

# How Soccer Players' Box Score Statistics Effect on their Rating and Market Value

**Haoyang Yan**

*Shanghaitech University    Shanghai    201210    China*

**ABSTRACT.** *This paper discusses the application of advance statistics in soccer. In order to valuate the performance and value of soccer players, we introduce more advance statistics and try to detect the effect on players' rating and market value from their box score statics.*

*In section 2 of this paper, we introduce an advance statistic - Box Rating (BR) - to valuate the performance of a soccer player in a match. The derivation of BR comes from the rating given by Whoscored regressed by 26 series of box score statistics (including goals, assists, passes and so on) of 2024 players in European top 5 leagues last 3 competition seasons (2017-2020). We then analyse the ridge trace of the regression coefficient to test the robustness of the regression.*

*In section 3 of this paper, to avoid the confounders like age, injury and commercial value which have unobserved effect to players' market value, we use 35 series of player ability ratings from FIFA20 as the Instrument Variable and do two stage least square (2SLS), getting the true effect players' box score statistics exerting on their market value.*

*In section 4 of this paper, we compute the correlation among ratings and market value, and do some hypotheses-test of our results. We also point out some limitation of players' statistics.*

**KEYWORDS:** *box score statistics, rating, market value*

## **1 Introduction**

Advance statistics are widely used in professional sports leagues to value players' performance. Based on detailed data recording and high-score matches, professional leagues of baseball, hockey and basketball have introduced plenty kinds of advanced statistics to value players' performance comprehensively. Adjusted Plus-Minus (APM) in National Hockey League (NHL) [1] [2] and Regularized Adjusted Plus-Minus (RAPM) in National Basketball Association

(NBA) [3] [4] [5] are two of the representative examples of advanced statistics. The common characteristic of APM and RAPM is that they are both derived on the basis of plus-minus, which means how many scores the team win/lose when the player is on the field, independent with what the player actually do on the field, but can generally accurately value how strong the player effect the match.

However, when applied to soccer, statistics from plus-minus like APM and RAPM face difficult problems to solve [6] [7]. Firstly, the score of soccer match is too small, which means the sample size is not large enough to insure a low-level stochasticity. Secondly, the first string of a soccer team is relatively fixed and there are only 3 substitution quota for a team in one match. These two factors make plus-minus statistics full of stochasticity and hard to reflect soccer players' real contribution to matches.

The advance statistics used in soccer are mostly one-sided. There has been developed advance statistic to predict expected goal which value offence chances and quality of shoots [8] [9], advance statistic to value players from their actions [10], advance statistic to value performance of goalkeeper [11], advance statistics to compare the team performance between competitions home and away [12].

There is still one way to value players' general performance on-field. To value the contribution of basketball players without large-scale sample and abundant plus-minus advance statistics, Daniel Myers develop the Box Plus/Minus (BPM), which is a advance statistic regressing to the Regularized Adjusted Plus Minus (RAPM). [13] Through the regression function, we can easily gain BPM of players with simple box score statistics, avoiding the limitation of immature statistics in previous time. This method can also be used in soccer. Although there is no meaningful plus-minus statistic like RAPM for soccer, we can use the rating from media to replace.

Valuing players' personal market value is a more complex problem. Not only on-field performance, but also other factors such as injury, age, personal background effect their market value [14] [15]. Thus, detecting the true effect players' box score statistics exerting on their market value is a challenging work for us.

## 2. BR - valuate players' performance by box score statistics

### 2.1 preprocessing and cleaning

In this part, we use web crawler to obtain the box score and rating data from whoscored and do some preprocessing and cleaning to these data. www.whoscored.com is a unique website and one of the fastest growing in the sports industry, specializing in the in-depth analysis of detailed soccer data.

Firstly, we filter out data of soccer players in the European top 5 league (England Premier League, Spanish La Liga, Italian Serie A, German Bundesliga,

French Ligue 1) last 3 competition seasons (2017-2020). Secondly, we filter out that a player played more than 5 matches in a competition season to avoid the error of small size sample. Thirdly, we match the players' name to the players' name in the dataset of FIFA2020, cleaning up players who cannot match (for players with the same name, we use his height and weight to confirm). After that, we get 4718 rows of data of 2024 different players in last 3 years.

In the vertical direction of the data, we reserve 1 column of players' average rating in this season and 11 columns of players' essential information, including name, season name, tournament name, team name, age, height, weight, position, appearance, substitute on, minutes played. After some preprocessing and cleaning, we reserve 26 columns of box score, which are all players' average box score statistics of every match in a season. For those statistics which have multiple items, such as total, won, lose, accurate, inaccurate, we only reserve the

total and rate to valuate players' yield and efficiency. What's more, we introduce the shot adjusted distance to estimate players' average shot distance, which

equals to  $\frac{1}{4} \text{shotsixyardboxrate} + 2 \text{shotpenaltyarearate} + 3 \text{shotoutboxrate}$ . Detailed information of these boxscore is listed on Table 1.

### 2.2 linear regression

In order to detect the relationship between rating and 26 series of box score statistics, we use multiple linear regression to regress the rating by box score statistics. The predicted value of rating,  $\hat{rating}$ , is introduced as Boxscore Rating (BR).

