

격자기반 분석을 통한 위치기반 소셜 미디어 데이터와 부동산 가격지수 간의 공간적 상관성 분석 연구

Analyzing Spatial Correlation between Location-Based Social Media Data and Real Estates Price Index through Rasterization

박우진* · 어승원** · 유기윤***
Park, Woo Jin · Eo, Seung Won · Yu, Ki Yun

요 旨

본 연구에서는 위치기반 소셜 미디어 데이터의 공간적 분포가 지역별 부동산 지수와 어떠한 공간적 관련성을 가지는지에 대해 알아보려고 한다. 두 데이터는 상이한 자료 형식을 가지고 있어, 이를 보완할 수 있는 방법론으로 본 연구에서는 격자 기반의 공간분석 방법을 적용하였다. 대상 데이터로는 2013년 8월 한 달간의 지오태그된 트윗 데이터와 행정구역별 주택가격지수(매매, 전세)를 이용하였으며, 공간적 범위는 서울과 수도권 일부를 포함하도록 하였다. 두 데이터 간의 상이한 공간적 단위를 고려하여 2,000m 단위의 격자망을 구성하고 이에 맞게 두 데이터를 격자 데이터 형태로 변환하였다. 변환된 두 데이터에 대하여 Hot spot 분석을 실시하여 공간적 분포를 시각적으로 비교하였으며, 공간시차를 고려한 이변량 공간적 상관계수를 측정함으로써 정량적 분석을 실시하였다. 시각적, 정량적 분석 결과, 서초구 지역이 트윗 데이터와 주택매매가격지수 데이터에서 공통적인 Hotspot 지역으로 탐색되었으나 주택전세가격지수 데이터와는 뚜렷한 공간적 상관성이 탐색되지 않았다.

핵심용어 : 소셜 미디어, 지오태그, 부동산 가격지수, 격자기반 분석, 핫스팟 분석, 공간적 상관성

Abstract

In this study, the spatial relevance between the regional housing price data and the spatial distribution of the location-based social media data is explored. The spatial analysis with rasterization was applied to this study, because the both data have a different form to analyze. The geo-tagged Twitter data had been collected for a month and the regional housing price index about sales and lease were used. The spatial range of both data includes Seoul and the some parts of the metropolitan area. 2,000m grid was constructed to consider the different spatial measure between two data, and they were combined into the constructed grids. The Hotspot Analysis was operated using the combined dataset to see the comparison of spatial distribution, and the bivariate spatial correlation coefficients between two data were measured for the quantitative analysis. The result of this study shows that Seocho-gu area is detected as a common hotspot of tweet and housing sales price index data. though the spatial relevance is not detected between tweet and housing lease price index data.

Keywords : Social Media, Geo-tag, Real Estate Price Index, Grid based Analysis, Hot Spot Analysis, Spatial Correlation

1. 서 론

1.1 연구배경 및 목적

최근 스마트폰의 사용량이 증가하면서 웹 플랫폼 및

모바일 플랫폼을 기반으로 한 소셜 미디어 산업이 급성장 하고 있으며 이에 따라 다양한 형태의 소셜 네트워크 서비스가 개발되고 있다. 특히 GPS 등의 위치정보 측정기술이 발달하면서 위치정보가 태그된 위치기반

Received: 2015.01.12, revised: 2015.02.23, accepted: 2015.03.09

* 서울대학교 환경정화기술 및 위해성평가 연구센터 연수연구원(Center of Environmental Remediation and Risk Assessment, Seoul National University, woojin1@snu.ac.kr)

** 서울대학교 대학원 건설환경공학부 석사과정(Department of Civil & Environmental Engineering, Seoul National University, esw1026@snu.ac.kr)

*** 교신저자 · 정회원 · 서울대학교 건설환경공학부 정교수(Corresponding author, Member, Department of Civil & Environmental Engineering, Seoul National University, kiyun@snu.ac.kr)

소셜 네트워크 서비스에 대한 개발 및 활용 사례가 점차 늘고 있다. 이러한 움직임에 힘입어 최근에 소셜 미디어 정보를 이용하여 기존의 통계적 정보와의 연관성을 찾기 위한 연구가 다양한 분야에서 진행되고 있다. 특히 구글 검색 추이를 이용하여 독감 유행을 예측하는 모델을 개발한 사례, 카드 결제정보의 공간적 분포 패턴을 이용하여 심야버스노선을 설계한 사례 등이 대표적이다.

