
공학 수치해석
Numerical Analysis for Engineers

박은천
단국대학교 건축대학

LECTURE NOTE
NOTE-2012-ver-20121211
11 December 2012

EUNCHURN PARK
(eunchurn@kmctech.co.kr)
Structrual Analysis & Dynamics Laboratory
Department of Architectural Engineering
Dankook University School of Architecture

Department of R&D Development
Korea Maintenance Co., LTD.
KM Tower 12F
130-5, Guro-5-dong, Gurogu, Seoul (152-842)
PHN: +82-2-830-7071 (236)
SEL: +82-10-4499-6420
FAX: +82-2-830-5256

공학 수치해석

Numerical Analysis for Engineers

박은천
단국대학교 건축대학

Dept. of Architectural Engineering, Dankook Univ. School of Architecture, NOTE-2012-ver-20121211

11 December 2012

Table of Contents

I 곡선적합		2.0.3 최적적합을 위한 조건	6
(Curve Fitting)	2	2.0.4 직선의 최소제곱적합	7
1 최소제곱회귀분석		2.0.5 선형회귀분석 오차의 정량화 . . .	7
(Lest Square Regression)	3	2.0.6 비선형 관계식의 선형화	8
1.1 통계학 기초	3	2.0.7 다항식 회귀분석	9
1.1.1 정규분포	4		
1.1.2 신뢰구간의 산출	5	II 부록	13
2 선형회귀분석	6	Appendices	13

Part I

곡선적합

(Curve Fitting)

공학에서 측정된 데이터를 사용하여 물리적 모델 혹은 요소들을 추정하는 경우가 대부분이다. 이 때 데이터는 연속적이기 보다는 이산적인 값으로 주어지는 경우가 많고, 이산적인 값 사이에 있는 임의의 점에서의 값을 추정해야 하는 경우가 있다. 또한 복잡한 함수를 단순화된 모델로 만들어야 하는 경우도 있다. 이 경우 임의의 수의 특정한 점들에서 그 함수값을 추적하고, 이들 계산값과 보간법을 이용하여 보다 간단한 형태의 함수로 표현할 수 있다. 이와 같은 응용들을 모두 일컬어 곡선적합(curve fitting)이라고 한다.

곡선적합은 데이터와 관련된 오차의 크기에 따라서 접근방법이 두 가지로 구별된다. 첫 번째 방법은 데이터가 상당한 크기의 오차 또는 "노이즈(noise)"를 내포하고 있을 경우로 데이터의 일반적인 경향을 나타내는 하나의 곡선을 유도해 내는 방법이다. 이 경우, 각각의 데이터값은 정확하지 않을 수 있기 때문에 유도되는 곡선이 모든 데이터점들을 통과하도록 노력할 필요는 없으나, 유도되는 곡선은 각 데이터점들의 경향을 따르도록 결정하여야 한다. 이러한 방식으로 데이터의 경향을 적절하게 표현하는 한 가지 방법으로 최소제곱회귀분석(least-square regression)이 있다. 두 번째 방법은 데이터값들이 상당히 정확하여 각각의 데이터점들을 직접 통과하는 하나의 곡선, 또는 일련의 곡선들을 매끄럽게 연결하는 방법이다. 이 경우의 데이터는 보통 도표에 있는 값들로서, 예를 들면 온도의 함수로 표현되는 물의 밀도에 대한 값, 또는 가스의 비열에 대한 값들을 말한다. 값이 알려진 두 점들 사이에 있는 임의의 점에서 함수값을 추정하는 방법을 보간법(interpolation)이라고 한다.

