

1

2

**ARTICLE TYPE**

3

4

# Forecasting basketball players' performance using sparse functional data<sup>†</sup>

5

6

7

8

9

G. Vinué\*<sup>1</sup> | I. Epifanio<sup>2</sup>

10

11

12

<sup>1</sup>Department of Statistics and O.R.,  
University of Valencia, 46100 Burjassot,  
Spain

<sup>2</sup>Dept. Matemàtiques and Institut de  
Matemàtiques i Aplicacions de Castelló,  
Campus del Riu Sec. Universitat Jaume I,  
71 Castelló, Spain

**Correspondence**

\*G. Vinué. Email: guillermo.vinue@uv.es

## Summary

Statistics and analytic methods are becoming increasingly important in basketball. In particular, predicting players' performance using past observations is a considerable challenge. The purpose of this study is to forecast the future behavior of basketball players. The available data are sparse functional data, which are very common in sports. So far, however, no forecasting method designed for sparse functional data has been used in sports. A methodology based on two methods to handle sparse and irregular data, together with the analogous method and functional archetypoid analysis is proposed. Results in comparison with traditional methods show that our approach is competitive and additionally provides prediction intervals. The methodology can also be used in other sports when sparse longitudinal data are available.

**KEYWORDS:**

Forecasting, Functional data analysis, Archetypal analysis, Functional sparse data, Basketball

## 1 | INTRODUCTION

The statistical analysis of sports is a fast-growing field. In particular, sports forecasting is a strongly expanding field. Proof of this is the increasing number of forecasting methods developed covering several sports and sportive topics [37]. Certainly, in professional sports, getting any notion in advance about the future performance of players can have important consequences on the roster composition in terms of renewing or buying players, or in terms of establishing levels of remuneration in function of that performance, for instance. One of the sports that has been revolutionized by data analytics is basketball. Basketball analytics started to attract attention with the publications by [24] and [17]. More recently, other papers and books have been released [22, 34, 32]. Technological advances have made it possible to collect more data than ever about what is happening on the field, requiring new methods of analysis. There is currently a need for innovative methods that exploit the full potential of the data and that make it possible to generate additional value for athletes and technical staff. Here too, one of the main challenges is to use past performance to predict future performance [32]. To address this open question, some forecasting methods have been developed. Following [19, Chapter 1.4], two main approaches can be distinguished based on the type of data used: time-series forecasting and cross-sectional forecasting. On the one hand, forecasting using data collected over time describes the likely outcome of the time series in the immediate future, based on knowledge of recent outcomes. On the other hand, cross-sectional forecasting methods use data collected at a single point in time. The goal here is to predict a target variable using some explanatory variables which are related to it.

Two well elaborated methods can be found using historical time data: (i) College Prospect Rating (CPR), which is a score assigned to college basketball players that attempts to estimate their NBA potential [32, 33]; (ii) Career-Arc Regression Model

<sup>†</sup>The data and software associated with this paper are available at <https://www.uv.es/vivogui/software>.

1

Vinué ET AL

2 Estimator with Local Optimization (CARMELO), which, for a player of interest, identifies similar players throughout NBA  
3 history and uses their careers to forecast the future player's activity [35].

4 Regarding cross-sectional models, a Weibull-Gamma with covariates timing model is proposed in [18] to predict the points  
5 scored by players over time. In this case, the time variable is years playing in NBA. Another interesting approach is presented in  
6 [30], where correlations and regression models are computed to figure out which foreign players will be successful in the NBA,  
7 by using their previous statistics in international competitions.

8 In addition to the effort of predicting individual performance, there have also been other approaches focusing on teams and  
9 other features of the game. Some models using simulation have been developed to forecast the outcome of a basketball match  
10 [38, 43]. A comparison between predictions based on NCAAB and NBA match data is discussed in [47]. A dynamic paired  
11 comparison model is described in [3] for the results of matches in two basketball and football tournaments. Furthermore, in [4]  
12 a process model is used with player tracking data for predicting possession outcomes.

13 We wish to consider a new perspective by using Functional Data Analysis (FDA) in sports. FDA is a relatively new branch  
14 of Statistics that analyses data drawn from continuous underlying processes, often time, i.e. a whole function is a datum. Let  
15 us assume that  $n$  smooth functions,  $x_1(t), \dots, x_n(t)$ , are observed, with the  $i$ -th function measured at  $t_{i1}, \dots, t_{in_i}$  points. In our  
16 study,  $x_i(t)$  represents the metric value of player  $i$  for a certain age  $t$ . An important point we would like to emphasize here is that  
17 the time component of the FDA approach we are considering will represent players' ages. As such, in this paper we propose  
18 different models for aging curves, which is a well-recognized and important topic within the more general area of forecasting  
19 player performance. As mentioned in [35], the most important attribute of all, in terms of determining a player's future career  
20 trajectory, is his/her age.

21 The goals of FDA coincide with those of any other branch of Statistics, and the classical summary statistics can be also  
22 defined, such as the mean function  $\bar{x}(t) = n^{-1} \sum_{i=1}^n x_i(t)$ , the variance function  $var_X(t) = (n-1)^{-1} \sum_{i=1}^n (x_i(t) - \bar{x}(t))^2$  and the  
23 covariance function  $cov_X(t_h, t_l) = (n-1)^{-1} \sum_{i=1}^n (x_i(t_h) - \bar{x}(t_h))(x_i(t_l) - \bar{x}(t_l))$ . An excellent overview of FDA is found in [28],  
24 while methodologies for studying functional data nonparametrically are found in [15]. [29] introduces related software and [27]  
25 presents some interesting applications in different fields. Other recent applications include [9] and [23]. In all these problems, a  
26 continuous function lies behind these data even though functions are sampled discretely at certain points. The FDA framework is  
27 highly flexible since the sampling time points do not have to be equally spaced and both the argument values and their cardinality  
28 can vary across cases. When functions are observed over a relatively sparse set of points, we have sparse functional data. An  
29 excellent survey on sparsely sampled functions is provided by [21].