이와 관련된 연구로, Mei et al.,(2006)의 연구에서는 웹 블로그들을 분석하여 태풍이나 아이팟 나노 출시 등의 이벤트가 발생하였을 때의 사회적 트렌드와 어떤 연관성이 있는지에 대한 연구를 시공간적 패턴 분석 기법을 통해 실시한 바 있다. Sakaki et al.,(2010)의 연구에서는 일본 지역 내에서의 지진의 발생과 트윗 데이터의 공간적 분포 패턴 간의 유사성 및 시간적 추이를 분석하여 지진이 발생할 위치를 예측하는 모델을 개발한 바 있다. Wu et al.,(2009)의 연구에서는 주택가격지수와 구글 검색 빈도와의 시계열적 상관성을 분석하여 주택 가격에 대한 예측 모델을 개발하고 검증한 바 있다. Lee et al.,(2013)의 연구에서는 소셜 미디어 데이터를 비롯한 빅데이터의 분석기법에 대해 통계 모델, 텍스트 마이닝, 심리분석 기법 등이 적용될 필요가 있다고 밝히고 있다. Dashti et al.,(2014)의 연구에서는 지진이나 태풍과 같은 자연재해로 인한 피해 정보를 취합하기 위해 소셜 미디어 데이터를 활용하는 방안에 대한 연구를 진행하였다.

이와 같이 위치기반 소셜 미디어 정보의 분포 패턴과 기존의 통계적 데이터 간의 공간적 상관성을 찾는 연구 사례는 점차 증가하는 추세에 있는 것으로 판단된다 (Mennis and Guo, 2009; Stefanidis et al., 2013). 이에 본 연구에서는 위치기반 소셜 미디어 정보의 공간적 분포와 부동산 가격 정보의 공간적 분포를 비교하여 두 데이터 간의 공간적 관련성을 파악하고자 하였으며, 이를 위해 격자기반 분석 방법론을 적용해보고자 한다.

1.2 연구의 범위 및 방법

본 연구의 대상 데이터는 지오태그 정보를 포함하고 있는 트윗(Tweet) 데이터를 주요 대상 데이터로 하였고 부동산 가격 정보로는 행정구역(시군구) 단위의 주택가격지수(매매지수, 전세지수)를 보조 데이터로 하였다. 두 대상 데이터의 시간적 범위는 2013년 8월 한 달간 수집된 자료로 한정하였고 공간적 범위는 서울특별시를 중심으로 하여 일산, 분당, 과천 등 주요 신도시들을 포함하는 수도권으로 하였다.

위치기반 소셜 미디어 정보와 부동산 가격 정보와

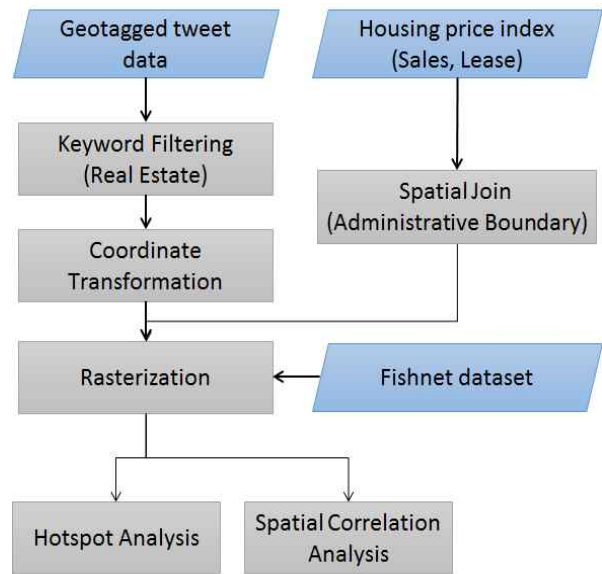


Figure 1. Workflow of the study

의 공간적 관련성을 분석하기 위한 방법론으로 본 연구에서는 두 자료의 격자 자료화(Rasterization)를 통한 데이터 통합, 각 대상 데이터의 공간적 분포를 시각적으로 확인하기 위한 Hotspot 분석, 그리고 공간시차를 고려한 상관관계 분석으로 구성된 일련의 방법론을 적용해 보았다.