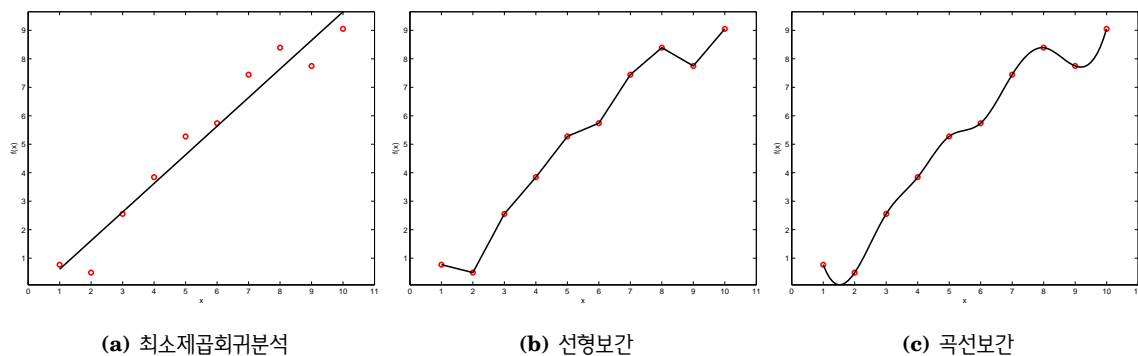


Figure 0.1: 10개의 점들을 최적의 곡선으로 보간하는 세가지 방법

경제공학에서 이윤도표나, 열역학에서의 증기표 그리고 토목공학에서 응력변형도 곡선등과 같이 도표화된 데이터로부터 중간값을 결정하는 방식을 경험하였다. 비록 널리 사영되는 공학적 물성치들이 대부분이 도표화 혹은 공식화되어 있음에도 불구하고 많은 물성치들은 이러한 편리한 형태로 얻을 수 없다. 자신이 직접 데이터를 측정하거나, 예측관계식을 개발하여야 하는 특별한 경우와 같이 새로운 분야가 종종 발생된다. 실험데이터를 보간할 때는 경향분석(trend analysis)이나 가상실험(hypothesis testing)과 같은 두 가지 형태의 응용문제와 마주치게 된다.

경향분석은 예측을 하기 위하여 데이터의 형(pattern)을 분석하는 과정을 나타낸다. 데이터가 높은 정확도로 측정된 경우 보간다항식을 이용할 수 있으나, 정확도가 높지 않은 데이터는 보편적으로 최소제곱회귀분석이 사용된다.

실험적 곡선적합의 두 번째의 공학적 적용은 가상실험이다. 이것은 현존하는 수학적 모델과 측정된 데이터를 비교하는

것이다. 만약 모델의 계수들을 알지 못한다면, 관측된 데이터에 최적으로 적합한 계수들을 결정하여야 한다. 한편, 만약 모델 계수의 추정값이 이미 확보된 경우에는 관측된 값과 모델의 예측값을 서로 비교하여 모델의 적합성을 검사하는 것이 필요할 것이다. 이 때 대안으로 제시되는 모델과 비교되어 실험적인 관측에 근거한 "최적"인 모델이 선택된다.

1 최소제곱회귀분석 (Lest Square Regression)

1.1 통계학 기초

공학적 연구 과정에서 어떤 특정한 양을 여러번 측정하였다고 가정하자.

6.495	6.595	6.615	6.635	6.485	6.555
6.665	6.505	6.435	6.625	6.715	6.655
6.755	6.625	6.715	6.575	6.655	6.605
6.565	6.515	6.555	6.395	6.775	6.685

Table 1: 철 구조물의 열팽창계수 [$\times 10^{-6} \text{in}/(\text{in} \cdot ^\circ \text{F})$] 측정

예로 표1는 철 구조물(steel structure)의 열팽창계수에 관한 24개의 측정된 값들을 나타내고 있다. 표에 있는 값들을 보면 데이터는 최소 6.395로부터 최대 6.775의 범위 내에서 제한된 분량의 정보만을 제공하고 있다. 적절한 통계법을 이용하여 데이터를 요약하는 것으로 데이터 집합의 특수한 성질에 관하여 보다 많은 정보를 제공할 수 있는 추가적인 관찰 결과를 얻을 수 있다. 이러한 분석적인 통계법들은 (1) 데이터 분포의 중심위치 (2) 데이터 집합의 분산 정도등을 나타내는데 주로 사용된다.

가장 보편적인 통계값은 산술평균(arithmetic mean)이다. 어떤 표본의 산술평균(\bar{y})은 다음과 같이 각각의 데이터값(y_i)들을 합하고, 이것을 데이터점의 수(n)으로 나눈 것으로 정의된다.

$$\bar{y} = \frac{\sum y_i}{n} \quad (1.1)$$

여기서 합(이후의 모든 합)은 $i = 1$ 부터 n 까지를 나타낸다. 또한, 가장 보편적인 표본의 분포도는 다음과 같이 정의되는 평균을 중심으로 하는 표준편차(standard deviation, s_y)이다.

$$s_y = \sqrt{\frac{S_t}{n-1}} \quad (1.2)$$

여기서 S_t 는 데이터점들과 평균 사이의 잔차를 제곱한 총합을 말한다. 즉 S_t 는 다음과 같다.

$$S_t = \sum (y_i - \bar{y})^2 \quad (1.3)$$

따라서 만약 개별적인 측정값들이 평균값 주변에 넓게 퍼져 있다면 S_t (결과적으로는 s_y)의 값은 커질 것이다. 만약에 측정값들이 모여 있다면 표준편차는 줄어든다. 분포는 표준편차의 제곱으로 다음과 같이 나타낼 수 있으며, 이것을 분산(variance, s_y^2)이라고 한다.

