





















KBL 2021-2022

플레이오프 승부 예측

22년도 1학기 주제분석

범주형자료분석팀 박지성 박지민 서희나 윤경선 이지윤

















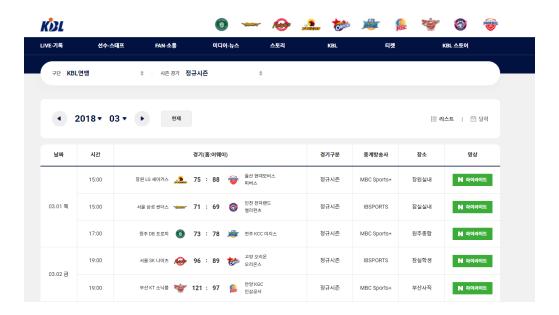






데이터 수집

KBL 홈페이지에서 2017-2022년 모든 경기의 경기/선수 데이터 크롤링





	Date	team	Home/Away	entry	player_num	name	Min	Pts	2PT	3PT	 DR	тот	DK	AST	то	Stl	BS	PF	FO	PP
1 2	2019.11.01	인천 전자랜드 엘리편 츠	Home		20	정영삼	4:25	2	1/1	0/0	 0	0	0	0	0	0	0	1	0	1
2 2	2019.11.01	인천 전자랜드 엘리편 츠	Home		1	민성주	16:34	4	2/3	0/1	 3	4	0	1	0	2	0	4	0	2
3 2	2019.11.01	인천 전자랜드 엘리편 츠	Home	•	30	박찬희	28:36	14	2/5	2/4	 7	7	0	5	3	0	0	3	2	1
4 2	2019.11.01	인천 전자랜드 엘리편 츠	Home		11	홍경기	0:0	0	0/0	0/0	 0	0	0	0	0	0	0	0	0	0
5 2	2019.11.01	인천 전자랜드 엘리편 츠	Home	•	6	자바위	27:30	18	1/2	5/5	 0	2	0	0	1	0	0	1	1	1
8 2	2019.11.30	전주 KCC 이지스	Away	•	43	이대성	33:15	24	0/2	7/16	 1	1	0	3	3	2	0	- 1	5	0
9 2	2019.11.30	전주 KCC 이지스	Away	•	2	송교창	35:58	7	2/6	0/3	 8	9	0	11	1	2	0	4	4	2
10 2	2019.11.30	전주 KCC 이지스	Away		9	최승육	5:30	0	0/0	0/0	 0	0	0	0	0	0	0	0	0	0
11 2	2019.11.30	전주 KCC 이지스	Away		5	유현준	12:15	0	0/0	0/0	 1	1	0	0	3	0	0	2	0	0

후..크롤링..쉽지않아..error..NA..그만..멈춰..























시즌의 독립성 검정



2*2*k 형태의 3차원 분할표의 독립성 검정

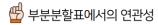
✓ 오즈비의 동질성 검정을 위해 고안된 카이제곱 검정법



즉, K개의 층에 대하여 <mark>동질연관성</mark>이 있는지 확인해보자

4 연관성 측도

오즈비 (Odds Ratio) | 3차원 분할표



동질 연관성(homogeneous association)

조건부 오즈비가 모두 같은 값을 가지는 경우 $(\theta_{XY(1)} = \theta_{XY(2)} = \cdots = \theta_{XY(K)})$

• 대칭적 : XY에 동질 연관성 존재 → YZ와 XZ간에도 동질연관성이 존재

동질연관성에 대한 자세한 내용은 범주 1주차 클린업 참고~























시즌의 독립성 검정



2*2*k 형태의 3차원 분할표의 독립성 검정

✓ 오즈비의 동질성 검정을 위해 고안된 카이제곱 검정법



즉, K개의 층에 대하여 <mark>동질연관성</mark>이 있는지 확인해보자



 H_0 : Homogeneity of odds ratio

각 k개의 층은 동일한 조건부 오즈비를 지닌다 (p-value < 0.05 기각)























시즌의 독립성 검정



파울

2점 파울 수치형 변수 → 범주형으로 변환



중간값보다 크면 1 작으면 0 부여



시즌(k)을 기준으로 <mark>파울</mark>(범주)과 <mark>승패</mark>에 관한 분할표 만들기 Breslow-Day test on Homogeneity of Odds Ratios

data: foul X-squared = 1.1448, df = 3 p-value = 0.7663

 H_0 : Homogeneity of odds ratio vs H_1 : not H_0

(p-value <0.05 기각)



각 시즌별로의 승패에 대한 오즈값이 동질적이다

→ 각 시즌에서 승패에 대한

파울의 영향은 유사!





















파생변수 - 기존 경기 데이터 활용

Q1) 왜 비율로 확인해보고자 하였나요?

단위가 득점인 경우 승/패에 대해 직접적인 지표로써 작용된다고 판단

승패 여부에 영향을 주는지 확인해보자!!

Q2) 왜 4가지의 경우로 더 세분화를 하였나요?