Fig. 1은 본 연구에서 적용한 연구흐름도를 나타낸다.

2. 대상데이터

본 연구의 대상 데이터는 소셜 미디어 데이터와 부동산 가격지수 데이터이다. 먼저, 소셜 미디어 데이터로는 지오태그 정보를 포함하고 있는 트윗 데이터를 주요 대상 데이터로 하였다. 트윗 데이터는 2013년 8월 한 달간 수집된 자료로 한정하였고 공간적 범위는 서울특별시를 중심으로 하여 일산, 분당, 과천 등 주요 신도시들을 포함하는 수도권으로 하였다. 즉, 트윗 데이터의 위치값에 대하여 사각형 형태의 공간질의(Spatial Query)를 적용하였는데, 범위는 좌상단 (176778.7, 472112.8), 우하단 (218778.7, 432112.8)으로 하였다.

부동산 가격 정보로는 한국감정원에서 구축, 제공하고 있는 행정구역(시군구) 단위의 주택가격지수(매매지수, 전세지수)를 보조 대상 데이터로 하였으며, 시간적 범위는 2013년 8월로 한정하고 공간적 범위는 트윗 데이터의 수집범위에 맞게 서울특별시(25개 구), 인천광역시(계양구, 부평구, 남동구의 3개 구), 경기도(고양시, 부천시, 안양시, 성남시 등) 일부지역을 포함하여 총 51개 시군구에 대한 주택가격지수를 수집하였다.

3. 격자기반 공간분석 방법론

3.1 데이터 공간자료화

본 연구에서는 위치 태그를 가진 트윗 데이터를 수집하기 위해, 트위터 공개 API¹⁾를 활용하여 위치 태그를 가진 트윗 데이터 중 대상 범위 내에 포함되는 트윗 데이터만을 자동으로 추출하였다. 키워드 필터링을 적용하였는데, 부동산 관련 키워드(부동산, 주택가격, 집값, 아파트 시세, 실거래가, 주택매매, 전세, 월세, 아파트 분양, 양도소득세 등)를 포함한 트윗 데이터들을 필터링하였다. 총 193,831개의 트윗 데이터 중 키워드 필터링을 통해 250개 데이터가 추출되었다. 트윗 데이터 지오태그 상의 위치정보는 경위도 좌표로 되어 있어 좌표 변환을 통해 경위도 좌표를 TM(Transverse Mercator) 좌표체계로 변환하였으며, 이를 위해 국토지리정보원의 NGIpro(Ver. 2.53) 소프트웨어를 활용하였다.

다음으로, 주택가격지수(매매지수, 전세지수) 데이터는 행정경계(시군구) 폴리곤 데이터를 활용하여 공간조인(Spatial Join)을 통해 공간자료화 하였다.

3.2 격자 데이터화(Rasterization)

트윗 데이터와 주택가격지수 데이터는 각각 포인트 데이터와 폴리곤 데이터 형태이기 때문에 데이터의 형식에도 차이가 있지만 공간적 세밀도에서도 차이를 보인다. 이러한 데이터 차이가 있어 두 데이터의 공간적 분포를 비교하기 위해서는 통합된 형태의 데이터를 구축할 필요가 있다(Liao et al., 2010). 이에 본 연구에서는 격자 데이터셋 형태로 두 데이터를 변환시킴으로써 공간적 분포의 비교분석과 공간적 상관성 분석을 보다 용이하게 하고자 하였다.

격자형 데이터셋의 구축을 위해 우선, ArcMap Toolbox 중 'Create Fishnet' 기능을 적용하여 격자망 데이터를 생성하였다. 격자망 데이터셋의 범위는 트윗 데이터를 추출한 범위와 동일하게 설정하였다. 적절한 격자의 크기를 도출하는데 있어서, 일반적인 방격분석(Quadrat Analysis)에서 사용하는 점 자료 수 대비 대상면적 비율²⁾, 대상지역 내 시군구 최소면적³⁾ 등을 중

합적으로 고려하여 2km×2km 격자크기를 적용하였다(Yu, 1998).