$$s_y^2 = \frac{\sum (y_i - \bar{y})^2}{n-1} \quad (1.4)$$

여기서 식(1.2)와 식(1.4)에 있는 분모가 $(n-1)$ 인데 이것을 자유도(degree of freedom)이라고 한다. 따라서 s_x 와 s_y 는 $(n-1)$ 자유도에 기초를 두고 있다. 데이터 분포를 수치화하는 데 이용되는 통계값으로는 분산계수(coefficient of variation; c.v.)가 있으며, 이것은 평균에 대한 표준편차의 비율이다. 이 값은 정규화된 분포도를 다음과 같이 백분율의 형태로 나타내고 있다.

$$\text{c.v.} = \frac{s_y}{\bar{y}} \times 100 \quad (1.5)$$

1.1.1 정규분포

데이터가 평균값 주위에 분포되어 있는 형태가 있다. 히스토그램(histogram)은 데이터의 분포에 대한 간단한 시각적 표현을 제공해 주며, 측정값을 구간별로 분류해서 구성하게 된다. 이때 측정단위는 수평축에 나타내고 각 구간의 발생 빈도는 수직축에 나타낸다. 만약 많은 수의 데이터를 갖고 있을 경우 이러한 히스토그램은 매끄러운 곡선으로 근사화될 수가 있다. Figure 1.1.1에 나타난 종 모양으로 된 대칭곡선은 이러한 근사화된 특성 곡선의 하나로서 정규분포(normal

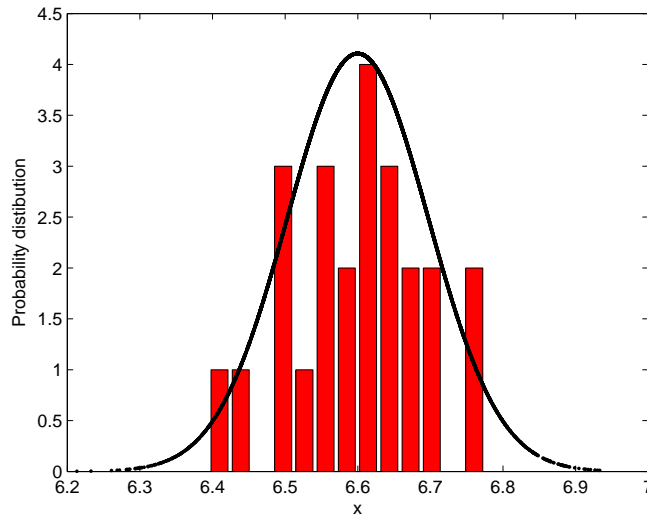


Figure 1.1: 정규분포

distribution)이라고 하며, 데이터점들의 수가 증가할 수록 히스토그램은 결국 이러한 정규분포에 접근하게 된다.

평균, 표준편차, 잔차의 제곱합 그리고 정규분포와 같은 개념들은 모두 공학문제와 큰 연관성을 갖고 있다. 매우 간단한 예로는 이들이 어떤 특정한 측정값에 대한 신뢰도를 정량화시키는데 사용될 수가 있다는 것이다. 만약 어떤 양이 정규적으로 분포되어 있다면 $\bar{y} - s_y$ 와 $\bar{y} + s_y$ 사이의 구간에는 전체 측정 개수의 약 68%를 포함하게 된다. 같은 방식으로 아래의 식들과 같이 정량적인 방법으로 물리적인 양을 가늠할 수 있다.

$$\bar{y} - s_y \leq \text{측정 개수의 68\%} \leq \bar{y} + s_y$$

$$\bar{y} - 1.96s_y \leq \text{측정 개수의 95\%} \leq \bar{y} + 1.96s_y$$

$$\bar{y} - 2.54s_y \leq \text{측정 개수의 99\%} \leq \bar{y} + 2.54s_y$$

예를들어 Table 1에서와 같이 주어진 열팽창계수 데이터($\bar{y} = 6.6$ 과 $s_y = 0.097133$)에 대하여 우리는 약 95%의 측정값들이 6.409619와 6.790381 사이에 존재한다고 할 수 있다. 만약에 어떤 사람이 7.35라는 값을 측정했다고 하면 우리는 그 측정값은 틀렸다고 의심하게 될 것이다. 이러한 측정값에 대한 신뢰도에 대한 평가는 다음절에서 다뤄진다.

