

2022 UROP 영역별 감정 사전 기반 분석 보고서

소속: 국민대학교 소프트웨어학부
학번: 20181703
이름: 평선호
지도교수: 박하명 교수
팀원: 윤상원 학우(20181651)

1. 서론

현재 인터넷상의 데이터가 급증함에 따라 사회 또는 기업에서 해당 데이터를 적절히 활용하여 서비스를 제공하고 있다. 데이터를 활용하는 분야 중에서 주위에서 제일 접하기 쉽고 상품을 구매할 때 참고를 많이 하는 리뷰데이터를 선정하였다.

리뷰데이터의 긍정,부정을 판별하고 상품별, 즉 영역별로 감성사전을 구축하는 과정을 통해 단순 긍정,부정을 평가하는 것이 아닌 상품의 소비방법에 맞는 감성사전을 구축하였다. 이를 통해 리뷰를 입력했을 때 리뷰 속 긍정,부정 단어들을 판별하여 별점예측을 하는 서비스를 제공하고자 한다.

감성사전 기반 감정 분석을 통해 리뷰 속 어휘들의 극성을 판별하고 이를 계산하여 점수화를 진행하였다.

2. 연구배경

2.1 연구데이터 수집

본 연구에서는 다양한 종류의 물품을 판매하는 전자상거래 사이트인 쿠팡(coupang) 리뷰데이터를 선정하여 영역별로 '가전 디지털', '생활 용품', '여성패션', '스포츠/레저', '완구/취미' 5가지 영역을 선정하여 상품 리뷰 데이터 크롤링을 진행하였다. 이때 상품 데이터 수집은 python의 selenium 라이브러리를 사용하여 크롤링을 통해 진행하였다. 이때 크롤링을 했을 때, 기존 사용자들이 입력한 별점을 바탕으로 데이터의 긍정,부정을 분류하였다. 리뷰 데이터의 약 80~90%정도가 별점 4,5점을 주면서 긍정적인 어휘를 바탕으로 리뷰를 작성하는 것을 확인하여 별점 4,5점을 긍정 리뷰로 판단하고 별점 1,2,3점의 리뷰들은 부정 리뷰로 판단하였다.

2.2 연구 절차

2.2.1 긍정 부정 분류 후 데이터 전처리

본 연구를 진행하는 과정 중, 크롤링한 리뷰데이터를 바탕으로 soynlp 오픈소스를 사용하여 한국어 분석을 진행했다. 크롤링한 리뷰데이터에는 기존 형식화 된 글들이 아닌 불용어가 포함되어 있는 리뷰데이터가 존재하였다. 감성 사전구축에 불필요하게 반복되는 이모티콘, 문구, 영어 리뷰, 오타 수정 등 리뷰 전처리 과정을 불필요한 요소들을 정리하였다.

2.2.2 데이터 태깅 및 어휘 통일 진행

한국어 특성상 하나의 어근에서 파생될 수 있는 많은 단어들이 존재한다. 이러한 한국어의 특징이 존재할 시 사전을 구축하는데 영향이 끼칠수 있으므로 어휘들을 토큰화를 진행하였다. 추가적으로 모든 어휘들을 형태소 분석을 통해 모든 어휘들에 형태소를 태깅을 진행하였다.

Ex) 맛있다, 맛있고, 맛있는 -> 맛있다 / Verb 로 통일을 하며 토큰화를 진행을 하였다.

2.2.3 기준 단어 선정

감성사전을 구축하기 위해 본 연구에서 긍정, 부정 리뷰 데이터속에서 기준 단어를 선정하였다. 해당 기준 단어들은 추후 리뷰 데이터들과 SO-PMI 계산을 통해 어휘들이 극성을 판별 할 수 있는 기준이 되어 주기때문에 기준단어라 불린다.

형태소만으로 극성을 판별하기 힘든 동사,부사 등은 제외하고 극성을 판별하기 쉬운 형용사를 기준으로 선정하여 리뷰데이터 속 형용사인 어휘들의 빈도 수를 계산하였다. 긍정 리뷰 속 형용사 어휘의 빈도수와 어휘를 쌍으로 계산하고, 부정 리뷰 속 형용사 어휘의 빈도수와 어휘를 쌍으로 계산하였다.