이러한 격자망 데이터셋에 트윗 데이터를 중첩시켜 각 격자 내에 포함된 트윗 데이터의 개수를 세어 격자의 속성정보로 입력하였다. 주택가격지수 역시 격자망 데이터셋과 중첩하여 각 격자에 대한 가격지수를 재계산 하여 입력하였다. 이때 한 격자 내에 다수의 행정경계 폴리곤이 만날 경우에는 각 폴리곤의 격자 내 포함 비율을 계산하여 면적 비율을 가중치로 한 주택가격지수의 가중평균합을 구하였다(Shin, 2004; Jung et al., 2009).

3.3 Hotspot 분석

본 연구에서는 소셜 미디어인 트윗 데이터와 주택가격지수 데이터의 전반적인 공간적 분포를 비교하기 위해 Hotspot 분석을 실시하였다. Hotspot 분석은 데이터셋의 각 사상(Feature)에 대해 Getis-Ord Gi* 통계량을 계산하고 z값과 p값을 이용하여 높은 값 혹은 낮은 값의 사상이 공간적으로 군집을 이루는지를 탐색해준다(Kim and Park, 2013). 본 연구에서는 ArcMap Toolbox 중, 'Hotspot Analysis' 도구를 이용하여 격자화된 트윗 데이터와 주택가격지수 데이터 각각에 대한 국지적 공간 클러스터링을 탐색하여 시각적으로 비교하였다.

3.4 공간적 상관성 분석

공간적 상관성 분석은 트윗 데이터와 주택가격지수 데이터의 공간적인 분포 패턴이 어느 정도의 유사성 또는 상관성을 가지는지를 수치적으로 파악하기 위한 과정이다. 이를 위해 본 연구에서는 이변량 공간 상관성 측도로, 공간시차를 적용한 피어슨 상관계수를 적용하였다(Lee, 2001). 이 상관계수는 일반적인 상관분석에서 쓰이는 피어슨 상관계수에 공간시차(Spatial Lag)를 적용한 것으로, 여기서 공간시차란 공간가중치 행렬(Spatial Weight Matrix)에 의해 정의된 이웃들의 가중평균합이다. 트윗 데이터와 주택가격지수의 격자 데이터에 대하여 식(1)과 같이 공간적 상관계수(Spatial Correlation Coefficient, SCC)를 계산하였다(Lee, 2014).

1) 트위터 공개 API: 소셜 네트워크 서비스 업체인 '트위터(Twitter)'에서 자사의 트윗 데이터의 활용을 지원하기 위해 개발자들에게 공개한 데이터 제공 인터페이스. 트윗, 사용자, 위치정보 등의 정보를 수집할 수 있음. <https://dev.twitter.com/overview/api> (2014년 12월 28일 방문)
 2) 점 자료 수 대비 대상면적 비율(A/N), A는 대상지역의 면적, N은 점의 수임. 본 연구에서는 A=1680km², N=250 이므로 이를 대입하면 격자면적은 6.72km² 이고, 정방형 격자인 경우 한 변 길이는 2.59km가 도출됨

3) 대상지역 내 행정구역 최소 면적이 10km²(중구)이고 이에 대한 제곱근은 3.16km 수준이므로 정방형 격자를 고려하면 격자의 한 변 크기는 이보다 크지 않아야 함

$$SCC_{AB} = \frac{\sum_i^n (\tilde{A}_i - \bar{\tilde{A}})(\tilde{B}_i - \bar{\tilde{B}})}{\sqrt{\sum_i^n (\tilde{A}_i - \bar{\tilde{A}})^2} \sqrt{\sum_i^n (\tilde{B}_i - \bar{\tilde{B}})^2}} \quad (1)$$

여기서 SCC_{AB} 는 A, B 데이터 간의 공간적 상관계수, i 는 각 격자, n 은 격자의 개수, \tilde{A}_i 는 i 번째 격자에서 A 데이터의 공간시차에 의한 격자값, \tilde{B}_i 는 i 번째 격자에서 B 데이터의 공간시차에 의한 값이다. 이에 대한 결과값은 -1과 1사이의 범위를 갖게 되며, 양수는 양의 상관성, 음수는 음의 상관성을 나타내고 절대값이 클수록 상관성의 정도가 크다는 것을 나타낸다(Lee, 2014).

