간의 종속성이 존재하는 경우라 할 수 있다. 이러한 경우는 기기의 대체를 통해서도 종속성이 제거되지 않는 경우이다. 따라서 다수의 고장 원인을 갖는 기기의 신뢰성 분석에 있어 연속되는 고장 원인간의 종속성을 고려하는 것이 보다 현실적인 모형이라 할 수 있다[8].

고장에 관한 자료는 크게 고장 시간에 대한 자료와 고장 원인에 대한 자료로 구분된다. 그러나 현실적인 상황에서 고장 원인에 관한 자료 수집에는 많은 비용이 요구되거나 원인의 식별 자체가 불가능한 경우가 많다[19]. 따라서 고장 원인에 대한 자료는 일반적으로 불완전(incomplete) 자료 형태로 주어진다고 할 수 있다. 이와 같이 고장 원인에 대한 완전한 정보가 주어지지 않은 경우에 고장 자료를 잠재적 자료(masked data)라 한다[6]. 경쟁적 위험 모형을 이용하여 잠재적인 자료를 대상으로 분석을 수행한 연구가 일부 제시되었으나 이를 모든 고장에 대해 고장 자료가 완전히 잠재적(fully masked)인 경우에 응용하는 것은 불가능하다고 할 수 있다[2,5,6,15,16,19]. 또한 분포 혼합 모형을 이용하기 위해서는 일반적으로 정상성(stationarity) 가정을 보장할 수 있는 고장 자료가 요구되는 데 현장에서 이러한 자료를 요구하는 것은 현실적으로 어려운 일이라 할 수 있다[1,8,10].

이와 같이 식별되지 않는 다수의 고장 원인을 갖는 수리 가능한 기기의 수명 분석을 수행하기 위해서는 고장 원인간의 종속성을 반영할 수 있는 새로운 모형과 고장 원인에

대한 자료가 제한적인 상황에서도 분석이 가능한 자료 분석 방법론의 개발이 요구된다고 할 수 있다.

계량 경제학(Econometrics)분야에서 최근 들어 활발히 연구가 진행되고 있는 마코프 전이 모형(이하 MSM)은 고장 원인간의 종속성과 고장 원인별 수명 분포를 동시에 표현할 수 있는 모형이다. MSM은 변동의 추세가 어떤 영향에 의해 한 시점에 급격히 변화하는 경우를 모형화하기 위해 도입된 방법론이다. 급격한 변동은 외부의 환경적 요인의 변화에 의해 발생된다고 볼 수 있으며, 이러한 요인들의 변화를 MSM에서는 상태(state) 또는 체계(regime)의 전이로 모형화한다. 또한 이러한 상태의 변화는 계량화가 불가능하기 때문에 관측이 불가능한 자료로 간주하며, 관측 가능한 자료를 대상으로 회귀 분석을 수행한다. 다만 상태의 전이를 마코프 연쇄(Markov chain)로 모형화하여 상태 전이가 일어날 때마다 회귀 분석 모형의 모수를 보정하는 방법으로 변동의 추이를 분석한다. MSM은 Quandt[18]에 의해 제안된 이후로 연구가 진행되지 않았다가 최근 들어 Hamilton에 의해 구체적인 방법론이 제안된 이후로 연구가 활발히 진행되고 있는 분야이다 [7,9,10].

본 논문에서는 고장 원인을 하나의 가상적인 상태(state)로 간주하여 고장 원인간의 종속성 및 고장 원인별 수명 분포의 특성을 MSM으로 모형화하고 상태 전이 확률 및 고장들의 MLE를 추정한다. 자료의 형태는 고장 원인이 식별되지 않은 일련의 고장 시간으로서 불완전 자료를 대상으로 한다. 불완전 자료하에서의 MLE 추정을 위하여 Demp-

ster 등[4]이 제시한 EM 원리를 이용하여 다항 시간 알고리즘을 제시한다. 또한 유의한 고장 원인의 개수를 결정하는 절차와 상태 전이간의 독립성을 검정하는 절차를 제안하고 이를 이용하여 통합된 분석 절차를 제시한다.

본 논문에서 제안된 모형은 기존의 분포 혼합 모형을 포함하는 보다 유연한(flexible) 모형이라 할 수 있다. 특히 이 모형은 기기의 수명 분포가 시간과 고장 회수의 함수로 표현됨으로써 고장 시간 자료에 대해 정상성 가정을 사용하지 않으며 또한 후속 고장에 의한 수명 분포 변화의 예측이 가능하다는 점에서 기존 모형과 차이가 있음을 알 수 있다.

본 논문의 구성은 다음과 같다. 먼저 2절에서는 논문에 사용되는 기호와 가정을 정리하였으며, 3절에서는 고장 원인에 대한 종속성 모형을 제시하였으며, 4절에서는 모수에 대한 추정 알고리즘을 제시하였으며, 5절에서는 상태수 결정 절차 및 상태 전이에 대한 독립성 검정 절차, 그리고 이를 이용한 통합적 분석 절차를 제시하였다. 마지막으로 6절에서는 다양한 예제를 통해 모형의 타당성을 검증해 보며 실용예를 제시한다.

2. 기호 및 가정

본 논문에서 사용되는 기호를 정리하면 다음과 같다.

N	총 고장 회수
k	상태의 수
y_n	$(n-1)$ 번째와 n 번째 고장간 시간 간격; $n=1, 2, \dots, N$

y_n	n 까지 관측된 자료 벡터 ; $y_n = (y_1, \dots, y_n)$
s_n	n 번째 고장을 지배하는 기기의 상태 변수 ; $n=1, 2, \dots, N$
s_n	상태 변수 벡터 ; $s_n = (s_{1n}, \dots, s_{kn})$
I	지시 변수 ; $s_n = i$ 이면 $I(s_n = i) = 1$, 아니면 $I(s_n = i) = 0$
$f_i(t)$	$s_n = i$ 일 때 y_n 이 따르는 확률 밀도 합 수 ; $i=1, 2, \dots, k$
$F_i(t)$	$s_n = i$ 일 때 y_n 의 분포 함수
λ_i	$f_i(t)$ 의 모수, 고장률
π_i	$f_i(t)$ 를 선택할 확률 ; 안정 상태(steady state) 확률
p_{ij}	상태가 i 에서 j 로 전이할 확률 ; 임의의 n 에 대해 $p(s_n = j s_{n-1} = i) = p_{ij}$, $i=1, \dots, k, j=1, \dots, k$
ρ_i	초기 상태 확률 ; $\rho_i = p(s_1 = i), i=1, 2, \dots, k$
λ	λ_i 들의 벡터 ; $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_k)$
ρ	ρ_i 들의 벡터 ; $\rho = (\rho_1, \dots, \rho_k)$
P	전이 확률들로 구성된 전이 행렬
p	p_{ij} 로 이루어진 벡터 ; $p = (p_{11}, \dots, p_{ij}, \dots, p_{kk})$, $i=1, \dots, k, j=1, \dots, k$
θ	전체 모수 벡터 ; $\theta = (\lambda, \rho, p)$
I_0	Fisher 정보 행렬(Fisher information matrix)

논문에서 사용되는 가정은 다음과 같다.

- ① 고장 원인에 따른 y_n 의 분포는 지수 분포이다. 즉, $F_i(t) = 1 - e^{-\lambda_i t}$
- ② 연속되는 고장 원인간에는 1차 종속성이 존재할 수 있다.
- ③ 기기는 수리 가능하다.
- ④ 존재 가능한 고장 원인의 개수는 유한이다.

- ⑤ 고장 원인에 대한 자료는 모든 고장에 대해 잠재적(masked)이다.

모든 고장에 대해 고장 원인 자료가 잠재적이지 않고 일부의 자료가 주어진 경우에도 본 논문에서 제시된 방법론을 사용할 수 있으나, 가정 ⑤는 고장 원인에 대한 자료가 더욱 제한적인 상황으로 일반화하기 위해 사용된 가정이다.

3. 마코프 전이 모형

일반적으로 수리 가능한 기기에서 기기의 수리가 부품 단위에서 수행될 경우에는 비록 수리 완료후 기기가 작동 상태로 복구되었다 할지라도 고장을 일으킨 부품중에 수리 작업 시 발견되지 않아 계속 고장 상태로 남는 부품이 존재할 수 있으며, 고장을 일으키지 않은 부품에는 수리 작업이 수행되지 않으므로 그때까지 열화가 진행된 상태에서 다시 작동을 시작하게 될 것이다. 따라서 고장 발생시 기기의 완전 대체나 완전 분해 점검이 수행되지 않는 한 고장 원인의 종속성이 존재한다고 할 수 있다. 또한 기기의 완전 대체가 이루어졌다 할지라도 환경적 요인과 같은 고장의 간접적 원인에 의한 종속성도 생각해 볼 수 있다. 예를 들어 기기의 고장이 전자 부품의 고장으로 인하여 발생된 경우에 설계상 전기 배선의 문제나 입력 전압의 불안정성 등이 간접적 고장 원인이 될 수 있는 데, 이는 기기의 대체만으로 완전히 고장 원인이 제거될 수 없는 상황이므로 이 경우에도 고장 원인의 종속성이 존재하는 경우라 할 수 있다.