The regression formula and result is as follow:

$$BR = \hat{rating} = \beta_0 + \beta_i * \text{boxscore} \quad (1)$$

$$i=1$$

Table 1: 26 series of box score statistics and their (new) regression coefficients

 $\beta_i$ , ( $\gamma_i$ ) to rating

$i$	$boxscore_i$	$\beta_i$	$\gamma_i$
1	goal	1.0223	1.0568
2	assist	0.6939	0.7434
3	key pass	0.1356	0.1399
4	turnover	(0.0584)	-
5	dispossessed	(0.0422)	-
6	clearance	0.0387	0.0797
7	save	0.2059	0.2312
8	tackle won	0.1021	0.1117
9	interception	0.0895	0.1306
10	shots total	0.0547	0.0440
11	shot adjusted distance	(0.0121)	-
12	shot on target rate	0.0637	0.0499
13	penalty taken	(0.5345)	(0.5595)
14	pass long	(0.0053)	-
15	pass long accurate rate	0.0776	0.1932
16	pass short	0.0044	-
17	pass short accurate rate	0.1597	-
18	dribble	0.0935	0.0536
19	dribble rate	0.0908	0.1212
20	duel aerial	0.0311	-
21	duel aerial rate	0.2022	-
22	goal own	(0.6545)	(0.5834)
23	foul given	0.0365	0.0346
24	foul committed	(0.0246)	-
25	yellow card	(0.1549)	(0.1738)
26	red card		

	(1.1921)	
	(1.2305) 0	
	constant	
	5.6716	
	5.8121	
$\overline{R^2}$	0.8915	0.8299
$R$	0.8909	0.8290

It is obvious that goals and assists are two of the most influential factors to players' rating. Besides, turnover dispossessed, penalty taken (this is reasonable because the goals' value become less when a player take penalty), goal own, yellow and red card do negative effect on rating.

The goodness of fit of the linear regression  $R^2 = 0.8915$  and the adjusted  $R$ -squared  $R^2 = 0.8909$  show good fitness, but we still need further test to certify this multiple linear regression model is suited to explain the relationship between rating and box score.

### 2.3 ridge regression

To estimate the autocorrelation among these box score statistics, we need to analyse the trace of their ridge regression. Before doing ridge regression, we linearly normalize all the box score into the range [0,1]. Then we use ridge regression

$$\beta(k) = (\text{boxscore}^T * \text{boxscore} + k * I)^{-1} * \text{boxscore}^T * \text{rating} \quad (2)$$

where the ridge coefficient  $k \in [0, 10]$ ,  $I$  is the identity matrix, and  $\beta(k)$  is the normalized regression coefficient under ridge coefficient  $k$ .

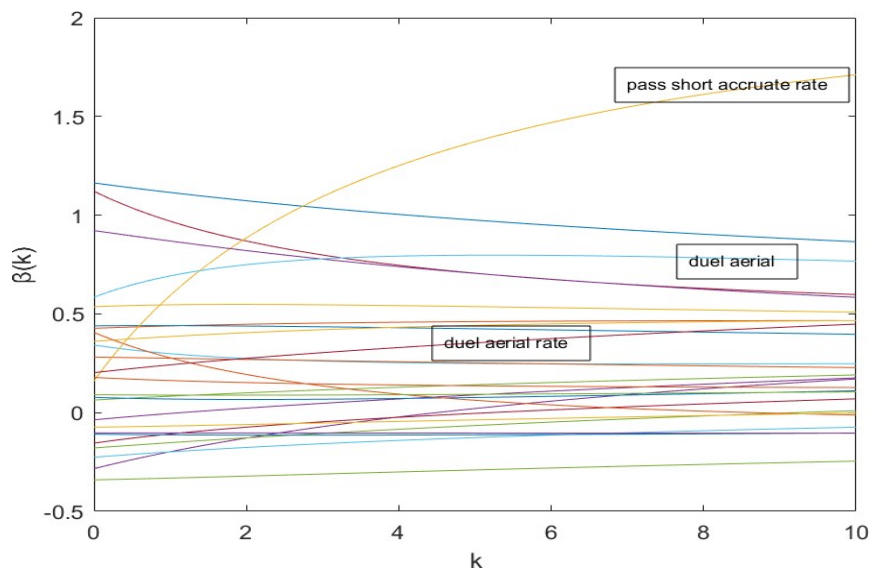
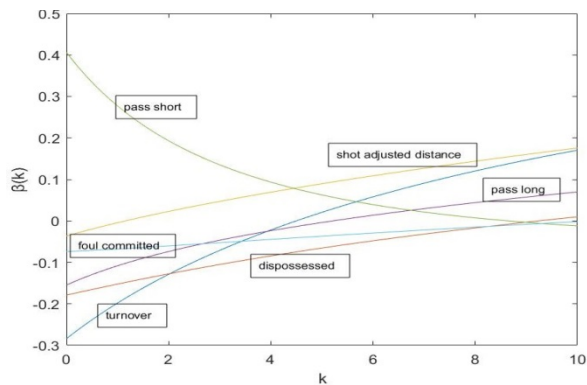


Figure 1: ridge trace of coefficients (after normalizing) (without  $\beta_0$ )

In the ridge trace figure, there are 3 series of regression coefficients - pass short accurate rate, duel aerial, duel aerial rate - are obviously abnormal, which increase when the ridge coefficient  $k$  increase.



What's more, the absolute values of 6 series of regression coefficients - turnover, dispossessed, shot adjusted distance, pass long, pass short, foul committed - have been less than 0.05 in range[0,10].

Figure 2: ridge trace of 6 series of regression coefficients less than 0.05 (after normalizing)

This abnormality implies that these series of boxscore may not satisfy the linear relation with the rating or not stable to be significant. So we delete these 9 series of box score statistics, using other 17 series to do linear regression. The new regression coefficients  $\gamma_i$  are showed in table 1, with  $R^2 = 0.8299$  and

$\bar{2}$

$R = 0.8290$ .

### 3. how players' box score statistics effect their market value

#### 3.1 choosing of instrument variable

In this section, we want to detect the true effect players' box score statistics exerting on their market value. However, there are many confounders correlated with both box score and market value, which may exaggerate or shrink this effect. Here are some examples, younger players are usually more valuable with similar boxscore and ability because they are more potential and may work for the team much longer; injury players do not appear much and are not in good body conditions, which means their box score may not be impressive, while their market value will not decrease that much; players whose position are forward may be well-known to spectators and have higher commercial value.