1.1.2 신뢰구간의 산출

앞 절에서와 같이 통계의 주목적은 집단으로부터 제한적으로 추출된 표본을 이용하여 그 집단의 특성을 파악하는데 있다. 생산되는 모든 구조물용 강재의 열팽창계수를 측정한다는 것이 불가능함은 명백하다. Table 1와 같이 표본추출에 근거하여 전체 집단의 특성을 규명할 수가 있다.

제한된 표본으로부터 모르는 집단의 물성을 "추론"해야 하기 때문에 이러한 노력을 '통계적인 추론'이라고 한다. 통계적 추론의 결과는 보통 집단에 대한 인자들의 추정값으로 주어지기 때문에 이러한 과정을 '추정'이라고도 한다.

기호 \bar{y} 및 s_y 등은 표본의 평균 및 표준편차를 나타내고, 기호 μ 및 σ 등은 집단의 평균 및 표준편차를 각각 나타낸다. 전자의 표현은 "산출(estimated)"된 평균 및 표준편차라고 하고, 후자의 표현은 종종 "참(true)"평균 및 표준편차라고 한다.

구간추정자(interval estimator)는 그 인자가 주어진 확률로 갖게 되는 값의 범위를 제시해 준다. 그러한 구간들은 한쪽이거나 양쪽인 것으로 설명될 수 있다. 한쪽 구간이라 함은 인자의 산출이 참값보다 작거나 또는 큰 신뢰도를 나타낸다. 반면에 양쪽 구간은 산출값이 항상 참값과 일치하게 되는 보다 일반적인 상황을 다루게 된다. 양쪽 구간이 보다 일반적이기 때문에 여기서는 이에 초점을 맞추어 설명한다. 양쪽 구간은 다음과 같이 기술 될 수 있다.

$$P\{L \leq \mu \leq U\} = 1 - \alpha \quad (1.6)$$

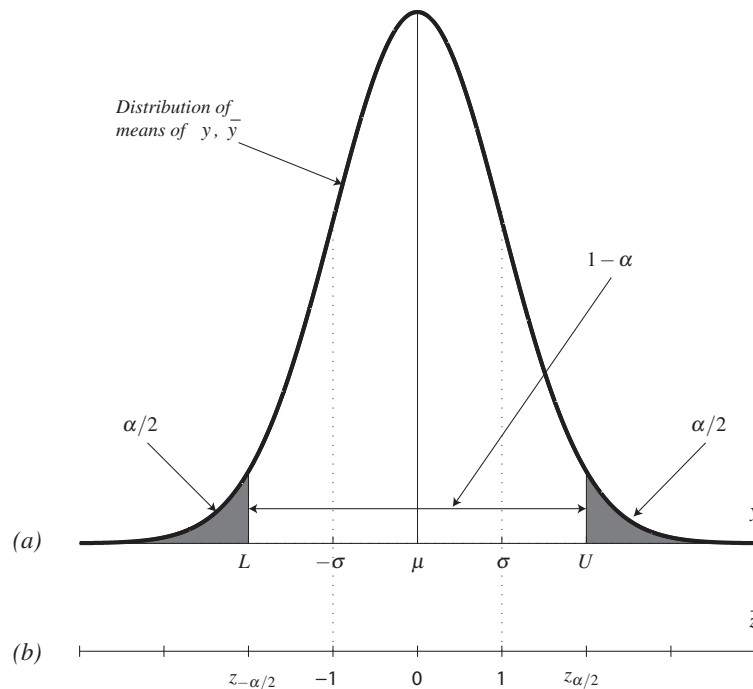


Figure 1.2: 양쪽 구간 신뢰도 범위

이것은 "y의 참평균 μ 가 L 과 U 사이의 구간에 놓일 수 있는 확률이 $(1 - \alpha)$ 이다"라고 읽는다. 여기서 α 값을 유의수준 (significant level)이라고 부르며, 따라서 신뢰구간 (confidence level)을 결정하는 문제는 L 과 U 를 구하는 것으로 축소된다. 필수조건은 아니지만 Figure 1.1.2과 같이 분포의 각 끝단 꼬리에 같은 크기로 $\alpha/2$ 씩 분포된 α 확률을 갖는 양쪽 구간을 관찰하는 것이 보편적이다.

y의 분포에 대한 참분산 σ^2 을 알수 없지만 알고있다고 가정하면 표본의 평균 \bar{y} 는 평균이 μ 이고 분산값이 σ^2/n 인

정규분포로부터 얻어질 수가 있다. Figure 1.1.2에 설명된 문제에서 실제적인 평균 μ 를 모르고 있다. 그러므로 표본의 평균 \bar{y} 와 관련해서 정규곡선이 정확하게 어디에 놓이게 되는지를 알지 못한다. 따라서 표준정규추정(standard normal estimate)로 새로운 값을 산출한다.