각 독립적인 득점 방법을 하나씩 고려해 봄으로써 결과에

강한 욕인을 죽는 용소가 있는지 파악하고자 함

자유투 득점

페인트 득점

페인트존 외 2점슛 득점

3점슛 득점















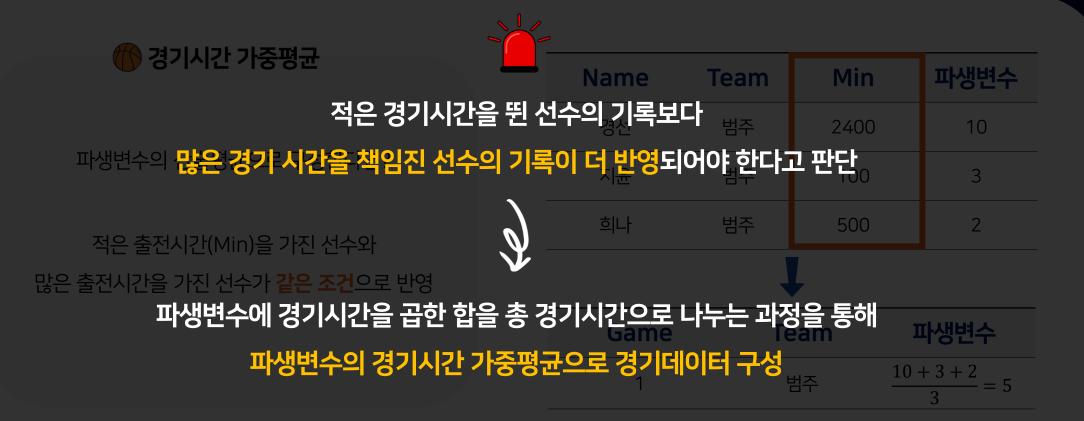








0. 경기시간 가중평균

























모델링 결과 NBA P/O 변수선택

차원 축소 (Dimension Reduction)

Feature의 개수가 증가할수록 차원의 저주 발생, 모델의 예측 신뢰도 감소

차원의 저주에 대한 자세한 내용은 데마팀 1주차 클린업 참고~

Method



변수 선택

분석 목적에 부합하는

소수의 예측 변수만을 선택

Ex) Stepwise Selection





변수 추출

상관관계가 높은 변수끼리 묶어

변환을 통해 새로운 변수 추출

Ex) 주성분 분석, 요인분석























상관관계 – 범주형 vs 범주형



Cramer's V Test

카이제곱 검정 통계량을 이용



Crammer 계수 계산



범주형 변수끼리의 **연관성 측도** 파악

	Home/Away	FB_p	ТО_р	SC_p	BC_p	연장여부
Home/Away	0.999	0.020	0.025	0.032	0.028	0.000
FB_p	0.020	1.000	0.181	0.045	0.074	0.018
TO_p	0.025	0.181	1.000	0.058	0.069	0.079
SC_p	0.032	0.045	0.058	1.000	0.057	0.078
BC_p	0.028	0.074	0.069	0.057	1.000	0.073
연장여부	0.000	0.018	0.079	0.078	0.073	0.996

→ '연장여부'와 'FB_p'에 대한 크래머 계수 값























모델링 NBA P/O 결과 변수선택

상관관계 – 순서형 vs 연속형



☆ 순서형 변수 vs 연속형 변수

순서형 변수를 Ordinal encoding



앞선 두 연속형 변수의 상관관계 파악을 위해 활용한 Spearman 상관분석 진행!

ТО_р		ТО_р
높음		3
	40	
낮음		1
보통		2























상관관계 분석 – 명목형 vs 연속형

Point Biserial correlation coefficient

이분화된 명목형 변수 0과 1로 Encoding



Encoding된 명목형 변수와 수치형 변수의 Pearson 상관계수로 상관관계 파악























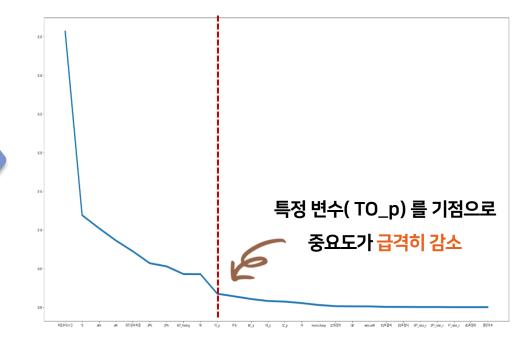
변수 선택 기준 (1) – Select K Best



모든 변수의 Feature Importance 파악

Select K best에서 K=27개로 설정하여 적합한 후 score점수 비교

	Feature	Score		
	득점우위시간	694.61		
	TS	231.86		
9개 -	eFG	199.00		
- "				
	TR	83.27		
	ТО_р	33.89		
		_		























