

A Space Model to Annual Rainfall in South Korea

Eui-Kyoo Lee¹⁾

Abstract

Spatial data are usually obtained at selected locations even though they are potentially available at all locations in a continuous region. Moreover the monitoring locations are clustered in some regions, sparse in other regions. One important goal of spatial data analysis is to predict unknown response values at any location throughout a region of interest. Thus, an appropriate space model should be set up and their estimates and predictions must be accompanied by measures of uncertainty. In this study we see that a space model proposed allows a best interpolation to annual rainfall data in South Korea.

Keywords : Kriging, Semivariogram, Nonlinear Least Squares, BLUP

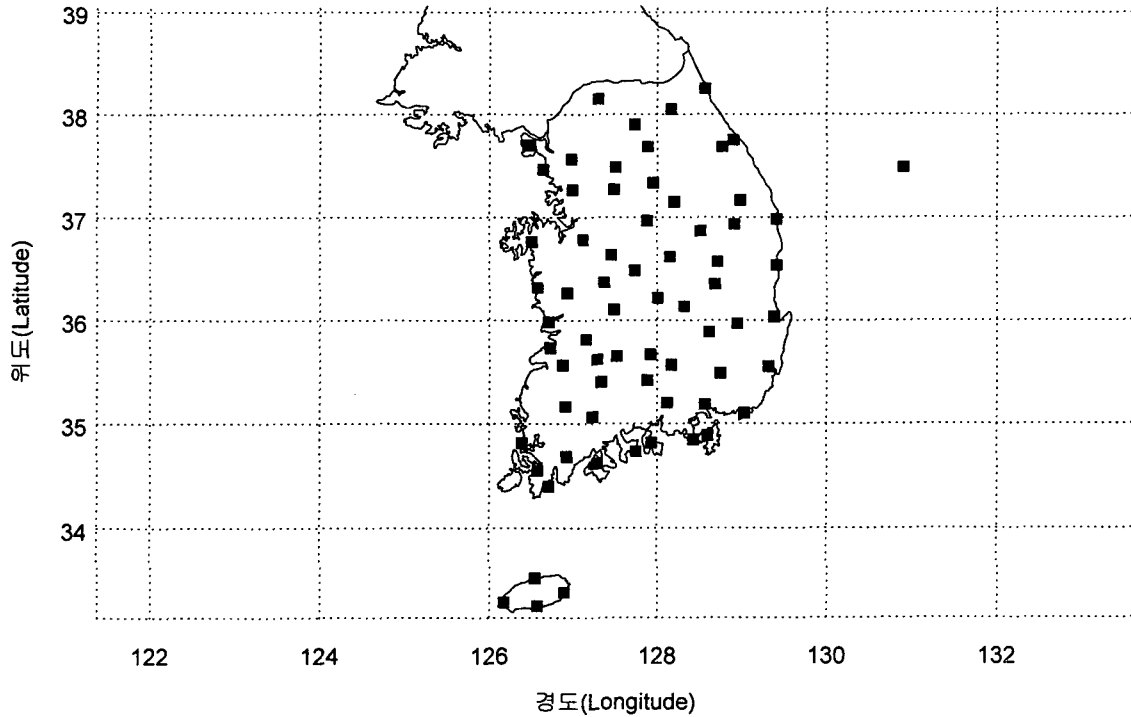
1. 서론

공간자료는 관측지점에서 얻어지는 자료를 말한다. 일반적으로 공간자료는 모든 지점에서 관측될 수 없을 뿐 아니라 불규칙적인 지점에서 자료가 수집된다. 이러한 자료는 각 지점이 서로 가까울수록 유사한 값을 가지리라 예상할 수 있다. 이는 시계열 자료가 종종 시간적으로 상관성을 갖는 것 처럼 공간상에서 얻어진 자료는 공간적인 상관성을 갖는다. 그러므로 공간자료분석은 이러한 공간적 상관성을 바탕으로 이루어져야 할 것이다. 즉 공간자료분석은 공간 변이도함수(semivariogram function) 또는 공분산함수(covariance function)를 근거로 하여 관측되지 않은 지점에서의 예측을 주목적으로 한다. 이미 공간자료에 대한 모형분석은 환경이나 기상과 같은 많은 분야에서 적용되고 있다. 본 논문에서는 남한지역의 평년 강우량 자료에 대한 공간모형을 설정하고 자료의 상관성을 효과적으로 이용하여 분석하고자 한다.

국내 전 지역에 흩어져 있는 68개 강우량 관측소의 위치가 [그림 1.1]에 표시되어 있고 이 공간상 위치들에서 얻어진 관측값은 강우량 평년값이다. 여기서 평년값이란 각 관측지점에서 1971년부터 2000년까지의 연도별 강우량을 평균한 값을 말한다. 일반적으로 우리는 보간법에 의해 관측되지 않은 지점의 강우량을 예측한다. 예를 들면 위도와 경도를 격자로 하여 이변량 보간법(bivariate interpolation)에 근거하여 추정을 하는데 이는 자료의 상관성을 충분히 이용하지 못하기 때문에 때때로 만족스럽지 못한 결과를 얻게 된다. 따라서 본 논문에서는 공간적 상관성을 최대한 이용하는 최적의 보간법(optimal interpolation)인 크리깅(Kriging) 추정법을 통하여 강우량 평년값에 대한 공간모형분석을 하고자 한다.