본 논문에서는 MSM을 사용하여 고장 원인의 1차 종속성을 모형화한다. 이 모형에서는 각 고장 간격 y_n 의 분포를 결정하는 상태 변수 s_n 이 존재한다. 이때 상태는 고장을 발생시킨 원인으로 정의한다. 모형을 정리하면 다음과 같다.

- ① 고장 시간 y_n 의 분포는 상태 변수 s_n 에 의해 결정된다. 즉, $s_n=i$ 일 때, $y_n \sim F_i(t)$
- ② 상태 변수 $\{s_n, n=1,2,\dots,N\}$ 는 마코프 특성을 갖는다. 즉, $\Pr(s_n=j | s_{n-1}=i, s_{n-2}, \dots, s_1) = \Pr(s_n=j | s_{n-1}=i) = p_{ij}$
- ③ 각 상태 변수 값에 따른 y_n 의 분포는 지수 분포이다. 즉, $F_i(t) = 1 - e^{-\lambda_i t}$
- ④ 가능한 기기의 고장 원인의 개수는 유한이다. 따라서 가능한 상태는 유한 마코프 연쇄 (finite Markov chain)를 이룬다.

MSM을 도시하면 다음 [그림 1]과 같다.

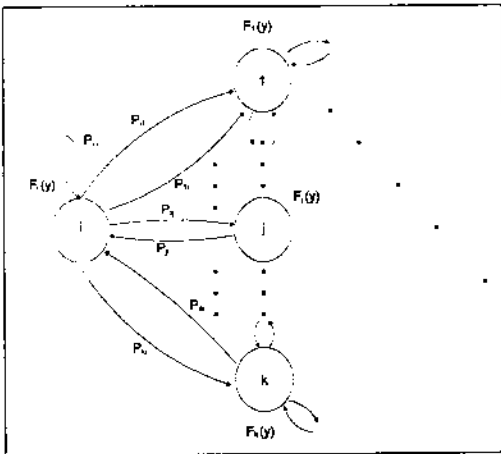


그림 1. 상태수가 k인 마코프 전이 모형

이때, 전체 모수 벡터 $\theta = (\lambda, \rho, p)$ 이다. n

회의 고장이 발생한 후 기기의 수명 분포 $F(n,t)$ 를 표현해 보면 다음 식 (1)과 같다.

$$F(n,t) = \sum_{i=1}^k \sum_{j=1}^k \rho_i P_{ij}^{(n)} F_j(t) \quad (1)$$

위에서 $P_{ij}^{(n)}$ 은 상태 i에서 n번 전이 후에 j 상태에 있을 확률을 의미하며, 이는 P행렬을 n회 곱했을 때 (i, j)성분을 의미한다. 이를 경쟁적 위험 모형 및 분포 혼합 모형과 비교해 보면 다음과 같다. 먼저 경쟁적 위험 모형에서는 직접적 고장이 고장 원인들간의 직렬 구조에서 발생한다는 가정을 사용하므로 이 모형을 이용하여 기기의 수명 분포를 표현하면 다음 식 (2)와 같다.

$$F(t) = 1 - e^{-(\sum_{i=1}^k \lambda_i)t} \quad (2)$$

분포 혼합 모형에서 기기 수명 분포는 다음 식 (3)과 같이 주어진다.

$$F(t) = \sum_{j=1}^k \pi_j F_j(t) \quad (3)$$

위에서 보는 바와 같이 기존 모형은 시간 t만의 함수이다. 기존의 모형에서는 연속되는 고장 원인에 대해 독립성 가정을 사용하므로 고장이 발생해도 항상 동일한 분포를 따르는 것으로 해석된다. 그러나 MSM을 사용하면 n개의 고장 시간 자료를 가지고 추정을 수행하여 구한 $F(n,t)$ 를 사용하여 n번 고장 이후의 기기 수명 분포 변화를 설명할 수 있다. 즉, 기존의 모형이 고장 회수에 대해 정적(static)인 모형인 반면, 본 모형은 시간과 기기의 고장 회수에 대한 함수로 기기의 수명 분포를 표현함으로써 동적(dynamic)인 모형

이라 할 수 있다.

MSM과 분포 혼합 모형의 수명 분포를 비교하면 [그림 2]와 같다. [그림 2]는 모수가 $(\lambda_1, \lambda_2, \rho_2, p_{11}, p_{22}) = (0.1, 0.01, 1, 0.01, 0.01)$ 인 경우에 대해 수명 분포 $F(n, t)$ 를 n 의 변화에 따라 도시한 그림이다. $n = \infty$ 인 경우는 분포 혼합 모형이 된다.

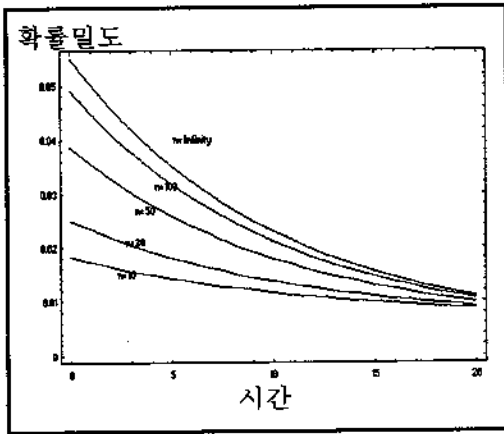


그림 2. 마코프 전이 모형과 분포 혼합 모형($n = \infty$)
 $(\lambda_1 = 0.1, \lambda_2 = 0.01, \rho_2 = 1, p_{11} = 0.01, p_{22} = 0.01)$

분포 혼합 모형은 MSM의 특수한 경우가 됨을 다음과 같이 설명할 수 있다. MSM에서 $\{s_n, n = 1, 2, \dots, N\}$ 가 에르고딕 마코프 연쇄(ergodic Markov chain)를 이룬다면 다음 식이 성립한다.

$$\lim_{n \rightarrow \infty} p_{ij}^{(n)} = p_j \quad (4)$$

위에서 p_j 는 분포 혼합 모형에서 π_j 와 동일하다. 따라서 분포 혼합 모형은 본 논문에서 제시된 MSM의 안정 상태, 즉 고장이 무수히 많이 발생한 후 전이된 분포로 설명될 수 있다. 위 관계를 이용하면 분포 혼합 모형에서

정상성을 보장하기 위해 추정을 위해 많은 수의 고장 시간 자료를 요구하는 가정을 사용하지 않고도 MSM을 이용한 추정을 수행한 후에 식 (4)를 이용하여 분포 혼합 모형에 대한 추정을 수행할 수 있음을 의미한다. 또한 모든 i, j, l 에 대해, $p_{ij} = p_{lj}$ 이면, MSM은 분포 혼합 모형과 동일해진다. 따라서 분포 혼합 모형은 MSM의 특수한 예라고 할 수 있다.

4. 모수 추정 절차

본 절에서는 앞에서 제시한 모형에서 모수 벡터 $\theta = (\lambda, \rho, p)$ 에 대한 추정 절차의 확립 및 그에 따른 알고리즘을 제시한다. 먼저 추정을 위해 사용되는 자료를 살펴보면 고장 시간 자료가 주어졌으며, 상태 s_n 에 대한 자료는 전혀 주어지지 않은 상황이므로, 고장 시간 자료는 완전하며, 상태에 대한 자료는 잠재적이라고 정의할 수 있다. 본 논문에서는 모수 추정 방법으로 최대 우도 추정 방법으로 모수를 추정하는 경우에 사용할 수 있는 수치적 계산 알고리즘으로는 Newton-Raphson 방법, Davidon-Fletcher-Powell 방법, EM 알고리즘 등이 널리 알려져 있다[10]. 특히 상태 벡터에 대해 정보가 주어지지 않기 때문에 본 모형의 추정 문제는 국부 최적해(local optimum)의 존재 여부, 경계 조건(boundary condition)에서 우도의 증가 방향 성등에 대한 예측이 불가능하므로 알고리즘의 둔감성(robustness)이 알고리즘 채택의 주요 결정 기준이 된다. 따라서 위 알고리즘에서 수치적으로 둔감(robust)하며, 다중의 국부 최적해(multiple local optimum)를 갖는 경

우에 있어서도 최적해에 빨리 수렴해가는 것이 알려진 EM 알고리즘을 채택한다[9][10].

