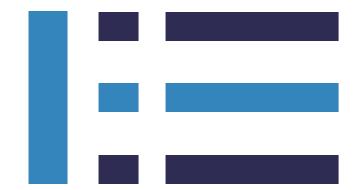




EPL 경기 예측 SOCIALMEDIA ANALYTICS PROJECT

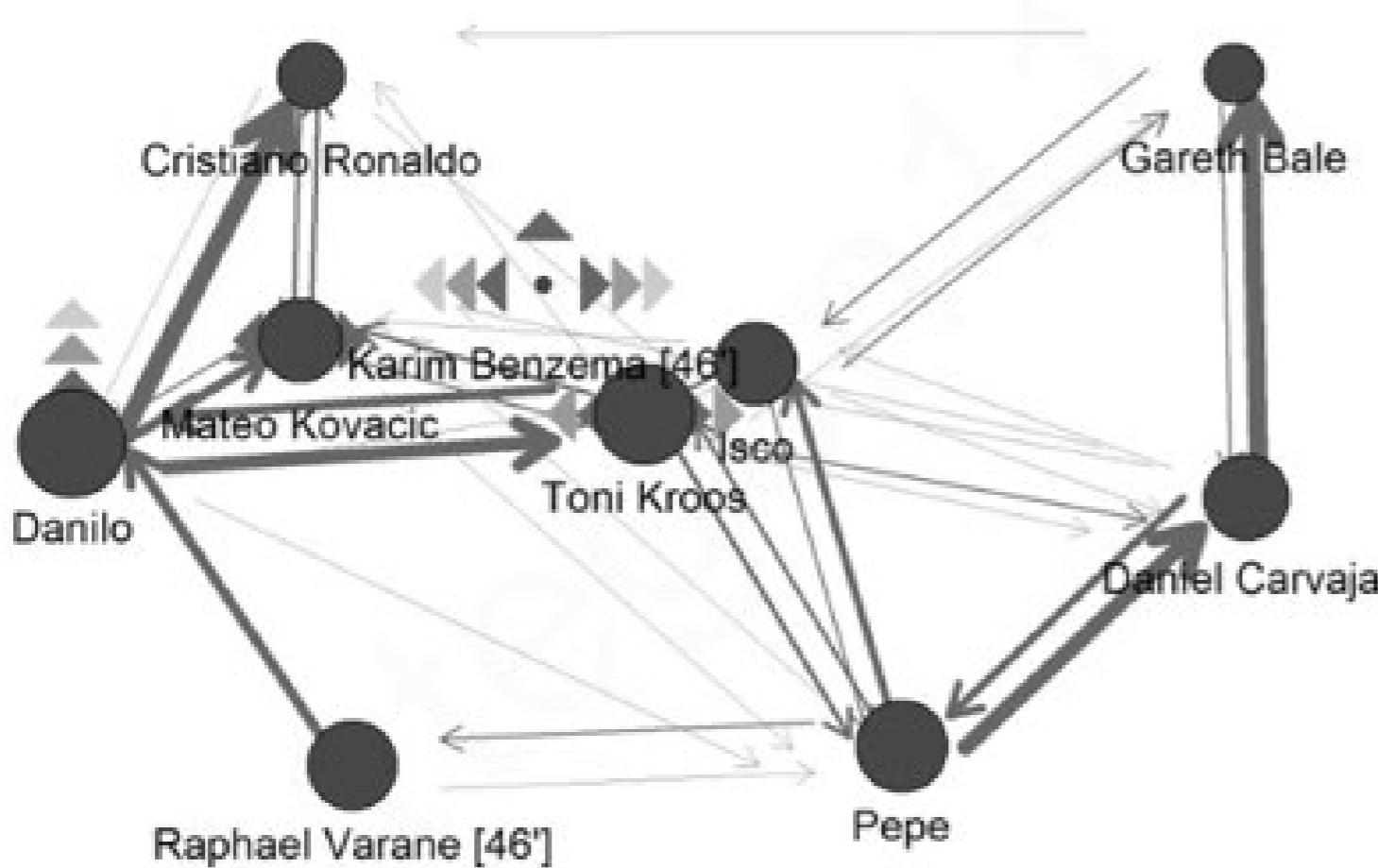
윤기태, 조일혜 황준선

PRESENTATION CONTENTS



- Topic
- Motivation
- Model Design
- Data
- Modeling & Result
- Conclusion

PROJECT TOPIC



[그림] 선수 개인의 성향을 데이터로 분석한 전술판

프로스포츠와 데이터의 융합

- 최근의 **프로스포츠는 데이터 사이언스 분야**와 융합하여 가파른 성장세를 보이고 있다.
- 프로스포츠 관련 데이터 활용은 선수의 훈련, 선수에 대한 통계, 경기에 대한 통계, 승패 예측 등이 있다.
- 때문에 프로스포츠와 관련된 데이터가 범람하고 있는데 이를 **승패에 대한 예측**에 이용해보기로 하였다.

PROJECT TOPIC

Home > 뉴스

슈퍼컴퓨터, 맨유 우승 예측.. 2위는 맨시티

기사입력 : 2016.08.11 기사보내기 : [트위터](#) [페이스북](#)

슈퍼컴퓨터를 이용하여
프로 스포츠의 **승패를 예측**하는 시대



PROJECT TOPIC



VS



PROJECT TOPIC

An aerial photograph of a large, modern stadium filled with spectators. The stadium is brightly lit from within, and the surrounding area is dark. The text "Who Win???" is overlaid in the center of the image.

Who Win???

MOTIVATION



순수한 궁금증

- | 0.02%의 희박한 확률로 우승
- | 예상치 못한 변수가 등장하는 프로스포츠에서 얼마나 정확하게 승패를 예측할 수 있을까?란 호기심을 자극함