1) Lecturing Professor, Department of Applied Statistics, Konkuk University, Seoul, 143-701, Korea
E-mail: ekyoolee@konkuk.ac.kr

다음 2장에서는 공간모형분석의 이론적 배경을 소개하고 3장에서 강우량 자료에 대한 모형의 추정과 분석을 다룬다. 또한 공간변이도를 달리한 모형들을 비교분석하며 4장에서 결론을 갖는다.



[그림 1.1] 국내 68개의 강우량 관측소

2. 공간모형의 배경

s 를 k 차원의 유클리드(Euclid) 공간상의 위치벡터라 하자. 그리고 $Z(s)$ 는 공간상 위치 s 에서 관측되는 확률변수이고 $z(s)$ 는 $Z(s)$ 의 실현값이라 하자. Cressie(1993)는 공간 확률과정 $Z(s)$ 를 대규모변동(large-scale variation), 소규모변동(small-scale variation), 미세규모변동(micro-scale variation)과 측정오차(measurement error) 네가지 구성요소로 분해한다. 즉,

$$Z(s) = \mu(s) + e_{ss}(s) + e_{ms}(s) + e_{me}(s). \quad (2.1)$$

대규모변동은 전지역적 추세를 나타내는 결정적 요소를 의미하며 소규모변동은 그 결정적 추세 주변의 확률적 변동을 말한다. 미세규모변동은 표집된 거리보다 더 작은 거리에서의 변동을 나타내며 주로 측정오차와 결합되어 표현하기도 한다.

시계열분석과 유사하게, 공간적 확률과정에 대한 모수추정과 추론을 위하여 2차 공간적 정상성과 같은 동질성 가정을 필요로 한다. 공간상 2차 정상성은 평균이 일정하며 공분산은 위치에 의존하지 않고 단지 거리 (그리고 경우에 따라서는 방향)에 의존함을 의미한다(Christensen 1991). 즉,

$$\begin{aligned} E[Z(\mathbf{s})] &= \mu \\ \text{Cov}[Z(\mathbf{s}), Z(\mathbf{s} + \mathbf{d})] &= \sigma(\mathbf{d}). \end{aligned} \quad (2.2)$$

이차 정상성보다 약한 가정인 내재적 정상성(intrinsic stationarity)은 평균이 일정하며 증가분의 분산이 위치에 의존하지 않음을 요구한다. 즉,

$$\begin{aligned} E[Z(\mathbf{s})] &= \mu \\ \text{Var}[Z(\mathbf{s}) - Z(\mathbf{s} + \mathbf{d})] &= 2\gamma(\mathbf{d}). \end{aligned} \quad (2.3)$$

일반적으로 공간모형분석은 이러한 정상성 가정하에서 분석한다.

2.1 변이도함수(Semivariogram function)와 크리깅(Kriging)

공간자료분석의 주요 목적은 관심지역에 걸쳐 관측되지 않은 지점에서의 관측값을 예측하는 것이다. 크리깅(Kriging)은 지리통계학(Geostatistics)에서의 최적의 예측방법(Matheron 1962)으로서 그 예측량은 관측된 값들의 가중선형결합이다. 여기서 가중치는 예측량이 비편향성과 최소예측분산을 갖도록 결정되어진다. 실제로 이 예측량은 일반선형모형에서의 최량선형불편예측량(Best Linear Unbiased Predictor; BLUP)과 일치한다(Goldberger 1962). 한편 가중치벡터는 변이도 또는 공분산함수로 표현될 수 있다(Cressie 1993). (이에 대해서는 본 논문의 부록에서 다루기로 한다.) 그러므로 변이도 또는 공분산함수는 알려져 있거나 추정되어야 한다.

변이도함수(semivariogram function)는 공간 상관성의 척도이며 거리(그리고 경우에 따라서는 방향)의 함수로서 다음과 같이 정의된다.

$$\gamma(\mathbf{d}) = \frac{1}{2} \text{Var}[Z(\mathbf{s}) - Z(\mathbf{s} + \mathbf{d})] \quad (2.4)$$

따라서

$$\begin{aligned} \gamma(\mathbf{d}) &= \sigma(0) - \sigma(\mathbf{d}) \\ &= [1 - \rho(\mathbf{d})]\sigma(0) \end{aligned} \quad (2.5)$$

가 성립한다. 여기서 $\rho(\mathbf{d})$ 는 상관함수이다. 즉, 변이도는 공분산함수나 상관함수의 관계식이 됨을 알 수 있다.

전통적인 변이도 추정량은 Matheron(1962)에 의해 다음과 같이 추정된다.

$$\hat{\gamma}(\mathbf{d}) = \frac{1}{2|N(\mathbf{d})|} \sum_{M(\mathbf{d})} [z(\mathbf{s}_i) - z(\mathbf{s}_j)]^2 \quad (2.6)$$

여기서 $|N(\mathbf{d})|$ 는 거리 d 의 차이를 갖는 가능한 지점들의 중복되지 않는 쌍의 수이다. 이는 평균이 일정하다는 가정 ($E[Z(\mathbf{s}_i) - Z(\mathbf{s}_j)] = 0$)하에서 $\gamma(\mathbf{d})$ 의 불편추정량이다.