본 절에서는 먼저 최대 우도 함수의 정의 및 구체적인 모수 추정 알고리즘을 제시한다. 앞 절에서 제시한 가정을 사용하여 완전 자료에 대한 대수 우도 함수를 표현해 보면 다음과 같다.

$$\begin{aligned} \log f(s_N, y_N; \theta) &= \sum_{i=1}^k I(s_1 = i) [\log f(y_1 | s_1 = i; \theta) + \log p_i] \\ &+ \sum_{n=2}^N \{ \sum_{j=1}^k I(s_n = j) \log f(y_n | s_n = j; \theta) \\ &+ \sum_{j=i=1}^k I(s_{n-1} = i, s_n = j) \log p_{ij} \} \quad (5) \end{aligned}$$

앞에서 설명한 바와 같이 $s_n = \{s_n, n = 1, \dots, N\}$ 는 잠재적 자료이기 때문에 기존의 최대 우도 방법을 사용하여 모수를 추정하기 위해서는 s_n 에 대해 식 (5)의 $\log f(s_N, y_N; \theta)$ 를 2^N 번 합을 구해 표본 우도 함수 $\log f(y_N; \theta)$ 를 구하고 이를 최대화해야 한다. 그러나 이는 현실적으로 불가능하므로 EM 알고리즘을 사용하여 모수 추정을 수행한다. EM알고리즘은 [그림 3]에 요약되어 있는 바와 같이 기대값 계산 단계(Expectation Step)와 최대화 단계(Maximization Step)를 교대로 수행함으로써 모수 θ 에 대한 추정을 수행하는 방법론이다. 이 계산을 위해서 먼저 조건부 대수 우도 함수의 기대값 $E[\log f(s_N, y_N; \theta) | y_N, \theta]$ 의 계산이 선행되어야 한다. 식 (5)로부터 $E[\log f(s_N, y_N; \theta) | y_N, \theta]$ 는 다음 식과 같이 유도된다.

$$\begin{aligned} E[\log f(s_N, y_N; \theta) | y_N, \theta] &= \sum_{i=1}^k \Pr(s_1 = i | y_N; \theta) [\log \lambda_i - \lambda_i y_1 + \log p_i] \end{aligned}$$

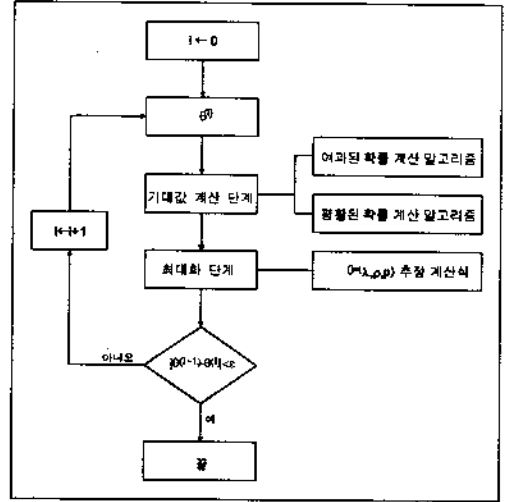


그림 3. 모수 추정 알고리즘의 구성

$$\begin{aligned} &+ \sum_{n=2}^N [\sum_{j=1}^k \Pr(s_n = j | y_N; \theta) [\log \lambda_j - \lambda_j y_n] \\ &+ \sum_{j=1}^k \sum_{i=1}^k \Pr(s_{n-1} = i, s_n = j | y_N; \theta) \log p_{ij}] \quad (6) \end{aligned}$$

알고리즘을 정리하면 다음과 같다.

MLE 추정 알고리즘

단계 1 $k \leftarrow 0, \theta^{(l)}$ 선택

단계 2 (여파된 확률, 평활된 확률 계산)

$i = 1, \dots, k, j = 1, \dots, k$ 에 대해

$$\Pr(s_n = i | y_N; \theta^{(l)}),$$

$$\Pr(s_{n-1} = i, s_n = j | y_N; \theta^{(l)}) \text{ 계산}$$

단계 3 (기대값 계산 단계)

단계 2의 결과를 대입하여

$$E[\log f(s_N, y_N; \theta) | y_N, \theta] \text{ 계산}$$

단계 4 (최대화 단계)

$$\theta^{(l+1)} = \arg \max_{\theta} E[\log f(s_N, y_N; \theta) | y_N, \theta^{(l)}]$$

계산

단계 5 (수렴성 확인)

$|\theta^{(l+1)} - \theta^{(l)}| < \epsilon$ 이면, $\theta = \theta^{(l+1)}$
 아니면, $l \leftarrow l+1$, 단계 2 로 간다.

위 알고리즘의 단계 2에서는 y_N 이 주어졌다는 조건하에서 상태와 상태의 전이에 대한 확률 $\Pr(s_n = i | y_N, \theta^{(l)})$ 와 $\Pr(s_{n-1} = i, s_n = j | y_N, \theta^{(l)})$ 의 계산이 필요한데 이를 평활된 확률(smoothed probability)이라 한다[7]. 평활된 확률을 구하기 위해서 여과된(filtered) 확률 $\Pr(s_n = i | y_n, \theta^{(l)})$ 와 $\Pr(s_{n-1} = i, s_n = j | y_n, \theta^{(l)})$ 의 계산이 선행되어야 한다. 여과된 확률 및 평활된 확률의 계산을 위한 알고리즘을 정리하면 다음과 같다.

여과된 확률의 계산 알고리즘

($\theta = \theta^{(l)}$ 하에서 계산)

단계 0

$$f(y_2 | y_1) = \sum_{s_1=1}^k \sum_{s_2=1}^k f(y_2 | s_2) \Pr(s_2 | s_1) \Pr(s_1) \Pr(s_2, s_1 | y_1) \\ = [f(y_2 | s_2) \Pr(s_2 | s_1) \Pr(s_1)] / f(y_2 | y_1)$$

$n = 3, \dots, N$ 에 대해 다음을 계산한다.

단계 1 $f(y_n, s_n, s_{n-1} | y_{n-1})$

$$= \sum_{s_{n-2}=1}^k f(y_n | s_n) \Pr(s_n | s_{n-1}) \Pr(s_{n-1}, s_{n-2} | y_{n-1})$$

단계 2 $f(y_n | y_{n-1}) = \sum_{s_n=1}^k \sum_{s_{n-1}=1}^k f(y_n, s_n, s_{n-1} | y_{n-1})$

단계 3 $\Pr(s_n, s_{n-1} | y_n) = f(y_n, s_n, s_{n-1} | y_{n-1}) / f(y_n | y_{n-1})$

특히 위에서 $f(y_n | y_{n-1})$ 과 $\Pr(s_n, s_{n-1} | y_n)$ 는 평활된 확률의 계산을 위해 사용되므로 기억할 필요가 있다. 이때 $f(y_n | y_{n-1})$ 은 $(1 \times (N-1))$ 벡터로, $\Pr(s_n, s_{n-1} | y_n)$ 은 $(k^2 \times (N-1))$ 행렬로 기억한다.

평활된 확률의 계산 알고리즘

($\theta = \theta^{(l)}$ 하에서 계산)

$n = N-2, \dots, 1$ 에 대해 다음을 계산한다.

단계 1 $f(y_N, \dots, y_{n+1}, s_N, \dots, s_{n+1} | y_n)$

$$= f(y_N | s_N) \Pr(s_N | s_{N-1}) \dots \Pr(s_{n+1} | s_n) \Pr(s_{n-1}, s_n | y_n)$$

단계 2 $f(y_N | y_n)$

$$= f(y_N | y_{N-1}) f(y_{N-1} | y_{N-2}) \dots f(y_{n+1} | y_n)$$

단계 3 $\Pr(s_{n-1}, s_n | y_N)$

$$= \sum_{s_n=1}^k \dots \sum_{s_{n-1}=1}^k f(y_N, \dots, y_{n+1}, s_N, \dots, s_{n-1} | y_n) / f(y_N | y_n)$$

위 알고리즘의 단계 1과 단계 2 계산에 여과된 확률 계산 알고리즘에서 저장된 결과가 쓰인다. 또한 n 을 증가하여 계산을 수행할 때도 직전에서 계산된 단계 2, 단계 3의 결과를 사용함으로써 저장 용량과 계산 회수를 상당히 줄일 수 있다. 기존에 제시된 알고리즘 [9], [14]와 비교해 보면, 기존의 알고리즘은 단계 3 계산을 위해서 주변(marginal) 분포 계산이 요구되며, 또한 n 에 따라 계산을 수행할 때 직전의 결과를 사용하지 않는다. 따라서 단계 2와 단계 3에서 불필요한 계산이 수반되나 본 논문에서 제시한 알고리즘은 후방(backward) 계산 방식을 사용함으로써 총 $3(N-2)$ 회의 계산만으로 평활된 확률 계산이 가능한 효율적 알고리즘을 제시하였다.