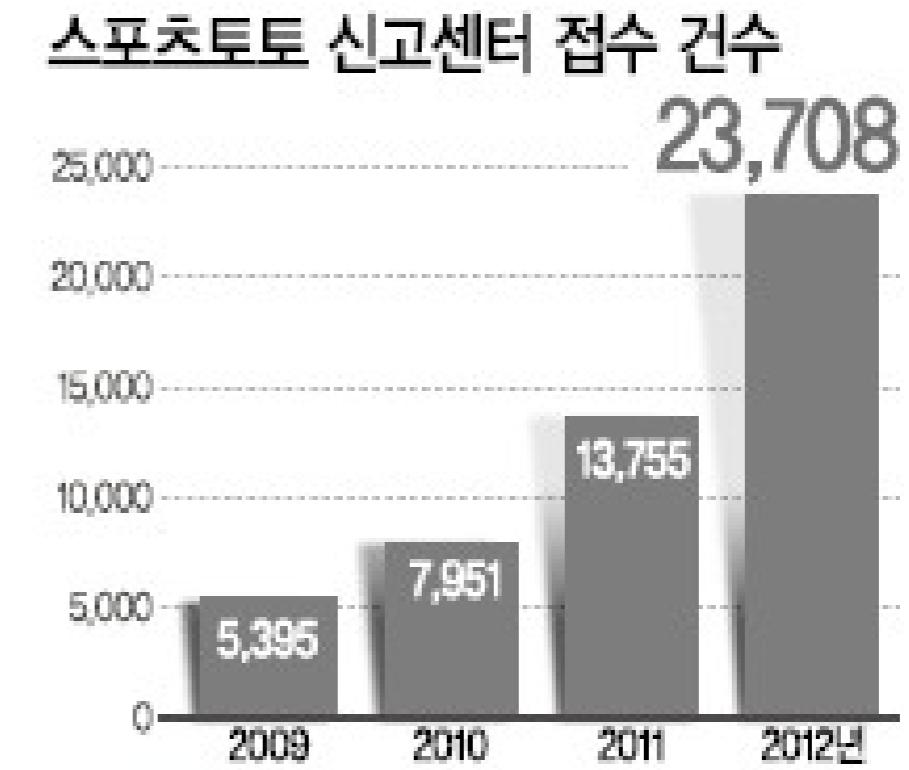


MOTIVATION



불법 도박 근절

최근 스포츠 승패 예측과 관련된 **불법 도박 문제**가 대두되고 있는데, 프로스포츠에 대한 **승패 예측 모델**이 정확해진다면 이 문제가 근절될 것이다.



사설 스포츠토토 운영 사이트 추정치(2012년)

구분	개수
대규모(가입자 2000~3000명)	10~15개
중규모(가입자 1000여명)	40~50개
소규모(가입자 500여명)	200~300개

MODEL DESIGN



선수 데이터 중,
예측에 영향을 미치는
중요하지 않은 변수로는
무엇이 있을까?



MODEL DESIGN

Name	Season	Apps	Mins	Goals	Assists	SpG	KeyP	Drb	Fouled	Off	Disp	UnsTch	Rating	Team.x	Meta	Team.y	Cost
Andy King	2014	16(8)	1374	2	0	0.6	0.2	0.4	0.6	0	0.4	0.7	6.73	Leicester	AM(C)	leicester-city	1700000
Andy King	2015	9(16)	1050	2	1	0.8	0.4	0.3	0.3	0	0.2	0.5	6.56	Leicester	AM(C)	leicester-city	2130000
Andy King	2016	6(4)	624	0	0	0.6	0.4	0.3	0.9	0.1	0.6	0.7	6.53	Leicester	AM(C)	leicester-city	2980000
Aruna Dindan	2009	18(1)	1479	8	2	2.5	1	1.3	3.1	0.4	2.2	0	7.04	Portsmouth	AM(C)	NA	NA
David Dunn	2009	20(3)	1720	9	2	1.7	1.6	1.2	1.7	0.5	1.9	0	7.11	Blackburn	AM(C)	NA	NA
David Dunn	2010	17(10)	1389	2	1	0.9	0.6	0.7	0.6	0	1.1	1	6.5	Blackburn	AM(C)	NA	NA
David Dunn	2011	21(5)	1613	2	1	1.5	0.7	0.6	1.1	0.1	0.8	0.9	6.63	Blackburn	AM(C)	NA	NA
Geovanni	2009	16(10)	1539	3	0	1.8	0.7	0.8	1.9	0.2	1.5	0	6.63	Hull	AM(C)	hull-city	3830000
Giorgos Kara	2012	20(5)	1556	1	0	0.8	0.6	0.7	2.2	0	1	0.3	6.69	Fulham	AM(C)	NA	NA
Jonathan Greening	2009	15(8)	1502	1	1	0.3	0.5	0.2	0.6	0	0.7	0	6.68	Fulham	AM(C)	west-bromwich	3490000
Ricardo Gardner	2009	11(10)	1107	1	1	0.7	0.7	0.4	0.8	0.2	1.2	0	6.61	Bolton	AM(C)	NA	NA
Robert Koren	2013	10(12)	1136	2	0	1.2	0.5	0.2	0.4	0.1	0.6	0.8	6.39	Hull	AM(C)	hull-city	1280000
Benik Afobe	2016	0(10)	176	0	1	0.9	0.8	0.1	0.2	0.6	0.8	0.3	6.21	Bournemouth	AM(C),FW	afc-bournemouth	7650000
Danny Ings	2014	35	3035	11	4	2.8	1	1.6	1.4	0.5	2.1	2.5	6.89	Burnley	AM(C),FW	burnley-fc	2550000
Dimitar Berbatov	2009	24(10)	2000	10	5	0	1.0	1.0	1.1	0.7	0.0	0	7.41	Leeds United	AM(C),FW	leeds-united	201517



→ 축구에 대한 직관적 판단을 통해, 과적합을 발생시키거나
의미 없는 계산을 하게 할 변수들을 미리 제거하였음.