앞에서 언급하였듯이 추정된 변이도함수(또는 공분산함수)는 크리깅 절차에서 최적의 가중치를 계산하기 위하여 사용되어 진다. 그런데 구하여진 공분산행렬이 양정치가 될 것을 보증하도록 이론적인 변이도함수를 추정된 변이도값들에 적합시킨다. 이론적인 변이도함수 중 스피어리컬(spherical) 또는 가우시안(Gaussian) 변이도함수가 대표적인데 스피어리컬 변이도함수는 다음과 같이 정의된다.

$$\gamma(d) = \begin{cases} 0 & , d=0 \\ \theta_1 + \theta_2 \left[1.5 \left(\frac{d}{\theta_3} \right) - 0.5 \left(\frac{d}{\theta_3} \right)^3 \right] & , 0 < d \leq \theta_3 \\ \theta_1 + \theta_2 & , d > \theta_3. \end{cases} \quad (2.7)$$

식(2.7)에서의 모수들은 각각 중요한 의미를 갖는다. 너겟(nugget) θ_1 은 관측지점간의 거리가 0으로 접근할 때의 변이도값의 극한값이다. 이것은 미세규모변동이나 측정오차를 의미한다. 실(sill)은 관측지점간의 거리가 증가할 때의 변이도값의 수렴값이다. θ_2 는 부분 실(partial sill)이라 하고 실(sill)에서 너겟값을 뺀 값이다. 범위(range) θ_3 는 관측값이 더 이상 상관되지 않는 관측지점들간의 거리이다. 대표적 변이도 중 또 하나인 가우시안 변이도함수는 다음과 같이 정의된다.

$$\gamma(d) = \begin{cases} 0 & , d=0 \\ \theta_1 + \theta_2 \left[1 - \exp \left\{ - \left(\frac{d}{\theta_3} \right)^2 \right\} \right] & , d > 0. \end{cases} \quad (2.8)$$

식(2.8)에서의 모수들은 범위 θ_3 를 제외하고 식(2.7)의 경우와 같이 해석된다. 가우시안 모형의 범위 θ_3 는 점근적 실을 가지므로 실질적 범위는 아니다. Journel 과 Huijbregts(1978)은 $\sqrt{3}\theta_3$ 를 실질적 범위로 간주한다. 이와같이 추정된 변이도함수 또는 공분산함수를 이용하여 각 관찰치의 가중치를 결정하고 이들의 선형결합으로 관측되지 않은 지점의 값을 추정하게 된다.

2.2 공간적 추세의 추정과 유니버설 크리깅(Universal Kriging)

시계열자료가 시계열추세와 관찰값간의 자기상관관계를 갖는 반면 공간자료는 공간적 추세와 공간적 상관관계가 존재할 수 있다. 공간상관성은 앞서 소개된 변이도함수를 적합시킴으로써 추정할 수 있다. 그런데 고전적인 변이도 추정량 $\hat{\gamma}(\mathbf{d})$ 는 평균이 위치에 상관없이 불변하다는 정상성 가정을 만족한다면 불편추정량이다. 즉 위치에 상관없이 평균이 일정하다면 변이도함수는 다음과 같이 정의되기 때문이다.

$$\begin{aligned} \gamma(\mathbf{d}) &= \frac{1}{2} \text{Var}[Z(\mathbf{s}) - Z(\mathbf{s} + \mathbf{d})] \\ &= \frac{1}{2} E[Z(\mathbf{s}) - Z(\mathbf{s} + \mathbf{d})]^2 \end{aligned} \quad (2.9)$$

그러나 평균이 일정하지 않다면,

$$\text{Var}[Z(s) - Z(s+d)] = E[Z(s) - Z(s+d)]^2 - \{E[Z(s)] - E[Z(s+d)]\}^2 \quad (2.10)$$

이므로 만약 공간적 추세가 존재한다면 표본변이도는 더이상 불편성을 만족하지 않는다. 따라서 (2.10)에서 볼수 있듯이 추세가 존재하는 경우에 이를 무시하고 추정된 변이도는 종종 거리가 멀어짐에 따라 증가하는 행태를 보이게 된다.

Cressie(1993)는 공간적 추세를 제거한 잔차를 이용하여 각 변이도(semivariogram)를 추정한 후 이를 이용하여 관측되지 않은 지점에서의 값을 예측하는데, 이는 시계열에서 추세를 제거한 후 상관도(correlogram)를 추정하는 것과 같은 맥락으로 볼 수 있다. 비정상 시계열자료에 대해 차분을 통하여 정상성을 만족시키는 ARIMA모형이 유용하듯이 추세가 제거된 값은 종종 이차 정상성을 갖는다(Matheron 1973). 그러므로 변이도함수를 추정하기전에 공간추세는 추정된 추세함수로 제거되어야 한다. 잔차의 공분산함수와 본래 자료의 공분산함수가 다름에도 불구하고 Kintanidis (1993)는 일반화 공분산함수 또는 변이도함수가 공간적 변동의 핵심부분을 포함함을 보였다. 게다가 Christensen(1990)은 크리깅 추정과 추정량의 표준오차는 원자료의 변이도함수 대신 일반화된 변이도함수를 사용할 때 불변함을 보였다. 그러므로 잔차의 변이도함수는 공간자료의 추정을 위해 사용될 수 있다.