30 As regards the forecast of functional time series, there is a body of research, such as [31, 2, 20], where functions are measured  
31 over a fine grid of points. However, only a few works deal with the problem of forecasting sparse functional data [11]. Notice  
32 that when functions are observed over a dense grid of time points, it is possible to fit a separate function for each case using  
33 any reasonable basis. Nevertheless, in the sparse case, this approach fails and the information from all functions must be used  
34 to fit each function. This is because each individual has been observed at different time points. Therefore, any fixed grid that is  
35 formed will contain many missing observations for each curve. A very good explanation of this issue can be found in [21].

36 Sports data are sparse and irregular. They are sparse because most players do not have a very long career in the same league.  
37 And they are irregular because each player's career lasts for a different length of time. Despite the fact that time series data or  
38 movement trajectories are very common in sports, FDA has been mostly used in sport biomechanics or medicine [14, 16]. To  
39 the best of our knowledge, there are only two references about sports analytics using FDA. In [44], FDA was introduced for the  
40 study of players' aging curves and both hypothesis testing and exploratory analysis were performed. [41] extended archetypoid  
41 analysis (ADA) for sparse functional data (see also [42, 13]), showing the potential of FDA in sports analytics. In particular,  
42 in [41, Section 5.2] it was demonstrated that advanced analysis with FDA reveals patterns in the players' trajectories over the  
43 years that could not be discovered if data were simply aggregated (averaged, for example). We take advantage of this fact and  
44 continue the work done in [41, Section 5.2] with a view to predict how players will evolve.

45 In this paper, we propose a methodology to predict player's performances using sparse functional data. Two metrics will be  
46 analyzed: Box Plus/Minus (BPM)<sup>1</sup> and Win Shares (WS)<sup>2</sup>. Analysis using BPM will allow us to establish a plausible comparison  
47 with CARMELO, while analysis with another variable such as WS will allow us to evaluate differences in career arcs. To that end,  
48 we will focus on two existing methods designed to handle sparse and irregular data: (i) Regularized Optimization for Prediction  
49 and Estimation with Sparse data (ROPES), originally developed by Alexander Dokumentov and Rob Hyndman [11, 10]; (ii)

50  
51  
52  
53  
54  
55  
56 <sup>1</sup><https://www.basketball-reference.com/about/bpm.html>

57 <sup>2</sup><https://www.basketball-reference.com/about/ws.html>

Principal components Analysis through Conditional Expectation (PACE), originally developed by Fang Yao, Hans-Georg Müller and Jane-Lin Wang [45].

Our methodology will also involve using the method of analogues based on functional archetypoid analysis (FADA), which will allow us to refine predictions for the players of interest and to achieve a more reliable forecast, in line with the expectations of basketball analysts. We will apply them to a very comprehensive database of NBA players. Results will be obtained using the R software [25].

As noted earlier, forecasting future performance is also very relevant to other sports (see for instance [1]). We would like to emphasize that our methodology can also be used in other sports when sparse longitudinal data are available. Data and R code (including a web application created with the R package **shiny** [5]) to reproduce the results can be freely accessed at <https://www.uv.es/vivogui/software>. The rest of the paper is organized as follows: Section 2 reviews ROPES, PACE, ADA and FADA. Section 3 will be concerned with the data and input variables used. Section 4 presents two analyses: (i) ROPES and PACE are compared with each other and with standard benchmarks; (ii) The reliability of ROPES predictions for current players using the method of analogues with FADA is shown. A comparison with CARMELO results is also provided and the implications of the archetypoid coefficients is discussed. The paper ends with some conclusions in Section 5. An appendix contains a validation study to choose an optimal blend of tuning parameters which ROPES depends on, which is crucial for good predictive activity. It also shows how this methodology can also be proposed for forecasting international players. The appendix is available as an online supplement.

## 2 | METHODOLOGY

### 2.1 | ROPES

The method ROPES (Regularized Optimization for Prediction and Estimation with Sparse data), proposed by Alexander Dokumentov and Rob Hyndman [11, 10], solves problems involving decomposing, smoothing and forecasting two-dimensional sparse data. In practical terms, where the aim is to interpolate and extrapolate the sparse longitudinal data, made up of  $n$  observations, and presented over the time dimension with  $m$  time points, the following optimization problem is solved:

$$\{(U, V)\} = \underset{U, V}{\operatorname{argmin}}(||W \odot (Y - UV^T)||^2 + ||U||^2 + ||\text{DIFF}_2(m, \lambda_2)V||^2 + ||\text{DIFF}_1(m, \lambda_1)V||^2 + ||\text{DIFF}_0(m, \lambda_0)V||^2) \quad (1)$$

where:

- $Y$  is an  $n \times m$  matrix.
- $U$  is an  $n \times k$  matrix of “scores” (“coefficients”),  $k = \min(n, m)$ .
- $V$  is an  $m \times k$  matrix of “features” (“shapes”).
- $\|\cdot\|$  is the Frobenius norm.
- $\odot$  is the element-wise matrix multiplication.
- $W$  is an  $n \times m$  “masking matrix” of weights.
- $\lambda_0, \lambda_1$  and  $\lambda_2$  are smoothing parameters.