MODEL DESIGN



구단과 감독 각각의 승점으로 승부 예측변수

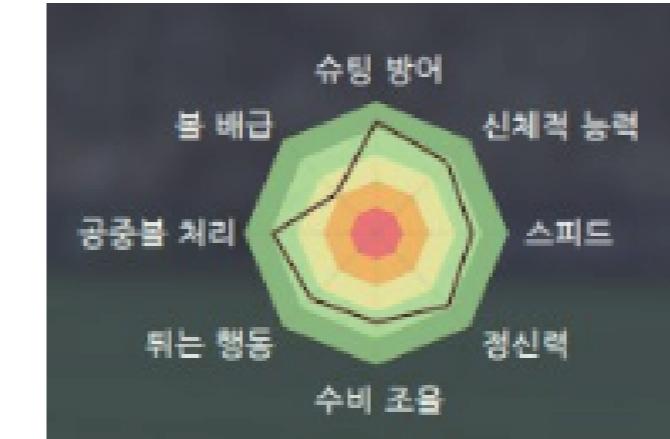


선수 데이터 변수를 어떤 방식으로 적용 할 것인가

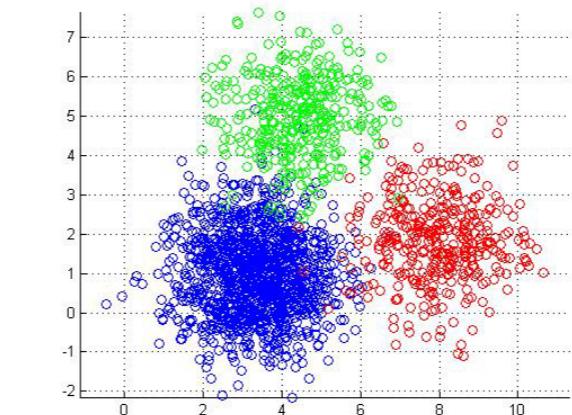
MODEL DESIGN



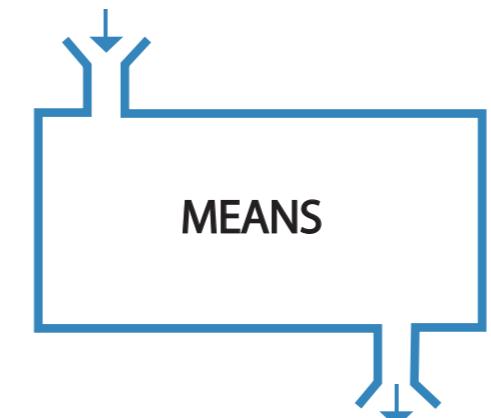
~~선수 데이터를 점수화~~



~~클러스터링 이용 선수 그룹핑~~



선수 스탯의 평균을 이용



DATA

This screenshot shows a football transfer website interface. At the top, there are tabs for NEWS, TRANSFERS & RUMOURS, MARKET VALUES, COMPETITIONS, FORUMS, MY TM, and LIVE. Below the tabs, a search bar and a user profile icon are visible. The main content area displays a table of player statistics. The columns include Player, Name, Age, Nationality (Nat.), Club(s), Market value, App., and various performance metrics like Goals, Assists, and Yellow cards. Each row contains a small thumbnail of the player's face and their club's logo.

This screenshot shows a Microsoft Excel spreadsheet titled "model_data.csv". The data consists of two main sections: "match_flag" and "player". The "match_flag" section contains 20 rows of data, each representing a match result (win, draw, or loss) at home or away. The columns include "match_flag", "home_away", "team", "match_accumulate", "coach", "match_accumulate_defense", "defense", "defensive", "defensive_defense", "offensive", "offensive_defense", "offensive_defensive", and "offensive_defensive_defense". The "player" section contains 20 rows of data, each representing a player's performance across various metrics. The columns include "Player", "Name", "Age", "Nat.", "Club(s)", "Market value", "App.", and numerous performance metrics such as Goals, Assists, and Yellow cards.

This screenshot shows a continuation of the Microsoft Excel spreadsheet from the previous image. It contains 20 more rows of data, labeled 11 through 30, corresponding to the "match_flag" and "player" sections. The data follows the same structure as the first part, providing detailed statistics for each match and player.

축구 통계 사이트에서
데이터 크롤링

데이터를 합쳐서
Input Data 생성

결측값 제거등의
데이터 전처리 과정

DATA



2009년부터의 EPL 데이터 수집

축구 통계 사이트에서
데이터 크롤링

구단 데이터 - 모든 경기의 승/무/패
해당 경기 선수 라인업, 경기 감독 정보

감독 데이터 - 모든 경기의 승/무/패
해당 시즌과 소속팀

선수 데이터 - 경기 관련 Stat 데이터
이적료와 소속팀, 평점과 해당 시즌

DATA

감독 데이터 - 모든 경기의 승/무/패
해당 시즌과 소속팀

구단 데이터 - 모든 경기의 승/무/패

해당 경기 선수 라인업, 경기 감독 정보



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	match	home_away	team	match_ahead	coach	match_ahead_coach	defensive_defense	defensive_defense	defensive_defense	defensive_defense	offensive_defense								
2	lose	home	FC Seoul	-0.9999	-1.13648	-0.91342	-1.18437	-0.94499	0.09032	0.443299	0.412697	0.223578	0.151358	0.299599	0.232399	-0.7784	-0.84443	-0.59	
3	win	home	FC Seoul	-1.0999	-1.13648	-0.91342	-1.18437	-0.27783	0.028916	1.143835	1.216785	-0.272789	0.386047	0.266599	0.620219	-0.89443	-0.18222	-1.16509	-0.45
4	win	away	FC Seoul	-0.70719	-0.92827	-1.70009	-1.0874	-0.32584	-0.46481	0.766688	0.069166	-0.60003	0.1954	0.003121	-0.4841	-0.41424	-0.01512	-1.14122	-0.476
5	draw	home	FC Seoul	-0.70719	-0.92827	-1.70009	-1.0874	-0.32584	-0.46481	0.766688	0.069166	-0.60003	0.1954	0.003121	-0.4841	-0.41424	-0.01512	-1.14122	-0.476
6	lose	home	FC Seoul	-0.26541	-0.87027	-0.85669	-0.6381	-0.59509	-0.40557	0.71269	0.272789	-0.46481	0.620219	0.386047	-0.272789	-0.45599	-0.17844	-1.16509	-0.583
7	draw	home	FC Seoul	-0.65811	-0.54657	-0.69409	-0.44526	-0.48166	-0.48777	0.419507	0.285847	-0.67126	0.792427	0.277216	-0.52368	-0.10004	-1.03628	-0.778	
8	lose	away	FC Seoul	-0.65811	-0.54657	-0.69409	-0.44526	-0.48166	-0.48777	0.419507	0.285847	-0.67126	0.792427	0.277216	-0.52368	-0.10004	-1.03628	-0.778	
9	lose	away	FC Seoul	-0.70719	-0.92827	-1.70009	-1.0874	-0.32584	-0.46481	0.766688	0.069166	-0.60003	0.1954	0.003121	-0.4841	-0.41424	-0.01512	-1.14122	-0.476
10	draw	away	FC Seoul	-0.70719	-0.92827	-1.70009	-1.0874	-0.32584	-0.46481	0.766688	0.069166	-0.60003	0.1954	0.003121	-0.4841	-0.41424	-0.01512	-1.14122	-0.476
11	lose	away	FC Seoul	-1.14898	-0.85887	-0.6794	-0.63417	-1.02458	-1.0804	-0.46027	1.07804	-0.359152	-0.36058	0.033243	-0.72391	-0.37185	-0.50027	-0.71785	-1.38
12	lose	home	FC Seoul	-1.14898	-0.85887	-0.6794	-0.63417	-1.02458	-1.0804	-0.46027	1.07804	-0.359152	-0.36058	0.033243	-0.72391	-0.37185	-0.50027	-0.71785	-1.38
13	win	away	FC Seoul	-1.10934	-0.95452	1.03214	0.848804	-1.18557	-1.31815	0.33919	-0.9392	-0.60044	-0.78773	-0.64854	0.144428	1.235673	0.615062	1.303583	0.287
14	win	away	FC Seoul	-1.0765424	0.954519	0.43262	1.043186	-1.18557	-1.31815	0.33919	-0.9392	-0.60044	-0.78773	-0.64854	0.144428	1.235673	0.615062	1.303583	0.287
15	win	away	FC Seoul	-1.0765424	0.954519	0.43262	1.043186	-1.18557	-1.31815	0.33919	-0.9392	-0.60044	-0.78773	-0.64854	0.144428	1.235673	0.615062	1.303583	0.287
16	draw	home	FC Seoul	-0.765704	0.954519	0.43262	1.043186	-1.18557	-1.31815	0.33919	-0.9392	-0.60044	-0.78773	-0.64854	0.144428	1.235673	0.615062	1.303583	0.287
17	win	home	FC Seoul	-0.7229	0.14779	0.296110	0.343229	-0.3189	0.047459	1.474224	-0.57771	0.246349	-0.87742	-0.20152	0.339666	-0.42427	2.233864	-0.40446	1.2902
18	win	home	FC Seoul	0.12729	0.164779	0.30471	0.352844	-0.38078	0.031072	1.245689	-0.33201	0.278834	-0.5862	0.245242	1.073266	2.259804	-0.41756	1.4596	
19	win	home	FC Seoul	0.12729	0.164779	0.30471	0.352844	-0.38078	0.031072	1.245689	-0.33201	0.278834	-0.5862	0.245242	1.073266	2.259804	-0.41756	1.4596	