그런데 어떤 추세가 존재한다면 변이도는 정확히 추정할 수 없듯이 어떤 공간적 상관성이 존재한다면 추세 또한 정확하게 추정할 수 없게 된다. 그러나 위와 같은 이유에서 우리는 추세를 추정하고 추세가 제거된 잔차의 변이도를 추정하게 된다. 이처럼 추세와 변이도를 동시에 고려하는 크리깅을 유니버설 크리깅(universal kriging)이라 한다.

3. 공간모형분석

이 장에서는 남한 지역내의 평년강우량에 대한 공간모형을 분석한다. 연구의 간결성을 위하여 68개의 관측지점에서 울릉도 1개지점과 제주도 4개 지점은 제외하였다. 강우량자료에 공간모형을 적합시키기 전에 공간상의 추세와 변동을 대략적으로 파악하기 위해 여러 그래픽을 이용한다. 남한지역에서의 평년 강우량을 이변량 보간법으로 추정한 값이 [그림 3.1]에 나타나 있다. 3차원 그림에서 볼 수 있듯이 강우량은 남해안과 강원도에서 크며 대구지역과 내륙지역에서 적은 값이 관측되어짐을 알 수 있다. 특히 대관령지점은 지역적 특성인 적설량의 영향으로 연평균값 1717mm로 주변지역에 비해 많은 강우량을 보이며 거제와 남해는 각각 1797mm와 1790mm으로 최다 관측되었다.

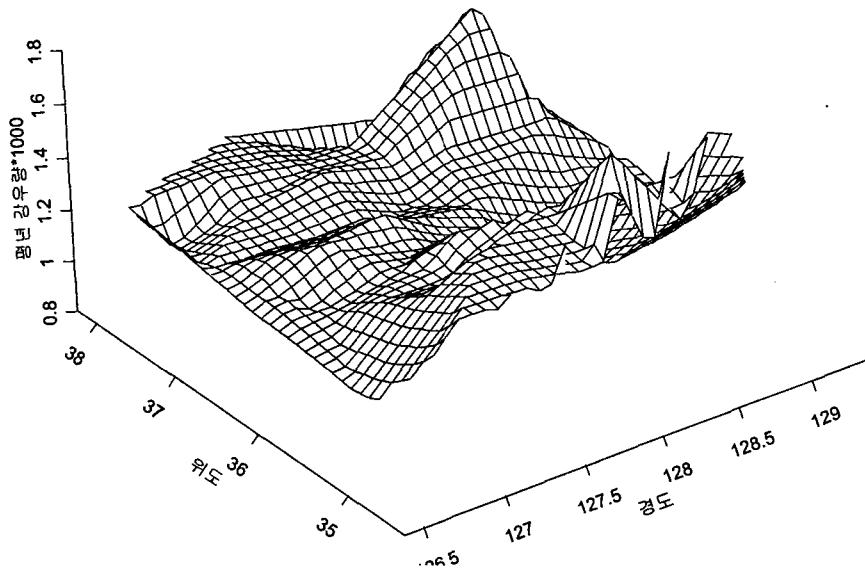
또한 강우량을 경도, 위도 각각에 대하여 산점도를 그려 추세를 보았는데 이차추세함수로 나타냈다. 때때로 공간추세를 비모수적 방법인 경도와 위도에 대한 국소회귀(local regression)로 그 공간추세를 평활하기도 하는데 여기서는 전자의 모수적 방법을 사용한다. 이제 모형을 설정하고 제안된 모형의 모수를 추정하며 관측되지 않은 지점에서의 평년강우량을 예측하고 비교한다.

3.1 공간 모형의 추정과 예측

강우량자료의 변동의 근거는 무엇인가? 이 경우 대규모변동은 전지역에 걸친 지리적 차이가 될 것이며 소규모변동은 지역적 환경차이가 변동의 원인이 될 수 있을 것이다. 따라서 $\mu(s)$ 는 공간 추세로서 경도와 위도의 이차함수이며 $e(s)$ 는 추세함수주변의 변동을 설명하는 확률적 요소로 미세규모변동과 측정오차를 포함한다. 이를 벡터로 나타내면

$$\begin{aligned} Z(s) &= \mu(s) + e(s) \\ &= X\beta + e(s) \end{aligned} \tag{3.1}$$

와 같이 표현할 수 있다. 여기서 $x' = (1, \text{경도}, \text{위도}, \text{경도}^2, \text{위도}^2, \text{경도} * \text{위도})$ 이며 각 항은 64×1 열벡터이다.



