



'Break The Frame'

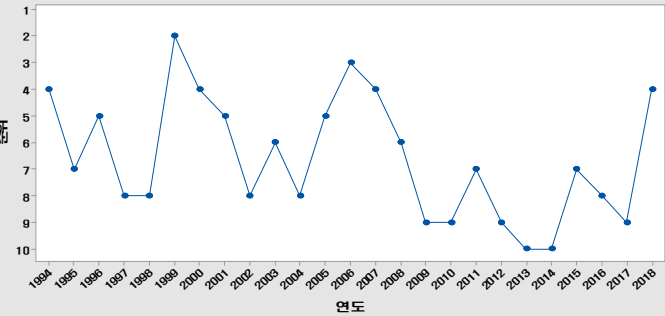
세이버메트릭스로 본 한화 이글스의 비상

호서대학교
응용통계학과
이정호

1. 배경 및 목적

Break The Frame 2018년 한화 이글스의 팀 슬로건입니다. 한화 이글스가 2018년도 프로야구 판을 세차게 흔들고 있습니다. 전년도 8위 팀이 올해는 전반기 3위로 마무리 하였습니다. 어떤 이유에서 높은 승률을 가지게 된 지 확인해보려 합니다. 그리고 어떤 변수들이 승률에 영향을 끼치는지 알아보겠습니다. 변수는 KBO 공식 홈페이지에서 투수와 타자 각각 15개를 설정했습니다. 변수를 확인해보면 방어율과 실점 혹은 타석과 타수 등 변수를 값이 다중 공선성이 존재하게 됩니다. 그러므로 주성분 분석을 시행하여 변수축약을 한 뒤, 회귀분석을 시행하고 승률 예측 모델을 만들어서 승률 예측해보겠습니다.

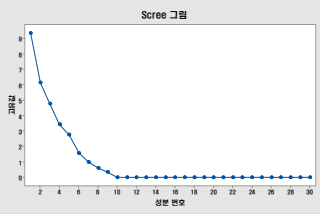
한화 이글스 역대 순위



2. 주성분 분석

Table with 23 columns (C1-C23) and 10 rows (1-10) showing principal component analysis results. Columns include variables like '1. 투수', '2. 타', '3. 타석', etc. Rows show loadings for components 1 through 10.

야구 데이터의 변수들은 측정단위가 서로 다르고, 값의 차이가 크기 때문에 표준화 변수를 사용해야 합니다. 또 주성분 분석을 실시합니다. 주성분 분석의 목적은 가장 적은 수의 주성분을 사용하여 분산의 최대량을 설명하는 것입니다. 변수의 척도가 서로 다르고 모든 변수에 동일한 가중치를 부여함으로써 올바른 결과를 얻기 위해 30개의 변수의 행렬 유형을 표준화된 상관계수로 선택하였습니다. 이 때 표준화된, 각 값에 평균을 빼준 뒤 표준편차로 나누어 주는 것 입니다.



주성분 분석을 통해서 나온 고유치와 누적 비율을 이용해서, 고유치 1 이상이고 누적 비율이 90% 이상의 주성분을 선택하였습니다. Scree 그림을 확인해보면 성분이 6-7으로 넘어갈 때 그래프가 급격히 변화합니다. 그러므로 성분 수를 6개로 선택하여 6개의 주성분을 얻었습니다.

Table showing correlation coefficients between variables and the first six principal components. Columns are labeled '고유 벡터' and rows are labeled with variables like '1. 투수', '2. 타', etc.

주성분을 6개로 보았을 때 첫 번째 주성분에서 득점, 타점, 총루타, OPS가 큰 양의 연관성이 있기 때문에 공격력이라 하였습니다. 두 번째 주성분은 홈런, 삼진이 큰 음의 연관성이 있기 때문에 투수력이라 하였습니다. 세 번째 주성분에서는 타수, 타석, 자책점, 이닝, 실점, 피안타가 큰 음의 연관성이 있기 때문에 체력이라 하였습니다. 네 번째 주성분은 볼넷허용, 보크, 홈런허용, 피안타, 볼넷이 상관이 있기 때문에 제구력이라 하였습니다. 다섯 번째 주성분은 도루실책, 도루가 큰 양의 상관이 있기 때문에 주력이라 하였습니다. 마지막 주성분은 폭투, 삼진(투수), 보크가 높음으로 정선력이라 하였습니다.

3. 변수 설명

Table explaining variables. Columns: 투수 (Pitcher), 설명 (Description), 타자 (Batter), 설명 (Description). Rows include variables like '방어율', '승', '세이브', '홀드', '이닝', '피안타', '피홈런', '볼넷', '피희생타', '삼진', '폭투', '보크', '실점', '자책점', 'WHIP'.

4. 회귀 분석

축약된 변수를 이용해서, 회귀분석을 하여 각 부분에 대한 적합 된 회귀모형을 만들고, y 값을 추정하여 순위를 예측해보겠습니다. 주성분 점수와 각 팀의 승률을 이용하여 회귀분석을 실시하기 위해서 먼저 각 주성분의 값에 표준화 값을 곱한 후 더하여 변수를 값을 얻었습니다. 적합 회귀모형을 실시하였습니다.

Table showing regression coefficients for variables C1 through C9. Columns: C1-C9, 공격력, 투수력, 체력, 제구력, 주력, 정선력, 승률, 적용치. Rows: 1. 투수, 2. 타, 3. 타석, 4. 내전, 5. LG, 6. 삼성, 7. 롯데, 8. KIA, 9. KT, 10. NC.

모형이 얼마나 잘 적합되었는지 확인하기 해보겠습니다. R-제곱 값을 보겠습니다. R-제곱은 회귀분석의 설명력을 보여준다고 볼 수 있습니다. 이 모형은 반응 변동의 약 96%를 설명하고 있습니다. 분산분석 값을 확인해보면 회귀의 F-값에 P-값이 0.05 이내로 유의함을 보여주고 있습니다. 즉 회귀분석의 결과에 유의성이 있다고 말할 수 있습니다.

회귀분석 결과를 확인하겠습니다. 반응모형과 각 양의 간의 연관성을 통계적으로 유의한지를 확인하기 위해 p값을 유의수준과 비교하여 귀무가설을 평가해 보겠습니다. 공격력을 제외하고 모두 값이 유의하지 않음을 볼 수 있습니다. 이는 공격력을 제외한 모든 양이 연관성이 없다는 것을 볼 수 있습니다.

5. 결론

순위에 영향을 미치는 변수 30개를 사용하여 주성분 분석에 의한 회귀분석을 하였습니다. 변수의 측정 단위가 다르기 때문에 변수를 표준화하였습니다. 하지만 유사한 항목이 많아서 가장 문제와 많은 변수에 의한 다중 공선성 문제에 의해서 주성분 분석을 하였습니다. 얻어진 주성분 분석을 통해 축약된 변수와 승률을 이용하여 적합 회귀 모형을 만들고, 값을 추정하여 2018년도 한국 프로야구 정규 시즌 순위를 비교하였습니다. 결과와 같이 한화 이글스의 승률이 예측값과 나와서 3위를 하여 10년 만에 가을 야구에 가서 플레이오프 한국시리즈 진출 나아가서 우승까지 노려봤으면 좋겠습니다. 한화 이글스 파이팅!