데이터를 합쳐서
Input Data 생성

취합한 데이터를 Model에 적용하기
위해 Input Data 형태로 정리

선수 데이터 - 경기 관련 Stat 데이터
이적료와 소속팀, 평점과 해당 시즌

DATA



전처리 과정

filter를 통해 데이터를 합치는 과정에서 중복되거나
에러가 나는 데이터를 정리, **분포조정**을 진행.

어쩔 수 없는 결측값에 대한 처리 : 일단 모델에 적용해 본 뒤, 결측값을 제거한 것과 무시한 것 등의 결과를 비교하여 최종적으로 처리함.

	A1	match_flag	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	match_flag	home_away_team_mean_ecumolatcoach	match_ecumolat_defensive	defensive_defensive	defensive_defensive	defensive_defensive	defensive_defensive	defensive_defensive	offensive_defensive											
2	lose	away	-1.09899	-0.91242	-1.18437	-0.65499	0.03022	0.74529	0.855478	0.023571	-0.056	0.129399	-0.32895	-0.7894	-0.57783	-0.84843	-0.59			
3	win	home	-1.09899	-1.18464	-0.91242	-1.18437	-0.27883	0.28916	0.144823	0.121678	-0.05	-0.388047	0.126982	-0.62219	-0.68843	-0.16232	-1.16509	-0.45		
4	win	away	-0.70719	-0.92827	-1.70001	-1.0874	-0.32354	-0.43111	0.766885	0.869169	-0.05	0.19354	0.803223	-0.4541	-0.41424	-0.02152	-0.14122	-0.476		
5	lose	home	-0.26541	-0.70272	-0.17559	-0.20172	-0.35669	-0.3613	0.359055	0.359055	-0.05	-0.71216	-0.109851	0.104079	-0.72345	-0.40559	-0.17656	-0.558		
6	lose	home	-0.65811	-0.54657	-0.69409	-1.4453	-0.1616	-0.4877	0.47	-0.268261	-0.16177	-0.0746	0.792427	-0.77212	-0.52989	-0.10004	-0.13632	-0.45		
7	win	away	-0.65811	-0.54657	-0.69409	-0.67032	-0.45637	-0.5317	0.267347	-0.5612	-0.08157	0.482519	-0.58775	-0.45066	-0.48732	-0.9243	-0.92			
8	lose	away	-0.70719	-0.54657	-0.1294	-0.92259	-0.76145	-0.375	-0.708089	0.115813	-0.17622	-0.19645	0.308868	-0.64764	-0.39958	-0.25345	-0.88811	-0.116		
9	lose	away	-1.14988	-0.85887	0.101951	0.153857	-0.2956	-0.58122	1.028991	-0.06191	0.476271	-0.42569	-0.04543	-0.79653	-0.25847	-0.57111	-0.63264	-0.476		
10	win	home	1.10934	0.45524	0.16207	0.247	-0.28168	-0.4844	-0.33156	-0.72861	-0.09469	-0.69572	-0.175597	0.393233	0.145697	0.856812	0.374676	-0.535		
11	win	away	1.10934	0.45524	0.16207	0.247	-0.28168	-0.4844	-0.33156	-0.72861	-0.09469	-0.69572	-0.175597	0.393233	0.145697	0.856812	0.374676	-0.535		
12	lose	home	0.1279	0.164778	0.236161	0.342823	-0.06669	0.774131	0.708493	-0.09344	0.105944	-0.17576	0.44521	-0.4912	-0.72153	1.785	-0.52535	0.973		
13	win	home	0.1279	0.164778	0.236161	0.342823	-0.06669	0.774131	0.708493	-0.09344	0.105944	-0.17576	0.44521	-0.4912	-0.72153	1.785	-0.52535	0.973		
14	win	away	0.1279	0.164778	0.236161	0.342823	-0.06669	0.774131	0.708493	-0.09344	0.105944	-0.17576	0.44521	-0.4912	-0.72153	1.785	-0.52535	0.973		
15	win	away	0.765424	0.541994	0.16207	0.041316	-1.42608	-0.17189	-0.17198	-0.17198	-0.14134	-0.41434	-0.80322	-0.144366	-0.138695	-0.91796	0.10464	-0.0124		
16	draw	home	0.765424	0.541994	0.16207	0.041316	-1.42608	-0.17189	-0.17198	-0.17198	-0.14134	-0.41434	-0.80322	-0.144366	-0.138695	-0.91796	0.10464	-0.0124		
17	win	away	0.1279	0.164778	0.236161	0.342823	-0.06669	0.774131	0.708493	-0.09344	0.105944	-0.17576	0.44521	-0.4912	-0.72153	1.785	-0.52535	0.973		
18	win	home	0.1279	0.164778	0.236161	0.342823	-0.06669	0.774131	0.708493	-0.09344	0.105944	-0.17576	0.44521	-0.4912	-0.72153	1.785	-0.52535	0.973		
19	win	away	0.1279	0.164778	0.30471	0.236162	-0.06870	0.031076	0.245589	-0.33307	0.287984	-0.3662	0.42542	0.173266	0.259084	0.254849	-0.47556	0.1456		
20	win	home	0.1279	0.164778	0.30471	0.236162	-0.06870	0.031076	0.245589	-0.33307	0.287984	-0.3662	0.42542	0.173266	0.259084	0.254849	-0.47556	0.1456		

