

이 문서는 Minitab 통계 소프트웨어의 보조 도구에서 사용되는 방법과 데이터 검사를 개발하기 위해 Minitab 통계 학자들이 실시한 연구에 대해 설명하는 전체 백서 중 하나입니다.

# 다중 회귀 분석

## 개요

보조 도구의 다중 회귀 분석 절차에서는 최소 제곱 추정 방법을 사용하여 예측 변수(X)가 최대 5개이고 연속 반응 변수가 1개인 선형 및 2차 모형을 적합합니다. 사용자가 모형 유형을 선택하고 보조 도구가 모형 항을 선택합니다. 이 문서에서는 보조 도구에서 회귀 모형을 선택하기 위해 사용하는 기준에 대해 설명합니다.

또한 Minitab에서는 유효한 회귀 모형을 얻는 데 중요한 여러 요인을 조사합니다. 첫째, 표본은 검정에 충분한 검정력을 제공하고 X와 Y 간 관계의 강도에 대해 충분히 정밀한 추정치를 제공할 수 있을 만큼 충분히 커야 합니다. 그 다음으로, 분석 결과에 영향을 미칠 수 있는 비정상적인 데이터를 식별하는 것이 중요합니다. 또한 오류 항이 정규 분포를 따른다는 가정도 고려하고 전체 모형의 가설 검정에 대한 비정규성의 영향을 평가합니다.

보조 도구에서는 이러한 요인에 따라 데이터에 대해 다음과 같은 검사를 자동으로 수행하고 검사 결과를 보고서 카드에 표시합니다.

- 데이터 양
- 비정상적인 데이터
- 정규성

이 문서에서는 이러한 요인이 실제 회귀 분석과 어떤 관계가 있는지 조사하고 보조 도구에서 해당 요인을 확인하기 위한 가이드라인을 정한 방법에 대해 설명합니다.

# 회귀 분석 방법

## 모형 선택

보조 도구의 회귀 분석에서는 연속형 반응 변수가 1개이고 예측 변수가 2-5개인 모형을 적합합니다. 예측 변수 중 1개는 범주형일 수 있습니다. 다음 두 가지 유형의 모형 중에서 하나를 선택할 수 있습니다.

- 선형:  $F(x) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$
- 2차:  $F(x) = \beta_0 + \sum_i \beta_i X_i + \sum_i \beta_{ii} X_i^2 + \sum_{i < j} \beta_{ij} X_i X_j$

보조 도구는 완전 선형 또는 2차 모형에서 모형 항을 선택합니다.

## 목적

Minitab에서는 모형 선택에 사용할 수 있는 여러 가지 방법을 조사하여 보조 도구에서 사용할 방법을 결정하고자 했습니다.

## 방법

Minitab에서는 후진, 전진, 단계별 등 세 가지 모형 선택 유형을 조사했습니다. 이러한 모형 선택 유형에는 Minitab에서 조사한 다음과 같은 여러 가지 옵션도 포함됩니다.

- 모형에 항을 추가하거나 제거하는 데 사용되는 기준
- 특정 항을 모형에 강제로 추가하거나 초기 모형에 포함시키는지 여부
- 모형의 계층 구조
- 모형 내 X 변수의 표준화

Minitab에서는 이러한 옵션을 검토하고 해당 절차의 결과에 대한 옵션의 영향을 확인하며 실무자들이 선호하는 방법을 고려했습니다.

## 결과

보조 도구에서 모형 항을 선택하기 위해 사용한 절차는 다음과 같습니다.

- 단계적 모형 선택이 사용됩니다. 잠재적 X 변수들은 종종 상관 관계가 있으므로 한 항의 효과가 모형에 포함되어 있는 다른 항에 따라 다릅니다. 단계적 선택은 한

단계에서 항을 추가하고 나중에 모형에 포함되어 있는 다른 항에 따라 제거할 수 있기 때문에 이러한 조건에서는 최선의 방법일 것입니다.