$E[\log f(s_N, y_N; \theta) | y_N, \theta]$ 를 최대화하는 모수 $\theta^{(l+1)} = (\lambda^{(l+1)}, \rho^{(l+1)}, p^{(l+1)})$ 를 구하기 위해서 각각의 모수 λ, ρ 들에 대해 $E[\log f(s_N, y_N; \theta) | y_N, \theta]$ 를 편미분하면 다음과 같은 계산식을 얻을 수 있다. (자세한 계산 과정은 부록 1 참조)

$\theta = (\lambda, \rho, \rho)$ 의 추정치 계산식

단계 1 $n = 2, \dots, N$ 에 대해 $\Pr(s_n | y_N; \theta^{(l)})$

계산

$$\Pr(s_n = i | y_N; \theta^{(l)}) = \sum_{s_{n-1}=1}^k \Pr(s_{n-1}, s_n = i | y_N; \theta^{(l)})$$

단계 2 $i = 1, \dots, k, j = 1, \dots, k$ 에 대해 다음

모수들의 추정치 계산

$$\hat{\lambda}_i^{(l+1)} = \frac{\sum_{n=1}^N \Pr(s_n = i | y_N; \theta^{(l)})}{\sum_{n=1}^N \{\Pr(s_n = i | y_N; \theta^{(l)}) \dots y_n\}}$$

$$\hat{\rho}_i^{(l+1)} = \Pr(s_n = 1 | y_N; \theta^{(l)})$$

$$\hat{p}_{ij}^{(l+1)} = \frac{\sum_{n=2}^N \Pr(s_{n-1} = i, s_n = j | y_N; \theta^{(l)})}{\sum_{n=2}^N \Pr(s_{n-1} = i | y_N; \theta^{(l)})}$$

5. 통계적 검정 절차

본 절에서는 앞 절에서 제안한 MSM에서 상태수 k 에 대한 결정을 위한 통계적 분석 절차에 대해 정리한다. 또한 3절에서 보인 바와 같이 MSM은 상태간의 종속성을 반영한 모형이며, 분포 혼합 모형은 상태의 전이가 현재 상태에 대해 독립임을 의미한다. 또한 분포 혼합 모형은 MSM의 특수한 예이다. 따라서 본 절에서는 추정된 MSM과 분포 혼합 모형의 동일성 검정을 위해 상태 전이간의 독립성 검정 절차를 개발하였다.

5.1 상태 수 결정 절차

상태 수 결정 절차에 대해 정리하면 아래와 같다. 상태수에 대한 결정은 먼저 상태 수 $k=2$ 에서 분석을 수행하고 귀무 가설 $H_0; \lambda_i = \lambda_j$ 가 채택되면, 상태수 $k=1$ 에서 $\theta = (\lambda)$ 를 추

정하며, 기각되면, 상태 수 $k \leftarrow k+1$ 에서 모두 θ 를 추정하고 $i \neq j (i = 1, \dots, k, j = 1, \dots, k)$ 에 대해 $k(k-1)/2$ 개의 귀무 가설 $H_0; \lambda_i = \lambda_j$ 에 대한 검정을 수행하는 방식이다. 검정 통계량 (test statistic)으로는 Wald 통계량을 사용한다[7]. [그림 4]는 상태 수 결정 절차를 이용하여 개발한 통합된 분석 절차를 제시하고 있다. 상태 수 결정 절차를 정리하면 다음과 같다.

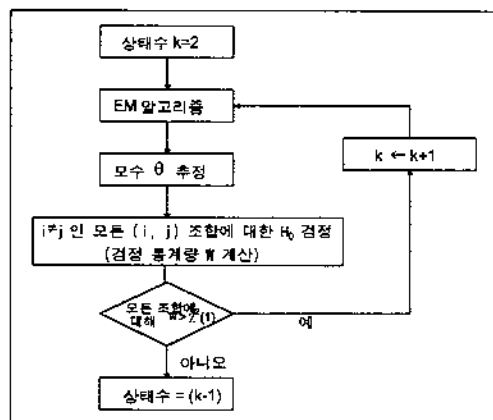


그림 4. 통합된 분석 절차

상태 수 결정에 대한 검정

가설 $H_0; \lambda_i = \lambda_j$

(단, $i \neq j, i = 1, \dots, k, j = 1, \dots, k$)

$H_1; \lambda_i \neq \lambda_j$

$$\text{통계량 } W(\text{Wald 통계량}) = \frac{(\hat{\lambda}_i - \hat{\lambda}_j)^2}{\text{Var}(\hat{\lambda}_i) + \text{Var}(\hat{\lambda}_j) - 2\text{Cov}(\hat{\lambda}_i, \hat{\lambda}_j)}$$

검정 모든 (i, j) 조합에서 $W < \chi^2_{\alpha}(1)$ 이면

→ 상태 수 $(k+1)$ 에서 추정 및 검정

아니면 → 상태 수를 $(k-1)$ 로 결정

위 검정 통계량 W 를 구하기 위해서 $\text{Var}(\hat{\lambda}_i), \text{Cov}(\hat{\lambda}_i, \hat{\lambda}_j)$ 값의 계산이 선행되어야 한다.

$Var(\hat{\lambda}_i)$ 와 $Cov(\hat{\lambda}_i, \hat{\lambda}_j)$ 에 관한 계산식에 대해 정리하면 다음과 같다.

Fisher 정보 행렬(Fisher information matrix)은 다음과 같이 정의할 수 있다.

$$I_0 = \begin{pmatrix} I_\lambda & 0 & 0 \\ 0 & I_p & 0 \\ 0 & 0 & I_p \end{pmatrix} \quad (7)$$

점근적 공분산 행렬(asymptotic covariance matrix) Σ 는 다음과 같다[16].

$$\Sigma = \begin{pmatrix} Var(\hat{\lambda}_i) & Cov(\hat{\lambda}_i, \hat{\lambda}_j) \\ Cov(\hat{\lambda}_i, \hat{\lambda}_j) & Var(\hat{\lambda}_j) \end{pmatrix} = I_\lambda^{-1} \quad (8)$$

위에서 I_λ 의 (i, j) 원소는 $E_{\theta} \{ (\frac{\partial \mathcal{L}}{\partial \lambda_i}) \cdot (\frac{\partial \mathcal{L}}{\partial \lambda_j}) \}$ 로 정의된다. $E_{\theta} \{ \cdot \}$ 는 { } 안의 식에 기대값을 취한 후에 θ 에 $\hat{\theta}$ 를 대입한 값을 의미한다. I_λ 의 각 원소에 대한 구체적인 계산식은 다음과 같다. (자세한 계산 과정은 부록 2 참조)

$$E_{\theta} \left\{ \left(\frac{\partial \mathcal{L}}{\partial \lambda_i} \right)^2 \right\} = \sum_{n=1}^N \Pr(s_n = i | y_n; \hat{\theta}) \cdot \frac{1}{\lambda_i^2} \quad (9)$$

$$E_{\theta} \left\{ \left(\frac{\partial \mathcal{L}}{\partial \lambda_i} \right) \left(\frac{\partial \mathcal{L}}{\partial \lambda_j} \right) \right\} = 0 \quad (10)$$

5.2 상태 전이에 대한 독립성 검정

앞에서 언급한 바와 같이 일단 상태 수 k 가 결정되면 부가적 검정으로 상태 전이간의 독립성 검정을 수행하여 분포 혼합 모형으로 볼 수 있는가에 대한 검정을 수행해 본다. 이때 귀무 가설로는 모든 j 에 대해 p_{ji} 가 동일하다는 가설을 사용한다. 그 의미는 상태의 전이가 현재 상태에 대해 독립적으로 이루어질 수 있으며, 이는 MSM이 정상 상태

(steady state)라는 해석을 할 수 있다. 앞 절에서 설명한 바와 같이 이 경우는 분포 혼합 모형으로 볼 수 있는 경우가 된다. 이에 대한 검정 절차는 다음과 같다.

상태 전이에 대한 독립성 검정

가설 H_0 : 모든 i 에 대해 $p_{1i} = p_{2i} = \dots = p_{ki}$

H_1 : 모든 i, j, l 에 대해 어느 하나의 $p_{ji} \neq p_{li} (l \neq j)$

통계량 $W_M = h^T(\hat{p}) [H^T \cdot I_p \cdot H]^{-1} h(\hat{p})$

검정 $W_M > \chi^2_{\alpha}((k-1)^2)$ 이면, 상태 전이간의 종속성이 존재한다.

아니면, 분포 혼합 모형으로 볼 수 있다.

위에서 \hat{p} 와 $h(\hat{p})$ 는 다음과 같이 정의된다.

$$\hat{p} = \begin{pmatrix} \hat{p}_{11} \\ \vdots \\ \hat{p}_{1(k-1)} \\ \hat{p}_{21} \\ \vdots \\ \hat{p}_{k(k-1)} \end{pmatrix} \quad (11)$$

$$h(\hat{p}) = \begin{pmatrix} \hat{p}_{11} - \hat{p}_{21} \\ \hat{p}_{11} - \hat{p}_{31} \\ \vdots \\ \hat{p}_{11} - \hat{p}_{k1} \\ \vdots \\ \hat{p}_{(k-1)(k-1)} - \hat{p}_{1(k-1)} \\ \vdots \\ \hat{p}_{(k-1)(k-1)} - \hat{p}_{(k-2)(k-1)} \end{pmatrix} \quad (12)$$

H 는 $h(\hat{p})$ 의 자코비안 행렬(Jacobian matrix)이다. 위에서 $h^T(\hat{p})$ 와 H^T 는 $h(\hat{p})$ 와 H 의 전치(transpose) 행렬을 의미한다. I_p 의 각 원소

에 대한 계산식은 다음과 같다.

$$E_0\left(\frac{\partial \mathcal{L}}{\partial p_{ij}}\right)^2 = \sum_{n=2}^N \Pr(s_{n-1}=i, s_n=j | y_w; \hat{\theta}) \frac{1}{p_{ij}^2} \quad (13)$$

$$E_0\left(\frac{\partial \mathcal{L}}{\partial p_{ij}}\right)\left(\frac{\partial \mathcal{L}}{\partial p_{i'j'}}\right) = 0 \quad (14)$$

단, $i \neq i'$ 또는 $j \neq j'$

6. 실험 예제

본 절에서는 다양한 실험 예제를 통해 본 논문에서 제안한 MSM에 대한 타당성 검증 및 제시된 추정 알고리즘과 통합된 분석 절차에 대한 실제 수행 예를 보이고 있다.