결측값 제거등의 데이터 전처리 과정

RESULT



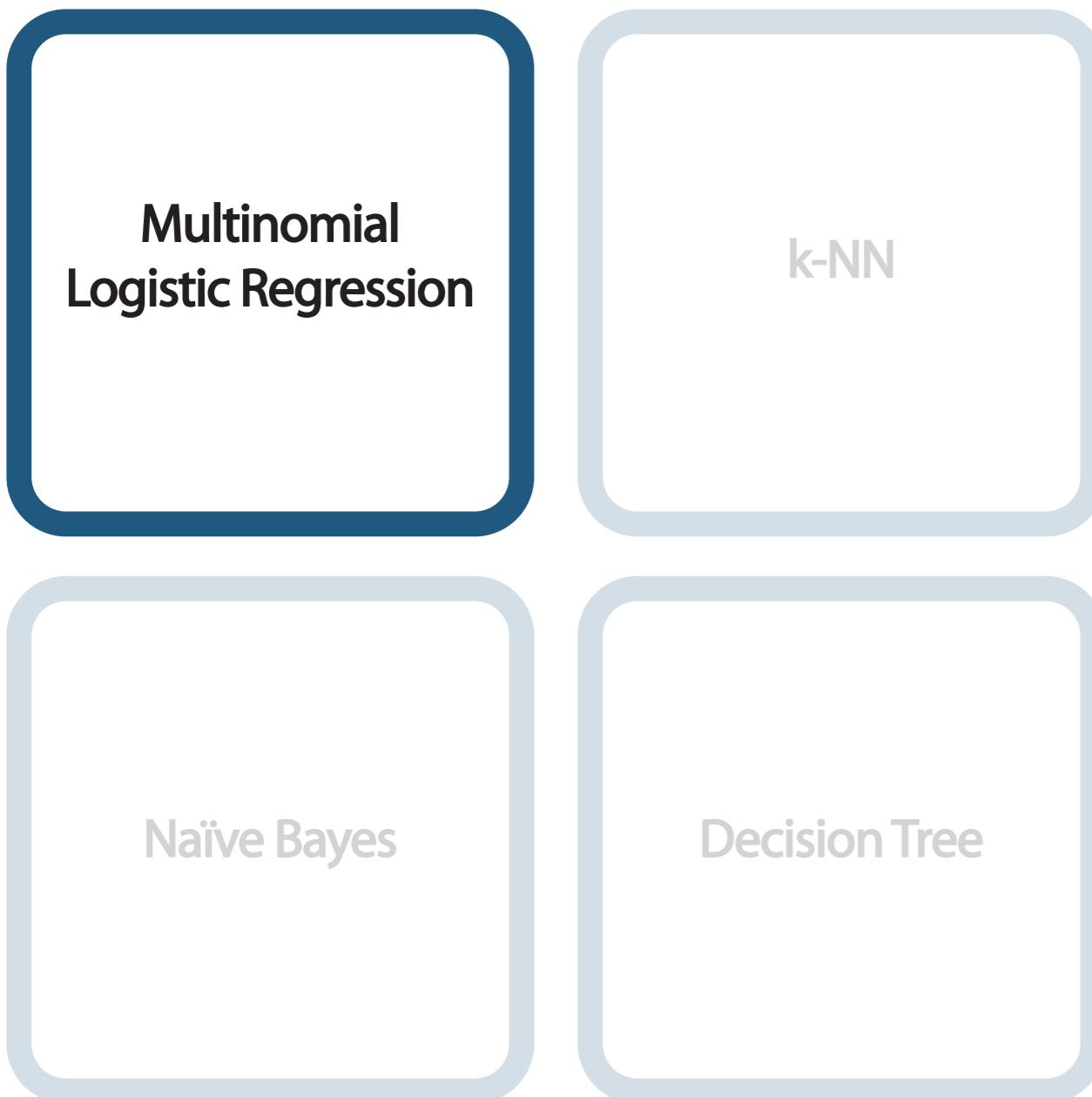
Multinomial
Logistic Regression

k-NN

Naïve Bayes

Decision Tree

RESULT



Logistic Regression

전역 탐색을 사용하려고 하였으나, 변수가 너무 많아서 노트북으로 감당하기 힘들어, **step 함수**를 이용하여 단계 선택 방식으로 **변수를 설정**하였음.

선택된 변수에 대해서 공선성 검증을 한 결과, 모두 10 이하로 나타났음.

RESULT

```
test_predicted_EPL_Im = predict(EPL_Im, type="class", newdata=test_data)
1-sum(test_data$match_flag != test_predicted_EPL_Im)/nrow(test_data)
```

```
## [1] 0.4900662
```

```
table(test_data$match_flag, test_predicted_EPL_Im)
```

```
##      test_predicted_EPL_Im
##      draw  lose  win
##  draw    8    7   26
##  lose    7   21   16
##  win   13    8   45
```

RESULT



Multinomial
Logistic Regression

k-NN

Naïve Bayes

Decision Tree



Naïve Bayes

rfe 함수를 이용하여 변수를 설정하였음.
laplace 조정을 시도한 결과, 아무런 의미가 없었음.

RESULT

```
EPL_nb = naiveBayes(formula = match_flag ~ player_rating+team_match_point+accumulated_team_match_point+offensive_goals+
    accumulated_coach_match_point+player_cost+passing_keyp+passing_assists+offensive_assists+offensive_spg+
    offensive_keyp+coach_match_point+passing_thrb+home_away+offensive_drb,
    data=train_data, laplace = 1)
test_predicted_EPL_nb = predict(EPL_nb, newdata = test_data)
1-sum(test_data$match_flag != test_predicted_EPL_nb)/nrow(test_data)
```

```
## [1] 0.4900662
```

```
table(test_data$match_flag, test_predicted_EPL_nb)
```

```
##      test_predicted_EPL_nb
##      draw  lose  win
##  draw   8   12   21
##  lose   6   26   12
##  win   12   14   40
```