In order to remove the interference by these confounders, we use the ability ratings of players in FIFA20 as the instrument variable, which is only correlated to players' performance on field, but not to any other factors. From the instrument variable and two stage least square, we can separate the endogenous and exogenous part of box score statistics, which reveal the true effect players' box score statistics exerting on their market value.

FIFA20 is one of the most popular soccer simulation video game published by Electronic Arts as part of the FIFA series. The ability rating in FIFA20 gives 35 different kinds of ratings of each player's ability, including pace, shooting, passing and so on.

### 3.2 first stage regression

In the first stage regression, we use 35 series of ability ratings to regress 26 series of box score statistics. Because the ratings of ability are mostly aggregated in range[60,100], we use the log-linear regression to fit the curve better, and regularize boxscore into range[1,e] and ability into range[0,1]. The coefficient  $\delta_{ij}$  and the goodness of fitting of the log-linear regression is showed in table 2.

$$\log(\widehat{boxscore}_j) = \delta_{0j} + \sum_{i=1}^{35} \delta_{ij} * ability_i \quad (3)$$

$$boxscore_j = e^{\delta_{0j}} * \sum_{i=1}^{35} \delta_{ij} * ability_i + u_j$$



Table 2: Regression coefficients  $\delta_{ij}$  between ability ratings and box score statistics

$\delta_{ij}$	goal assist key-	tu rn -	dispo- tackle- shot-	clear- inter- shots-	save h ty- ot	s h ty- ot	penal pass-	p as s-	p as s-	p as s-	dribble dribble- duel-	d u el-	g o el	fo o al	fo ul -	yellow-	red-
	Pass v e r	o ssess ance	Won ception Total adjusted-on-	Tak en long	long short	sh or t-	rate	Aerial Own	Given -Card	Card							
					distan	acc	acc	rate	ted								
					ce	urat	urat										
					tar	e-	e-										
					ge	rat	rat										
					t-	e	e										
					ra												
					te												

pace 0. 0. 0. 0. 0. 0. (0. (0 0. 0. 0. 0. 0. (0.0. 0. 0. 0. 0. (0.(0.0. 0. (0 (0 0.  
 1 15 0 1 0 2 42 .0 07 13 3 3 0 36 0 18 02 21 0 06 26 0 0 .2 4 0  
 7 6 7 9 0 ) 1) 5 3 6 ) 6 3 ) ) 3 5 0) 0) 2  
 accele(0. (0 (0.0. (0.(0.0. (0 0. (0 (0.(0.(0.0. (0. (0 0. 0. (0.(0.(0.0. 0. 0. 0. (0  
 ration 03 .0 02 0 02 14 1 .0 03 .0 14 05 07 0 11 .1 01 06 10 11 11 0 0 07 19 .0  
 ) 8) ) 2 ) ) 5 5) 1) ) ) ) 4 ) 6) ) ) ) 0 2 1)  
 sprint (0. 0. 0. (0.(0.(0.0. (0 (0 (0 (0.(0.(0.0. (0. (0 (0 0. 0. 0. 0. (0.(0. (0 0. 0.  
 \_spee 05 05 0 05 01 02 2 .0 .1 .0 21 13 02 2 02 .0 .0 01 1 0 2 03 14 .0 04 0  
 d ) 0 ) ) ) 3 7) 2) 4) ) ) ) 2 ) 2) 0) 0 4 5 ) ) 9) 3  
 shooti 0. 0. 0. (0.(0.0. (1. (0 (0 0. 0. 0. 0. (1.(0. 0. (0 (0 0. 0. (0.(0.(0. (0 (0 (0  
 ng 3 22 0 12 13 0 04 .0 .2 31 6 5 0 73 35 06 .2 .2 3 2 82 11 30 .2 .2 .3  
 0 6 ) ) 4 ) 5) 0) 1 0 7 ) ) 2) 9) 6 6 ) ) ) 1) 8) 7)  
 positi 0. 0. 0. 0. 0. (0.0. 0. (0 (0 (0.0. 0. (0.(0. (0 (0 0. (0.0. 0. (0.0. 0. 0. (0  
 oning 0 03 0 2 1 34 0 06 .1 .0 07 0 0 12 04 .1 .0 00 23 0 0 04 0 19 06 .0  
 3 7 1 5 ) 8 1) 3) ) 1 7 ) ) 5) 8) ) 2 7 ) 7 3)  
 ftnish 0. (0 (0.0. 0. (0.0. (0 (0 (0 (0.(0.(0.0. 0. (0 0. 0. (0.(0.0. 0. 0. 0. 0. 0. 0. 0.  
 ing 0 .0 06 2 1 14 4 .0 .0 .0 30 02 03 6 1 .2 05 18 19 03 2 0 1 09 02 0  
 2 7) ) 0 4 ) 6 6) 6) 1) ) ) ) 3 1 1) ) ) 6 4 7 9  
 shot\_ (0. (0 (0.(0.(0.0. (0 (0 (0 0. (0.(0.0. 0. (0 0. 0. (0.(0.0. 0. (0. (0 0. 0.  
 power 13 .1 08 04 06 07 1 .0 .0 .0 0 02 08 3 0 .0 04 00 17 06 1 0 06 .0 01 0  
 ) 3) ) ) ) ) 9 2) 5) 7) 1 ) ) 5 1 9) ) ) 4 8 ) 2) 6  
 long\_ (0. (0 0. 0. 0. (0.0. (0 0. 0. 0. (0.(0.0. 0. (0 0. 0. 0. (0.0. (0.0. 0. 0. 0. 0.  
 shots 08 .0 0 0 0 0 07 2 .0 03 04 0 24 01 3 1 .0 04 06 0 13 1 03 0 12 19 1  
 ) 6) 2 4 5 ) 2 3) 2 ) ) 6 3 5) 2 ) 8 ) 7 2  
 volley 0. 0. 0. 0. 0. (0.0. (0 (0 0. (0.(0.(0.0. 0. 0. 0. (0 (0.(0.0. 0. 0. 0. 0. 0. 0. 0.  
 s 0 01 0 0 0 09 0 .0 .0 02 18 05 04 0 0 02 02 .0 14 03 0 0 0 02 00 0  
 6 3 1 1 ) 8 4) 3) ) ) ) 7 9 3) ) ) 7 3 4 2  
 penalt 0. (0 0. 0. 0. 0. 0. 0. 0. 0. (0 0. 0. 0. 0. (0.0. (0 (0 (0.0. 0. 0. 0. 0. 0. 0. 0.  
 ies 0 .0 0 0 0 0 0 02 01 .0 0 0 2 1 00 02 .0 .0 03 0 0 0 0 08 00 0  
 3 2) 4 3 3 2 3 2) 3 4 5 1 ) 0) 3) ) 3 1 1 7 2  
 passin (0. 0. 0. (0.(0.(0.0. 0. 0. (0 0. 0. (0.0. 0. 0. 0. (0 (0.(0.(0.0. (0. (0 (0 (0  
 g 10 00 1 28 10 24 2 05 03 .0 0 1 11 0 3 70 27 .1 67 13 24 1 02 .3 .3 .0  
 ) 1 ) ) ) 9 8) 9 6 ) 5 6 4) ) ) ) 8 ) 3) 5) 4)  
 vision (0. 0. 0. 0. 0. (0.(0.0. (0 0. (0.(0.(0.(0.(0. (0 (0 0. 0. 0. 0. (0.(0. (0 (0 (0  
 04 07 0 0 0 16 01 00 .1 00 05 16 02 04 02 .0 .0 05 2 0 0 09 10 .0 .0 .0