특정 단어의 경우 긍정, 부정의 뜻을 가지는 것이 아니라, 단순히 자주 등장 하는 단어가 기준단어로 선정되는 것을 방지해야한다. 이를 해결하고자 긍정리뷰속에서 구한 기준단어들의 집합과 부정리뷰속에서 구한 기준단어들의 집합을 서로 빼주며 빈도수 문제를 해결하였다.

2.2.4 단어 극성 판별

리뷰데이터의 단어의 극성을 판별하기 위해 이전 과정에서 기준단어를 바탕으로 SO-PMI과정을 진행한 다. PMI(Pointwise Mutual Information)란? 두 단어 간의 유사성을 분석하는 지표로 활용되는 방법 이다.

$$PMI(w_1, w_2) = \log_2 \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$$

여기서 w_1 과 w_2 는 분석하고자 하는 두 단어를 나타내며 $PMI(w_1, w_2)$ 는 한 문서내에서 w_1 과 w_2 가 함께 출현할 확률을 w_1 과 w_2 가 각각 출현할 확률의 곱으로 나눈 값이다. 따라서 $PMI(w_1, w_2)$ 는 두 단어

가 같은 문서에서 나타날 확률을 나타내며 값이 클수록 두단어 유사성, 즉 극성이 비슷한 것으로 판별 할 수 있다. 하지만 단어의 출현 빈도수만으로 감성극성 판별은 단어의 극성강도를 정확히 판별하기 어려운 문제가 있을 수 있고, 기준단어의 선정에 따라 단어극성 판별의 편차가 크게 달라지는 문제점이 발생한다. 이러한 문제점을 보완하고자 보다 극성을 정확히 판별하는 방법인 SO-PMI(Semantic Orientation From Pointwise Mutual Information)방식을 사용하였다.

$$SO-PMI(w) = \sum_{pw \in PW} PMI(w, pw) - \sum_{nw \in NW} PMI(w, nw)$$

여기서 PW는 긍정기준단어의 집합, NW는 부정기준단어의 집합을 의미한다. SO-PMI(w)는 단어 w와 긍정단어 집합과의 PMI합에서 부정단어 집합과의 PMI합을 뺀 결과이다.

PMI를 계산하기 위해 soynlp의 모듈 중 하나인 pmi_func을 불러와 PMI계산을 진행하였다. 본 연구에서 감성 사전에 등재되는 단어는 형용사와 부사단어만을 활용하여 구축하였다. 따라서 전체 리뷰에서 등장하는 형용사와 부사 단어들을 추출하여 긍정 기준단어집합과 부정 기준단어집합과의 PMI를 계산한 다음 이를 빼면서 단어들의 SO-PMI를 계산한다.

2.2.5 감성사전 구축

전체 데이터에서 형용사와 부사들을 기준단어들과 PMI를 계산한 다음 이를 활용하여 감성사전 구축을 진행하였다. 점수화된 단어들의 집합들을 SO-PMI 값에 따라 ‘강한 긍정’, ‘약한 긍정’, ‘중립’, ‘약한 부정’, ‘강한 부정’ 5종류로 분류하였다. SO-PMI 점수를 기준으로 상위 20%가 ‘강한 긍정’, 상위 20~40%가 ‘약한 긍정’, 상위 40~60%가 ‘중립’, 상위 60~80% ‘약한 부정’, 상위 80~100%가 ‘강한 부정’으로 분류하였다.

2.2.6 감성 사전을 활용한 평점 예측

테스트 데이터 리뷰 집합의 감성 점수를 측정한다. 해당 리뷰의 물품이 속한 영역의 감성 사전을 사용하여 ‘강한 긍정’ 단어는 +2, ‘약한 긍정’ 단어는 +1, ‘중립’ 단어는 0, ‘약한 부정’ 단어는 -1, ‘강한 부정’ 단어는 -2점으로 계산한다. 이때 해당 사전을 json형식으로 파일화 하여 리뷰를 입력하면 json형식의 사전에 등재된 단어들을 점수로 변환하여 모두 합한 값을 단어의 수로 나누어 최대값이 2와 최소값이 -2로 나오게 고정한다.

```

"data": [
  {
    "index": 0,
    "text": "좋다",
    "score": 2
  },
  {
    "index": 1,
    "text": "편하다",
    "score": 2
  },
  {
    "index": 2,
    "text": "고급스럽다",
    "score": 2
  },
  {
    "index": 3,
    "text": "맛있었다",
    "score": 2
  },
  {
    "index": 4,
    "text": "깔끔하다",
    "score": 2
  },
  {
    "index": 5,
    "text": "예쁘다",
    "score": 2
  },
  {
    "index": 6,
    "text": "가볍다",
    "score": 2
  }
]