먼저 실험 1, 2, 3은 각각 자료 개수, 고장률 비, 전이 확률에 따른 추정의 정확도에 대한 실험이며, 실험 4는 상태 수 결정 절차에 대한 정확도를 알아보기 위해 수행된 실험이다. 실험 5는 본 논문에서 제시된 방법론을 이용하여 [3, ch 1, p6]에 제시된 냉방 설비의 고장 자료에 대한 분석을 수행하여, 실제 수행 예를 보이고 있다.

6.1 실험 1 (자료 개수의 영향, 추정의 일관성 검증)

(1) 입력 자료

상태 수 ; 2

고장률 ; $\lambda_1 = 0.1$ $\lambda_2 = 0.01$

(고장률 비 ; 10)

초기 확률 ; $\rho_1 = 0$ $\rho_2 = 1.0$

전이 확률 ; $p_{11} = 0.1$ $p_{22} = 0.1$

자료 개수 ; 10, 20, 30, 40, 50, 60, 70, 80, 90, 100

(2) 실험 결과

총 10 가지 경우에 대해, 각 100회씩 모의 실험을 수행하고, 그 자료를 이용하여 추정을 수행하였다. 자세한 실험 결과는 [표 1]에 정리되어 있다. [그림 5]는 변동 계수 (coefficient of variation : COV)를 척도로 하여 추정에 대한 정확도를 나타내고 있다. 그림에서 보이는 바와 같이 추정에 사용된 자료 개수가 커짐에 따라 COV는 작아지고 있음을 알 수 있다. 따라서 자료 개수에 따른 추정의 일관성(consistency)이 존재한다고 할 수 있다.

6.2 실험 2 (고장률 비의 영향)

(1) 입력 자료

상태 수 ; 2

고장률 비(λ_1/λ_2) ; 2, 4, 8, 16, 32, 64, 128, 256, 512 ($\lambda_1 = 0.1$ 로 고정)

초기 확률 ; $\rho_1 = 0$ $\rho_2 = 1.0$

전이 확률 ; $p_{11} = 0.1$ $p_{22} = 0.1$

자료 개수 ; 50 개

(2) 실험 결과

총 9 가지 경우에 대해 각 100회씩 모의 실험을 수행하고, 그 자료를 이용하여 추정을 수행하였다. 결과는 [표 2]에 정리되어 있다. [그림 6]은 고장률 비에 따른 모수 추정의 정확도를 95% 신뢰 구간을 척도로 하여 나타내고 있다. 그림에서 보는 바와 같이 신뢰 구간 안에 모수의 참값이 포함되어 있음을 알 수 있으며, 고장률의 비가 증가할수록 신뢰 구간은 추정에 사용되는 자료수가 동일하다 하더라도 보다 정밀하게 추정됨을 알 수 있다. 그림에서 보이는 바와 같이 λ_1 의 추

표 1. 자료 개수의 영향에 대한 실험 (실험 1)

($\lambda_1 = 0.1, \lambda_2 = 0.01, \rho_2 = 1.0, \rho_{11} = \rho_{22} = 0.1$)

자료 개수	$\hat{\lambda}_1$		$\hat{\lambda}_2$		$\hat{\rho}_2$		$\hat{\rho}_{11}$		$\hat{\rho}_{22}$	
	mean	COV	mean	COV	mean	COV	mean	COV	mean	COV
10	0.1231	0.5824	0.0121	0.6325	0.99	0.01	0.1386	1.6745	0.0879	1.7284
20	0.1114	0.4177	0.0111	0.4586	1.00	0.00	0.1503	1.4045	0.0939	1.4861
30	0.1110	0.3666	0.0102	0.2645	1.00	0.00	0.1318	1.1821	0.0941	1.0985
40	0.1047	0.2743	0.0102	0.2276	1.00	0.00	0.1196	1.1926	0.0947	1.0759
50	0.1088	0.2574	0.0102	0.1877	1.00	0.00	0.1167	1.066	0.1032	0.9351
60	0.1049	0.2582	0.0102	0.2052	1.00	0.00	0.1134	1.0368	0.0984	0.8281
70	0.1019	0.2044	0.0102	0.1834	1.00	0.00	0.1153	0.9813	0.0921	0.8165
80	0.1016	0.1986	0.0100	0.1635	1.00	0.00	0.1099	0.9489	0.0965	0.7449
90	0.1022	0.1872	0.0100	0.1623	1.00	0.00	0.1047	0.8987	0.0961	0.7117
100	0.1016	0.1680	0.0100	0.1665	1.00	0.00	0.1078	0.8165	0.1060	0.6589

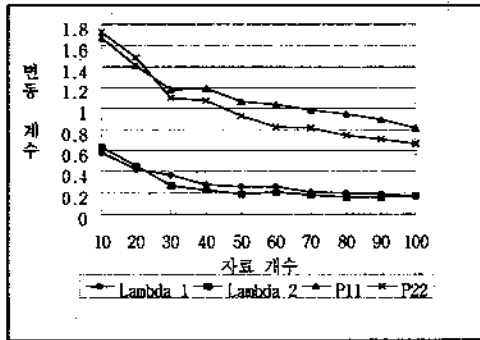


그림 5. 자료수에 따른 추정치의 변동 계수(실험 1)

$\lambda_1=0.1, \lambda_2=0.01, \rho_2=1, \rho_{11}=0.1, \rho_{22}=0.1$

정 신뢰 구간은 대략 고장률 비가 8이상에서 안정화되며, λ_2 에 대해서는 고장률 비의 증가에 따라 급격히 안정화됨을 나타내고 있다. 또한 ρ_2 에 대해서는 고장률 비가 2인 경우를 제외하고는 모두 참값으로 추정되며, ρ_{11} 과 ρ_{22} 에 대해서는 고장률 비가 16이상인 경우에

대해서 신뢰 구간이 안정됨을 알 수 있다. 또한 COV를 척도로 해석해도 똑같은 결과를 보임을 [표 2]를 통해 알 수 있다.

6.3 실험 3 (전이 유형의 영향)

(1) 입력 자료

상태 수 : 2

고장률 : $\lambda_1=1 \lambda_2=0.1$ (고장률 비 : 10)

초기 확률 : $\rho_1=0 \rho_2=1$

전이 확률 : ① $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ ② $P = \begin{pmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{pmatrix}$

③ $P = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix}$ ④ $P = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix}$

자료 개수 : 50개

(2) 실험 결과

위 4가지 경우에 대해 100회의 모의 실험을 수행하고 그 자료를 이용하여 자세한 결

표 2. 고장률 비의 영향에 대한 실험 (실험 2)

($\lambda_1 = 0.1, \rho_2 = 1.0, p_{11} = p_{22} = 0.1, N = 50$)

고장률 비 (λ_1/λ_2)	$\hat{\lambda}_1$		$\hat{\lambda}_2$		$\hat{\rho}_2$		\hat{p}_{11}		\hat{p}_{22}	
	mean	std	mean	std	mean	std	mean	std	mean	std
2	0.1186	0.1046	0.0493	0.0145	0.93	0.2564	0.1907	0.2506	0.1164	0.1935
4	0.1069	0.0251	0.0250	0.0041	1.00	0.0000	0.1449	0.1619	0.1060	0.1273
8	0.1053	0.0199	0.0127	0.0019	1.00	0.0000	0.1139	0.1050	0.1015	0.0739
16	0.1043	0.0177	0.0064	0.0009	1.00	0.0000	0.1031	0.0802	0.0988	0.0544
32	0.1034	0.0158	0.0032	0.0004	1.00	0.0000	0.0995	0.0692	0.0963	0.0473
64	0.1025	0.0146	0.0016	0.0002	1.00	0.0000	0.0963	0.0590	0.0950	0.0439
128	0.1027	0.0147	0.0008	0.0001	1.00	0.0000	0.0947	0.0524	0.0955	0.0431
256	0.1030	0.0149	0.0004	6 E-5	1.00	0.0000	0.0944	0.049	0.0960	0.0416
512	0.1031	0.0152	0.0002	3 E-5	1.00	0.0000	0.0949	0.0479	0.0959	0.0407