4. 적용 및 결과

4.1 데이터 격자화 및 통합

Fig. 2는 트윗 데이터와 시군구 행정구역경계, 그리고 격자망 데이터셋을 크기별로 중첩시킨 결과를 나타낸다.

Fig. 3, Fig. 4는 각각 주택매매가격지수와 주택전세가격지수의 격자화된 데이터의 공간적 분포를 나타낸다.

Fig. 3에 따르면 주택매매가격지수의 경우 과천시 지역이 높은 값을 나타내고 있으며, 파주시와 고양시 일

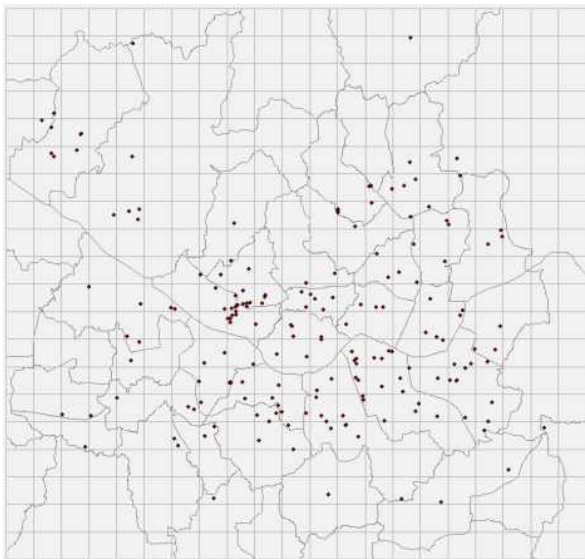


Figure 2. Tweets(Point), administrative boundary (Polygon) and fishnet dataset(Grid) for the test area

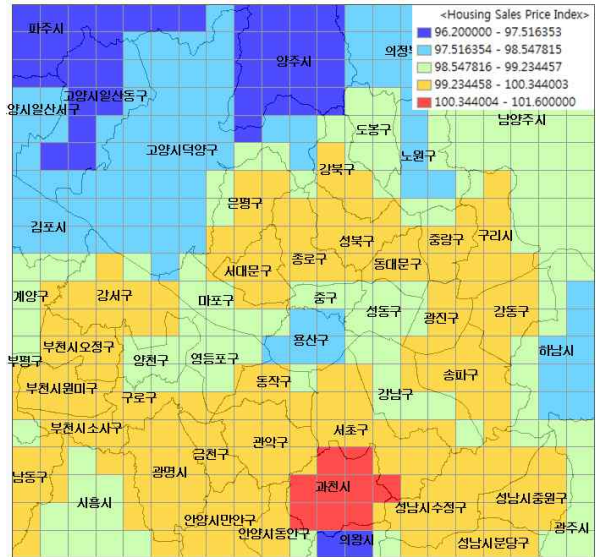


Figure 3. Grid data of housing sales price index

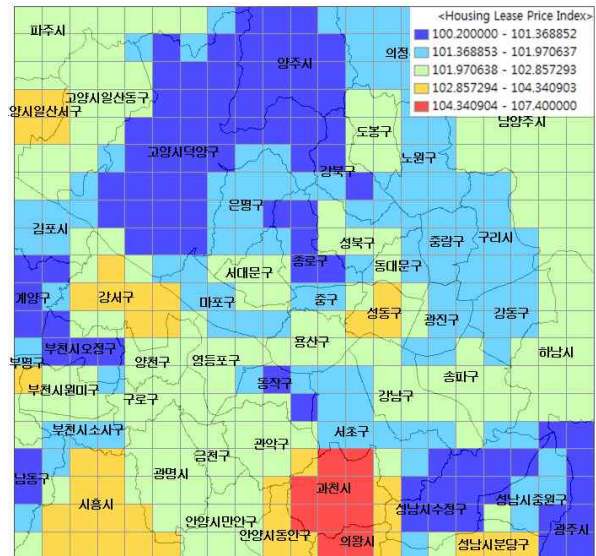


Figure 4. Grid data of housing lease price index

산동구, 양주시 지역 등이 낮은 값을 나타내고 있다. 또한 Fig. 4에 따르면 주택전세가격지수의 경우 과천시와 의왕시 지역이 높은 값을 나타내고 있으며, 고양시 덕양구, 양주시, 성남시 수정구, 종로구 등 지역이 낮은 값을 나타내고 있다.