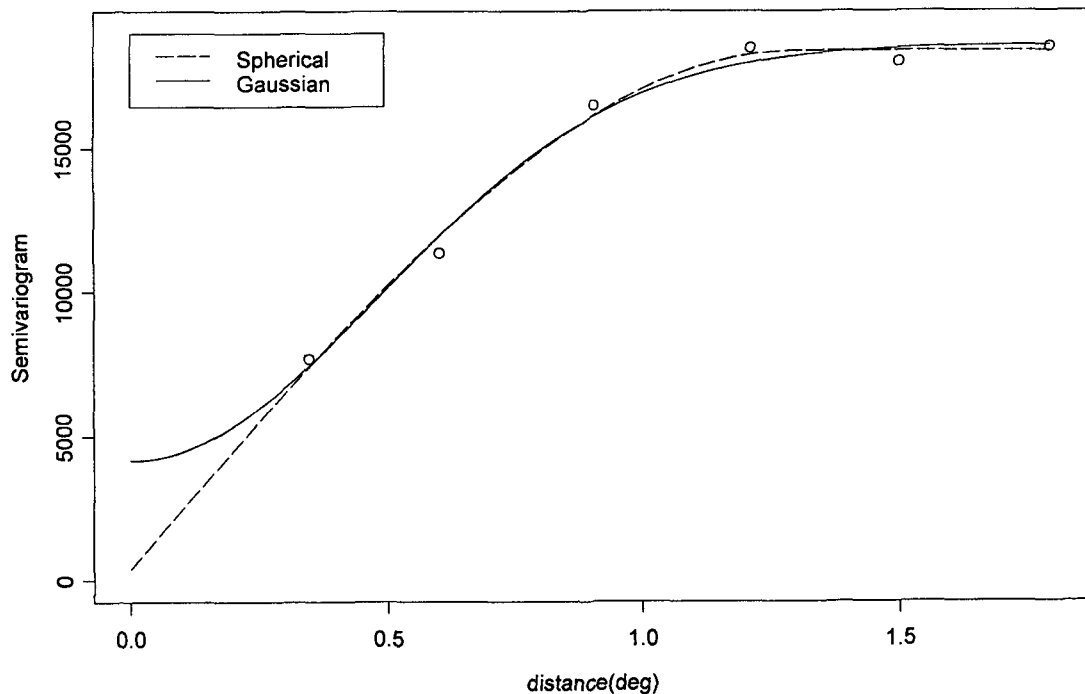
[그림 3.1] 이변량 보간법에 의한 평년강우량의 3차원 그림

이제 공간적 확률요소를 추정하기 위하여 이차추세함수를 적합후 관측값에서 이를 제거한 잔차들을 이용하여 구한 추정변이도값들을 [그림 3.2]의 각 빈점들로 나타냈다. 지역에서 최대거리의 반인 거리범위는 1.8도(약 200km)를 사용한다. 변이도값은 작은 거리에서 작은 값을 갖다가 거리가 점차 멀어질수록 그 값이 커지고 있다. 또한 방향성을 고려할 때 방향에 크게 의존하지 않았으며 구간수를 조정함에 따라 다소 값이 변하기는 하지만 전체적인 행태는 유사함을 확인하였다. 특히 각 구간에서 가능한 쌍들이 적어도 30개 이상 존재할 때 우리는 그 변이도값을 신뢰할 수 있는데 계산된 값들은 모두 이를 만족한다.

이제 이 추정된 변이도값을 이론적인 변이도함수로 적합시킨다. Journel과 Huijbregts(1978)은 이론적인 변이도함수로 선형(linear), 구형(spherical), 정규분포형(Gaussian), 지수형(exponential), 파동형(wave)등의 여러가지 변이도모형들을 제시하고 있다. 그런데 [그림 3.2]에 나타난 추정된 변이도값들은 초기구간에서는 빠르게 증가하다 일정 거리이후에는 천천히 수렴해가는 행태를 띤다. 따라서 스피어리컬(spherical) 또는 가우시안(Gaussian) 변이도모형이 변이도값의 변화행태를 적절하게 설명할 수 있을 것이다.

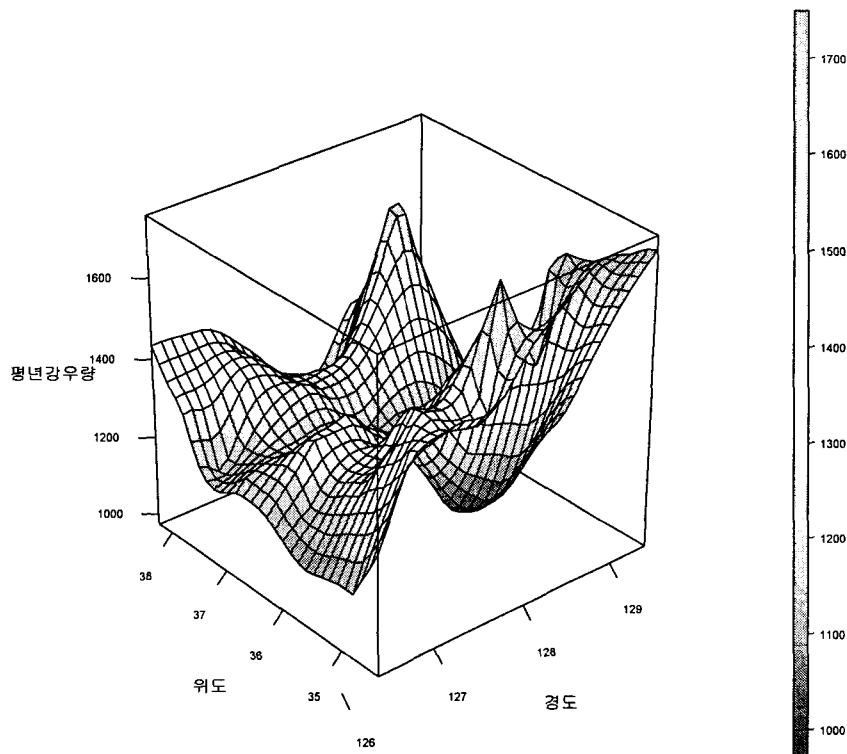
이들 비선형모형의 모수추정은 다양한 방법으로 수행되어질 수 있는데 여기서는 추정된 변이도값과 이론적인 변이도함수값과의 잔차제곱합을 최소화하는 방법에 의해 S+를 이용하여 모수를 추정하였다. 스피어리컬함수의 경우 $\hat{\theta}_1$ 은 374, $\hat{\theta}_2$ 은 18010, $\hat{\theta}_3$ 는 1.302도로 추정되었으며 가우시안의 경우는 각각 $\hat{\theta}_1=4155$, $\hat{\theta}_2=14439$, $\hat{\theta}_3=0.683$ 이다. (실제적 범위는 $\sqrt{3}(0.683)=1.18$).