- 모형의 계층 구조는 각 단계에서 유지되며 동일한 단계에서 여러 항을 모형에 추가할 수 있습니다. 예를 들어, 가장 유의한 항이  $X_1^2$ 인 경우 이 항은  $X_1$ 과 함께 추가되며  $X_1$ 이 유의한지 여부는 관계 없습니다. 계층 구조는 모형을 표준화된 단위에서 표준화되지 않은 단위로 변환할 수 있기 때문에 바람직합니다. 또한 계층 구조의 경우 모든 단계에서 모형에 여러 항을 추가할 수 있기 때문에 연관된 선형 항과 반응 간에 강력한 관계가 없는 경우에도 중요한 제곱 또는 교호작용 항을 식별할 수 있습니다.
- 항은  $\alpha = 0.10$ 을 기준으로 모형에 추가되거나 제거됩니다.  $\alpha = 0.10$ 을 사용하면 절차가  $\alpha = 0.15$ 를 사용하는 중심 Minitab의 단계적 절차보다 더 선택적이 됩니다.
- 모형 항 선택을 위해 예측 변수는 평균을 빼고 표준 편차로 나눠 표준화됩니다. 최종 모형은 표준화되지 않은 X의 단위로 표시됩니다. X를 표준화하면 선형 및 제곱 항 간의 상관 관계 대부분이 제거되어 불필요하게 고차항을 추가할 확률이 감소합니다.

# 데이터 검사

## 데이터 양

검정력은 귀무 가설이 거짓일 때 귀무 가설을 기각할 확률과 관계가 있습니다. 회귀 분석의 경우 귀무 가설은 X와 Y 간에 관계가 없다는 것입니다. 데이터 집합이 너무 작은 경우 검정의 검정력이 X와 Y 간에 실제로 존재하는 관계를 탐지하기에 적절하지 않을 수도 있습니다. 따라서 데이터 집합은 실제로 중요한 관계를 높은 확률로 탐지할 수 있을 만큼 충분히 커야 합니다.

### 목적

Minitab에서는 데이터 양이 X와 Y 간의 관계 및  $R^2_{\text{수정}}$ 의 정밀도에 대한 전체 F-검정의 검정력, X와 Y 간 관계의 강도 추정치에 미치는 영향을 확인하고자 했습니다. 이 정보는 데이터 집합이 데이터에서 관측된 관계의 강도가 관계의 기본 강도를 나타낸다고 확신할 수 있을 만큼 충분히 큰지 확인하는 데 중요합니다.  $R^2_{\text{수정}}$ 에 대한 자세한 내용은 부록 A를 참조하십시오.

### 방법


Minitab에서는 비슷한 방법을 사용하여 권장되는 표본 크기를 정하고 단순 회귀 분석에 사용했습니다.  $R^2_{\text{수정}}$  값의 변동성을 조사하여  $R^2_{\text{수정}}$ 을  $\rho^2_{\text{수정}}$ 에 가깝게 유지하기 위한 표본 크기를 정했습니다. 또한 권장되는 표본 크기가 Y와 X 변수 간 관계의 강도가 상당히 약한 경우에도 적절한 검정력을 제공한다는 것을 확인했습니다. 계산에 대한 자세한 내용은 부록 B를 참조하십시오.

### 결과

단순 회귀 분석과 마찬가지로,  $R^2_{\text{수정}}$ 의 관측치가  $\rho^2_{\text{수정}}$ 의 0.20 내에 들어간다고 90% 확신할 수 있을 만큼 충분히 큰 표본을 사용하는 것이 좋습니다. 모형에 항을 추가하면 필요한 표본 크기가 증가합니다. 따라서 각 모형 크기에 필요한 표본 크기를 계산했습니다. 권장되는 크기는 가장 가까운 5의 배수로 올립니다. 예를 들어, 모형에 상수 외에 선형 항 4개, 교호작용 항 3개, 제곱 항 1개 등 8개의 계수가 있는 경우 기준을 충족하기 위해 필요한 최소 표본 크기  $n = 49$ 입니다. 보조 도구는 이 값을 권장되는 표본 크기  $n = 50$ 으로 올립니다. 항의 수를 기준으로 권장되는 표본 크기에 대한 자세한 내용은 부록 B를 참조하십시오.