```

```

{
  "index": 118,
  "text": "힘들다",
  "score": -2
},
{
  "index": 119,
  "text": "좁다",
  "score": -2
},
{
  "index": 120,
  "text": "달달하다",
  "score": -2
},
{
  "index": 121,
  "text": "너무하다",
  "score": -2
},
{
  "index": 122,
  "text": "작다",
  "score": -2
},
{
  "index": 123,
  "text": "없다",
  "score": -2
},
{
  "index": 124,
  "text": "짧다",
  "score": -2
},

```

왼쪽 이미지가 사전에서 긍정 단어로 선정된 단어들의 일부이며, 오른쪽 이미지가 사전에서 부정 단어로 선정된 단어의 일부이다.

이후 평점 분포에 알맞게 감정 점수를 평점으로 변환한다. 본 연구를 진행함에 있어 수집한 리뷰 데이터를 분석 하였을 때, 리뷰의 대부분이 5점과 4점으로 분포하여있고, 3점 이하의 점수를 갖는 리뷰의 비중이 상당히 작다는 것을 확인 할 수 있었다. 5점과 4점의 분포를 계산하여 평점 분포를 평점 예측에 활용하였다. 훈련 데이터를 통해 만들어진 평점 비율을 앞서 구한 감성 점수들에 적용해 평점으로 변환한다.