RESULT

Multinomial
Logistic Regression

Naïve Bayes

k-NN

Decision Tree

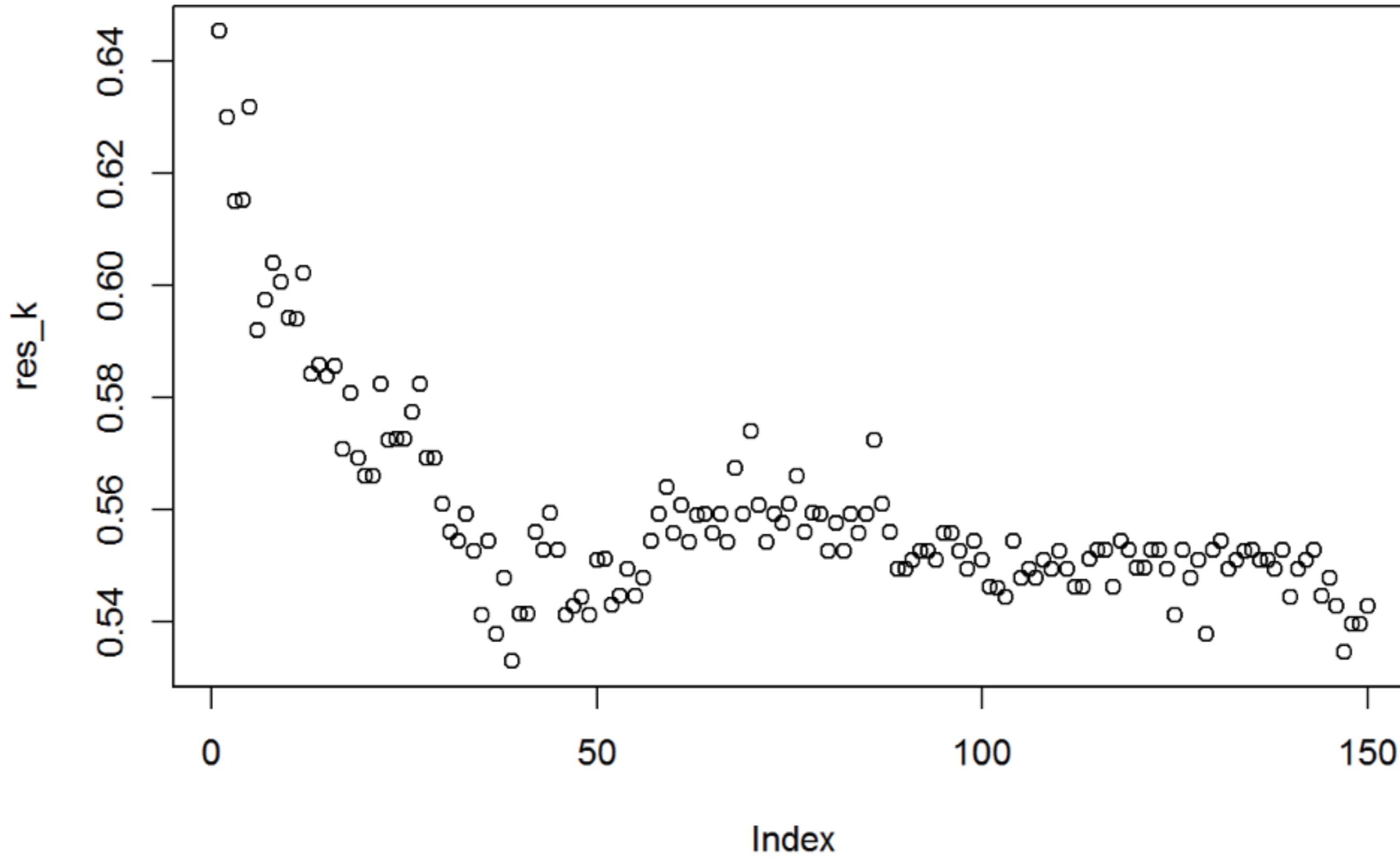


k-NN

?forward.search() 함수를 이용하여 변수를 설정하였음.

k값은 10-fold Cross Validation을 이용하여 Overfitting을 해결하였음.

RESULT



RESULT

```
fit_k = which(res_k == min(res_k))

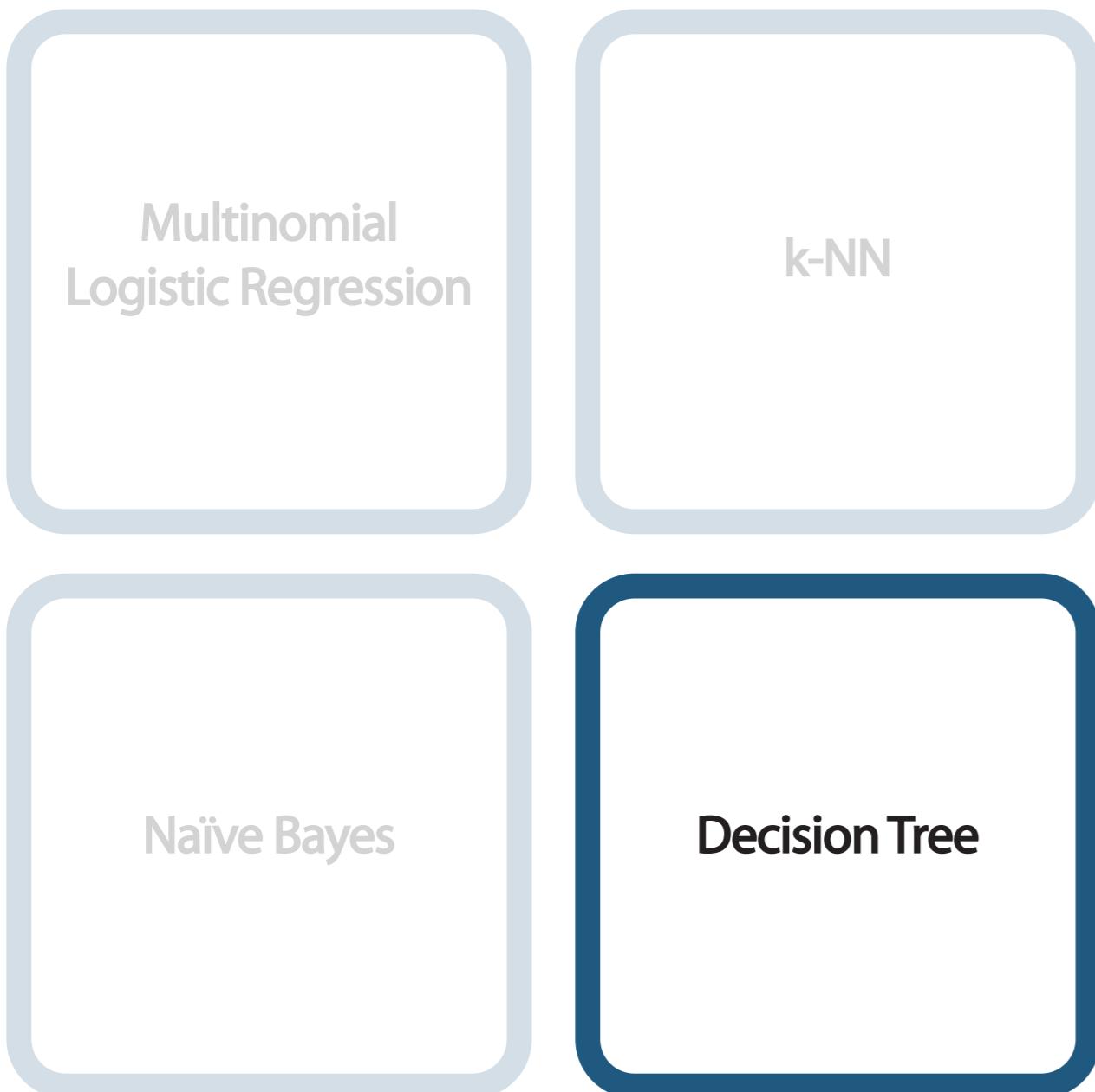
test_predicted_EPL_knn = knn(train = train_data, test = test_data, cl = train_data_cl, fit_k)
1-mean(test_data_cl != test_predicted_EPL_knn)
```

```
## [1] 0.5165563
```

```
table(test_data_cl, test_predicted_EPL_knn)
```

```
##          test_predicted_EPL_knn
## test_data_cl draw lose win
##          draw   6   9  26
##          lose   4  19  21
##          win    6   7  53
```