and where  $\text{DIFF}_i(m, \lambda)$  represents the discrete  $i$  times derivative operator multiplied by the scalar  $\lambda$ . In particular,  $\text{DIFF}_0(m, \lambda)$  is the identity matrix  $m \times m$  multiplied by the scalar  $\lambda$ ;  $\text{DIFF}_1(m, \lambda)$  is the matrix  $(m - 1) \times m$ , with  $-\lambda$  values in the main diagonal,  $\lambda$  values in the following upper diagonal and 0 otherwise;  $\text{DIFF}_2(m, \lambda)$  is the matrix  $(m - 2) \times m$ , with  $\lambda$  values in the main diagonal,  $-2\lambda$  values in the following upper diagonal,  $\lambda$  values in the following upper diagonal, and 0 otherwise.

ROPEs is equivalent to maximum likelihood estimation with partially observed data, which allows the calculation of confidence and prediction intervals. They are estimated using a Monte-Carlo style method. The original two sources [11, 10] should be referred to for all the specific details.

4

Vinué ET AL

## 2.2 | PACE

Functional Principal Components Analysis (FPCA) is a common tool to reduce the dimension of data when the observations are random curves. The usual computational methods for FPCA based on function discretization or basis function expansions are inefficient when data with only a few repeated and sufficiently irregularly spaced measurements per subject are available. As a quick reminder, when functions are measured over a fine grid of time points, a separate function for each individual can be used. In the sparse case, however, the information from all functions must be used to fit each function.

A version of FPCA, in which the FPC scores are framed as conditional expectations, was developed by Fang Yao, Hans-Georg Müller and Jane-Lin Wang to overcome this issue [45]. This method was referred to as Principal components Analysis through Conditional Expectation (PACE) for sparse and irregular longitudinal data. In practice, the prediction for the trajectory  $X_i(t)$  for the  $i$ th subject, using the first  $p$   $\phi_q$  eigenfunctions, is:

$$\hat{X}_i^p(t) = \hat{\mu} + \sum_{q=1}^p \hat{\xi}_{iq} \hat{\phi}_q(t) \quad (2)$$

where  $\hat{\mu}$  is the estimate of the mean function  $E(X(t)) = \mu(t)$  and  $\hat{\xi}_{iq}$  are the FPC scores. PACE and its implementation in the R library **fdapace** ([7]) use local smoothing techniques to estimate the mean and covariance functions of the trajectories, specifically a local weighted bilinear smoother is used for estimating the covariance. Generalized Cross Validation is used for bandwidth choice, which is the default method for the FPCA function in the R library **fdapace** (default parameters are considered; for example, 10 folds and a Gaussian kernel are used). The number of components  $p$  is determined using the Fraction-of-Variance-Explained threshold (0.9999 by default) computed during the SVD of the fitted covariance function.

The eigenfunctions  $\hat{\phi}_q(t)$  and the number  $p$  are estimated with the training set, and they are used in the estimation of the scores for the test set. This is the procedure we will follow in Section 4.1. With the scores and the estimated eigenfunctions, we obtain an approximation of the trajectories and they can be used to predict unobserved portions of the functions. [45] also explains the construction of asymptotic pointwise confidence intervals for individual trajectories and asymptotic simultaneous confidence bands.

## 2.3 | ADA

Archetypoid analysis (ADA) was presented in [42] and is an extension of archetypal analysis defined by [6] (see [8, 26] for other derived methodologies). In ADA, archetypes correspond to real observations (the so-called archetypoids). Let  $\mathbf{X}$  be an  $n \times p$  matrix of real numbers representing a multivariate data set with  $n$  observations and  $p$  variables. For a given  $g$ , the objective of ADA is to find a  $g \times p$  matrix  $\mathbf{Z}$  that characterizes the archetypal patterns in the data. In ADA, the optimization problem is formulated as follows:

$$RSS = \sum_{i=1}^n \|\mathbf{x}_i - \sum_{j=1}^g \alpha_{ij} \mathbf{z}_j\|^2 = \sum_{i=1}^n \|\mathbf{x}_i - \sum_{j=1}^g \alpha_{ij} \sum_{l=1}^n \beta_{jl} \mathbf{x}_l\|^2, \quad (3)$$

under the constraints

$$1) \sum_{j=1}^g \alpha_{ij} = 1 \text{ with } \alpha_{ij} \geq 0 \text{ and } i = 1, \dots, n \text{ and}$$

$$2) \sum_{l=1}^n \beta_{jl} = 1 \text{ with } \beta_{jl} \in \{0, 1\} \text{ and } j = 1, \dots, g \text{ i.e., } \beta_{jl} = 1 \text{ for one and only one } l \text{ and } \beta_{jl} = 0 \text{ otherwise.}$$

Archetypoids are computed with the R package **Anthropometry** [40].

### 2.3.1 | ADA for sparse data with FDA

ADA was defined for functions in [13], where it was shown that functional archetypoids can be obtained as in the multivariate case if the functions are expressed in an orthonormal basis, simply by applying ADA to the basis coefficients. When functions are measured over some sparse points, we have sparse functional data.

The basic idea of functional archetypoid analysis (FADA) is as follows. Based on the Karhunen-Loëve expansion, the functions are approximated as in Eq. 2. Because the eigenfunctions are orthonormal, to obtain FADA we can apply ADA to the  $n \times p$  matrix  $\mathbf{X}$ , with the scores (the coefficients in the Karhunen-Loëve basis).