변수 선택 기준 (2) - KS 검정

Kolmogorov Smirnov 검정

주어진 두 표본의 분포가 일치하는지의 여부를 검정할 때 이용

 H_0 : 비교하는 두 분포가 동질적이다

P - value <0.05 → 비교하는 분포 간 이질성 존재



만약 요인 별로 승/패에 대한 두 분포가 **이질적임**을 확인

→ 승/패를 결정짓는 중요한 요인으로 작용함을 알 수 있음























변수 선택

부제 : 변수듀스 101



두 가지의 기준에 대해 <mark>공통적</mark>으로 중요하다고 판단된 변수 최종적으로 선택



② KS test , 카이제곱 검정을 통한 승/패에 따른 이질성 확인



선택된 변수

득점우위시간, TS, eFG, Effi, 최다연속득점, 3P%, 2P%, AST_Rating, TR













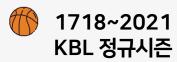




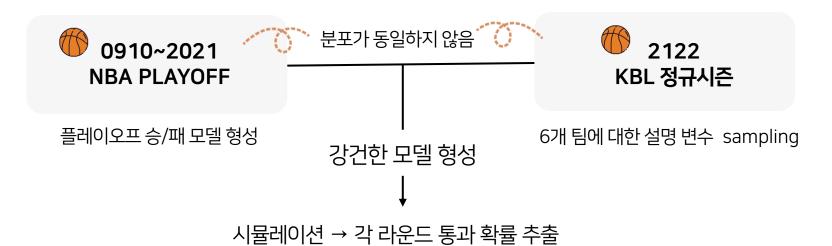




예측 흐름



경기 승리에 영향을 주는 7개 변수 선택→ 고정!























샘플링 (Sampling)

① 2122 KBL 팀 별 7개 설명변수에 대한 샘플링 진행

Sampling이란?

무작위로 표본을 추출하는 방법



진행방식 (목표: 경기 데이터(X)값 생성)

① 후보 분포 정의

'norm','t', 'f', 'chi', 'cosine', 'alpha', 'beta', 'gamma', 'dgamma', 'dweibull', 'maxwell', 'pareto', 'fisk'

② Fitter 함수를 사용하여 각 설명변수 별로 SSE 최소화하는 후보 분포 찾기

③ 각 설명 변수에 적합 된 분포에서 필요한 경기 수 만큼 random하게 하나의 값 선택















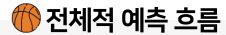








예측 흐름



시뮬레이터를 이용

대진표를 토대로 N번의 시뮬레이터 결과로 플레이오프 결과 예측



1번의 시뮬레이션마다 주어진 대진표에 따라 6강부터 4강, 결승까지 진출팀 예측























예측 흐름 - 시뮬레이션



팀/라운드 별 전체 승리 대비 팀 승리 지표



각 **팀 별 라운드 통과 확률** 계산 가능!

	Team	4강진출진출	결승진출확률	우승 확률
0	안양KGC	0.675	0.418	0.310
1	울산현대모비스	0.653	0.124	0.023
2	서울SK	1.000	0.816	0.328
3	대구한국가스공사	0.325	0.105	0.048
4	고양오리온	0.347	0.060	0.007
5	수원KT	1.000	0.477	0.284























모델링 결과 NBA P/O 변수선택

예측 흐름 - 평가지표



6 평가지표

실제 결승에 진출한 두 팀의 4강 진출 시 결승 진출 조건부 확률의 조화 평균



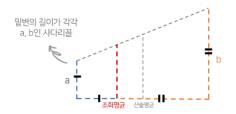
한 쪽 팀의 결승 진출 조건부 확률만 잘 예측하는 모델 패널티 부여

P(결승 진출 | 4강 진출)

1 혼동행렬

⑥ F1-Score | 조화평균

불균형한 데이터가 주어졌을 때도 보다 정확한 성능 파악



조화평균을 기하학적으로 접근해보면.

밑변의 길이와 동일한 거리에 떨어진 지점에서 빗변으로의 높이가 곧 조화평균!

조화 평균의 자세한 내용은 범주팀 3주차 클린업 참고!





















SVC w. linear kernel

SVC with linear kernel

하이퍼 파라미터 튜닝

데이터셋의 구분이 선형적인 관계를C값을 강낮게 튜닝하여 이상치를

일부 허용한다면을 통해 Auxiliary Dataset에 대한 문제를 해결 기대 최적의 결정 경계로 선형적인 경기를 생성 (Why? 선형적인 SYM모델)





C: 경계를 결정할 때 이상치를 허용하는 정도

C가 낮은 값을 가진 정확도 간 차이가 크지 않아 과적합을 이상치를 많이 허용 과적합 방지 방지하기 위해 C=0.3으로 채택























SVC(Support Vector Classifier) – 결과 해석



SVC Linear (C: 0.3)

2122KBL 정규리그 1,2위 서울, 수원 2122KBL PO 결승 진출팀 서울, 안양



전반적으로 해당 세 팀의 승리지분 높음 실제 경기와 유사하게 잘 예측