RESULT

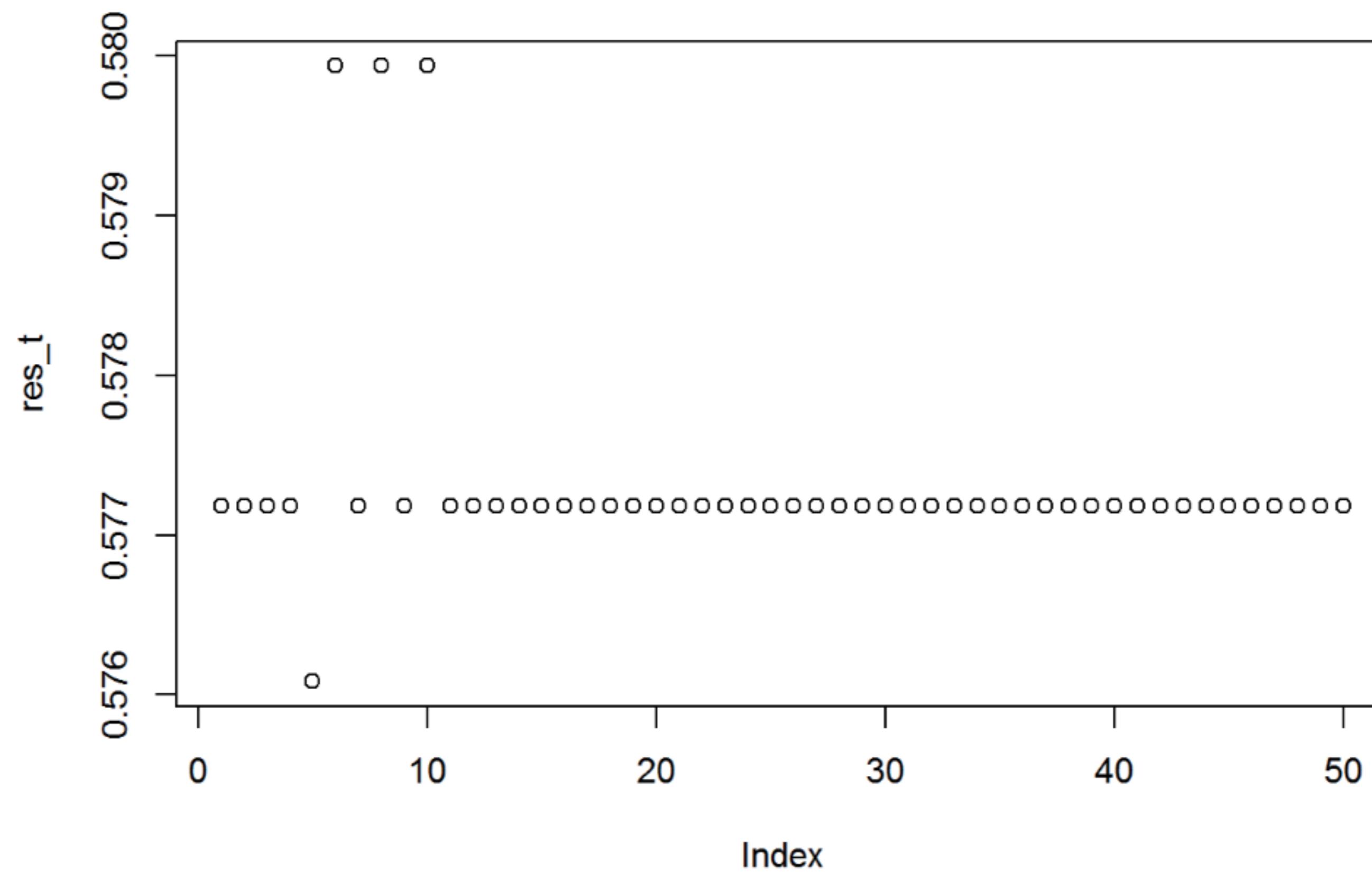


Decision Tree

?forward.search() 함수를 이용하여 변수를 설정하였음.

1~50 횟수의 Boosting을 실행한 결과, 의미가 있었음.

RESULT



RESULT

```
fit_t = which(res_t == min(res_t))

train_EPL_tree <- C5.0(formula = f, data=train_data, trials = fit_t)
test_predicted_EPL_tree = predict(train_EPL_tree, newdata = test_data)
1-sum(test_data$match_flag != test_predicted_EPL_tree)/nrow(test_data)
```

```
## [1] 0.5364238
```

```
table(test_data$match_flag, test_predicted_EPL_tree)
```

```
##          test_predicted_EPL_tree
##          draw  lose  win
##  draw     0   13   28
##  lose     0   24   20
##  win      0    9   57
```

CONCLUSION



Multinomial
Logistic Regression

k-NN

Naïve Bayes

Decision Tree



결과값이 가장 높은 Decision Tree로 결정하였다.