### 3 | DATA

We have used the R package **ballr** [12] to obtain the total advanced statistics for each NBA player from the 1973-1974 season to the latest season, 2017-2018, including the player's age on February 1st of that season (the convention in <https://www.basketball-reference.com/> is to provide players' age at the start of February 1st of every season). From the total set of statistics, we have focused on Box Plus/Minus (BPM) and Win Shares (WS).

BPM is a box score-based variable for estimating basketball players' quality and contribution to the team. It takes into account both the players' statistics and the team's overall performance. The final value enables us to evaluate the player's performance relative to the league average. BPM is a per-100-possession statistic and its scale is as follows: 0 is the league average, +5 means that the player has contributed 5 more points than an average player over 100 possessions, -2 is replacement level, and -5 is very bad. Replacement level players are those who replace a roster spot for short-term contracts, so they are not normal members of a team's rotation. We have chosen BPM because it was created to only use the information that is available historically. This means that BPM only takes into account those stats that have always been available in the box-scores. It does not consider new stats derived from play-by-play data or from tracking data. According to the website where BPM is explained<sup>3</sup>, "*it is possible to create a better stat than BPM for measuring players, but difficult to make a better one that can also be used historically*", which fits perfectly with the goal of our method. BPM is available from the 1973-1974 season.

We have chosen a second metric, which is also widely used: Win Shares (WS). It also has the advantage of taking the surrounding team into account. In particular, WS is a player statistic that distributes the team's success among the team players. It is calculated using player, team and league statistics. The sum of all the players' WS in a given team will be approximately equal to that team's total wins for the season. A player with negative WS means that the player took away wins that the teammates had generated.

The reason for analyzing two variables is to investigate differences in career arcs for different aspects of skill. This allows us to highlight the power of our approach and could be of interest to basketball fans/analysts. Any other statistic can be chosen.

We have removed the observations with fewer than five games played. They were related to very extreme BPM values, such as -86.7 for Gheorghe Muresan in 1998-1999<sup>4</sup> or -49.3 for Mindaugas Kuzminskas in 2017-2018<sup>5</sup>.

Fig. 1 illustrates the type of data we are working with. It shows the observations of certain players, whose values are represented as connected points. Players' ages will represent the time points to be used by our methodology. The initial range of ages in the database went from 18 to 44 years old. However, there were only a few measurements between ages 41 and 44, related to a few long-lasting players, so we have removed them. The age range finally considered is from 18 to 40 years old, i.e., there are 23 time points. In total there are 3075 players.

### 4 | RESULTS

Section 4.1 contains a comparison analysis between ROPES, PACE and two benchmark methods. Section 4.2 discusses the implications of the archetypoid coefficients. Section 4.3 specifies the type of projections obtained. Section 4.4 presents an interactive web application.

#### 4.1 | Comparison with other methods

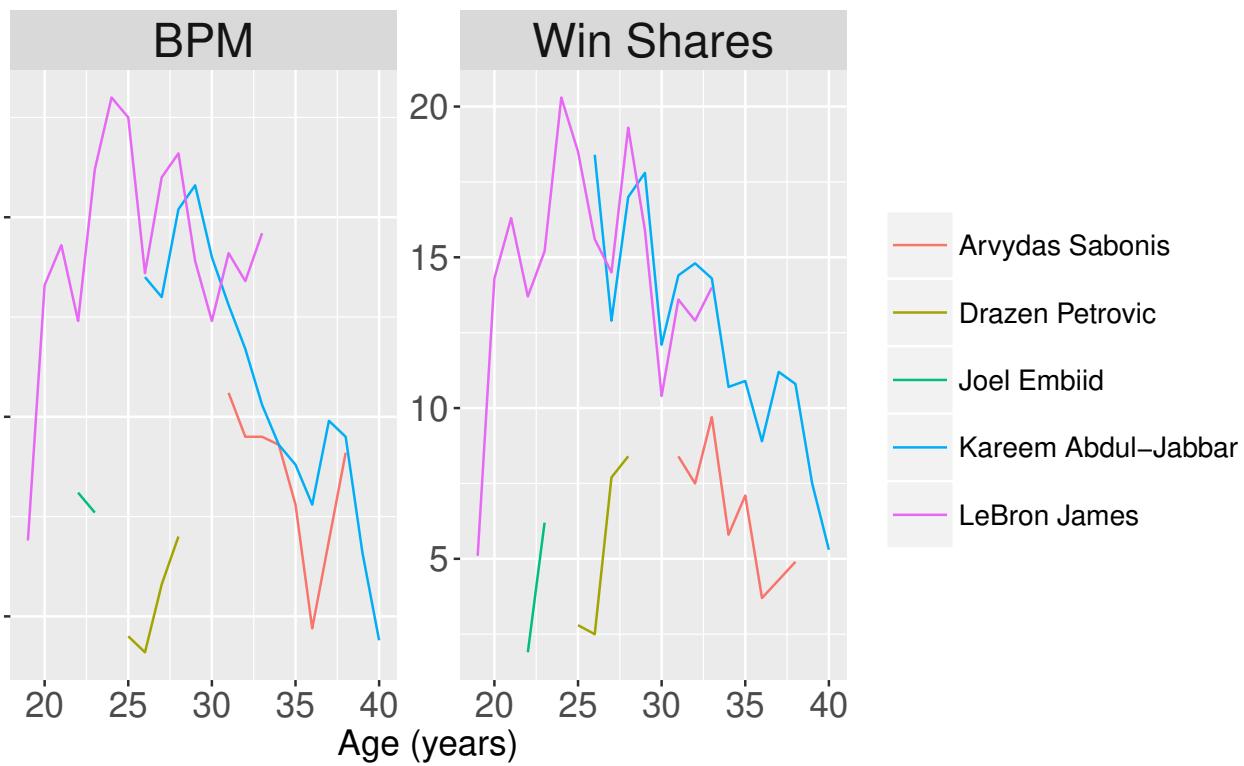
In order to evaluate the usefulness of ROPES and PACE, we carry out a comparison with each other and with two benchmark methods, such as the average method and the naïve method. In the average method, the forecast of the next value is the mean of the previous values. In the naïve option, the forecast is the value of the last observation. They are two common simple alternatives to more advanced techniques [19, Section 2.3].