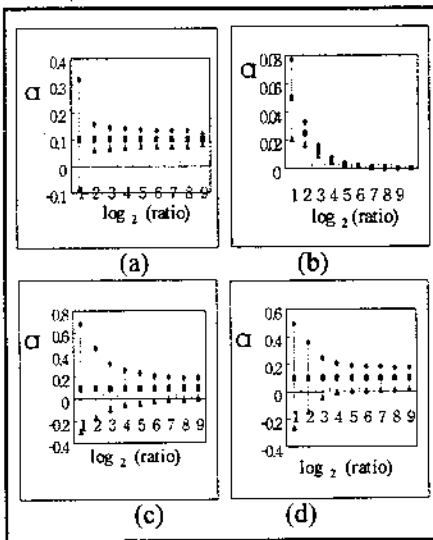


그림 6. 추정치에 대한 95% 신뢰구간(실험 2)
 (a) λ_1 에 대한 신뢰 구간 (b) λ_2 에 대한 신뢰 구간
 (c) p_{11} 에 대한 신뢰 구간 (d) p_{22} 에 대한 신뢰 구간

과는 [표 3]에 정리되어 있다. 척도로는 상대 오차(relative error)를 사용하였는데, 이

유는 ①의 경우에 COV가 정의되지 않기 때문이다. [그림 7]은 상대 오차를 척도로 하여 추정의 정확도를 나타내고 있다. 그림에서 보이는 바와 같이 $p_{11} = 0.5, p_{22} = 0.5$ 의 경우에 오차가 가장 크며, 전이가 빈번하거나 작은 상태 전이 확률을 갖는 상태 구별이 뚜렷한 경우에는 오차가 작음을 알 수 있다.

6.4 실험 4 (상태 수 결정 절차에 대한 3원 배치 실험)

(1) 입력 자료

상태 수 ; 2

고장률 ; ① $\lambda_1 = 1, \lambda_2 = 0.1$ (고장률 비 : 10)

② $\lambda_1 = 1, \lambda_2 = 0.5$ (고장률 비 : 2)

초기 확률 ; $p_{11} = 0, p_{22} = 1$

전이 확률 ; ① $P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ ② $P = \begin{pmatrix} 0.1 & 0.9 \\ 0.9 & 0.1 \end{pmatrix}$

③ $P = \begin{pmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{pmatrix}$ ④ $P = \begin{pmatrix} 0.9 & 0.1 \\ 0.1 & 0.9 \end{pmatrix}$

표 3. 전이 확률의 영향에 대한 실험 (실험 3)

($\lambda_1 = 0.1, \lambda_2 = 0.01, \rho_2 = 1.0, N = 50, R.E = |\text{참값} - \text{추정값}| / \text{참값}$)

전이 확률(p_{11}, p_{22})	$\hat{\lambda}_1$			$\hat{\lambda}_2$			$\hat{\rho}_2$			\hat{p}_{11}			\hat{p}_{22}		
	mean	std	R.E	mean	std	R.E	mean	std	R.E	mean	std	R.E	mean	std	R.E
(0, 0)	0.1054	0.0211	0.0536	0.0101	0.002	0.0082	1.00	0.00	0.00	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
(0.1, 0.1)	0.1088	0.0313	0.0878	0.0100	0.0019	0.0064	1.00	0.00	0.00	0.1167	0.1244	0.0167	0.1032	0.0965	0.0321
(0.5, 0.5)	0.1158	0.0475	0.1579	0.0104	0.0022	0.0417	1.00	0.00	0.00	0.4561	0.1834	0.0878	0.4735	0.1900	0.0530
(0.9, 0.9)	0.1097	0.0617	0.0973	0.0102	0.0023	0.0248	0.99	0.10	0.01	0.8556	0.1428	0.0493	0.8637	0.0905	0.0404

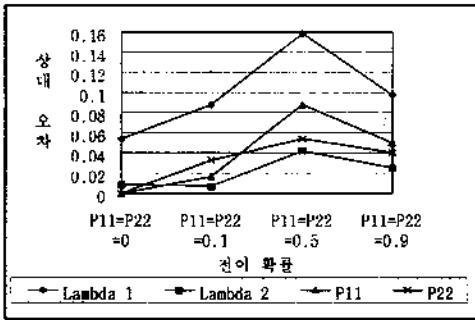


그림 7. 전이 확률에 따른 추정치의 상대 오차(실험 3)
($\lambda_1 = 0.1, \lambda_2 = 0.01, \rho_2 = 1, N = 50$)

$$\begin{aligned} \textcircled{5} P &= \begin{pmatrix} 0.9 & 0.1 \\ 0.5 & 0.5 \end{pmatrix} & \textcircled{6} P &= \begin{pmatrix} 0.9 & 0.1 \\ 0.9 & 0.1 \end{pmatrix} \\ \textcircled{7} P &= \begin{pmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{pmatrix} \end{aligned}$$

자료 개수 : ① 10개 ② 50개 ③ 100개
유의 수준 : 1%, 5%

(2) 실험 결과

유의 수준에 따라 위 42가지 경우에 대해 100회의 모의 실험을 수행하고, 그 자료를 이용하여 각각의 경우에 대해 상태 수 결정 절차를 수행하였다. 자세한 결과는 [표 4]에 정리되어 있다. 이때 척도로는 상태를 1로 추정할 비율, 즉 오차율을 사용하였다. [그림

8]에서 보는 바와 같이 분산 분석 결과 고장을 비, 전이 확률, 자료 개수가 모두 상태 수 결정에 영향을 미치는 것으로 나타났다. [그림 8]에서 자료 개수 10일 때와 고장률 비가 2일 때의 오차율은 실험 1, 2의 결과를 이용하면, 추정의 부정확성에 기인한 것으로 해석할 수 있다.

특히 [그림 8]은 $p_{11} = 0.1, p_{22} = 0.9$ 의 경우에 오차율이 가장 큰 것으로 나타내고 있다. 이 경우는 입력 모형의 특성상 초기에 상태 2에서 시작된 과정(process)이 상태 1로의 전이가 잘 일어나지 않는 경우로서 모의실험에 있어 고장률 λ_1 을 따르는 고장 시간은 잘 발생하지 않게 된다. 따라서 발생한 자료를 토대로 보면 고장률의 차이가 크지 않은 것으로 해석될 수 있어 상태 수 결정 절차는 정확한 결과를 보이고 있다고 할 수 있다. 또한 [표 4]를 보면 유의 수준이 1%인 경우가 5% 유의 수준을 사용한 경우에 비해 오차율이 모두 크게 나타나고 있으며, 오차율의 차이는 고장률 비가 2인 경우에 더 크게 나타나고 있다. 또한 그림을 보면 유의 수준이 변해도 오차율 곡선이 비슷한 양상을 보이고 있다. 이는 유의 수준에 대한 상태 수 결정 절차의 일관성을 보여 준다고 할 수 있다.

표 4. 상태수 결정에 대한 오차율 (실험 4)

($k = 2, \rho_2 = 1.0$)

(a) 유의 수준 = 1%

전이 확률 (p_{11}, p_{22})	고장률 (λ_1, λ_2)			
	(1, 0.1)		(1, 0.5)	
	자료	개수	자료	개수
(0, 0)	10	1.00	10	1.00
	50	0.00	50	0.62
	100	0.00	100	0.19
(0.1, 0.1)	10	1.00	10	1.00
	50	0.00	50	0.55
	100	0.00	100	0.29
(0.5, 0.5)	10	0.95	10	1.00
	50	0.00	50	0.54
	100	0.00	100	0.42
(0.9, 0.9)	10	0.97	10	1.00
	50	0.05	50	0.66
	100	0.00	100	0.31
(0.9, 0.5)	10	0.86	10	0.99
	50	0.01	50	0.52
	100	0.00	100	0.43
(0.9, 0.1)	10	0.72	10	0.99
	50	0.04	50	0.58
	100	0.00	100	0.48
(0.1, 0.9)	10	1.00	10	1.00
	50	0.62	50	0.73
	100	0.33	100	0.60

(b) 유의 수준 = 5%

전이 확률 (p_{11}, p_{22})	고장률 (λ_1, λ_2)			
	(1, 0.1)		(1, 0.5)	
	자료	개수	자료	개수
(0, 0)	10	0.42	10	1.00
	50	0.00	50	0.33
	100	0.00	100	0.07
(0.1, 0.1)	10	0.60	10	0.94
	50	0.00	50	0.37
	100	0.00	100	0.15
(0.5, 0.5)	10	0.63	10	0.94
	50	0.00	50	0.37
	100	0.00	100	0.19
(0.9, 0.9)	10	0.75	10	0.95
	50	0.01	50	0.46
	100	0.00	100	0.22
(0.9, 0.5)	10	0.34	10	0.89
	50	0.01	50	0.41
	100	0.00	100	0.34
(0.9, 0.1)	10	0.31	10	0.85
	50	0.02	50	0.43
	100	0.00	100	0.40
(0.1, 0.9)	10	0.96	10	0.95
	50	0.40	50	0.50
	100	0.10	100	0.43