) 2 7 9 ) ) 0) ) ) ) ) ) 6) 4) 0 2 1 ) ) 3) 9) 0)  
crossi (0. 0. 0. 0. (0. 0. (0. 0. 0. (0 (0. (0. 0. 0. (0. (0 (0 0. (0. (0. 0. (0. (0. (0 0. (0  
ng 08 17 2 0 03 1 06 10 08 .0 03 08 0 0 37 .1 .1 08 07 03 1 01 01 .0 06 .0  
) 0 5 ) 5 ) 8) ) ) 2 3 ) 7) 2) ) ) 0 ) ) 4) 2)  
freeki 0. 0. 0. 0. 0. 0. (0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
ck 0 13 1 0 0 0 02 03 09 13 0 0 0 1 0 06 01 01 1 02 0 0 0 01 15 0  
9 7 1 1 6 ) 5 1 5 1 6 0 ) 2 0 8 1  
short 0. (0 (0. 0. 0. (0. (0. 0. 0. (0 0. 0. (0. (0. (0 (0 (0 0. 0. (0. (0. 0. 0. 0. 0.  
\_pass 0 .1 05 0 0 13 23 06 03 .1 2 1 04 35 01 .0 .0 .1 3 0 04 07 02 13 08 0  
2 2) ) 2 1 ) ) 1) 3 2 ) ) ) 3) 2) 3) 9 1 ) ) ) 2  
long\_ (0. 0. (0. (0. 0. 0. (0. (0 0. 0. (0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
pass 03 01 02 05 0 1 05 .0 06 06 03 0 0 2 1 .0 02 09 2 05 0 0 0 02 05 0  
) ) ) 1 0 ) 4) ) 2 2 0 2 3) 0 ) 5 1 0 1  
curve 0. 0. 0. (0. (0. 0. 0. (0 (0 (0 0. 0. 0. 0. 0. (0. (0 (0 0. 0. (0. 0. (0. (0. (0 (0 0.  
0 04 0 02 05 0 0 .0 .0 .0 0 0 0 0 04 .0 .0 01 0 02 0 04 10 .0 .0 0  
2 4 ) ) 7 1 9) 8) 1) 1 0 1 6 ) 4) 0) 2 ) 3 ) ) 8) 8) 3  
ovr\_dr (0. (0 (0. 0. (0. (0. (0 (0 (0 (0. (0. 0. 0. (0 0. 0. (0. (0. (0. (0. 0. 0. 0. 0. 0.  
ibbeli 10 .1 20 0 03 08 08 .1 .0 .0 24 54 03 8 0 .4 00 02 06 16 37 20 0 36 77 0  
ng ) 5) ) 1 ) ) ) 4) 7) 7) ) ) ) 2 5 1) ) ) ) ) 3 6  
agilit (0. (0 (0. (0. 0. 0. (0. 0. (0 (0 0. (0. (0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
y 07 .1 14 01 0 0 04 05 .0 .0 0 07 02 14 0 .0 00 .0 0 01 09 0 1 04 04 .0  
) 1) ) ) 0 0 ) 1) 6) 8 ) ) ) 7 6) 5) 2 ) ) 4 1 1)  
balan (0. 0. 0. (0. 0. (0. (0. 0. (0 (0 0. 0. (0. (0. (0 (0 0. 0. (0. (0. 0. 0. 0. 0. 0. 0. 0.  
ce 03 01 0 01 0 11 07 05 .0 .0 1 1 01 13 06 .0 .0 00 0 08 20 00 0 .0 .1 .0  
) 5 ) 6 ) ) 2) 5) 5 1 ) ) ) 3) 1) 6 ) ) ) 1 5) 3) 1)  
reacti 0. 0. 0. (0. (0. 0. 0. (0 0. 0. (0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
ons 4 45 2 07 05 4 0 .1 16 33 08 2 2 2 0 40 01 02 08 1 04 0 0 .2 .1 .0  
9 8 ) ) 6 1 0) ) 6 0 0 6 ) 4 ) 6 0 3) 5) 4)  
ball\_ 0. 0. 0. 0. 0. 0. (0. (0 0. 0. 0. 0. (0. (0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
contr 0 06 0 0 0 0 07 .0 01 18 2 0 02 16 2 14 .0 09 1 0 0 1 1 07 01 .0  
ol 1 6 5 4 1 ) 4) 4 5 ) ) 3 0) 6 3 9 2 1 1)  
dribbl (0. 0. 0. 0. 0. 0. (0. 0. 0. 0. (0. 0. 0. (0. (0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
ing 03 04 1 2 3 15 0 23 03 .0 1 3 05 54 0 12 .0 37 1 0 1 02 1 .0 .3 .0  
) 4 2 7 ) 3 1) 9 1 ) ) 1 3) 0 2 7 ) 9 1) 9) 0)  
comp 0. 0. 0. (0. (0. 0. 0. (0 (0 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.