적합된 스피어리컬과 가우시안 변이도함수들이 [그림 3.2]에서 추정된 변이도값들 위에 적합되었다. 적합된 두함수는 초기구간을 제외하고 그리 다르지 않지만 종종 스피어리컬이 가우시안모형보다 선호된다. 그것은 초기 구간에서 직선형태로 나타나 너겟이 상대적으로 작게 추정되기 때문이다(Gunst 1995). 요약하면 거리에 따라서 상관성은 약해지다 약 150km정도의 거리이상에서는 상관성은 없는 것으로 간주한다. 즉 그 이상 떨어진 자료는 예측시 영향을 주지 않는다는 것이다.



[그림 3.2] 추정 변이도값과 적합된 변이도함수

유니버설 크리깅은 추세가 존재하는 공간모형에서의 최적의 보간법으로서 위에서 적합된 변이도함수를 이용하여 관측되지 않은 지점에서의 강우량을 예측한다. [그림 3.3]는 S+를 이용하여 (Kaluzny 1998) 예측하였는데 이것은 이차함수로 추정된 추세에다 추정된 스피어리컬 변이도함수에 의한 확률적 변동부분으로 조정한 결과를 보여준다. [그림 3.1]의 이변량 보간법에 의한 3차원 그림과의 차이는 공간모형이 모든 공간자료를 이용하여 공간적 상관성을 추정하고 이를 바탕으로 가중치를 달리하는 반면 이변량 보간법에 의한 모형은 공간 상관성을 충분히 이용하지 못하기 때문이다.