Minitab에서는 또한 권장되는 표본 크기가 충분한 검정력을 제공한다는 것을 확인했습니다. 또한 상당히 약한 관계( $\rho_{\text{수정}}^2 = 0.25$ )의 경우 검정력은 일반적으로 약 80% 이상입니다. 따라서 보조 도구에서 권장되는 표본 크기를 사용하면 관계의 강도 추정 시 적절한 검정력과 정밀도를 얻을 수 있습니다.

이러한 결과를 바탕으로 데이터 양을 확인하는 경우 보조 도구의 보고서 카드에는 다음과 같은 정보가 표시됩니다.

상태	조건
	<p><b>표본 크기 &lt; 권장되는 크기</b></p> <p>표본 크기가 관계의 강도에 대해 매우 정밀한 추정치를 제공하기에 충분히 크지 않습니다. R-제곱 및 수정된 R-제곱과 같은 측정값이 크게 달라질 수 있습니다. 이 크기의 모형에 대해 정밀한 추정치를 얻으려면 더 큰 표본을 사용해야 합니다.</p> <p><b>표본 크기 &gt;= 권장되는 크기</b></p> <p>표본이 관계의 강도에 대한 매우 정밀한 추정치를 얻기에 충분히 큼니다.</p>

## 비정상적인 데이터

보조 도구 회귀 분석 절차에서는 비정상적인 데이터를 표준화 잔차 또는 레버리지 값이 큰 관측치로 정의합니다. 이 방법은 일반적으로 회귀 분석에서 비정상적인 데이터를 식별하기 위해 사용됩니다(Neter et al., 1996). 비정상적인 데이터가 결과에 중대한 영향을 미칠 수 있기 때문에 유효한 분석을 위해 데이터를 수정해야 할 수도 있습니다. 그러나 공정의 본래 변동에 따라 비정상적인 데이터가 발생할 수도 있습니다. 따라서 이러한 데이터 점을 처리하는 방법을 정하려면 비정상적인 동작의 원인을 식별하는 것이 중요합니다.

### 목적

Minitab에서는 데이터 점이 비정상적이라는 신호를 보내기 위해 표준화 잔차 및 레버리지 값이 얼마나 커야 하는지 확인하고자 했습니다.

### 방법

Minitab에서는 Minitab의 표준 회귀 분석 절차(통계분석 > 회귀 분석 > 회귀 분석)를 바탕으로 비정상적인 관측치를 식별하는 지침을 개발했습니다.

## 결과

### 표준화 잔차



표준화 잔차는 잔차  $e_i$ 를 해당 표준 편차의 추정치로 나눈 값과 같습니다. 일반적으로, 관측치는 표준화 잔차의 절대값이 2보다 큰 경우 비정상적인 것으로 간주됩니다. 그러나 이 지침은 약간 보수적입니다. 모든 관측치의 약 5%가 우연히 이 기준을 충족할 것으로 예상됩니다(오차가 정규 분포를 따르는 경우). 따라서 관측치가 실제로 비정상적인지 확인하기 위해서는 비정상적인 동작의 원인을 조사하는 것이 중요합니다.

### 레버리지 값

레버리지 값은 관측치의 X 값에만 관련이 있으며 Y 값에는 종속되지 않습니다. 관측치는 레버리지 값이 모형 계수의 수(p)를 관측치 수(n)로 나눈 값의 3배보다 크면 비정상적인 것으로 간주됩니다. 또한 일부 교과서에서는  $\frac{2 \times p}{n}$ 를 사용하지만 이 값이 일반적으로 기준 값으로 사용됩니다(Neter et al., 1996).

데이터에 높은 레버리지 점이 포함되어 있는 경우 이 점이 데이터를 적합하기 위해 선택된 모형에 불필요한 영향을 미치는지 여부를 확인하십시오. 예를 들어, 극단적인 X 값 하나 때문에 선형 모형 대신 2차 모형을 선택하게 될 수 있습니다. 2차 모형의 관측된 곡면성이 공정에 대한 이해와 일관성이 있는지 여부를 확인해야 합니다. 일관성이 없는 경우 더 간단한 모형을 데이터에 적합하거나 추가 데이터를 수집하여 공정을 더 철저히 조사하십시오.