4.2 Hotspot 분석 결과

Fig. 5, Fig. 6, Fig. 7은 각각 격자화된 트윗 데이터와 주택가격지수 데이터에 대한 Hotspot 분석 결과를 나타낸다.

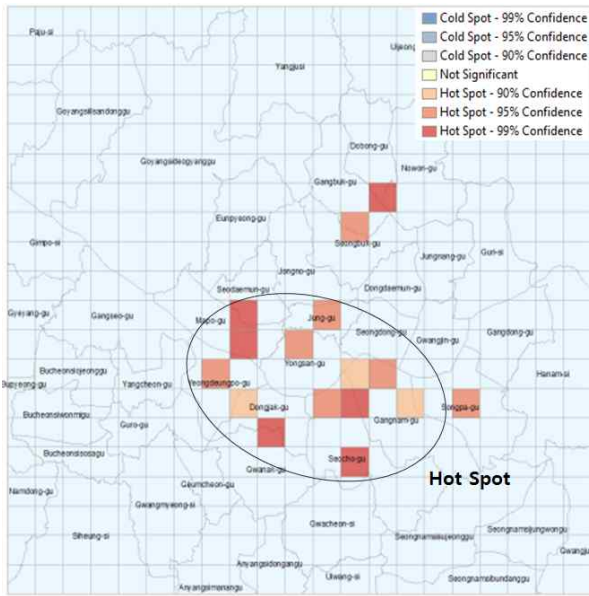


Figure 5. Results of hotspot analysis using tweet data

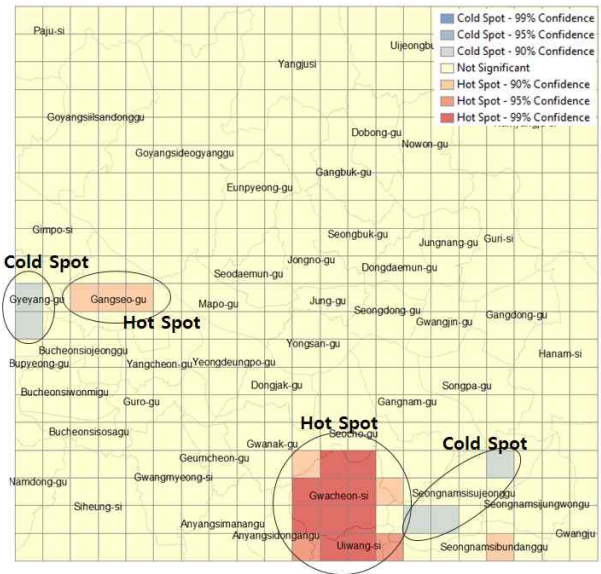


Figure 7. Results of hotspot analysis using housing lease price index

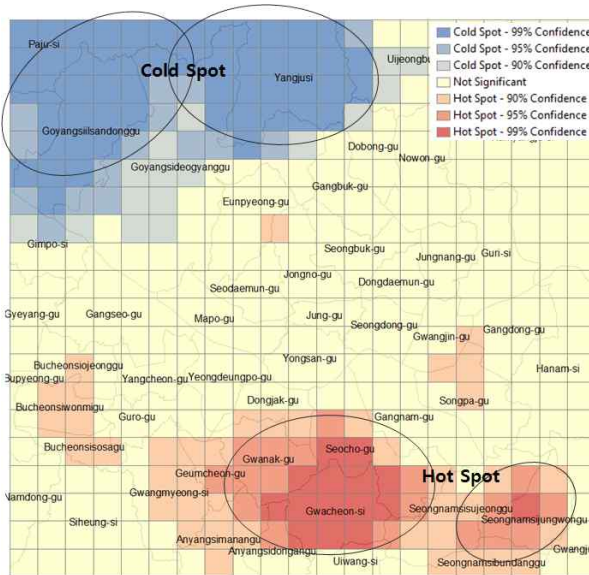


Figure 6. Results of hotspot analysis using housing sales price index

Fig. 5, Fig. 6, Fig. 7을 통해 세 데이터에 대한 Hotspot 지역과 Coldspot 지역을 비교 분석해본 결과는 다음과 같다. 트윗 데이터에 대한 Hotspot 지역은 마포구와 서대문구, 강북구, 서초구 등 지역에서 탐색되었다. 주택매매가격지수 데이터에 대한 Hotspot 지역은 과천시를 중심으로 하여 관악구, 서초구, 성남시 중원구 등의 지역에서 가장 두드러지게 나타났으며, 부천시 지역에서도 약한 Hotspot 지역이 탐색되었다. 고양시 지역은 전체적으로 Coldspot 지역으로 탐색되었다.

주택전세가격지수 데이터에 대한 Hotspot 지역은 매매 가격지수와 마찬가지로 과천시가 강한 Hotspot 지역으로 나타났으며, 안양시 동안구까지 Hotspot 지역이 나타났다.

결과적으로, 서초구는 트윗 데이터와 주택매매가격지수 데이터에서 공통적인 Hotspot 지역으로 나타났으며, 주택전세지수 데이터와 트윗 데이터 간에는 뚜렷한 공간적 상관성이 나타나지 않았다.

4.3 공간적 상관성 분석 결과