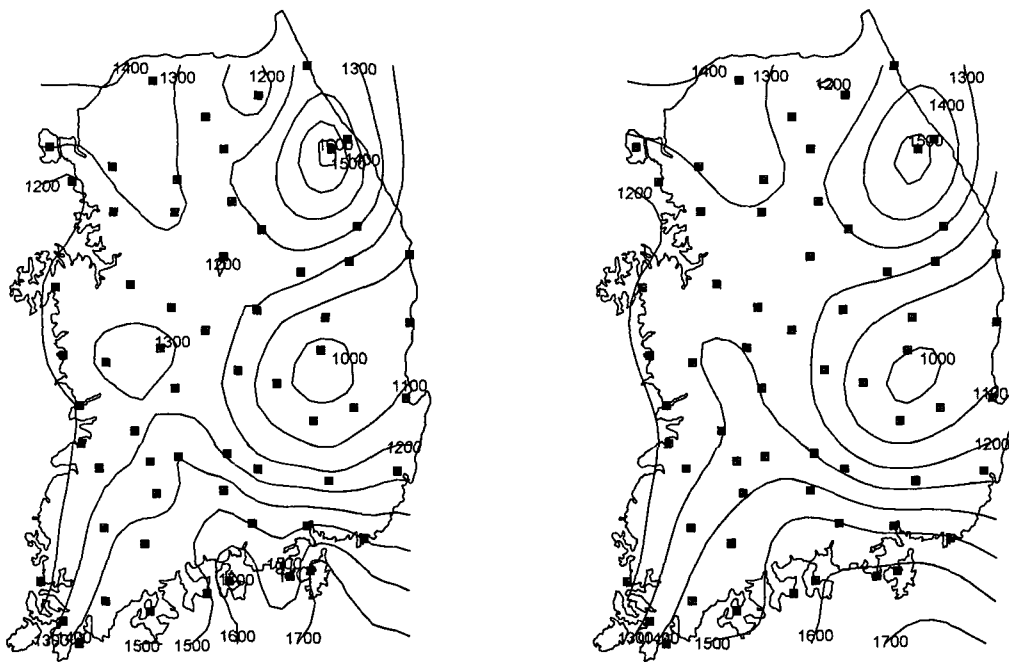


[그림 3.3] 스피어리컬 변이도함수에 근거한 예측값의 3차원 그림

3.2 스피어리컬과 가우시안 변이도함수 적합시의 예측값 비교

[그림 3.4]는 스피어리컬 변이도에 의한 모형의 예측값과 가우시안 변이도에 의한 모형의 예측값을 비교하기 위하여 등고선도로 표현한 것이다. 등고선도는 2차원으로 표현한 그림으로서 이어진 선은 같은 수준의 강우량을 나타내며 선의 간격이 좁아질수록 수준값은 가파르게 변화함을 의미한다. 짧은 거리에서는 가능한 쌍들이 존재하지 않아 변이도의 초기구간에서는 그 상관성의 정도를 측정하기 어렵다. [그림 3.2]에서 볼 수 있듯이 스피어리컬 변이도함수는 초기구간에서 급격히 상관성이 떨어지는 반면 가우시안 변이도함수의 경우는 초기구간에서 상관성이 완만히 떨어져 너겟이 스피어리컬 변이도함수에 비해 크다. 따라서 스피어리컬 변이도함수에 의한 모형에서는 예

측시 짧은 거리의 관측값이 구체적으로 반영되어 나타난다. 그러나 가우시안 함수의 경우 짧은 거리에서는 정보가 불확실하여 덜 구체적으로 반영되어 예측됨을 알 수 있다. 그런데 강수량 평년값은 30년간에 측정된 강수량의 평균값이므로 측정오차는 무시될 수 있으며 따라서 짧은 거리의 관측값이 보다 구체적으로 반영되도록 하는 것이 적절할 것이다. 그러므로 가우시안 변이도함수의 적합보다는 스피어리컬 변이도함수를 선택하여 분석하는 것이 타당할 것이다.



[그림 3.4] 스피어리컬(왼쪽)과 가우시안(오른쪽) 변이도함수 적합시 예측값의 등고선도

4. 결론

자료가 공간상에서 관측되어졌다면 이러한 공간자료는 종종 공간적으로 상관되어진다. 그러므로 관측되어지지 않은 지점에서의 예측은 공간적 상관성을 바탕으로 이루어져야 할 것이다. 따라서 본 논문에서 남한지역내의 강수량 평년값에 대한 예측모형으로서 공간상 이차추세함수와 그 추세 주변의 확률적 변동구조를 갖는 모형으로 설정하였다. 그리고 강수량 평년값 자료를 이용한 공간적 상관성을 스피어리컬 변이도모형으로 식별하고 추정하였다. 예측값들은 최소예측오차를 갖도록 계산된 가중치에 의한 관측값들의 가중선형결합으로 최적의 보간법이 됨을 확인하였다. 끝으로 강수량자료나 기온자료는 자료수집시 각 관측지점에서 뿐만아니라 각 시점에 따른 시계열 자료를 동시에 갖는다. 즉 월별 평년 강수량이나 연도별 평년 강수량 자료에 대한 모형을 시공간적으로 동시에 분석한다면 주어진 정보를 보다 더 효율적으로 이용하여 분석할 수 있을 것이다.

References

- [1] Christensen, R. C. (1991), *Linear Models for Multivariate, Time Series and Spatial Data*, New York: Springer-Verlag.
- [2] Christensen, R. C. (1990), The Equivalence of Predictions from Universal Kriging and Intrinsic Random-Function Kriging, *Mathematical Geology*, 22, 655-664.
- [3] Cressie, N. (1993). *Statistics for Spatial Data*, revised ed. New York: John Wiley & Sons.
- [4] Goldberger (1962), Best Linear Unbiased Prediction in the Generalized Linear Regression Model. *Journal of the American Statistical Association*, 57, 369-375.
- [5] Gunst, R. F. (1995), Estimating Spatial Correlations from Spatial-Temporal Meteorological Data, *Journal of Climate*, 8, 2454-2470.
- [6] Journel, A.G. and Huijbregts, C.J. (1978), *Mining Geostatistics*, New York: Academic Press.
- [7] Kaluzny, S.P., Vega, S.C., Cardoso, T.P. and Shelly, A.A. (1998), *S+ Spatial Stats. User's Manual for Windows and UNIX*, York: MathSoft.
- [8] Kitanidis, P. K. (1993), Generalized Covariance Functions in Estimation, *Mathematical Geology*, 25, 525-540.
- [9] Matheron, G. (1962), *Traite de Geostatistique Appliquee, Tome I. Memories du Bureau de Recherches Geologiques et Minières*, No. 14, Paris: Editions Technip.

[2003년 5월 접수, 2003년 8월 채택]

<부 록>

최량선형불편예측량(BLUP)과 크리깅(Kriging)

다음과 같은 공간모형을 고려하자.

$$\begin{aligned} Z(\mathbf{s}) &= \mu(\mathbf{s}) + e(\mathbf{s}) \\ &= \mathbf{x}'(\mathbf{s})\boldsymbol{\beta} + e(\mathbf{s}), \quad \mathbf{s} = \mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_n. \end{aligned}$$

여기서

$$\boldsymbol{\Sigma} = [\text{Cov} \{ Z(\mathbf{s}_i), Z(\mathbf{s}_j) \}]_{n \times n}.$$