비정상적인 데이터를 확인하는 경우 보조 도구의 보고서 카드에는 다음과 같은 상태가 표시됩니다.

상태	조건
	비정상적인 데이터 점이 없습니다.
	하나 이상의 큰 표준화 잔차 또는 하나 이상의 높은 레버리지 점이 있습니다.

## 정규성

회귀 분석에서 일반적인 가정은 랜덤 오차( $\epsilon$ )가 정규 분포를 따른다는 것입니다. 정규성 가정은 계수의 추정치( $\beta$ )에 대한 가설 검정을 수행하는 경우 중요합니다. 랜덤 오차가 정규 분포를 따르지 않는 경우에도 표본이 충분히 크면 검정 결과를 일반적으로 신뢰할 수 있습니다.

## 목적

Minitab에서는 정규 분포를 바탕으로 신뢰할 수 있는 결과를 제공하기 위해 필요한 표본 크기를 확인하고자 했습니다. 실제 검정 결과가 검정의 목표 유의 수준(알파 또는 제1종 오류율)과 얼마나 가깝게 일치했는지, 즉 여러 정규 분포에 대해 검정이 예상된 것보다 더 자주 또는 덜 자주 정규성의 귀무 가설을 잘못 기각했는지 여부를 확인하고자 했습니다.

## 방법



제1종 오류율을 추정하기 위해 정규 분포에서 크게 벗어난 치우친 분포, 두꺼운 꼬리를 갖는 분포 및 가는 꼬리를 갖는 분포를 사용하여 여러 시뮬레이션을 수행했습니다. 15의 표본 크기를 사용하여 시뮬레이션을 수행했습니다. 그리고 여러 모형에 대해 전체 F-검정을 조사했습니다.

각 조건에 대해 10,000번의 검정을 수행했습니다. Minitab에서는 각 검정에 대해 귀무 가설이 참이 되도록 랜덤 데이터를 생성했습니다. 그런 다음 0.10의 목표 유의 수준을 사용하여 검정을 수행했습니다. 10,000번 중에서 검정이 귀무 가설을 실제로 기각한 횟수를 집계하고 이 비율을 목표 유의 수준과 비교했습니다. 검정이 제대로 수행되는 경우 제1종 오류율은 목표 유의 수준에 매우 가깝습니다. 시뮬레이션에 대한 자세한 내용은 부록 C를 참조하십시오.

## 결과

전체 F-검정 모두에 대해 통계적으로 유의한 결과를 찾을 확률은 다른 비정규 분포와 크게 다르지 않습니다. 제1종 오류율은 모두 0.08820과 0.11850 사이로, 목표 유의 수준 0.10에 상당히 가깝습니다.

표본 크기가 비교적 작아도 검정이 제대로 수행되기 때문에 보조 도구에서는 데이터의 정규성을 검정하지 않습니다. 대신, 보조 도구에서는 표본 크기를 확인하고 표본이 15보다 작은 경우 알려줍니다. 보조 도구의 보고서 카드에는 회귀 분석에 대해 다음과 같은 상태가 표시됩니다.

상태	조건
	표본 크기가 15 이상이므로 정규성은 문제되지 않습니다.
	표본 크기가 15보다 작으므로 정규성이 문제될 수 있습니다. p-값을 해석할 때는 주의해야 합니다. 표본이 작으면 p-값의 정확성이 비정규 잔차 오류의 영향을 받기 쉽습니다.

# 참고 문헌

Neter, J., Kutner, M.H., Nachtsheim, C.J., & Wasserman, W. (1996). *Applied linear statistical models*. Chicago: Irwin.



# 부록 A: 모형 및 통계량

예측 변수 X와 반응 변수 Y의 관계를 보여주는 회귀 모형 형식은 다음과 같습니다.

$$Y = f(X) + \varepsilon$$

여기서 함수  $f(X)$ 는 주어진 X에 대한 Y의 기대값(평균)을 나타냅니다.

보조 도구에서는 두 가지 형식의 함수  $f(X)$ 를 선택할 수 있습니다.