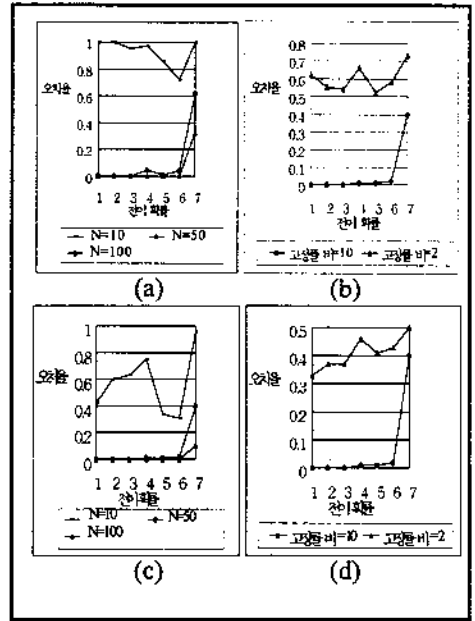


그림 8. 상태수 결정에 대한 오차(실험 4)

- (a) 전이 확률과 자료수에 따른 오차율
(고장률 비=10, $\alpha=1\%$)
- (b) 전이 확률과 고장률 비에 따른 오차율
($N=50, \alpha=1\%$)
- (c) 전이 확률과 자료수에 따른 오차율
(고장률 비=10, $\alpha=5\%$)
- (d) 전이 확률과 고장률 비에 따른 오차율
($N=50, \alpha=5\%$)

6.5 실험 5 (냉방 설비 자료 분석)

실험 5에서는 [3, ch 1, p6] 에 제시된 보잉 720 비행기의 냉방 설비에 대한 고장 시간 자료를 이용하여 본 논문에서 제시한 방법론을 사용하여 분석한 실용예를 보이고 있다. 분석을 통해 보면 자료 4, 6, 7, 12번에 대해서는 상태 수 2로 추정을 하였으며, 나머지 자료에 대해서는 하나의 지수 분포로 추정되었다. [표 5]는 상태 수가 2로 추정된 경우에 대해 추정 결과를 정리하였다. 결과

표 5. 냉방 설비 고장 자료에 대한 추정 결과

자료 번호	자료 개수	상태수	$\hat{\lambda}_1$	$\hat{\lambda}_2$	$\hat{\rho}_1$	$\hat{\rho}_{11}$	$\hat{\rho}_{22}$
4	15	2	0.027702	0.004329	1.00	0.34210	0.000005
6	30	2	0.075641	0.012694	1.00	0.39960	0.739400
7	27	2	0.009138	0.024494	0.00	0.08991	0.000005
12	12	2	0.063810	0.005796	1.00	0.00000	0.187800

를 보면 7번 자료에 대해서는 고장률이 0.009138인 지수 분포를 따르는 고장 시간이 발생하면 이후의 고장 시간은 고장률이 0.024494인 지수 분포를 따르며 이러한 전이를 반복한다고 해석할 수 있다. 7번 자료의 경우는 어떤 외부적 영향이 작용하여 수명 향상을 수반하는 수리 작업과 수명 감소를 가져오는 수리 작업이 교대로 이루어졌다고 말할 수 있다. 12번 자료의 경우도 거의 교대로 2개의 지수 분포를 따라 고장 시간이 발생했다고 해석할 수 있으며, 4번 자료는 고장률 0.004329인 지수 분포를 따라 고장 시간이 발생하면 다음에는 반드시 다른 고장률을 가지는 분포를 따라 고장 시간이 발생하며 이후에는 약 0.66의 확률로 전이가 발생할 수 있다고 해석할 수 있다.

7. 결론

본 논문은 잘 식별되지 않는 다수의 고장 원인을 갖는 기기의 신뢰성 분석을 위해 연속되는 고장 원인간의 종속성 및 고장 원인별 분포 특성 차를 모형화할 수 있는 신뢰성 모형 및 이에 따른 통계적 분석 절차의 개발을 수행하였다. 또한 MSM을 이용하여 기존의 신뢰성 함수와는 달리 기기의 수명 분포

를 고장 회수와 고장 시간의 함수로 표현함으로써 고장 회수에 동적인 기기의 수명 함수를 제시하였다. 분석 방법에 있어서는 N개의 자료가 존재하는 경우에 대해 평활된 확률 계산이 총 $3(N-2)$ 회의 계산만으로 가능한 효율적 알고리즘을 제시하였다. 또한 상태 수 결정 및 전이에 대한 독립성 검정을 위한 통계적 검정 절차를 개발하여 제시함으로써 통합적인 분석 절차를 수립하였다. 이를 이용하여 다양한 예제를 통해 타당성을 검증하고 실제 수행 결과를 제시하였다. 본 논문을 이용하면 고장 원인에 대한 부가적인 자료를 요구하지 않고서도 고장 원인에 따른 분포 특성의 차이 및 종속성 모형화가 가능하며, 또한 다항 시간의 복잡도를 가지는 계산만으로 모수의 추정이 가능하다고 할 수 있다. 모형면에 있어서는 분포 혼합 모형을 포함하는 보다 유연한 모형을 제시하였다고 할 수 있다. 따라서 본 논문에 제시된 방법론은 다수의 특정한 고장 원인이 존재하고 그 식별이 어려우며 고장 자료의 분포 특성이 IID가 아닌 분포 혼합적 특성을 가진 경우의 분석에 쉽게 응용될 수 있는 방법론이라 할 수 있다.

참 고 문 헌

- [1] Cheng, S. W., Fu, J. C., and Sinha, S. K., "An empirical procedure for estimating the parameters of a mixed exponential life testing model", IEEE Trans. Reliability, vol R-34, April, 1985, pp 60-64
- [2] Cox, D. R., "The analysis of exponentially distributed life-times with two types of failure", J. R. Statistical Soc., B, vol. 21, 1959, pp 411-421
- [3] Cox, D. R., and P. A. W. Lewis, The Statistical Analysis of Series of Events, John Wiley & Sons, 1966
- [4] Dempster, A. P., Laird, N. M., and Rubin, D. B., "Maximum likelihood from incomplete data via the EM algorithm", J. R. Statistical Soc., B, vol 39, 1977, pp1-3
- [5] Dinse, G. E., "Nonparametric estimation for partially-incomplete time and types of failure data", Biometrics, vol. 38, 1982, pp 417-431
- [6] Doganaksoy, N., "Interval estimation from censored & masked system-failure data", IEEE Trans. Reliability, vol R-40, Aug., 1991, pp 280-285
- [7] Engel, C., and Hamilton, J. D., "Long swings in the dollar: Are they in the data and do markets know it?", The American Economic Review, vol 80, Sep., 1990, pp 689-713
- [8] Hahn, G. J., and Meeker, W. Q., "Pitfalls and practical considerations in product life analysis, parts I and II", Journal of Quality Technology, vol. 14, 1982, pp 144-152 and 177-185
- [9] Hamilton, J. D., "Analysis of time series subject to changes in regime", Journal of Econometrics, vol 45, 1990, pp 39-70
- [10] Hamilton, J. D., Time Series Analysis, Princeton University Press, 1994
- [11] Kalbfleish, J. D., and Prentice, R. L., The Statistical Analysis of Failure Time Data, John Wiley and Sons, 1980
- [12] Kalbfleish, J. D., and Lawless, J. F., "Estimation of Reliability in Field-Performance studies", Technometrics, vol. 30, Nov., 1988, pp 365-378
- [13] Lawless, J. F., "Statistical Methods in Reliability", Technometrics, vol.25, Nov., 1983, pp 305-316
- [14] Lee, J. H., Nonstationary Markov Switching Models of Exchange Rates ; The Pound-Dollar Exchange Rate, Ph. D Dissertation, University of Pennsylvania, 1991
- [15] Miyakawa, M., "Analysis of incomplete data in competing risk model", IEEE Trans. Reliability, vol R-33, Oct., 1984, pp 293-296
- [16] Nelson, W., Applied life data analysis, John Wiley & Sons, 1982
- [17] Prentice, R. L., Williams, B. J., and Peterson, A. V., "On the regression analysis of multivariate failure time data", Biometrika, vol 68, 1981, pp373-379
- [18] Quandt, R. E., "The estimation of the parameters of a linear regression system

obeying two separate regimes”, Journal of the American Statistical Association, vol 55, 1958, pp 875-880

- [19] Usher, J. S., and Hodgson, T. J., “Maximum likelihood analysis of component reliability using masked system life-

test data”, IEEE Trans. Reliability, vol R-37, Dec., 1988, pp 550-555

95년 8월 최초 접수, 95년 11월 최종 수정

부록 1

최대 우도 추정치 \hat{p}_{ij} 를 구하기 위해 다음 식을 전개한다.