3. 연구결과

<p>5,배송빠르고 2,택배가 엉망이네용 저희집 밑에층에 말도없이 놔두고가고 5,아주좋아요 바지 정말 좋아서2개 더 구매했어요 이가격에 대박입니다. 바느질이 조금 엉성하긴 하지만 편하고 가성비 최고예요. 2,선물용으로 빨리 받아서 전달했어야 하는 상품이었는데 머그컵만 와서 당황했습니다. 전화했더니 바로주신다했지만 배송도 누락되어있었네요.. 확인안하고 바로 선물했으면 큰일날뻔했네요..이렇게 5,민트색상 예뻐요. 옆 손잡이는 거는 용도로도 사용되네요 ㅎㅎ 2,비추합니다 계란 뒤집을 때 완전 불편해요 πππ 코팅도 문어내고 보기엔 예쁘고 실용적으로 보였는데 생각보다 진짜 별로입니다. 1,주문을 11월6에 시켰는데 11월16일에 배송이 왔네요 ㅎㅎㅎ 여기 회사측과는 전화도 안되고 아무런 연락을 받을수가 없으니 답답하신 분들은 다른곳에서 사시는거 추천드립니다 2,넉넉한 길이로 주문했는데도 안 맞네요 별로예요 2,보풀이 계속 때처럼 나오다가 지금은 안나네요~ 2,110인데 전문속옷브랜드 위생팬티105보다 작은듯해요. 불편해요. 밴딩부분이 다 신축성없는 일반실로 되어있어 빅사이즈임에도 빅사이즈같지않아요. 입고벗을때 편하게 밴딩부분이 늘어나고 입으 5,사이즈도 딱이고 귀엽고 넘 좋아요 ㅎㅎ 2,베이지 색 구매했는데 약간 살색에 가까워요 2,화면별인가봐요;; 노란컬러가 돋보여요;; 저렴한맛에 그냥 씁니다 2,별루 1.3.5.7.9.11.13.15.17.19.21.23.25.27.29.31.33.35.37.39.41.43.45.47.49.51.53.55.57.59.61.63.65.67.69.71.73.75.77.79.81.83.85.87.89.91.93.95.97.99.101.103.105.107.109.111.113.115.117.119.121.123.125.127.129.131.133.135.137.139.141.143.145.147.149.151.153.155.157.159.161.163.165.167.169.171.173.175.177.179.181.183.185.187.189.191.193.195.197.199.201.203.205.207.209.211.213.215.217.219.221.223.225.227.229.231.233.235.237.239.241.243.245.247.249.251.253.255.257.259.261.263.265.267.269.271.273.275.277.279.281.283.285.287.289.291.293.295.297.299.301.303.305.307.309.311.313.315.317.319.321.323.325.327.329.331.333.335.337.339.341.343.345.347.349.351.353.355.357.359.361.363.365.367.369.371.373.375.377.379.381.383.385.387.389.391.393.395.397.399.401.403.405.407.409.411.413.415.417.419.421.423.425.427.429.431.433.435.437.439.441.443.445.447.449.451.453.455.457.459.461.463.465.467.469.471.473.475.477.479.481.483.485.487.489.491.493.495.497.499.501.503.505.507.509.511.513.515.517.519.521.523.525.527.529.531.533.535.537.539.541.543.545.547.549.551.553.555.557.559.561.563.565.567.569.571.573.575.577.579.581.583.585.587.589.591.593.595.597.599.601.603.605.607.609.611.613.615.617.619.621.623.625.627.629.631.633.635.637.639.641.643.645.647.649.651.653.655.657.659.661.663.665.667.669.671.673.675.677.679.681.683.685.687.689.691.693.695.697.699.701.703.705.707.709.711.713.715.717.719.721.723.725.727.729.731.733.735.737.739.741.743.745.747.749.751.753.755.757.759.761.763.765.767.769.771.773.775.777.779.781.783.785.787.789.791.793.795.797.799.801.803.805.807.809.811.813.815.817.819.821.823.825.827.829.831.833.835.837.839.841.843.845.847.849.851.853.855.857.859.861.863.865.867.869.871.873.875.877.879.881.883.885.887.889.891.893.895.897.899.901.903.905.907.909.911.913.915.917.919.921.923.925.927.929.931.933.935.937.939.941.943.945.947.949.951.953.955.957.959.961.963.965.967.969.971.973.975.977.979.981.983.985.987.989.991.993.995.997.999.1001.1003.1005.1007.1009.1011.1013.1015.1017.1019.1021.1023.1025.1027.1029.1031.1033.1035.1037.1039.1041.1043.1045.1047.1049.1051.1053.1055.1057.1059.1061.1063.1065.1067.1069.1071.1073.1075.1077.1079.1081.1083.1085.1087.1089.1091.1093.1095.1097.1099.1101.1103.1105.1107.1109.1111.1113.1115.1117.1119.1121.1123.1125.1127.1129.1131.1133.1135.1137.1139.1141.1143.1145.1147.1149.1151.1153.1155.1157.1159.1161.1163.1165.1167.1169.1171.1173.1175.1177.1179.1181.1183.1185.1187.1189.1191.1193.1195.1197.1199.1201.1203.1205.1207.1209.1211.1213.1215.1217.1219.1221.1223