그리고 관측되지 않는 지점에서의 반응변수 $Z(\mathbf{s}_0)$ 는 같은 모형을 따른다. 즉,

$$Z(\mathbf{s}_0) = \mathbf{x}'(\mathbf{s}_0)\boldsymbol{\beta} + e(\mathbf{s}_0).$$

그리고

$$\boldsymbol{\Sigma}_{z_0} = [\text{Cov} \{ Z(\mathbf{s}_i), Z(\mathbf{s}_0) \}]_{n \times 1}.$$

위와 같은 모형에서 최량선형불편예측량(best linear unbiased predictor; BLUP) (Goldberger 1962)은 선형성 ($\widehat{Z}(s_0) = w'Z$)과 불편성($E[\widehat{Z}(s_0)] = E[Z(s_0)] \Rightarrow X'w = x(s_0)$)을 만족하고 예측분산($\text{Var}[\widehat{Z}(s_0) - Z(s_0)]$)을 최소로 하는 예측량인데 이것은 라그랑즈 승수법(Lagrangian multiplier method)에 의하여 아래와 같이 구하여진다.

먼저 예측량은 관찰값의 선형결합이므로 다음과 같다.

$$\widehat{Z}(s_0) = w'Z$$

따라서 예측오차의 분산

$$\begin{aligned} \text{Var}\{Z(s_0) - \widehat{Z}(s_0)\} &= \text{Var}\{Z(s_0) - w'Z\} \\ &= \sigma_{z(s_0)}^2 - 2w'\Sigma_{z_0} + w'\Sigma w. \end{aligned}$$

이다. 그런데 불편성을 만족하는 조건하에서 예측오차의 분산을 최소로 하는 w 를 찾기 위해 다음의 라그랑즈함수를 고려하자.

$$L(w, \lambda) = \sigma_{z(s_0)}^2 - 2w'\Sigma_{z_0} + w'\Sigma w - 2\lambda'(X'w - x(s_0))$$

여기서 $L(w, \lambda)$ 를 최소로 하는 w, λ 는 다음과 같은 선형방정식을 푸는 것과 같다.

$$\begin{bmatrix} \Sigma & X \\ X' & \Phi \end{bmatrix} \begin{bmatrix} w \\ -\lambda \end{bmatrix} = \begin{bmatrix} \Sigma_{z_0} \\ x(s_0) \end{bmatrix}$$

즉,

$$\begin{aligned} \begin{bmatrix} w \\ -\lambda \end{bmatrix} &= \begin{bmatrix} \Sigma & X \\ X' & \Phi \end{bmatrix}^{-1} \begin{bmatrix} \Sigma_{z_0} \\ x(s_0) \end{bmatrix} \\ &= \begin{bmatrix} \Sigma^{-1}\{I - X(X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}\} & \Sigma^{-1}X(X'\Sigma^{-1}X)^{-1} \\ (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1} & -(X'\Sigma^{-1}X)^{-1} \end{bmatrix} \begin{bmatrix} \Sigma_{z_0} \\ x(s_0) \end{bmatrix} \end{aligned}$$

가 된다. 따라서 최량선형불편예측량(BLUE)은 다음과 같이 표현할 수 있다.

$$\begin{aligned} \widehat{Z}(s_0) &= x'(s_0)\widehat{\beta} + \Sigma'_{z_0}\Sigma^{-1}(Z - X\widehat{\beta}) \\ \widehat{\beta} &= (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}Z. \end{aligned}$$

그러므로 BLUP는 잔차의 선형결합으로 조정된 일반화 선형예측량이다. 이러한 조정은 관찰값 내의 오차와 예측하고자하는 지점에서의 오차로 부터의 정보를 사용한다. 또한 예측오차의 분산은 변이도함수로 다음과 같이 표현될 수 있다.

$$\begin{aligned}\text{Var}\{Z(s_0) - \widehat{Z}(s_0)\} &= \text{Var}\{Z(s_0) - \mathbf{w}'\mathbf{Z}\} \\ &= 2\mathbf{w}'\mathbf{\Gamma}_{z_0} - \mathbf{w}'\mathbf{\Gamma}\mathbf{w}.\end{aligned}$$

따라서 선형방정식은 다음과 같다.

$$\begin{bmatrix} \mathbf{\Gamma} & \mathbf{X} \\ \mathbf{X}' & \mathbf{\Phi} \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ \boldsymbol{\xi} \end{bmatrix} = \begin{bmatrix} \mathbf{\Gamma}_{z_0} \\ \mathbf{x}(s_0) \end{bmatrix}.$$

여기서 $\mathbf{\Gamma}$ 와 $\mathbf{\Gamma}_{z_0}$ 는 아래와 같이 정의된다.

$$\begin{aligned}\mathbf{\Gamma} &= [\gamma(s_i - s_j)]_{n \times n}, \\ \mathbf{\Gamma}_{z_0} &= [\gamma(s_i - s_0)]_{n \times 1}.\end{aligned}$$

그러므로 가중치 \mathbf{w} 는 공분산함수 $\boldsymbol{\Sigma}$ 와 $\boldsymbol{\Sigma}_{z_0}$ 대신 변이도함수 $\mathbf{\Gamma}$ 와 $\mathbf{\Gamma}_{z_0}$ 으로 대체하여 표현될 수 있다.