$$\frac{\partial \log f(y_{N_i}, s_{N_i}; \theta)}{\partial p_{ij}} = \sum_{n=1}^N I(s_{n-1} = i, s_n = j) \cdot \frac{1}{p_{ij}} \tag{A.1}$$

다음과 같은 함수를 정의한다.

$$\begin{aligned} Q(\theta^{(l+1)} | y_{N_i}; \theta^{(l)}) &= \sum_{s_N=1}^k \cdot \sum_{s_1=1}^k \log f(y_{N_i}, s_{N_i}; \theta^{(l+1)}) \cdot f(y_{N_i}, s_{N_i}; \theta^{(l)}) \\ &= \int_s \log f(y_{N_i}, s_{N_i}; \theta^{(l+1)}) \cdot f(y_{N_i}, s_{N_i}; \theta^{(l)}) \\ &= Q \end{aligned} \tag{A.2}$$

위에서 l 은 EM 알고리즘에서 l 번째 순환에서 계산된 결과임을 의미한다. Hamilton[9]은 다음 정리를 증명한 바 있다.

$$\begin{aligned} \frac{\partial Q(\theta^{(l+1)} | y_{N_i}; \theta^{(l)})}{\partial \theta^{(l+1)}} \Big|_{\theta^{(l+1)} = \hat{\theta}^{(l)}} = 0 \text{ 이면,} \\ \frac{\partial f(y_{N_i}; \theta)}{\partial \theta} \Big|_{\theta = \hat{\theta}^{(l)}} = 0 \end{aligned} \tag{A.3}$$

이 된다.

위 정리에 의해

$$\begin{aligned} \frac{\partial Q}{\partial p_{ij}^{(l+1)}} &= \int_s \frac{\partial f(y_{N_i}, s_{N_i}; \theta^{(l+1)})}{\partial p_{ij}^{(l+1)}} \cdot f(y_{N_i}, s_{N_i}; \theta^{(l)}) \\ &= \int_s [p_{ij}^{(l+1)}]^{-1} \sum_{n=2}^T I(s_{n-1} = i, s_n = j) \cdot f(y_{N_i}, s_{N_i}; \theta^{(l)}) \end{aligned} \tag{A.4}$$

이 식에서

$$\int_s I(s_{n-1}=i, s_n=j) \cdot \Pr(y_{N}, s_{N}; \theta^{(l)}) = \Pr(s_{n-1}=i, s_n=j | y_{N}; \theta^{(l)}) \cdot f(y_{N}; \theta^{(l)})$$

이므로, 위 식 (A.4)는 다음과 같이 정의할 수 있다.

$$\frac{\partial Q}{\partial p_{ij}^{(l+1)}} = [p_{ij}^{(l+1)}]^{-1} \sum_{n=2}^N \Pr(s_{n-1}=i, s_n=j | y_{N}; \theta^{(l)}) \cdot f(y_{N}; \theta^{(l)}) \quad (A.5)$$

EM 알고리즘에서 p_{ij} 에 대한 추정 절차는 위 식 (A.4)에 주어진 식을 만족하는 $\hat{p}_{ij}^{(l)}$ 을 찾는 과정이므로 결국 다음 식의 해를 구하는 과정이다.

$$\begin{cases} \text{Max} & Q \\ \text{s.t} & \sum_{j=1}^k p_{ij}^{(l+1)} = 1 \end{cases} \quad (A.6)$$

위 식 (A.6)에서 라그랑지안 승수(Lagrangian Multiplier)를 μ_i 라면 1차 조건식(first order condition)에 의해 다음 식을 구할 수 있다.

$$\frac{\partial Q}{\partial p_{ij}^{(l+1)}} = \mu_i \quad (A.7)$$

위 식 (A.5)와 식 (A.7)을 결합하여 다음 식을 얻을 수 있다.

$$\sum_{n=2}^N \Pr(s_{n-1}=i, s_n=j | y_{N}; \theta^{(l)}) = p_{ij}^{(l+1)} \mu_i / f(y_{N}; \theta^{(l)}) \quad (A.8)$$

다음은 j 에 대해 합을 취하면,

$$\sum_{n=2}^N \sum_{j=1}^k \Pr(s_{n-1}=i, s_n=j | y_{N}; \theta^{(l)}) = \sum_{j=1}^k p_{ij}^{(l+1)} \mu_i / f(y_{N}; \theta^{(l)}) \quad (A.9)$$

위 식 (A.9)에 의해 $\sum_{j=1}^k p_{ij}^{(l+1)} = 1$ 이므로 다음 식을 얻는다.

$$\sum_{n=2}^N \Pr(s_{n-1}=i | y_{N}; \theta^{(l)}) = \mu_i / f(y_{N}; \theta^{(l)}) \quad (A.10)$$

위 식 (A.8)과 (A.10)을 결합하여 $\mu_i / f(y_{N}; \theta^{(l)})$ 를 소거하면 \hat{p}_{ij} 의 추정치 계산식을 다음과 같이 구할 수 있다.

$$\hat{p}_{ij}^{(l+1)} = \sum_{n=2}^N \Pr(s_{n-1}=i, s_n=j | y_{N}; \theta^{(l)}) / \sum_{n=2}^N \Pr(s_{n-1}=i | y_{N}; \theta^{(l)}) \quad (A.11)$$

부록 2

먼저 식 (5)를 \mathcal{L} 라 정의하면,

$$\begin{aligned} \left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)^2 &= \left(\sum_{n=1}^N I(s_n=i) \left(\frac{1}{\lambda_i} - y_n\right)\right)^2 \\ &= \sum_{n=1}^N I(s_n=i) \left(\frac{1}{\lambda_i} - y_n\right)^2 + CPT \end{aligned} \quad (\text{A.12})$$

위에서 CPT 는 제품에 의해 생기는 제품합 이외의 항(cross product term)을 의미한다. 위 식 (A.12)에 기대값을 취하고 $\theta = \hat{\theta}$ 을 대입하여 $E_0\left(\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)^2\right)$ 을 구해 보자. 이때, $E_0(CPT) = 0$ 이 된다. 왜냐하면 지시 변수에 의해 s_n 값이 특정한 값을 가진다는 조건하에서는 y_n 들은 조건부 독립이므로 $E_0(CPT) = 0$ 이 된다. 따라서, $E_0\left(\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)^2\right)$ 은 다음 식 (A.13)과 같다.

$$\begin{aligned} E_0\left(\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)^2\right) &= \sum_{n=1}^N \Pr(s_n=i | y_n; \hat{\theta}) E_0\left(\frac{1}{\lambda_i} - y_n\right)^2 \\ &= \sum_{n=1}^N \Pr(s_n=i | y_n; \hat{\theta}) \cdot \frac{1}{\lambda_i^2} \end{aligned} \quad (\text{A.13})$$

이때, $i=1, \dots, k$ 이다.

이제 $E_0\left(\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)\left(\frac{\partial \mathcal{L}}{\partial \lambda_j}\right)\right)$ 를 구해 보자. 먼저 $\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)\left(\frac{\partial \mathcal{L}}{\partial \lambda_j}\right)$ 는 다음과 같다.

$$\begin{aligned} \left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)\left(\frac{\partial \mathcal{L}}{\partial \lambda_j}\right) &= \left(\sum_{n=1}^N I(s_n=i) \left(\frac{1}{\lambda_i} - y_n\right)\right) \cdot \left(\sum_{n=1}^N I(s_n=j) \left(\frac{1}{\lambda_j} - y_n\right)\right) \\ &= \sum_{n=1}^N I(s_n=i) I(s_n=j) \left(\frac{1}{\lambda_i} - y_n\right) \left(\frac{1}{\lambda_j} - y_n\right) + \text{앞 식 이외의 } CPT \end{aligned} \quad (\text{A.14})$$

위 식 (A.14)에서 n 이 동일한 CPT 항은 지시 변수 값이 동시에 서로 다른 값(i, j)를 가질 수 없기 때문에 0이 된다. 그리고, s_n 값이 특정한 값을 가진다면 y_n 들은 조건부 독립이 된다. 이러한 성질을 이용하여 $E_0\left(\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)\left(\frac{\partial \mathcal{L}}{\partial \lambda_j}\right)\right)$ 를 계산해 보면 다음 식 (A.15)와 같다.

$$\begin{aligned} E_0\left(\left(\frac{\partial \mathcal{L}}{\partial \lambda_i}\right)\left(\frac{\partial \mathcal{L}}{\partial \lambda_j}\right)\right) &= E_0\left\{\sum_{n=1}^N \sum_{n'=1}^N I(s_n=i) I(s_{n'}=j) \left(\frac{1}{\lambda_i} - y_n\right) \left(\frac{1}{\lambda_j} - y_{n'}\right)\right\} \\ &= \sum_{n=1}^N \sum_{n'=1}^N \Pr(s_n=i) \Pr(s_{n'}=j) E_0\left\{\left(\frac{1}{\lambda_i} - y_n\right) \left(\frac{1}{\lambda_j} - y_{n'}\right)\right\} \\ &= 0 \end{aligned} \quad (\text{A.15})$$