$$\bar{z} = \frac{\bar{y} - \mu}{\sigma/\sqrt{n}} \quad (1.7)$$

이것은 \bar{y} 와 μ 사이의 정규화된 거리를 나타낸다. 통계학적 이론에 의하면 이 값은 평균 0이고, 분산값이 1인 정규분포가 되어야 한다. 더욱이 \bar{z} 가 Figure 1.1.2에서 면적표시가 되지 않은 영역에 떨어질 확률은 $(1 - \alpha)$ 이어야 한다. 그러므로 α 의 확률로써 다음과 같은 표현이 가능하다.

$$\begin{aligned} \frac{\bar{y} - \mu}{\sigma/\sqrt{n}} &< -z_{\alpha/2} \\ \frac{\bar{y} - \mu}{\sigma/\sqrt{n}} &> z_{\alpha/2} \end{aligned}$$

$z_{\alpha/2}$ 는 표준정규 임의변수(standard normal random variable)이다. 이 값은 $(1 - \alpha)$ 확률에 걸쳐 있으며, 평균을 전후해서 무차원화된 좌표축을 따라 측정된 거리를 나타낸다. 이 값은 각 상용프로그램으로도 구할 수 있다. 예를들어 $\alpha = 0.05$ 즉 95%의 확률에 걸쳐구간 $z_{\alpha/2}$ 는 약 1.96과 같다. 이런 결과는 확률 $(1 - \alpha)$ 로써 다음과 같이 쓸 수 있다.

$$L \leq \mu \leq U$$

여기서, L 및 U 는 각각 다음과 같이 정의된다.

$$\begin{aligned} L &= \bar{y} - \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \\ U &= \bar{y} + \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \end{aligned}$$

2 선형회귀분석

최소제곱 근사의 간단한 방법으로 관측치에 직선으로 적합시키는 것이다.

$$y = a_0 + a_1x + e \quad (2.1)$$

여기서, a_0 는 절편(intercept), a_1 는 기울기(slope) 그리고 e 는 관측값과 모델값의 차이로 오차(error)이다. 오차는 식 (2.1)를 변형하여

$$e = y - a_0 - a_1x \quad (2.2)$$

2.0.3 최적적합을 위한 조건

(1) 데이터를 통과하는 최적의 직선(best fit)을 구하는 방법은 모든 주어진 데이터에 대한 오차의 합을 최소화 시키는 것이다.

$$\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - a_0 - a_1x_i) \quad (2.3)$$

문제점?

(2) 오차의 절대값의 합을 최소화 하는 방법

$$\sum_{i=1}^n |e_i| = \sum_{i=1}^n |y_i - a_0 - a_1x_i| \quad (2.4)$$

문제점?

(3) 최소-최대(minimax) 판별조건은 직선으로부터 떨어진 각 점들의 최대 변위가 최소가 되도록 선택하는 것이다.

위에서 언급한 방법들의 단점을 극복하기 위하여 측정된 y 와 선형 모델을 이용하여 계산된 y 사이의 잔차에 대한 제곱의 합을 최소화하는 방법이 고안됨.

$$\begin{aligned}
 S_r &= \sum_{i=1}^n e_i^2 \\
 &= \sum_{i=1}^n (y_{i,measured} - y_{i,model})^2 \\
 &= \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2
 \end{aligned} \tag{2.5}$$

2.0.4 직선의 최소제곱적합

a_0 와 a_1 의 값을 결정하기 위하여 식(2.5)은 각각의 계수에 대하여 편미분을 취한다.

$$\begin{aligned}
 \frac{\partial S_r}{\partial a_0} &= -2 \sum (y_i - a_0 - a_1 x_i) \\
 \frac{\partial S_r}{\partial a_1} &= -2 \sum [(y_i - a_0 - a_1 x_i) x_i] \\
 \sum y_i - \sum a_0 - \sum a_1 x_i &= 0 \\
 \sum y_i x_i - \sum a_0 x_i - \sum a_1 x_i^2 &= 0
 \end{aligned}$$

$\sum a_0 = na_0$ 이므로 이 식은 a_0 과 a_1 에 대한 2원 1차 연립방정식으로 주어진다.

$$na_0 + (\sum x_i) a_1 = \sum y_i \tag{2.6}$$

$$(\sum x_i) a_0 + (\sum x_i^2) a_1 = \sum x_i y_i \tag{2.7}$$

이들을 정규 방정식(normal equation)이라고 한다. 이들을 연립방정식으로 풀면 a_1 은 다음과 같이 된다.

$$a_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} \tag{2.8}$$

이 결과를 식(2.6)에 대입해서 a_0 를 구하면

$$a_0 = \bar{y} - a_1 \bar{x} \tag{2.9}$$

여기서 \bar{y} 와 \bar{x} 는 각각 y 와 x 의 평균이다.