모형 유형	$f(X)$
선형	$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$
2차	$\beta_0 + \sum_i \beta_i X_i + \sum_i \beta_{ii} X_i^2 + \sum_{i < j} \beta_{ij} X_i X_j$

계수  $\beta$ 의 값은 알려져 있지 않으며 데이터로부터 추정해야 합니다. 추정 방법은 최소 제곱으로 표본 내 잔차 제곱합을 최소화합니다.

$$\min \sum_{i=1}^n (Y_i - \hat{f}(X_i))^2.$$

잔차는 관측된 반응  $Y_i$ 과 추정된 계수를 바탕으로 한 적합치  $\hat{f}(X_i)$  간의 차이입니다. 이 제곱합의 최소값이 주어진 모형에 대한 SSE(오차 제곱합)입니다.

## 전체 F-검정

이 방법은 전체 모형(선형 또는 2차)의 검정입니다. 선택된 회귀 함수  $f(X)$  형식에 대해 다음 사항을 검정합니다.

$H_0$ :  $f(X)$  가 일정함

$H_1$ :  $f(X)$  가 일정하지 않음

## 수정된 $R^2$

수정된  $R^2$  ( $R^2_{\text{수정}}$ )은 모형에서 반응의 변동성 중 X로 인한 부분을 측정합니다. X와 Y 간의 관측된 관계를 측정하는 데는 두 가지 일반적인 방법이 있습니다.

$$R^2 = 1 - \frac{SSE}{SSTO}$$

및

$$R^2_{\text{수정}} = 1 - \frac{SSE/(n-p)}{SSTO/(n-1)}$$

설명

$$SSTO = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

SSTO는 총 제곱합으로, 전체 평균  $\bar{Y}$ 에 대한 반응의 변동을 측정하며, SSE는 회귀 함수  $f(X)$ 에 대한 변동을 측정합니다.  $R^2_{\text{수정}}$ 에서의 수정은 전체 모형 내 계수의 수( $p$ )에 대한 것으로,  $\varepsilon$ 의 분산을 추정할 수 있도록  $n - p$  자유도가 남게 됩니다.  $R^2$ 은 모형에 계수가 추가되는 경우 감소하지 않습니다. 그러나  $R^2_{\text{수정}}$ 은 수정으로 인해 계수를 추가해도 모형이 개선되지 않는 경우 감소할 수 있습니다. 따라서 모형에 항을 추가해도 반응의 추가 분산이 설명되지 않는 경우  $R^2_{\text{수정}}$ 은 감소하며, 추가 항이 유용하지 않다는 것을 나타냅니다. 따라서 수정된 측정값을 사용하여 다른 크기의 모형을 비교해야 합니다.

## F-검정과 $R^2_{\text{수정}}$ 의 관계

전체 모형의 검정에 대한 F-통계량은  $R^2_{\text{수정}}$  계산에도 사용되는 SSE 및 SSTO의 측면에서 표현할 수 있습니다.

$$F = \frac{(SSTO - SSE)/(p-1)}{SSE/(n-p)}$$
$$= 1 + \left(\frac{n-1}{p-1}\right) \frac{R^2_{\text{수정}}}{1-R^2_{\text{수정}}}$$

위의 공식을 보면 F-통계량이  $R^2_{\text{수정}}$ 의 증가 함수임을 알 수 있습니다.

따라서 검정은  $R^2_{\text{수정}}$ 가 검정의 유의 수준( $\alpha$ )에 의해 지정된 특정 값을 초과한 경우에만  $H_0$ 을 기각합니다.

# 부록 B: 데이터 양

이 항목에서는 관측치의 수  $n$ 이 전체 모형 검정의 검정력 및 모형 강도의 추정치  $R^2_{\text{수정}}$ 의 정밀도에 미치는 영향에 대해 설명합니다.