In order to check the performance of all the methods, we have applied them to the test set of 385 players (see the validation study of the appendix to know how this set is created). Table 1 reports an extract of the results. It contains the following information for all players in the 2017-2018 season: (i) their age; (ii) their actual BPM value; (iii) the predictions with ROPES (using the optimal  $\lambda$  combination from the validation study of the appendix), PACE and the simple methods; (iv) the squared

<sup>3</sup><https://www.basketball-reference.com/about/bpm.html>

<sup>4</sup>[https://www.basketball-reference.com/players/m/muresgh01.html#all\\_advanced](https://www.basketball-reference.com/players/m/muresgh01.html#all_advanced)

<sup>5</sup>[https://www.basketball-reference.com/players/k/kuzmimi01.html#all\\_advanced](https://www.basketball-reference.com/players/k/kuzmimi01.html#all_advanced)

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

**FIGURE 1** : BPM and Win Shares values for some players present in the database, at their corresponding ages (colors in the online version). We have chosen these players because of two reasons: Firstly, we wanted to represent the values of well-known players. Sabonis and Petrovic are two of the best international players of all time. Abdul-Jabbar and James are two of the best players in history. Embiid is one of the most promising players nowadays. Secondly, we wanted to fill the entire range of ages, both with short and long careers.

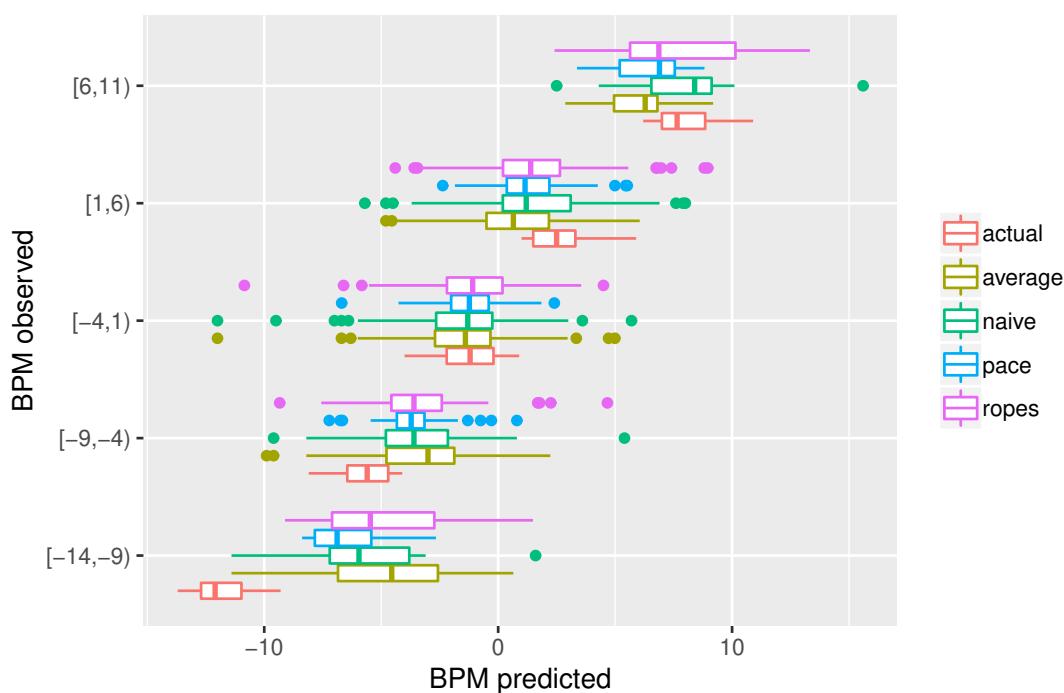
**TABLE 1** : Actual and predicted BPM values for the 2017-2018 season with ROPES, PACE, the average method and the naïve method, for the test set of players who played during the 2017-2018 season and at least one season before. The difference between actual and predicted values and MSE are also provided. MSE is highlighted in bold. Extract of the results.