osure 1 09 0 07 04 0 02 .1 .1 0220 1 0 03 1 25 11 05 3 08 0 02 0 .2 .1 0  
2 8 ) ) 9 ) 3) 1) ) 2 7 ) 6 1 ) 7 ) 1 2) 4) 0  
defen (0. (0 0. (0. (0.0. 0. 0. 0. (0 (0. (0.0. 1. (0. (0 0. (0 1. (0.0. 0. (0. (0 0. 0.  
se 25 .1 0 36 25 2 4 34 75 .3 47 16 0 0 03 .0 09 .0 6 37 6 2 19 .0 25 3  
) 5) 8 ) ) 8 2 4) ) ) 5 9 ) 7) 4) 0 ) 8 3 ) 2) 3  
interc (0. 0. (0. (0. (0. (0.0. 0. 0. (0 0. (0. (0.0. 0. 0. (0 (0.0. (0. (0. (0.0. 0. (0  
eption 00 03 03 03 07 04 04 10 12 .0 0 06 03 13 1 19 06 .1 29 0 03 03 02 03 03 .0  
s ) ) ) ) ) ) 0) 2 ) ) ) 8 0) ) 3 ) ) ) 6)  
headi 0. (0 (0.0. (0.0. (0. (0 (0 0. (0.0. 0. (0. (0.0. (0 (0 (0.0. 0. 0. (0. 0. (0 (0  
ng 1 .0 14 0 06 2 09 .0 .0 08 06 1 0 10 08 03 .0 .1 20 3 0 0 03 01 .0 .0  
6 5) ) 5 ) 3 ) 8) 1) ) 3 0 ) ) 4) 2) ) 2 3 3 ) 1) 1)  
aware (0. (0 (0. (0. (0.0. (0. (0 (0 0. 0. 0. (0. (0. (0.0. 0. (0 (0.0. (0. (0.0. (0 (0 (0  
nesss 05 .0 06 03 04 0 15 .0 .1 02 1 0 02 29 01 05 01 .0 42 0 18 05 0 .0 .1 .0  
) 3) ) ) ) 2 ) 6) 3) 7 3 ) ) ) 9) ) 5 ) ) 3 6) 2) 9)  
stand 0. 0. 0. 0. 0. (0. (0.0. (0 0. 0. 0. (0. (0. (0.0. (0 0. (0.0. (0. (0.0. 0. (0 (0  
\_tack 0 05 0 2 2 04 15 02 .1 04 2 0 11 35 03 19 .0 10 49 2 18 07 0 03 .1 .1  
le 0 1 0 3 ) ) 7) 4 5 ) ) ) 1) ) 0 ) ) 6 0) 1)  
slide (0. (0 (0. (0. (0.0. (0. 0. 0. (0 0. (0. (0. (0. (0 (0 (0 (0. (0. (0. (0.0. 0. 0. (0  
\_tack 04 .0 09 11 12 0 04 15 10 .0 0 06 00 03 13 .0 .0 .0 14 02 02 00 0 05 08 .0  
le ) 4) ) ) ) 9 ) 8) 8 ) ) ) ) 0) 0) 2) ) ) ) 3 3)  
physi 0. 0. 0. 0. 0. 0. (0. 0. 0. 0. 0. 0. 0. (0. (0. (0 (0 0. 0. 0. 0. 0. 0. 0. 0. 0.  
cal 1 07 0 2 1 5 39 04 14 24 5 3 0 34 07 .1 .2 05 3 5 0 0 1 35 09 0  
6 1 3 4 1 ) 6 0 6 ) ) 2) 6) 8 4 2 0 0 5  
jumpi (0. (0 (0. (0. (0.0. 0. (0 0. (0 (0. (0.0. (0. (0 (0 (0 (0.0. 0. 0. (0. (0 (0 (0  
ng 05 .1 10 08 12 0 0 .0 03 .0 06 06 00 0 05 .0 .0 .1 07 0 1 0 10 .0 .0 .0  
) 0) ) ) ) 8 4 3) 4) ) ) ) 5 ) 4) 1) 2) ) 9 2 1 ) 0) 2) 1)  
stami 0. 0. 0. 0. 0. (0.0. 0. 0. 0. (0. (0.0. 0. (0. 0. 0. 0. (0. (0.0. (0.0. 0. 0. (0  
na 0 12 1 1 1 21 1 33 15 08 11 13 0 1 00 18 04 04 10 03 0 02 1 21 13 .0  
4 9 5 1 ) 2 ) ) 3 3 ) ) ) 4 ) 2 6)  
streng (0. (0 (0. (0.0. (0.0. (0 (0 (0 (0. (0.0. 0. (0 0. 0. (0. (0.0. (0. (0. (0 (0 (0  
th 10 .1 11 07 0 17 0 .1 .0 .1 15 15 03 0 0 .1 06 01 06 15 0 01 04 .0 .1 .0  
) 2) ) ) 2 ) 7 2) 8) 0) ) ) ) 6 6 0) ) ) 7 ) ) 7) 0) 2)  
aggre (0. (0 (0.0. 0. (0.0. 0. (0 (0 (0. (0.0. 0. (0. (0 (0 (0 (0. (0. (0. (0.0. 0. 0. 0.  
ssion 02 .0 07 0 0 16 0 02 .1 .0 12 09 0 0 04 .1 .0 .0 10 06 05 03 1 21 31 0  
) 7) ) 3 2 ) 7 4) 5) ) ) 0 2 ) 2) 1) 2) ) ) ) 2 6