2.0.5 선형회귀분석 오차의 정량화

최소제곱법으로 얻은 직선은 점들의 경향을 나타내는 "최적"의 유일한 직선이라 할 수 있다. 잔차들이 계산되는 방식을 자세히 분석해보면 보간법의 여러 다른 성질들을 발견할 수 있다. 식(2.5)에서 정의된 제곱합을 상기해보면

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2 \tag{2.10}$$

여기서 식(1.3)과 식(2.10)이 유사함을 보인다. 즉 (1) 데이터와 직선과의 차이는 데이터의 전체 범위에 걸쳐서 유사한 크기를 갖고, (2) 직선을 중심으로 한 데이터점들의 분포는 정규분포를 이룬다. 만약 이들 판별조건이 만족된다면 최소제

곱회귀분석은 a_0 과 a_1 를 구하는 가장 좋은방법이라고 할 수 있다.(Draper and Smith, 1981). 이것을 "최대우도법칙 (maximum likelihood principle)"이라 한다. 회귀분석직선에 대한 표준편차는 다음과 같이 결정된다.

$$s_{y/x} = \sqrt{\frac{S_r}{n-2}} \quad (2.11)$$

식(1.3)의 데이터의 표준편차와 오차의 제곱합 식(2.10)을 이용하여 상대적인 오차로 정규화하면

$$r^2 = \frac{S_t - S_r}{S_t} \quad (2.12)$$

여기서, r^2 를 결정계수(coefficient of determination), r 을 상관계수(correlation coefficient)라 한다. $S_r = 0$, $r = r^2 = 1$ 인 경우, 완전한 적합이므로 데이터는 회귀분석직선에 의하여 100%만족한다. 컴퓨터 계산을 위해 식(2.12)를 다음 식과 같이 사용하는 것이 편리하다.

$$r = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}} \quad (2.13)$$

Algorithm 2.1 선형회귀분석 알고리즘

```

function REGRESSION( $x, y, n, a_1, a_0, s_{y/x}, r^2$ )
     $sumx = 0$   $sumxy = 0$   $st = 0$ 
     $sumy = 0$   $sumx2 = 0$   $sr = 0$ 
    for  $i = 1, n$  do
         $sumx = sumx + x_i$ 
         $sumy = sumy + y_i$ 
         $sumxy = sumxy + x_i * y_i$ 
         $sumx2 = sumx2 + x_i * x_i$ 
    end for
     $x_m = sumx / n$ 
     $y_m = sumy / n$ 
     $a_1 = (n * sumxy - sumx * sumy) / (n * sumx2 - sumx * sumx)$ 
     $a_0 = y_m - a_1 * x_m$ 
    for  $i = 1, n$  do
         $S_t = S_t + (y_i - y_m)^2$ 
         $S_r = S_r + (y_i - a_1 x_i - a_0)^2$ 
    end for
     $s_{y/x} = \sqrt{S_r / (n - 2)}$ 
     $r^2 = (S_t - S_r) / S_t$ 
    return  $a_1, a_0, s_{y/x}, r^2$ 
end function

```

2.0.6 비선형 관계식의 선형화

선형식으로 회귀분석 할 수 있는 비선형 방정식의 형태는 다음과 같은 지수모델로 변형할 수 있다.

$$y = \alpha_1 e^{\beta_1 x} \quad (2.14)$$

여기서 α_1, β_1 은 상수이다. 이 모델은 공학의 여러분야에서 증가하는 양($\beta_1 > 0$), 또는 감소하는 양($\beta_1 < 0$)이 그 자신의 크기에 직접 비례하는 특성을 나타내는 모델이다. 예로 인구 증가모델, 또는 방사능 감소등이다. 또한 역방정식형태를 예를 들 수 있다.

$$y = \alpha_2 x^{\beta_2} \quad (2.15)$$

여기서 α_2, β_2 은 상수이다. 또한 포화성장률 방정식(saturation-growth-rate equation)이 있다.

$$y = \alpha_3 \frac{x}{\beta_3 + x} \quad (2.16)$$

여기서 α_3, β_3 은 상수이다. 이 모델은 제한된 조건하의 인구 성장모델이나 x 가 증가함에 따라 성장이 정지 즉, 포화상태에 이르는 비선형관계식을 나타내고 있다. 이러한 방정식모델은 간단한 조작으로 선형화가 가능하기 때문에 단순 선형회귀 분석으로 보간식을 구할 수 있다. 예를들어 식(2.14)의 양반여 자연로그를 취하면,

$$\begin{aligned} \ln y &= \ln \alpha_1 + \beta_1 x \ln e \\ &= \ln \alpha_1 + \beta_1 x \end{aligned} \quad (2.17)$$