Player	Age	BPM	ROPEs (900,10,10)		PACE		Average		Naïve	
			BPM <sub>pr</sub>	(BPM <sub>pr</sub> - BPM) <sup>2</sup>	BPM <sub>pr</sub>	(BPM <sub>pr</sub> - BPM) <sup>2</sup>	BPM <sub>pr</sub>	(BPM <sub>pr</sub> - BPM) <sup>2</sup>	BPM <sub>pr</sub>	(BPM <sub>pr</sub> - BPM) <sup>2</sup>
Aaron Brooks	33	-4.3	-4.14	0.03	-3.32	0.96	-2.21	4.37	-4.6	0.09
...	...	...	...	...	...	...	...	...	...	...
Josh Richardson	24	1.4	0.66	0.55	0.56	0.71	0.40	1.00	0.2	1.44
...	...	...	...	...	...	...	...	...	...	...
Zaza Pachulia	33	0.8	1.00	0.04	0.54	0.07	-0.75	2.40	2.7	3.61
Mean (MSE)		-0.92	-0.61	<b>6.73</b>	-0.87	<b>3.24</b>	-1.15	<b>7.59</b>	-0.92	<b>7.11</b>
Sd		3.31	3.01	16.28	2.35	7.6	2.72	15.41	3.22	16.53
Mean ± Sd		(2.39,4.23)	(2.4, 3.62)		(1.48,3.22)		(-3.87,1.57)		(-4.14,2.3)	

difference between predictions and actual values  $(\text{BPM}_{pr} - \text{BPM})^2$ ; (v) the resulting total MSE (highlighted in bold). PACE obtains the smallest MSE, followed by ROPES. It is interesting to note that the mean BPM obtained with the naïve method is practically the same as the actual one (both rounded to  $-0.92$ ). This result is most probably because over- and under-predictions cancel each other out.

Fig. 2 displays the boxplots for the actual BPM values together with the BPM predictions for each method in different intervals.

The intervals  $[-14, -9]$  and  $[-9, -4]$  refer to players with a very bad performance (according to the BPM scale). We see in both cases that the predictions are far away from the true values. All the methods give a conservative forecast for such extreme values. PACE is the method that provides the closest results in these two intervals. ROPES is close to PACE in  $[-9, -4]$ . In the



**FIGURE 2** : Boxplots for the actual BPM values together with the BPM predictions for each method in different intervals.

interval  $[-4, 1]$ , the four methods show similar values with respect to the actual ones. In the interval  $[1, 6)$ , ROPES gives the most similar predictions with respect to the actual ones. In the interval  $[6, 11)$ , again ROPES and PACE give the most accurate predictions. Remarkably, the naïve method shows outliers in all the intervals. Fig. 2 is showing that PACE and ROPES tend to overestimate bad players and underestimate good ones. This can be very helpful for team managers in their search for future stars (this additional interpretation was suggested by one of the referees).

Overall, PACE is the method that performs best. ROPES is able to beat the simple benchmark methods, showing an improvement with respect to them. The main drawback of the current PACE implementation is the lack of prediction intervals. The main goal of this paper is to draw attention to the added value of using an FDA approach to forecast players' performance, which has not been done so far. Therefore, even though PACE should give somewhat more accurate predictions than ROPES, in next Section we will use ROPES to forecast future players' activity because it does provide prediction intervals. Prediction intervals are very helpful and important because they express how much uncertainty is associated with the forecast.

## 4.2 | Implications of the archetypoid coefficients

We have analyzed a total of 8 players from different teams, namely Devin Booker, Clint Capela, Joel Embiid, Nikola Jokic, Tyus Jones, Zach LaVine, Donovan Mitchell and Jayson Tatum. They are representing several career status. Embiid, Jokic, Mitchell and Tatum are already established figures (especially Embiid and Jokic). Booker and Capela are a step below the super stars but they are also very good players. Something similar could be said about Tyus Jones, who is constantly improving his skills. Finally, LaVine is an offensive specialist.

In a first attempt to compute predictions using all the players of the data set, we realized that the ROPES method had some pull towards the mean of the entire sample (like the other methods discussed in Section 4.1 but not as strong as them). This gave unrealistic performance predictions for both the best and most promising players. Therefore, in order to refine predictions, it is much more suitable to use the so-called “method of analogues”. The idea is to find players related to the one of interest and then use their documented activity to obtain the predictions. We know how other players already performed, so we can use their information to gain an approximate idea about the future performances of others. The method of analogues has been used for years in fields such as climatology [46] and epidemiology [39]. Recently, an R package has been released that contains analogue methods for palaeoecology [36]. The CARMELO method is also based on this scheme.

8

Vinué ET AL

In order to find related players, we use archetypoid analysis (ADA, see [42] for theoretical details). ADA searches for extreme observations (the so-called archetypoids) to describe the frontiers of the data. In this technique, the BPM (and WS) function of a player is approximated by a mixture of archetypoids, which are themselves functions of boundary players (outstanding - positive or negative- performers). Archetypoids are specific players and the  $\alpha$  coefficients represent how much each archetypoid contributes to the approximation of each individual. The most comparable archetypoid should be the one corresponding to the largest value of the  $\alpha$  coefficients for the player of interest.

We choose the number of archetypoids for each metric following the screeplot explained by [42]. Five are selected for BPM and four for WS.

The archetypoids for the BPM metric are (their career BPM shown in brackets): Devin Gray (-8.4), Darryl Dawkins (-2.52), Diamond Stone (-24.1), Eddy Curry (-6.5) and LeBron James (9.21).

LeBron James is the representative of super star players. He is one of the best players in history. This is in line with the expected results since James has achieved the highest BPM values.

Darryl Dawkins represents the replacement level players (as a reminder, -2 is replacement level). Dawkins had a long NBA career. He was selected with the fifth pick in the 1975 NBA draft and played for 14 seasons, where he averaged double figures in scoring in nine of them. He lead the league in fouls committed in three seasons. In his case, his performance does not fit exactly with the replacement level description, but his average BPM does.