const	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)	(0.0)
ant	18.1	10.0	0.03	2.3	11.10	1.3	0.09	3.3	.088	.03	0.06	0.00	14.19	0.0	0.0	0.0	0.0	0.0	0.0
	)7)	)4)	)05)	)175)	)693)	)7)399214													
R2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	4.30	5.5	5.7	9.56	64.58	6.2	1.7	1.58	59.50	2.5	5.0	2.37	20.0	0.0	0.0	0.0	0.0	0.0	0.0
	7.3	8.3	3.1	6.9	9.4	5.5	2.3	6.4	3.6	4.3	2.0	3.2	0.0	0.0	0.0	0.0	0.0	0.0	0.0
R2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	4.29	5.5	5.7	9.56	63.58	6.2	1.7	1.58	58.50	2.5	5.0	2.36	20.0	0.0	0.0	0.0	0.0	0.0	0.0
	6.3	8.3	3.1	6.9	9.4	5.5	2.3	6.4	4.1	2.5	4.3	2.0	3.2	0.0	0.0	0.0	0.0	0.0	0.0

Some of the coefficients are easy to explain, like the high positive correlation between the box score of pass-short (-accruate-rate) and the ability of passing. Some are not that comprehensible, like the high positive correlation between the box score of pass-long and the ability of defense. That is because players with higher defending ability are usually posited in backfield and have more chance to give long pass.

Some of these regression do not have significant goodness of fit, especially the regression to the box score statistics of penalty-taken, goal-own, yellow-card and red-card. Probably because the sample is too small or these box score statistics are mostly happened accidental.

### 3.3second stage regression

In this part, we use web crawler to obtain players' market value from transfer market. www.transfermarkt.co.uk is a website which has footballing information, such as scores, results, statistics, transfer news, and fixtures. Its valuation to player's market value is widely accepted among social media. The data is showed in Great Britain Pound £.

Same as the process we did in section 2.1, we match names of players who has the market value data to names already in our dataset. 3896 rows are matced and others are cleaned.

We use the coefficients  $\delta_{ij}$  in section 3.2 to predict the boxes<sup>^</sup>core. Then we do the second stage regression, which regress the predicting boxes<sup>^</sup>core to the market value in equation 4, where  $\zeta_j$  is the regression coefficient and  $(\sum_{j=1}^{26} \zeta_j * u_j + \epsilon)$  is the disturbing term. For comparision, we also do a linear regression without instrument variable, where  $\eta_j$  is the regression coefficient and  $v$  is the disturbing term in equation 5.

$$\begin{aligned}
 value &= \zeta_0 + \sum_{j=1}^6 \zeta_j * boxscore_j + \varepsilon \\
 &= \zeta_0 + \sum_{j=1}^6 \zeta_j * box\hat{score}_j + \left( \sum_{j=1}^6 \zeta_j * u_j + \varepsilon \right)
 \end{aligned} \tag{4}$$

$$value = \eta_0 + \sum_{j=1}^6 \eta_j * boxscore_j + v \tag{5}$$

Table 3: (two stage) regression coefficient ( $\zeta_j$ ),  $\eta_j$   
from box score statistics to market value

$j$	$boxscore_j$	$\zeta_j$	$\eta_j$
1	goal	9.57E+0 7	3.69E+0 7
2	assist	2.63E+0 8	2.48E+0 7
3	key pass	- 3.23E+0 7	- 2.33E+0 6
4	turnover	5.14E+0 6	- 1.15E+0 6
5	dispossessed	1.75E+0 7	- 7.29E+0 5
6	clearance	1.47E+0 7	- 1.82E+0 5
7	save	4.44E+0 7	4.09E+0 5
8	tackle won	5.74E+0 6	1.17E+ 05
9	interception	- 2.47E+0 7	- 2.03E+0 6
10	shots total	4.42E+0 6	3.12E+0 6
11	shot adjusted distance	7.98E+0 6	- 1.84E+0 6
12	shot on target rate	4.71E+0 7	6.98E+0 5
13	penalty taken	- 2.03E+0 8	- 2.80E+0 7
14	pass long	- 3.57E+0 6	- 1.45E+0 5
15	pass long accurate rate	- 3.45E+0 7	3.57E+0 6
16	pass short	1.11E+ 7	3.73E+0

6		06	5
1	pass short	-	1.81E+0
7	accurate rate	3.63E+0	7
1	dribble	-	3.65E+0
8		3.15E+0	6
1	dribble rate	7.83E+0	3.53E+0
9		7	6
2	duel aerial	-	-
0		1.07E+0	4.06E+0
2	duel aerial rate	1.92E+0	7.38E+0
1		7	6
2	goal own	6.68E+0	-
2		8	1.18E+0
2	foul given	-	-
3		1.94E+0	1.15E+0
2	foul committed	3.48E+0	-
4		7	3.10E+0
2	yellow card	-	-
5		4.97E+0	3.84E+0
2	red card	-	-
6		6.53E+0	1.37E+0
0	constant	-	-
		5.48E+0	1.87E+0
		7	7
$R^2$		0.3230	0.3556
$R^z$		0.3183	0.3511

In table 3, some of the coefficients  $\zeta_j$  (more specifically, penalty taken, goal own and red card) are abnormally large due to their small sample. What's more, coefficients of two stage regression  $\zeta_j$  are mostly larger than coefficients  $\eta_j$  without instrument variable. This reveal that we may have underestimated the effect players' box score statistics exerting on their market value because of those confounders. Giving an example, a player's market value will increase about £100,000,000 (that's the market value of Lionel Messi, one of the most famous soccer player in the world) if he score one more goal in every competition.