따라서 x 에 대한 $\ln y$ 의 그림은 β_1 인 기울기와 $\ln \alpha_1$ 인 절편을 갖는다. 마찬가지로 식(2.15)은 양변에 상용로그를 취해 선형화할 수 있다.

$$\log y = \beta_2 \log x + \log \alpha_2 \quad (2.18)$$

식(2.16)는 역수를 취하여 다음과 같이 선형화될 수 있다.

$$\frac{1}{y} = \frac{\beta_3}{\alpha_3} \frac{1}{x} + \frac{1}{\alpha_3} \quad (2.19)$$

예제 역 방정식의 선형화

데이터의 log 변환을 이용하여 Table 2에 있는 데이터를 식(2.15)으로 나타내어라.

x	y	$\log x$	$\log y$
1	0.5	0	-0.301
2	1.7	0.301	0.226
3	3.4	0.477	0.534
4	5.7	0.602	0.753
5	8.4	0.699	0.922

Table 2: 역방정식으로 적합되는 데이터

2.0.7 다항식 회귀분석

2차다항식의 모델 방정식

$$y = a_0 + a_1 x + a_2 x^2 + e \quad (2.20)$$

잔차의 제곱합은 다음식과 같다.

$$S_r = \sum_{i=1}^n (y_i - a_0 - a_1 x_i - a_2 x_i^2)^2 \quad (2.21)$$

각 미지계수에 대해 편미분을 수행하면,

$$\begin{aligned}\frac{\partial S_r}{\partial a_0} &= -2 \sum (y_i - a_0 - a_1 x_i - a_2 x_i^2) \\ \frac{\partial S_r}{\partial a_1} &= -2 \sum x_i (y_i - a_0 - a_1 x_i - a_2 x_i^2) \\ \frac{\partial S_r}{\partial a_2} &= -2 \sum x_i^2 (y_i - a_0 - a_1 x_i - a_2 x_i^2)\end{aligned}$$

S_r 이 최소값을 갖는 경우 값을 0으로 두고 정규방정식으로 정리할 수 있다.

$$\begin{aligned}na_0 + (\sum x_i) a_1 + (\sum x_i^2) a_2 &= \sum y_i \\ (\sum x_i) a_0 + (\sum x_i^2) a_1 + (\sum x_i^3) a_2 &= \sum x_i y_i \\ (\sum x_i^2) a_0 + (\sum x_i^3) a_1 + (\sum x_i^4) a_2 &= \sum x_i^2 y_i\end{aligned}$$

행렬로 표시하면 다음과 같다.

$$\begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \begin{Bmatrix} a_0 \\ a_1 \\ a_2 \end{Bmatrix} = \begin{Bmatrix} \sum y_i \\ \sum x_i y_i \\ \sum x_i^2 y_i \end{Bmatrix} \quad (2.22)$$

즉 2차 다항식을 이용한 선형회귀문제는 쉽게 n 차 다항식으로 확장될 수 있다. n 개의 데이터를 m 차 다항식으로 적합을 시키는 과정을 행렬법으로 접근해보자. 우선 1차식의 경우

$$\begin{aligned}y_1 &= a_0 + a_1 x_1 + e_1 \\ y_2 &= a_0 + a_1 x_2 + e_2 \\ y_3 &= a_0 + a_1 x_3 + e_3 \\ &\vdots \\ y_n &= a_0 + a_1 x_n + e_n\end{aligned}$$

행렬식으로 작성하면

$$\begin{Bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{Bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{Bmatrix} a_0 \\ a_1 \end{Bmatrix} + \begin{Bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_n \end{Bmatrix} \quad (2.23)$$

$$\mathbf{Y} = \mathbf{z} \cdot \mathbf{A} + \mathbf{e} \quad (2.24)$$

m 차 다항식의 경우

$$\begin{Bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{Bmatrix} = \begin{bmatrix} z_{11} & \cdots & z_{1m} \\ \vdots & & \vdots \\ z_{n1} & \cdots & z_{nm} \end{bmatrix} \begin{Bmatrix} a_0 \\ \vdots \\ a_m \end{Bmatrix} + \begin{Bmatrix} e_1 \\ e_2 \\ e_3 \\ \vdots \\ e_n \end{Bmatrix} \quad (2.25)$$

$$\mathbf{Y} = \mathbf{z} \cdot \mathbf{A} + \mathbf{e} \quad (2.26)$$

잔차는

$$\mathbf{e} = \mathbf{Y} - \mathbf{z} \cdot \mathbf{A} \quad (2.27)$$

잔차의 제곱합은

$$S_r = \mathbf{e}^\top \mathbf{e} \quad (2.28)$$

$$= (\mathbf{Y} - \mathbf{z} \cdot \mathbf{A})^\top (\mathbf{Y} - \mathbf{z} \cdot \mathbf{A}) \quad (2.29)$$