.1225.1227.1229.1231.1233.1235.1237.1239.1241.1243.1245.1247.1249.1251.1253.1255.1257.1259.1261.1263.1265.1267.1269.1271.1273.1275.1277.1279.1281.1283.1285.1287.1289.1291.1293.1295.1297.1299.1301.1303.1305.1307.1309.1311.1313.1315.1317.1319.1321.1323.1325.1327.1329.1331.1333.1335.1337.1339.1341.1343.1345.1347.1349.1351.1353.1355.1357.1359.1361.1363.1365.1367.1369.1371.1373.1375.1377.1379.1381.1383.1385.1387.1389.1391.1393.1395.1397.1399.1401.1403.1405.1407.1409.1411.1413.1415.1417.1419.1421.1423.1425.1427.1429.1431.1433.1435.1437.1439.1441.1443.1445.1447.1449.1451.1453.1455.1457.1459.1461.1463.1465.1467.1469.1471.1473.1475.1477.1479.1481.1483.1485.1487.1489.1491.1493.1495.1497.1499.1501.1503.1505.1507.1509.1511.1513.1515.1517.1519.1521.1523.1525.1527.1529.1531.1533.1535.1537.1539.1541.1543.1545.1547.1549.1551.1553.1555.1557.1559.1561.1563.1565.1567.1569.1571.1573.1575.1577.1579.1581.1583.1585.1587.1589.1591.1593.1595.1597.1599.1601.1603.1605.1607.1609.1611.1613.1615.1617.1619.1621.1623.1625.1627.1629.1631.1633.1635.1637.1639.1641.1643.1645.1647.1649.1651.1653.1655.1657.1659.1661.1663.1665.1667.1669.1671.1673.1675.1677.1679.1681.1683.1685.1687.1689.1691.1693.1695.1697.1699.1701.1703.1705.1707.1709.1711.1713.1715.1717.1719.1721.1723.1725.1727.1729.1731.1733.1735.1737.1739.1741.1743.1745.1747.1749.1751.1753.1755.1757.1759.1761.1763.1765.1767.1769.1771.1773.1775.1777.1779.1781.1783.1785.1787.1789.1791.1793.1795.1797.1799.1801.1803.1805.1807.1809.1811.1813.1815.1817.1819.1821.1823.1825.1827.1829.1831.1833.1835.1837.1839.1841.1843.1845.1847.1849.1851.1853.1855.1857.1859.1861.1863.1865.1867.1869.1871.1873.1875.1877.1879.1881.1883.1885.1887.1889.1891.1893.1895.1897.1899.1901.1903.1905.1907.1909.1911.1913.1915.1917.1919.1921.1923.1925.1927.1929.1931.1933.1935.1937.1939.1941.1943.1945.1947.1949.1951.1953.1955.1957.1959.1961.1963.1965.1967.1969.1971.1973.1975.1977.1979.1981.1983.1985.1987.1989.1991.1993.1995.1997.1999.2001.2003.2005.2007.2009.2011.2013.2015.2017.2019.2021.2023.2025.2027.2029.2031.2033.2035.2037.2039.2041.2043.2045.2047.2049.2051.2053.2055.2057.2059.2061.2063.2065.2067.2069.2071.2073.2075.2077.2079.2081.2083.2085.2087.2089.2091.2093.2095.2097.2099.2101.2103.2105.2107.2109.2111.2113.2115.2117.2119.2121.2123.2125.2127.2129.2131.2133.2135.2137.2139.2141.2143.2145.2147.2149.2151.2153.2155.2157.2159.2161.2163.2165.2167.2169.2171.2173.2175.2177.2179.2181.2183.2185.2187.2189.2191.2193.2195.2197.2199.2201.2203.2205.2207.2209.2211.2213.2215.2217.2219.2221.2223.2225.2227.2229.2231.2233.2235.2237.2239.2241.2243.2245.2247.2249.2251.2253.2255.2257.2259.2261.2263.2265.2267.2269.2271.2273.2275.2277.2279.2281.2283.2285.2287.2289.2291.2293.2295.2297.2299.2301.2303.2305.2307.2309.2311.2313.2315.2317.2319.2321.2323.2325.2327.2329.2331.2333.2335.2337.2339.2341.2343.2345.2347.2349.2351.2353.2355.2357.2359.2361.2363.2365.2367.2369.2371.2373.2375.2377.2379.2381.2383.2385.2387.2389.2391.2393.2395.2397.2399.2401.2403.2405.2407.2409.2411.2413.2415.2417.2419.2421.2423.2425.2427.2429.2431.2433.2435.2437.2439.2441.2443.2445.2447.2449.2451.2453.2455.2457.2459.2461.2463.2465.2467.2469.2471.2473.2475.2477.2479.2481.2483.2485.2487.2489.2491.2493.2495.2497.2499.