위의 3.4절에서 언급한대로 트윗 데이터와 주택가격지수 데이터의 공간적 상관성을 분석하기 위해 격자화된 트윗 데이터와 매매지수, 전세지수 데이터 간의 SCC를 측정하였다. 아래의 표는 측정된 SCC값을 정리한 표이다(Table 1).

Table 1에 따르면 전체적으로 주택가격지수 데이터의 공간적 분포는 트윗 데이터의 공간적 분포와 음의 상관성을 띄는 것을 확인할 수 있고 특히 주택전세가격지수보다는 주택매매가격지수가 음의 상관성이 더 뚜렷한 것으로 나타났다.

Table 1. SCC between tweet data and housing price index

SCC	Tweet
Housing sales price index	-0.2736
Housing lease price index	-0.0699

5. 결 론

본 연구에서는 위치기반 소셜 미디어 데이터의 공간적 분포가 지역별 부동산 지수와 어떠한 공간적 관련성을 가지는지에 대해 알아보고자 하였다. 두 데이터는 상이한 자료 형식을 가지고 있어, 이를 보완할 수 있는 방법론으로 본 연구에서는 격자 기반의 공간분석 방법을 적용하고자 하였다. 대상 데이터로는 서울과 수도권 일부 지역에 대해 2013년 8월 한 달간의 지오테그된 트위터 데이터와 행정구역별 주택가격지수(매매, 전세)를 이용하였다. 두 데이터 간의 상이한 공간적 단위를 고려하여 2,000m 단위의 격자망을 구성하고 이에 맞게 두 데이터를 격자화하였다.

통합된 두 데이터에 대하여 Hotspot 분석을 실시하여 공간적 분포를 시각적으로 비교하였으며, 두 데이터 간의 공간적 상관계수를 측정함으로써 정량적 분석을 실시하였다. 분석 결과, 서초구 지역에서 트위터 데이터와 주택매매가격지수 데이터에 대한 공통적인 Hotspot 지역이 탐색되었다. 그러나 전체적으로 주택가격지수 데이터의 공간적 분포는 트위터 데이터의 공간적 분포와 음의 상관성을 띄는 것을 확인할 수 있고 특히 주택전세가격지수보다는 주택매매가격지수가 음의 상관성이 더 뚜렷한 것으로 나타났다.