잔차의 제곱합을 최소화하기 위해 편미분을 취하면,

$$\frac{\partial \{\mathbf{e}^\top \mathbf{e}\}}{\partial \mathbf{A}} = 0 \quad (2.30)$$

결론적으로 벡터미분을 계산하여 구하면,

$$\mathbf{A} = \left[\mathbf{z}^\top \mathbf{z} \right]^{-1} \mathbf{z}^\top \mathbf{Y} \quad (2.31)$$

벡터미분

먼저 $x_1^2 + x_2^2$ 를 벡터 $\begin{bmatrix} x_1 & x_2 \end{bmatrix}^\top$ 에 대해 미분을 하면 다음과 같은 $n \times 1$ 벡터 $\begin{bmatrix} 2x_1 & 2x_2 \end{bmatrix}^\top$ 가 만들어진다. $1 \times m$ 행렬을 $n \times 1$ 벡터로 미분하면 $n \times m$ 행렬이 만들어진다. 즉, scalar를 벡터에 대해 미분하면 벡터의 각 요소로 한번씩 scalar를 미분하여 요소가 생성되는 미분한 벡터와 크기가 같은 벡터가 만들어진다. $1 \times m$ 행벡터를 $n \times 1$ 열벡터로 미분하면 $n \times m$ 행렬이 만들어진다.

$$\mathbf{A} = \begin{bmatrix} x_1^2 + x_2^2 & 2x_1 & 2x_1^3 + x_2 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\frac{\partial \mathbf{A}}{\partial \mathbf{x}} = \begin{bmatrix} 2x_1 & 2 & 6x_1^2 \\ 2x_2 & 0 & 1 \end{bmatrix}$$

따라서 벡터미분은 다음과 같은 성질을 갖는다.

- $$\frac{\partial (\mathbf{x}^\top \mathbf{A})}{\partial \mathbf{x}} = \mathbf{A}$$
- $$\frac{\partial (\mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}^\top$$
- $$\frac{\partial (\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{A} \mathbf{x} + \mathbf{A}^\top \mathbf{x}$$

벡터미분을 통하여 식(2.31)을 구해보자 식(2.29)을 전개하면

$$\begin{aligned} \mathbf{e}^\top \mathbf{e} &= (\mathbf{Y} - \mathbf{z} \cdot \mathbf{A})^\top (\mathbf{Y} - \mathbf{z} \cdot \mathbf{A}) \\ &= (\mathbf{Y}^\top - \mathbf{A}^\top \mathbf{z}^\top) (\mathbf{Y} - \mathbf{z} \cdot \mathbf{A}) \\ &= \mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{z} \mathbf{A} - \mathbf{A}^\top \mathbf{z}^\top \mathbf{Y} + \mathbf{A}^\top \mathbf{z}^\top \mathbf{z} \mathbf{A} \end{aligned}$$

즉, 식(2.31)은 다음 식의 전개를 통해 구할 수 있다.

$$\frac{\partial \{\mathbf{e}^\top \mathbf{e}\}}{\partial \mathbf{A}} = \frac{\{\mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{z} \mathbf{A} - \mathbf{A}^\top \mathbf{z}^\top \mathbf{Y} + \mathbf{A}^\top \mathbf{z}^\top \mathbf{z} \mathbf{A}\}}{\partial \mathbf{A}} \quad (2.32)$$

$$= -\mathbf{z}^\top \mathbf{Y} - \mathbf{z}^\top \mathbf{Y} + \mathbf{z}^\top \mathbf{z} \mathbf{A} + \mathbf{z}^\top \mathbf{z} \mathbf{A} \quad (2.33)$$

$$= -2\mathbf{z}^\top \mathbf{Y} + 2\mathbf{z}^\top \mathbf{z} \mathbf{A} \quad (2.34)$$

$$= 0 \quad (2.35)$$

$$\therefore \mathbf{A} = [\mathbf{z}^\top \mathbf{z}]^{-1} \mathbf{z}^\top \mathbf{Y} \quad (2.36)$$

Part II

부록

Appendices