Eddy Curry, Devin Gray and Diamond Stone are representatives of players with a short-term career or with overall poor performance. Eddy Curry was selected fourth overall in the 2001 NBA draft and had a long NBA career. He led the NBA in field goal percentage in the 2002-2003 season but he did not really meet the expectations that his talent was indicating. Devin Gray had an irrelevant NBA career, playing a total of 27 games in two NBA seasons. Diamond Stone only played seven games in the NBA. From the basketball point of view, Devin Gray and Diamond Stone are related to the “bad” archetypical player profile. Both players played very little in the NBA. However, from the mathematical point of view, they are not exactly representing the same profile, since Devin Gray has a BPM of -8.4, while Diamond Stone has a BPM of -24.1 (the third worst BPM value in our database, the range of BPM values goes from -31.5 to 21). Therefore, Stone is representing even a much more extreme bad pattern than Gray. We acknowledge that the five archetypoids computed for BPM are not exactly representing all the players’ typologies available in the database. In some cases, a greater number of archetypoids is needed to capture other players’ profiles. Even though we have determined the optimal number of archetypoids with the screeplot, in a real situation it would be up to the analyst to decide how many representative cases to consider.

Regarding the WS metric, the archetypoids are (their career WS shown in brackets): Steve Burtt (0), Ben Wallace (5.84), Otis Birdsong (4.03) and LeBron James (14.6). Again, as expected, James is the representative of super star players. The fact that James is selected as the “best” archetypoid in both metrics is indicating how he can excel in many aspects of the game.

Otis Birdsong and Ben Wallace represent very good players. Otis Birdsong played twelve NBA seasons and appeared in four NBA All-Star Games. He was selected with the second pick of the 1977 NBA draft. Ben Wallace was very good at grabbing rebounds and blocking opponent shots. He won the NBA Defensive Player of the Year Award four times and won a championship with the Pistons in 2004.

Steve Burtt represents ordinary players. He played 101 games in four NBA seasons between 1984-1985 and 1992-1993.

Table 2 shows the  $\alpha$  values for the 8 players selected for the BPM and WS archetypoids.

As mentioned before, in ADA each datum is expressed as a mixture of actual observations (archetypoids). In particular, the  $\alpha$  coefficients of each player are of great utility because they allow us to determine the composition of each player according to the archetypoid players, and to establish a clustering of similar players [41]. A discussion of the implications of the resulting archetypoid coefficients is given next. The BPM composition of Embiid is as follows. Embiid’s profile matches 49% of James’, 23% of Gray’s, 15% of Dawkins’ and 13% of Curry’s, so this is reflecting the fact that Embiid is on his way to become a super star like James is, but he still has some room of improvement. Regarding the WS composition, Embiid’s profile is 42% explained by Burtt’s, 29% by Birdsong’s, 18% by Wallace’s and 11% by James’. In this case, the Embiid’s room for development is even more evident. The archetypoid coefficients for Jokic are quite impressive. His highest  $\alpha$  is with James in both variables. His BPM profile matches 68% of James’ BPM profile (the highest similarity to James in this analysis by far) and his WS profile matches 43% of James’ WS profile (only Tatum has a close value) and 35% of Birdsong’s. This high similarity with respect to James implies that Jokic’s performance is already very good. **Proof of his remarkable activity is that he has received his first All-Star and All-NBA First Team selections in 2018-19 season. In addition, his twelve triple-doubles ranked second on the season behind only Russell Westbrook (34).** Capela is another player worth highlighting. His highest BPM coefficient is also with James, though his profile still matches 33% of Stones’, which indicates some shortcomings in his performance. Regarding

**TABLE 2** : Similarity of the 8 players selected to the BPM and WS archetypoids according to the  $\alpha$  coefficients.

Player	BPM Archetypoids				
	D. Gray	D. Dawkins	D. Stone	E. Curry	L. James
D. Booker	0.32	0.12	0.07	0.13	0.36
C. Capela	0.00	0.08	0.33	0.10	0.49
J. Embiid	0.23	0.15	0.00	0.13	0.49
N. Jokic	0.06	0.19	0.00	0.07	0.68
T. Jones	0.31	0.14	0.07	0.12	0.36
Z. LaVine	0.39	0.15	0.05	0.11	0.30
D. Mitchell	0.37	0.12	0.00	0.12	0.39
J. Tatum	0.42	0.06	0.00	0.09	0.43
Player	WS Archetypoids				
	S. Burtt	B. Wallace	O. Birdsong	L. James	
D. Booker	0.70	0.12	0.08	0.10	
C. Capela	0.00	0.14	0.76	0.10	
J. Embiid	0.42	0.18	0.29	0.11	
N. Jokic	0.19	0.03	0.35	0.43	
T. Jones	0.58	0.12	0.24	0.06	
Z. LaVine	0.77	0.10	0.07	0.06	
D. Mitchell	0.59	0.08	0.12	0.21	
J. Tatum	0.59	0.06	0.00	0.35	

his WS composition, his profile matches 76% of Birdsong's, 14% of Wallace's and 10% of James', which indicates a remarkable team productivity. On the other hand, the BPM profiles for Booker, Jones, Tatum and Mitchell are as close to James' as to Gray's. In terms of James' coefficient, it is very difficult to approach his excellence, so we shouldn't expect Booker, Jones, Tatum and Mitchell to achieve such an incredible level. However, the Gray's coefficient is also indicating that their performance is not very far away from an average player yet, so they still need to make further efforts to display a difference from competitors. Their WS profiles reaffirm this claim. Finally, the LaVine's BPM and WS profiles are most similar to Gray's and Burtt's profiles, respectively, so he is not showing an outstanding performance at all.