이는 주택가격지수의 공간적 분포를 설명하는데 있어서 트위터 데이터를 활용하였을 때, 일부 지역에 대해서는 유사한 분포를 탐지할 수 있으나 전체적인 경향을 비교하였을 때에는 상관성이 높지 않은 것으로 해석할 수 있다. 따라서 소셜 미디어 데이터를 활용하여 주택가격지수의 패턴을 분석, 예측하기 위해서는 보다 다양한 요소를 고려할 필요가 있을 것으로 판단된다.

향후에는 보다 넓은 대상지역과 시간적 범위, 그리고 다양한 분야에 대해 소셜 미디어 데이터와 부동산 관련 통계데이터를 수집, 분석함으로써 보다 다양하고 심도 깊은 분석방법들을 적용하는 것이 필요할 것으로 판단된다. 또한, 공간적 상관성뿐만 아니라 시간적 상관성까지 같이 고려한다면 보다 유의미한 관련성을 도출할 수 있을 것으로 판단된다. 이러한 노력들은 향후 소셜 미디어 정보를 이용한 다양한 공간 빅데이터 분석 모델을 개발하는 데 있어서 유용하게 활용될 수 있을 것으로 예상된다.

감사의 글

본 연구는 ‘국토교통부 국토공간정보연구사업 국토공간정보의 빅데이터 관리, 분석 및 서비스 플랫폼 기

술개발(14NSIP-B091011-01)과제’의 연구비 지원에 의해 연구되었습니다.

References

1. Dashti, S., Palen, L., Heris, M. P., Anderson, K. M., Anderson, S., and Anderson, T. J., 2014, Supporting disaster reconnaissance with social media data: a design-oriented case study of the 2013 colorado floods, Proceedings of the 11th International ISCRAM Conference.
2. Jung, D., Kim, S., and Kim, K., 2009, The central place analysis with the characteristics of the distribution of the land price using GIS, Journal of the Korean Society for Geospatial Information System, Vol. 17, No. 3, pp. 420-421.
3. Kim, G., and Park, G., 2013, Hot spot analysis on forest carbon stocks using getis-ord spatial statistics, Proceedings of 2012 Summer Conference, Korea Forest Society, pp. 420-421.
4. Lee, B., Lim, J., and Yoo, J., 2013, Utilization of social media analysis using big data, Journal of the Korea Contents Association, Vol. 13, No. 2, pp. 211-219.
5. Lee, S., 2001, Developing a bivariate spatial association measure: an integration of pearson's r and moran's i, Journal of Geographical Systems, Vol. 3, No. 4, pp. 369-385.
6. Lee, Y., 2014, A study on detection methodology for influential areas in social network using spatial statistical analysis methods, Journal of the Korean Society for Geospatial Information System, Vol. 22, No. 4, pp. 21-30.
7. Liao, S., and Bai Y., 2010, A new grid-cell-based method for error evaluation of vector-to-raster conversion, Computational Geosciences, Vol. 14, No. 4, pp. 539-549.
8. Mei, Q., Liu, C., Su, H., and Zhai, C., 2006, A probabilistic approach to spatiotemporal theme pattern mining on weblogs, Proceedings of the 15th international conference on World Wide Web, ACM.
9. Mennis, J., and Guo, D., 2009, Spatial data mining and geographic knowledge discovery—an introduction, Computers, Environment and Urban Systems, Vol. 33, No. 6, pp. 403-408.
10. Sakaki, T., Okazaki, M., and Matsuo, Y., 2010, Earthquake shakes twitter users: real-time event

- detection by social sensors, Proceedings of the 19th international conference on World Wide Web, ACM.
11. Shin, J., 2004, Research on areal interpolation methods and error measurement techniques for reorganizing incompatible regional data units, Journal of the Korean Association of Regional Geographers, Vol. 10, No. 2, pp. 389-406.
 12. Stefanidis, A., Crooks, A., and Radzikowski, J., 2013, Harvesting ambient geospatial information from social media feeds, GeoJournal, Vol. 78, No. 2, pp. 319-338.
 13. Wu, L., and Brynjolfsson, E., 2009, The future of prediction: how google searches foreshadow housing prices and sales, NBER Conference Technological Progress & Productivity Measurement, WISE, ICIS.
 14. Yu, K., 1998, Generalization of point feature in digital map through point pattern analysis, Journal of GIS Association of Korea, Vol.6, No. 1, pp. 11-23.