2501.2503.2505.2507.2509.2511.2513.2515.2517.2519.2521.2523.2525.2527.2529.2531.2533.2535.2537.2539.2541.2543.2545.2547.2549.2551.2553.2555.2557.2559.2561.2563.2565.2567.2569.2571.2573.2575.2577.2579.2581.2583.2585.2587.2589.2591.2593.2595.2597.2599.2601.2603.2605.2607.2609.2611.2613.2615.2617.2619.2621.2623.2625.2627.2629.2631.2633.2635.2637.2639.2641.2643.2645.2647.2649.2651.2653.2655.2657.2659.2661.2663.2665.2667.2669.2671.2673.2675.2677.2679.2681.2683.2685.2687.2689.2691.2693.2695.2697.2699.2701.2703.2705.2707.2709.2711.2713.2715.2717.2719.2721.2723.2725.2727.2729.2731.2733.2735.2737.2739.2741.2743.2745.2747.2749.2751.2753.2755.2757.2759.2761.2763.2765.2767.2769.2771.2773.2775.2777.2779.2781.2783.2785.2787.2789.2791.2793.2795.2797.2799.2801.2803.2805.2807.2809.2811.2813.2815.2817.2819.2821.2823.2825.2827.2829.2831.2833.2835.2837.2839.2841.2843.2845.2847.2849.2851.2853.2855.2857.2859.2861.2863.2865.2867.2869.2871.2873.2875.2877.2879.2881.2883.2885.2887.2889.2891.2893.2895.2897.2899.2901.2903.2905.2907.2909.2911.2913.2915.2917.2919.2921.2923.2925.2927.2929.2931.2933.2935.2937.2939.2941.2943.2945.2947.2949.2951.2953.2955.2957.2959.2961.2963.2965.2967.2969.2971.2973.2975.2977.2979.2981.2983.2985.2987.2989.2991.2993.2995.2997.2999.3001.3003.3005.3007.3009.3011.3013.3015.3017.3019.3021.3023.3025.3027.3029.3031.3033.3035.3037.3039.3041.3043.3045.3047.3049.3051.3053.3055.3057.3059.3061.3063.3065.3067.3069.3071.3073.3075.3077.3079.3081.3083.3085.3087.3089.3091.3093.3095.3097.3099.3101.3103.3105.3107.3109.3111.3113.3115.3117.3119.3121.3123.3125.3127.3129.3131.3133.3135.3137.3139.3141.3143.3145.3147.3149.3151.3153.3155.3157.3159.3161.3163.3165.3167.3169.3171.3173.3175.3177.3179.3181.3183.3185.3187.3189.3191.3193.3195.3197.3199.3201.3203.3205.3207.3209.3211.3213.3215.3217.3219.3221.3223.3225.3227.3229.3231.3233.3235.3237.3239.3241.3243.3245.3247.3249.3251.3253.3255.3257.3259.3261.3263.3265.3267.3269.3271.3273.3275.3277.3279.3281.3283.3285.3287.3289.3291.3293.3295.3297.3299.3301.3303.3305.3307.3309.3311.3313.3315.3317.3319.3321.3323.3325.3327.3329.3331.3333.3335.3337.3339.3341.3343.3345.3347.3349.3351.3353.3355.3357.3359.3361.3363.3365.3367.3369.3371.3373.3375.3377.3379.3381.3383.3385.3387.3389.3391.3393.3395.3397.3399.3401.3403.3405.3407.3409.3411.3413.3415.3417.3419.3421.3423.3425.3427.3429.3431.3433.3435.3437.3439.3441.3443.3445.3447.3449.3451.3453.3455.3457.3459.3461.3463.3465.3467.3469.3471.3473.3475.3477.3479.3481.3483.3485.3487.3489.3491.3493.3495.3497.3499.3501.3503.3505.3507.3509.3511.3513.3515.3517.3519.3521.3523.3525.3527.3529.3531.3533.3535.3537.3539.3541.3543.3545.3547.3549.3551.3553.3555.3557.3559.3561.3563.3565.3567.3569.3571.3573.3575.3577.3579.3581.3583.3585.3587.3589.3591.3593.3595.3597.3599.3601.3603.3605.3607.3609.3611.3613.3615.3617.3619.3621.3623.3625.3627.3629.3631.3633.3635.3637.3639.3641.3643.3645.3647.3649.3651.3653.3655.3657.3659.3661.3663.3665.3667.3669.3671.3673.3675.3677.3679.3681.3683.3685.3687.3689.3691.3693.3695.3697.3699.3701.3703.3705.3707.3709.3711.3713.3715.3717.3719.3721.3723.3725.3727.3729.3731.3733.3735.3737.3739.3741.3743.3745.3747.3749.3751.3753.3755.3757.3759.3761.3763.3765.3767.3769.3771.3773.3775.3777