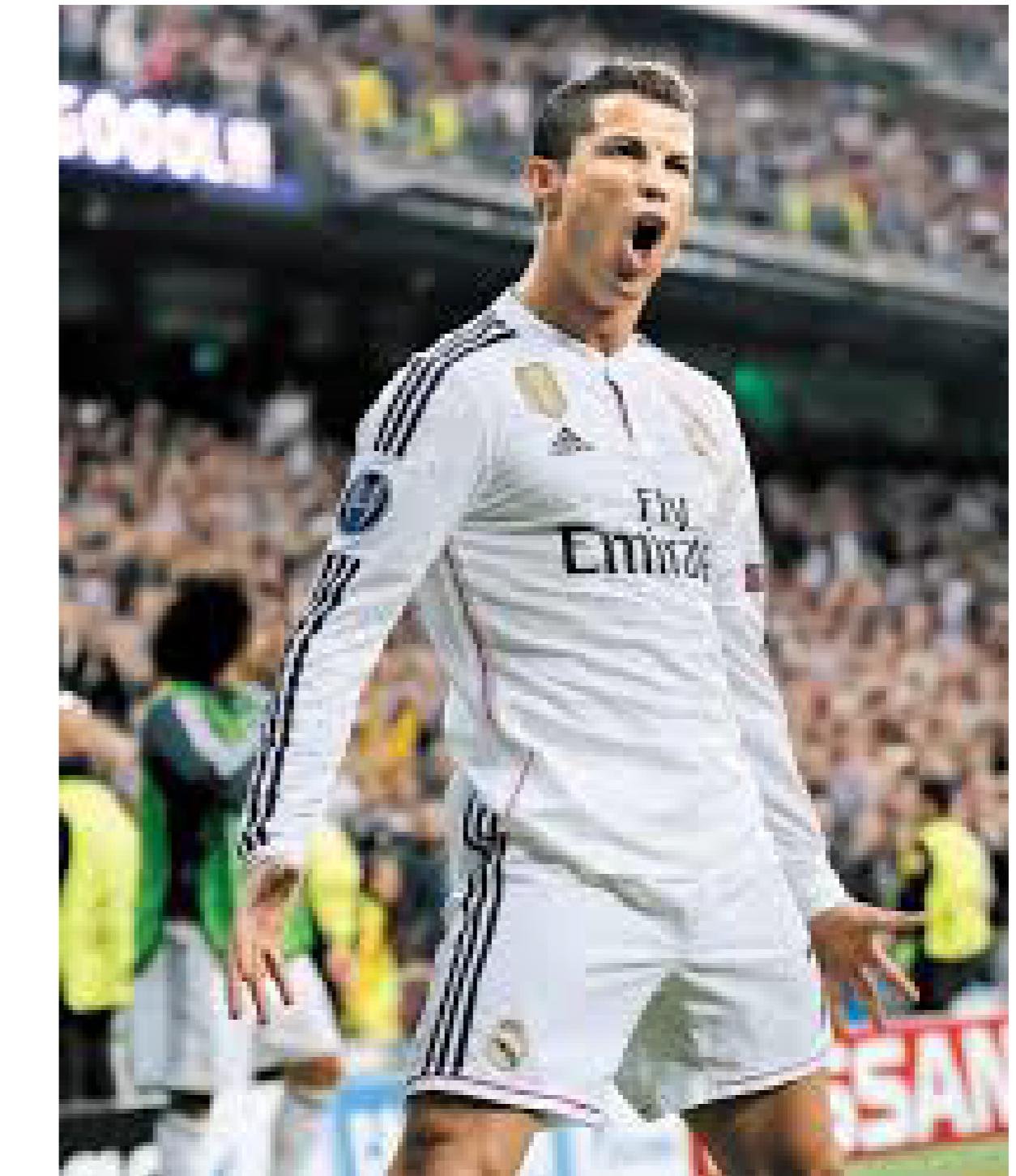
CONCLUSION



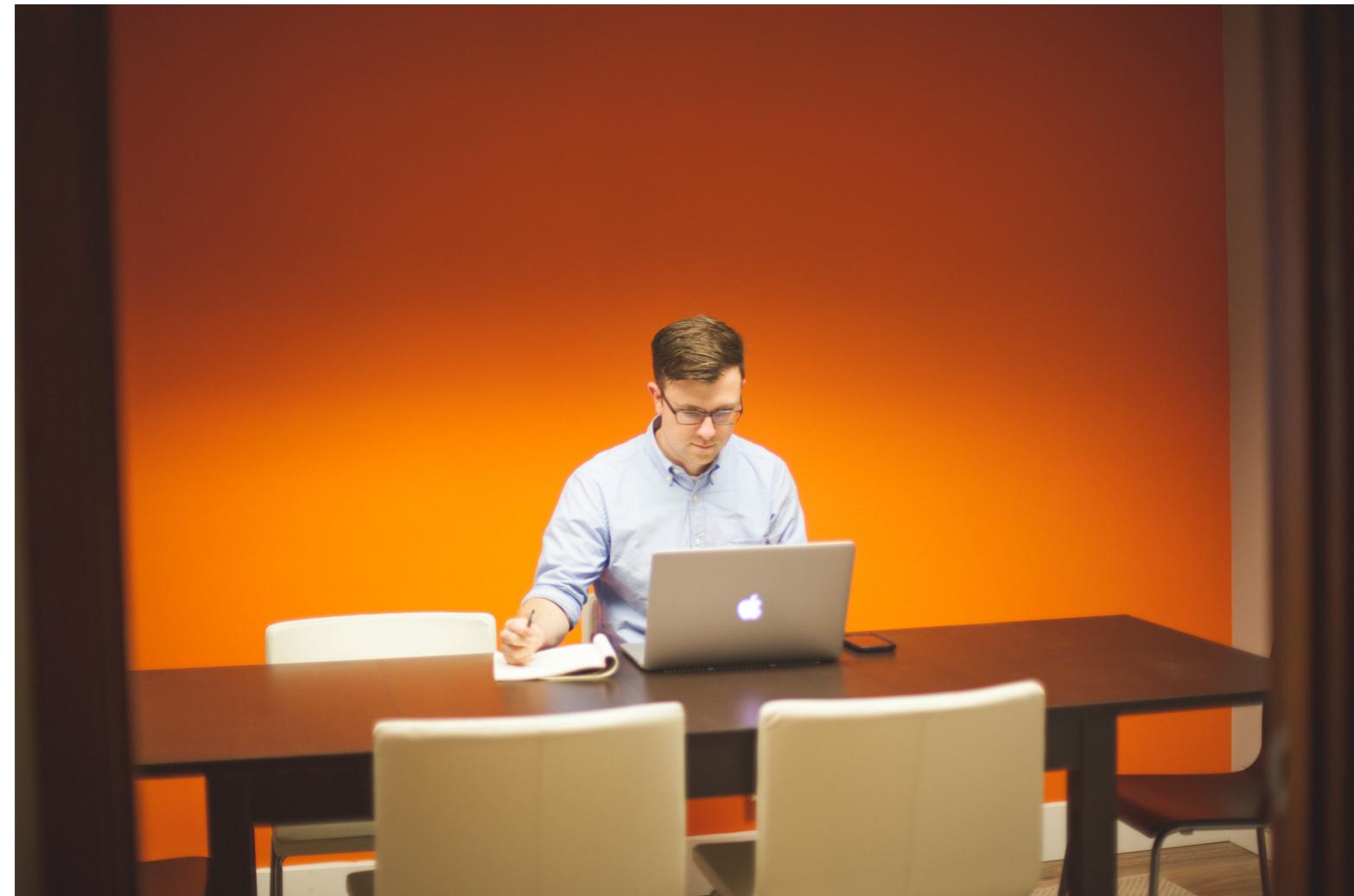
중요 변수

Decision Tree의 경우, [Home/Away](#)와 [Player Rating](#) 두 가지의 변수가 예측의 상관 관계가 가장 유의했다.

다른 모형들의 경우도 마찬가지로 구단이나 감독의 데이터보다는 [선수의 시장가치나 경기 스탯](#)에 많은 영향을 받는 것으로 나타났다.



CONCLUSION



한계점 및 소감

데이터를 어떻게 준비하고 처리할 지에 대해서 어려움이 있었다.

참조 문헌이나 자료가 전무한 모형을 만드는 작업이었기 때문에, 모형을 디자인하는 과정에서 난항을 겪었다.

Thanks!