#### 4.some property and hypothesis

##### 4.1correlation among ratings and value

For each of the 3896 rows, we extract 4 series of data - the average rating of the player this year given by whoscored, the Boxscore Rating (BR) derived from the box score statistics, the overall ability rating given by FIFA20 (that's the weighted average ability rating of players) and the market value given by transfermarket. We compute the correlation coefficients among these 4 series of data.

correlation coeficient rating	rati ng	Boxscore Rating	overall ability	market value
	1.0 00	0.944	0.518	0.489
Boxscore Rating	0.9 44	1.000	0.474	0.462
overall ability	0.5 18	0.474	1.000	0.573
market value	0.4 89	0.462	0.573	1.000

Those correlation coefficients reveal that compared with rating and boxscore in the field, the ability of players is more highly correlated to their market value, probably because rating and boxscore can not reflect all the effect a player exert to the competition and the team.

##### 3.2 are younger players more valuable?

In section 3.1 we refer that younger players are usually more valuable with similar boxscore and ability because they are more potential and may work for the team much longer. In this part we want to test whether this hypothesis is convincing.

We regress the overall ability rating and the age of players to their market value with multiple linear regression, then we structure the T-statistic of the regression coefficient  $\alpha_2$  from age to value.

$$\hat{value} = \alpha_0 + \alpha_1 * overallability + \alpha_2 * age(6)$$

$$\frac{\hat{\alpha}_2}{\alpha_2} =$$

$$T = \frac{\hat{\beta}_c - \beta_c}{\hat{\sigma}_e \sqrt{c_{22}}} \sim t(n-3) \quad (7)$$

where

$$c_{22} = \frac{1}{n-3} \sum_{i=1}^n (value_i - \hat{value}_i)^2$$

$c_{22}$  is the (3,3) element of matrix  $(X^T X)^{-1}$

X is the matrix of independent variable [1, overallability, age]

$$Se = \sum_{i=1}^n (value_i - \hat{value}_i)^2$$

Because of the large sample, we can approximately regard the t-distribution as gauss distribution with mean 0 and standard deviation 1.

	constant	overallability	age	$R^2$
$\alpha$	-1.28E+08	2.64E+06	-2.35E+06	0.5632
T-statistic			-40.33	

The result shows that a player's market value decrease about 2 million pounds when he get one year older while his ability stay stable. Due to the large absolute value of the T-statistic, we have more than 99.99% confidence ( $p < 0.0001$ ) to believe that there exists a negative correlation between players' market value and age.

### 3.3 limitation on valuating goalkeeper

In some aspect, the rating and box score statistics exist significant limitation on valuating players. This becomes apparently when valuating players from different positions, especially goalkeeper. To find the limitation, we divide players into two groups, 303 goalkeepers and 4415 non-goalkeepers, using Z-test to detect the difference

of rating and Boxscore Rating between these two group of players.

$$Z_{S_{gk}}^{rating} = \frac{rating_{gk} - rating_{ngk}}{\sqrt{\frac{S_{gk}^2}{n_{gk}} + \frac{S_{ngk}^2}{n_{ngk}}}} = 4.798 \quad (8)$$

$$Z_{BR} = \frac{BR_{gk} - BR_{ngk}}{\sqrt{\frac{S_{gk}^2}{n_{gk}} + \frac{S_{ngk}^2}{n_{ngk}}}} = -7.086$$

From the Z-statistic we have more than 99.99% confidence ( $p < 0.0001$ ) to believe that goalkeepers' rating and Boxscore Rating is less than non-goalkeepers'. For Boxscore Rating the difference is larger. As we all know, goalkeepers always activate in penalty area and seldom have box score statistics. This means valuating players via box score statistics is not fair for goalkeepers.

## References

- [1] Brian Macdonald. Adjusted plus-minus for nhl players using ridge regression with goals, shots, fenwick, and corsi. *Journal of Quantitative Analysis in Sports*, 8(3), 2012.
- [2] Brian Macdonald. A regression-based adjusted plus-minus statistic for nhl players. *Journal of Quantitative Analysis in Sports*, 7(3), 2011.
- [3] Shankar Ghimire, Justin A Ehrlich, and Shane D Sanders. Measuring individual worker output in a complementary team setting: Does regularized adjusted plus minus isolate individual nba player contributions? *PloS one*, 15(8):e0237920, 2020.
- [4] Jeremy Mertz, L Donald Hoover, Jean Marie Burke, David Bellar, M Lani Jones, Briana Leitzelar, and W Lawrence Judge. Ranking the greatest nba players: A sport metrics analysis. *International Journal of Performance Analysis in Sport*, 16(3):737–759, 2016.
- [5] Jeremias Engelmann. A new player evaluation technique for players of the national basketball association (nba). In *Proceedings of the MIT Sloan Sports Analytics Conference*, 2011.
- [6] Tarak Kharrat, Ian G McHale, and Javier López Peña. Plus–minus player ratings for soccer. *European Journal of Operational Research*, 283(2):726–736, 2020.
- [7] Francesca Matano, Lee F Richardson, Taylor Pospisil, Collin Eubanks, and Jining Qin. Augmenting adjusted plus-minus in soccer with fifa ratings. *arXiv preprint arXiv:1810.08032*, 2018.
- [8] Masoomeh Zamani, Mahmood Fathy, and Amin Sadri. A low cost algorithm for expected goal events detection in broadcast soccer video. *International Journal of Digital Content Technology and its Applications*, 4(8):118–125, 2010.
- [9] Andrew Rowlinson et al. Football shot quality: Visualizing the quality of soccer/football shots. 2020.
- [10] Guiliang Liu, Yudong Luo, Oliver Schulte, and Tarak Kharrat. Deep soccer analytics: learning an action-value function for evaluating soccer players. *Data Mining and Knowledge Discovery*, pages 1–29, 2020.
- [11] Paul Power, Aditya Cherukumudi, Sujoy Ganguly, Felix Wei, Long Sha, Jennifer Hobbs, Hector Ruiz, and Patrick Lucey. Trading places—simulating goalkeeper performance using spatial & body-pose data, 2019.