관계의 강도를 정량화하기 위해 표본 통계량  $R^2_{\text{수정}}$ 의 모집단 통계량으로 새로운 양인  $\rho^2_{\text{수정}}$ 을 소개합니다. 다음 사항을 기억하십시오.

$$R^2_{\text{수정}} = 1 - \frac{SSE/(n-p)}{SSTO/(n-1)}$$

따라서 다음과 같이 정의합니다.

$$\rho^2_{\text{수정}} = 1 - \frac{E(SSE|X)/(n-p)}{E(SSTO|X)/(n-1)}$$

연산자  $E(\cdot|X)$ 는 기대값 또는 주어진  $X$  값에 대한 랜덤 변수의 평균을 나타냅니다. 올바른 모형이 독립적으로 동일하게 분포된  $\varepsilon$ 이 있는  $Y = f(X) + \varepsilon$ 라고 가정하면 다음과 같은 결과를 얻게 됩니다.

$$\begin{aligned} \frac{E(SSE|X)}{n-p} &= \sigma^2 = \text{Var}(\varepsilon) \\ \frac{E(SSTO|X)}{n-1} &= \sum_{i=1}^n \frac{(f(X_i) - \bar{f})^2}{(n-1)} + \sigma^2 \end{aligned}$$

설명  $\bar{f} = \frac{1}{n} \sum_{i=1}^n f(X_i)$ .

따라서

$$\rho^2_{\text{수정}} = \frac{\sum_{i=1}^n (f(X_i) - \bar{f})^2 / (n-1)}{\sum_{i=1}^n (f(X_i) - \bar{f})^2 / (n-1) + \sigma^2}$$

## 전체 모형 유의성

전체 모형의 통계적 유의성을 검정하는 경우, Minitab에서는 랜덤 오차  $\varepsilon$ 이 독립적이며 정규 분포를 따른다고 가정합니다. 따라서  $Y$ 의 평균이 일정한 ( $f(X) = \beta_0$ )이라는 귀무 가설 하에서 F-검정 통계량은  $F(p-1, n-p)$  분포를 따릅니다. 대립 가설 하에서는 F-통계량이 비중심 모수를 가진 비중심  $F(p-1, n-p, \theta)$  분포를 따릅니다.

$$\theta = \sum_{i=1}^n (f(X_i) - \bar{f})^2 / \sigma^2$$

$$= \frac{(n-1)\rho_{\text{수정}}^2}{1-\rho_{\text{수정}}^2}$$

H<sub>0</sub>을 기각하는 확률은 비중심 모수와 함께 증가하며, 비중심 모수는 n 및  $\rho_{\text{수정}}^2$  에서 모두 증가합니다.

## 관계의 강도

단순 회귀 분석에서 이미 확인한 것처럼 데이터에 통계적으로 유의한 관계가 있어도 X와 Y 간에 반드시 강력한 관계가 있는 것은 아닙니다. 따라서 많은 사용자들이 관계가 실제로 얼마나 강력한지 알려주는  $R_{\text{수정}}^2$  과 같은 표시자를 찾습니다.  $R_{\text{수정}}^2$  을  $\rho_{\text{수정}}^2$  의 추정치로 고려하는 경우 추정치가 참  $\rho_{\text{수정}}^2$  값에 상당히 가깝다고 확신할 수 있습니다.

가능한 각 모형 크기에 대해 Minitab에서는 10% 이하의 확률로 절대 차이  $|R_{\text{수정}}^2 - \rho_{\text{수정}}^2|$  가 0.20보다 큰 n의 최소값을 식별하여 허용 가능한 표본 크기에 대한 적절한 임계값을 확인했습니다. 이는  $\rho_{\text{수정}}^2$  의 참 값과 관계가 없습니다. 권장되는 표본 크기 n(T)는 아래 표에 요약되어 있으며, 여기서 T는 모형에서 상수 계수와 다른 계수의 수입니다.

T	n(T)
1-3	40
4-6	45
7-8	50
9-11	55
12-14	60
15-18	65
19-21	70
22-24	75
25-27	80
28-31	85
32-34	90

T	n(T)
35-38	95
39-41	100
42-45	105
46-48	110
49-52	115
53-56	120
57-59	125
60-63	130
64-67	135
68-70	140
71-73	145

Minitab에서는 상당히 약한  $\rho_{\text{수정}}^2 = 0.25$ 의 값에 대해 모형의 전체 F-검정의 검정력을 평가하여 권장되는 표본 크기에서 검정력이 충분하다는 것을 확인했습니다. 아래 표의 모형 크기는 각 n(T) 값에 대한 최악의 경우를 나타냅니다. 동일한 n(T)에서 모형이 작을수록 검정력이 더 큼니다.