.3779.3781.3783.3785.3787.3789.3791.3793.3795.3797.3799.3801.3803.3805.3807.3809.3811.3813.3815.3817.3819.3821.3823.3825.3827.3829.3831.3833.3835.3837.3839.3841.3843.3845.3847.3849.3851.3853.3855.3857.3859.3861.3863.3865.3867.3869.3871.3873.3875.3877.3879.3881.3883.3885.3887.3889.3891.3893.3895.3897.3899.3901.3903.3905.3907.3909.3911.3913.3915.3917.3919.3921.3923.3925.3927.3929.3931.3933.3935.3937.3939.3941.3943.3945.3947.3949.3951.3953.3955.3957.3959.3961.3963.3965.3967.3969.3971.3973.3975.3977.3979.3981.3983.3985.3987.3989.3991.3993.3995.3997.3999.4001.4003.4005.4007.4009.4011.4013.4015.4017.4019.4021.4023.4025.4027.4029.4031.4033.4035.4037.4039.4041.4043.4045.4047.4049.4051.4053.4055.4057.4059.4061.4063.4065.4067.4069.4071.4073.4075.4077.4079.4081.4083.4085.4087.4089.4091.4093.4095.4097.4099.4101.4103.4105.4107.4109.4111.4113.4115.4117.4119.4121.4123.4125.4127.4129.4131.4133.4135.4137.4139.4141.4143.4145.4147.4149.4151.4153.4155.4157.4159.4161.4163.4165.4167.4169.4171.4173.4175.4177.4179.4181.4183.4185.4187.4189.4191.4193.4195.4197.4199.4201.4203.4205.4207.4209.4211.4213.4215.4217.4219.4221.4223.4225.4227.4229.4231.4233.4235.4237.4239.4241.4243.4245.4247.4249.4251.4253.4255.4257.4259.4261.4263.4265.4267.4269.4271.4273.4275.4277.4279.4281.4283.4285.4287.4289.4291.4293.4295.4297.4299.4301.4303.4305.4307.4309.4311.4313.4315.4317.4319.4321.4323.4325.4327.4329.4331.4333.4335.4337.4339.4341.4343.4345.4347.4349.4351.4353.4355.4357.4359.4361.4363.4365.4367.4369.4371.4373.4375.4377.4379.4381.4383.4385.4387.4389.4391.4393.4395.4397.4399.4401.4403.4405.4407.4409.4411.4413.4415.4417.4419.4421.4423.4425.4427.4429.4431.4433.4435.4437.4439.4441.4443.4445.4447.4449.4451.4453.4455.4457.4459.4461.4463.4465.4467.4469.4471.4473.4475.4477.4479.4481.4483.4485.4487.4489.4491.4493.4495.4497.4499.4501.4503.4505.4507.4509.4511.4513.4515.4517.4519.4521.4523.4525.4527.4529.4531.4533.4535.4537.4539.4541.4543.4545.4547.4549.4551.4553.4555.4557.4559.4561.4563.4565.4567.4569.4571.4573.4575.4577.4579.4581.4583.4585.4587.4589.4591.4593.4595.4597.4599.4601.4603.4605.4607.4609.4611.4613.4615.4617.4619.4621.4623.4625.4627.4629.4631.4633.4635.4637.4639.4641.4643.4645.4647.4649.4651.4653.4655.4657.4659.4661.4663.4665.4667.4669.4671.4673.4675.4677.4679.4681.4683.4685.4687.4689.4691.4693.4695.4697.4699.4701.4703.4705.4707.4709.4711.4713.4715.4717.4719.4721.4723.4725.4727.4729.4731.4733.4735.4737.4739.4741.4743.4745.4747.4749.4751.4753.4755.4757.4759.4761.4763.4765.4767.4769.4771.4773.4775.4777.4779.4781.4783.4785.4787.4789.4791.4793.4795.4797.4799.4801.4803.4805.4807.4809.4811.4813.4815.4817.4819.4821.4823.4825.4827.4829.4831.4833.4835.4837.4839.4841.4843.4845.4847.4849.4851.4853.4855.4857.4859.4861.4863.4865.4867.4869.4871.4873.4875.4877.4879.4881.4883.4885.4887.4889.4891.4893.4895.4897.4899.4901.4903.4905.4907.4909.4911.4913.4915.4917.4919.4921.4923.4925.4927.4929.4931.4933.4935.4937.4939.4941.4943.4945.4947.4949.4951.4953.4955.4957.4959.4961.4963.4965.4967.4969.4971.497</p>
--

```

{
  "index":0,
  "text":"좋다",
  "score":2
},
{
  "index":1,
  "text":"편하다",
  "score":2
},
{
  "index":2,
  "text":"고급스럽다",
  "score":2
},
{
  "index":3,
  "text":"만족하다",
  "score":2
},