- [12] V Barnett and S Hilditch. The effect of an artificial pitch surface on home team performance in football (soccer). *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 156(1):39–50, 1993.

### appendix: code

```
1 rating = xlsread ('WhoScoredFinal7. xlsx ', 'D3:D4720');
2 boxes = xlsread ('WhoScoredFinal7. xlsx ', 'O3:AO4720');
3 beta = regress ( rating , boxes);
4 BR = boxes * beta;
5 5
6 avgr = mean( rating );
7 SSR = sum((BR - avgr).^2);
8 SST = sum(( rating - avgr).^2);
9 R2 = SSR/SST;
10 10 ajR2 = 1 - (1 - R2)*(4718 - 1)/(4718 - 27 - 1);
11 11
12 12 ma=[];
13 13 for i = 1:27
14 ma = [ma, max(boxes(: , i))];
15 end
16 mam = repmat(ma, 4718 , 1);
17 norboxes = boxes./mam;
18 18
19 19 E = eye (27);
20 20 betak =[];
21 21 for k = 0:0.01:10
22 j = (norboxes .* norboxes + k * E)^(-1) * (norboxes ') * rating ;
23 betak = [betak , j];
24 end
25 25 for i = 1:26
26 26 plot (0:0.01:10 , betak ( i , :))
27 hold on
28 end
29 saveas (1 , '01. jpg ')
30 hold o f f
31 31
32 32 thr = zeros (1 , 27);
33 33 for k = 1:27
```

```
34 34 for j = 1:1001
35 35   if abs(betak(k, j)) < 0.05
36 36     thr (k) = 1;
37 end
38 end
39 end
40 40 for i = 1:26
41 41   if thr ( i )==1
42 42     plot (0:0.01:10 , betak( i ,:))
43 hold on
44 end
45 end
46 saveas (1 , '02. jpg ')
47 hold o f f
48 48
49 49   ajboxs = boxs;
50 50   ajboxs (:,[4 ,5 ,11 ,14 ,16 ,17 ,20 ,21 ,24]) = 0;
51 gamma = regress ( rating , ajboxs );
52 nBR = ajboxs * gamma;
53 nSSR = sum((nBR - avgr ).^2);
54 SST = sum(( rating - avgr ).^2);
55 nR2 = nSSR/SST;
56 56 najR2 = 1 - (1 -nR2)*(4718 -1)/(4718 -27 -1);
57 57
58 abi = xlsread ('WhoScoredFinal7. xlsx ','AR3:CA4720');
59 abi = abi /100;
60 reguboxs = log (norboxs*(exp(1) -1)+1);
61 delta = [];
62 62 for i = 1:26
63 delta = [ delta , regress (reguboxs(:, i ) , abi )];
64 end
65 reb = abi * delta ;
66 66 Rm = [];
67 67 ajRm = [];
68 68 for i = 1:26
69 avgnb = mean(reguboxs(:, i ));
70 SSR = sum(( reb(:, i ) - avgnb).^2);
71 SST = sum((reguboxs(:, i ) - avgnb).^2);
72 Rm = [Rm,SSR./SST];
73 73   ajRm = [ajRm,1 - (1 -SSR./SST)*(4718 -1)/(4718 -27 -1)];
74 74 end
75 75
```

```
76 ratingnew = xlsread ('WhoScoredFinal8. xlsx ', 'D3:D3898');
77 BRnew = xlsread ('WhoScoredFinal8. xlsx ', 'C3:C3898');
78 overall = xlsread ('WhoScoredFinal8. xlsx ', 'AQ3:AQ3898');
79 age = xlsread ('WhoScoredFinal8. xlsx ', 'H3:H3898');
80 boxesnew = xlsread ('WhoScoredFinal8. xlsx ', 'O3:AO3898');
81 value = xlsread ('WhoScoredFinal8. xlsx ', 'CC3:CC3898');
82 abinew = xlsread ('WhoScoredFinal8. xlsx ', 'AR3:CA3898');
83 abinew = abinew/100;
84 mamm = repmat(ma,3896,1);
85 85 mamm(:,27) = [];
86 preboxsnw = ((exp(abinew * delta )-1)/(exp(1)-1)).* mamm;
87 preboxsnw = [preboxsnw,ones (3896,1)];
88 zeta = regress ( value ,preboxsnw);
89 eta = regress ( value ,boxsnw);
90 90
91 prevalue1 = preboxsnw * zeta ;
92 prevalue2 = boxsnw * eta ;
93 avgv = mean( value );
94 SSR1 = sum(( prevalue1 - avgv).^2);
95 SSR2 = sum(( prevalue2 - avgv).^2);
96 SST = sum(( value - avgv).^2);
97 R21 = SSR1/SST;
98 R22 = SSR2/SST;
99 99 ajR21 = 1 - (1 -R21)*(3896 -1)/(3896 -27 -1);
100 100 ajR22 = 1 - (1 -R22)*(3896 -1)/(3896 -27 -1);
101 101
102 cor = [ratingnew ,BRnew, overall , value ];
103 corcoe = zeros (4);
104 for i = 1:4
105 for j = 1:4
106 corcoe ( i , j ) = corr ( cor (:, i ), cor (:, j ));
107 end
108 end
109
110 X = [ones (3896,1), overall ,age];
111 111 XX = (X'*X)^(-1);
112 112 c22 = XX(3,3);
113 alpha = regress ( value ,X);
114 a2 = alpha (3);
115 prevalue = X * alpha ;
116 Se = sum(( prevalue - avgv).^2);
117 R20 = Se/SST;
```