T	n(T)	검정력 $\rho_{\text{수정}}^2 = 0.25$
3	40	0.902791
6	45	0.854611
8	50	0.850675
11	55	0.831818
14	60	0.820592
18	65	0.798003

T	n(T)	검정력 $\rho_{\text{수정}}^2 = 0.25$
21	70	0.796425
24	75	0.796911
27	80	0.798856
31	85	0.789861
34	90	0.794367
38	95	0.788625
41	100	0.794511
45	105	0.790864
48	110	0.797487
52	115	0.795250
56	120	0.793698
59	125	0.800982
63	130	0.800230
67	135	0.799906
69	140	0.814664

# 부록 C: 정규성

보조 도구에서 사용되는 회귀 모형 형식은 모두 다음과 같습니다.

$$Y = f(X) + \varepsilon$$

랜덤 항  $\varepsilon$ 에 대한 일반적인 가정은 항들이 독립적이며 평균 0 및 일반 분산  $\sigma^2$ 을 갖는 동일하게 분포된 정규 랜덤 변수라는 것입니다.  $\beta$  모수의 최소 제곱 추정치는  $\varepsilon$ 이 정규 분포를 따른다는 가정을 무시하는 경우에도 여전히 최적의 선형 불편 추정치입니다. 정규성 가정은  $f(X)$ 에 대한 가설 검정에서처럼 이러한 추정치에 확률을 추가하려는 경우에만 중요합니다.

Minitab에서는 정규성 가정을 바탕으로 회귀 분석의 결과를 신뢰할 수 있는  $n$ 의 크기를 확인하고자 했습니다. 시뮬레이션을 수행하여 다양한 비정규 오차 분포 하에 가설 검정의 제1종 오류율을 조사했습니다.

표 1에는 10,000번의 시뮬레이션 중에서 세 가지 모형의 다양한  $\varepsilon$ 의 분포에 대해  $\alpha = 0.10$ 에서 전체 F-검정이 유의한 비율이 나와 있습니다. 이 시뮬레이션에서는 X와 Y 간에 관계가 없다는 귀무 가설이 참이었습니다. X 값은 Minitab의 RANDOM 명령어에 의해 다변량 정규 변수로 생성되었습니다. 모든 검정에 대해  $n=15$ 의 표본 크기를 사용했습니다. 모든 모형에는 5개의 연속형 예측 변수가 포함되었습니다. 첫 번째 모형은 X 변수가 모두 5개인 선형 모형이었습니다. 두 번째 모형에는 선형 및 제곱 항이 모두 있었습니다. 세 번째 모형에는 모든 선형 항과 7개의 이원 교호작용이 있었습니다.

**표 1** 비정규 분포의 경우  $n=15$ 인 전체 F-검정에 대한 제1종 오류율

분포	선형	선형 + 제곱	선형 + 7개의 교호작용
정규	0.09910	0.10270	0.10060
t(3)	0.09840	0.11850	0.11800
t(5)	0.09980	0.10010	0.10430
Laplace	0.09260	0.09400	0.09650
균등	0.10630	0.10080	0.09480
Beta(3, 3)	0.09980	0.10120	0.10020
지수	0.08820	0.09500	0.09960

분포	선형	선형 + 제곱	선형 + 7개의 교호작용
Chi(3)	0.09890	0.11400	0.10970
Chi(5)	0.09730	0.10590	0.10330
Chi(10)	0.10150	0.09930	0.10360
Beta(8, 1)	0.09870	0.10230	0.10490

시뮬레이션 결과를 보면 모든 오차 분포에 대해 통계적으로 유의한 결과를 찾을 확률이 명목 값 0.10과 크게 다르지 않습니다. 관측되는 제1종 오류율은 모두 0.08820과 0.11850 사이입니다.

© 2015, 2017 Minitab Inc. All rights reserved.

Minitab®, Quality. Analysis. Results.® and the Minitab® logo are all registered trademarks of Minitab, Inc., in the United States and other countries. See [minitab.com/legal/trademarks](http://minitab.com/legal/trademarks) for more information.