```

```

{
  "index":131,
  "text":"비싸다",
  "score":-2
},
{
  "index":132,
  "text":"기대하다",
  "score":-2
},
{
  "index":133,
  "text":"심하다",
  "score":-2
},
{
  "index":134,
  "text":"안 좋다",
  "score":-2
},
{
  "index":135,
  "text":"불편하다",
  "score":-2
},

```

위 2개는 감성사전의 긍정부분과 부정 부분의 일부분이다.

```

사전에 단어가 없는 경우 결과가 None으로 나타납니다!!!
종료하시려면 #을 입력해주세요!!!
-2:매우 부정, -1:부정, 0:중립, 1:긍정, 2:매우 긍정

리뷰를 입력하세요 :
[ '가볍다/Adjective', '좋다/Adjective', './Punctuation', '음선/Noun', '해/Josa', '피우치/Noun', '도/Josa', '있다/Adjective', '일마니/Noun', '좋다/Adjective', '아쉽다/Adjective', './Punctuation', '따로/Adverb', '구매/Noun', '히다/Verb', './Punctuation', '자다/Verb', '사용/Noun', '하디/Verb', './Punctuation' ]
[ '가볍다', '좋다', '있다', '좋다', '아쉽다' ]
이름 : 가볍다
극성 : 2
이름 : 좋다
극성 : 2
이름 : 있다
극성 : 0
이름 : 좋다
극성 : 2
이름 : 아쉽다
극성 : -2
total_score : 4
평균 : 0.8
예측한 별점은 4점입니다.

```

이후 리뷰를 입력하면 해당 리뷰를 토큰화를 진행한 다음 해당 어휘도에서 감성사전에 존재하는 극성 어휘들을 판별하여 계산한 후, 별점을 예측하게 한다.

4. 결론

본 연구에서는 SO-PMI를 활용해 영역별 감성 사전을 구축하고, 평점을 예측하는 연구를 수행하였다. 연구 결과 본 연구에서 구축한 영역별 감성 사전으로 평점을 예측하는 경우, 예측한 평점과 실제 평점이 유사하게 결과가 나왔다.

그러나, 본 연구를 진행하며 느낀 한계점은 다음과 같다.

쿠팡에 존재하는 리뷰의 90% 이상이 4점 또는 5점 리뷰이다보니, 부정적인 리뷰의 절대적인 수가 부족했다. 따라서 리뷰 비율의 분포를 바탕으로 평점을 예측하다 보니 불안정한 점수측정이 가끔씩 결과로 이어졌다. 또한 직접 데이터를 크롤링을 통해 가져오다보니 데이터의 양의 수가 절대적으로 적어 감성사전을 구축하는데 있어 감성사전에 들어가는 단어의 수가 부족하였다. 그리고 한국어의 특성상 예를들자

면 ‘하지 않다’와 ‘못하다’를 토큰화를 진행하면 ‘하다’와 ‘않다’, ‘못’과 ‘하다’ 이런식으로 구분되는데 이때 SO-PMI를 측정하는 부분에 있어서 ‘하다’라는 단어가 부정리뷰데이터에서 높은 점수가 측정되어 이후 별점을 예측할때 문장속 ‘하다’라는 단어가 있으면 이를 부정 단어로 취급을 하는 경우가 발생하였다. 만약 한국어가 아닌 영어 데이터로 감성사전을 구축하였다면 이러한 문제가 발생하지 않았을 것 같다. 이를 개선하기 위해서는 ‘못’이나 ‘않다’를 갖고 있는 단어를 아예 반대말로 바꿔 치환한다면 해당 문제점을 바로 잡을 수 있을 것 같다. 예를들어 ‘행복하다’를 ‘불행하다’ 라는 단어로 치환을 하는 방식이다.

또한 쿠팡 리뷰데이터 등 이커머스 데이터가 아닌 좀 더 객관적인 지표를 얻을 수 있는 음식점과 호텔 등 소비자에게 직접적으로 노출되지 않는 업종의 리뷰에서는 부정적인 리뷰 데이터를 보다 쉽게 얻을 수 있을거라 생각한다.