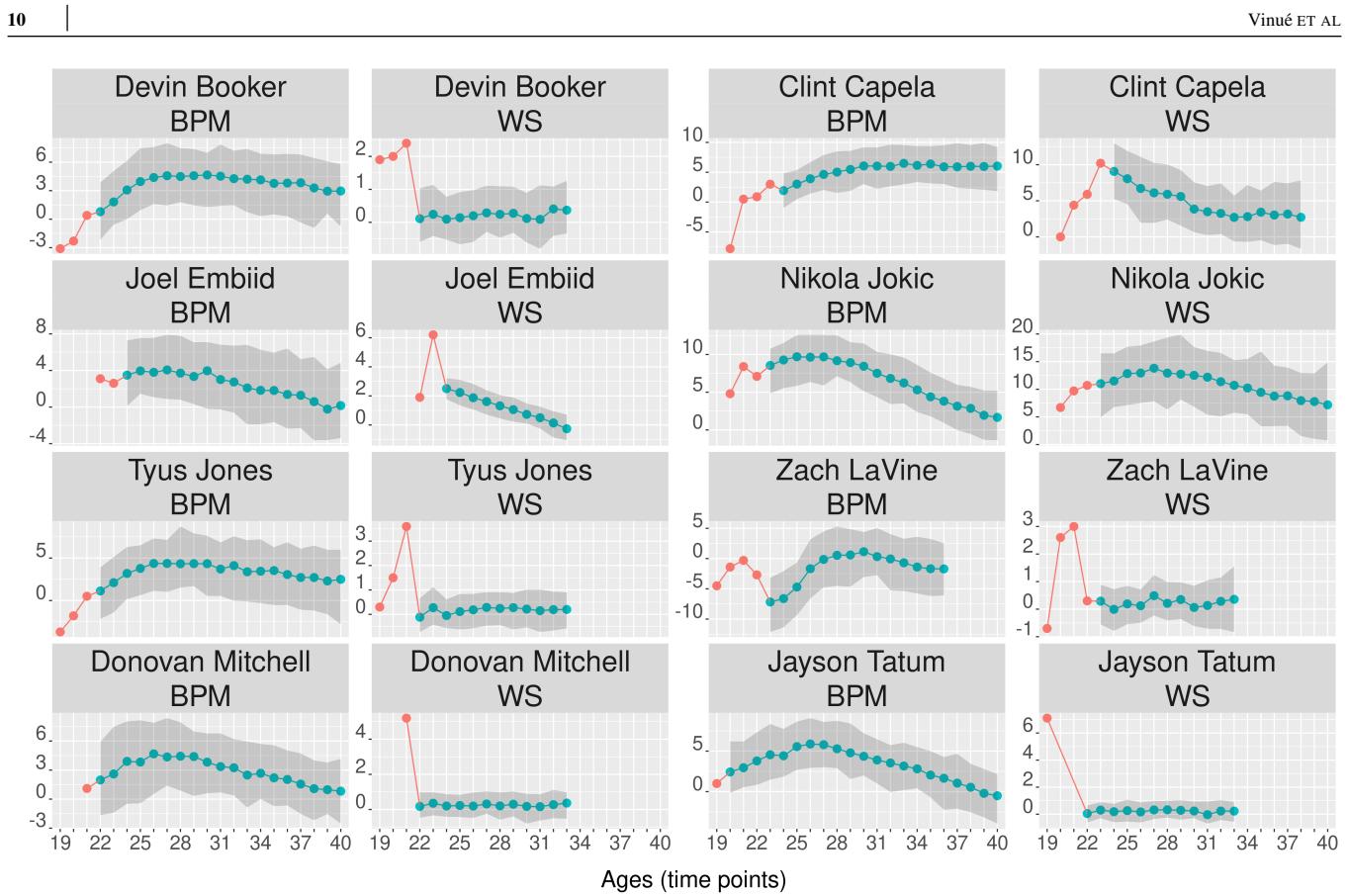
### 4.3 | Projections of future performance with ROPES and the method of analogues. Comparison with CARMELO

In order to select the cluster of analogous players, we first choose the archetypoid with the highest  $\alpha$ . As an illustration, Embiid's greatest similarity for BPM is with James. Then, the group of BPM similar players to Embiid is made up of James, together with the other players whose largest  $\alpha$  coefficient is also for James and who have an  $\alpha$  value greater than Embiid's  $\alpha$ . We will use this cluster to forecast the Embiid's future career arc. Current stars such as Chris Paul or Kevin Durant and stars of previous seasons such as Michael Jordan or Charles Barkley belong to this set. Likewise for WS.

The ROPES algorithm (with the lambda combination obtained in the validation study) is used to obtain  $p$ -forecasting intervals, where  $p = 0.05$  is the selected significance level. We will discuss the ROPES predictions with the ones that the CARMELO 2018-2019 version provides<sup>6</sup>. CARMELO is a basketball forecasting system released in the 2015-2016 season. Successive versions present some improvements [35]. To the best of our knowledge, it is the only publicly available projection system to compare our approach against. For each player of interest, CARMELO computes the similarity scores between that player and all historical players. To that end, it uses a number of statistics and players' attributes and a version of a nearest neighbor algorithm. The Wins Above Replacement (WAR) metric is computed for all historical players with a positive similarity score. The forecast is given by averaging these WAR values.

WAR reflects a combination of a player's projected playing time and his projected productivity while on the court. Productivity is measured by a blend of two-thirds Real Plus-Minus (RPM) and one-third Box Plus/Minus (BPM). BPM was solely used

<sup>6</sup><https://projects.fivethirtyeight.com/carmelo/>



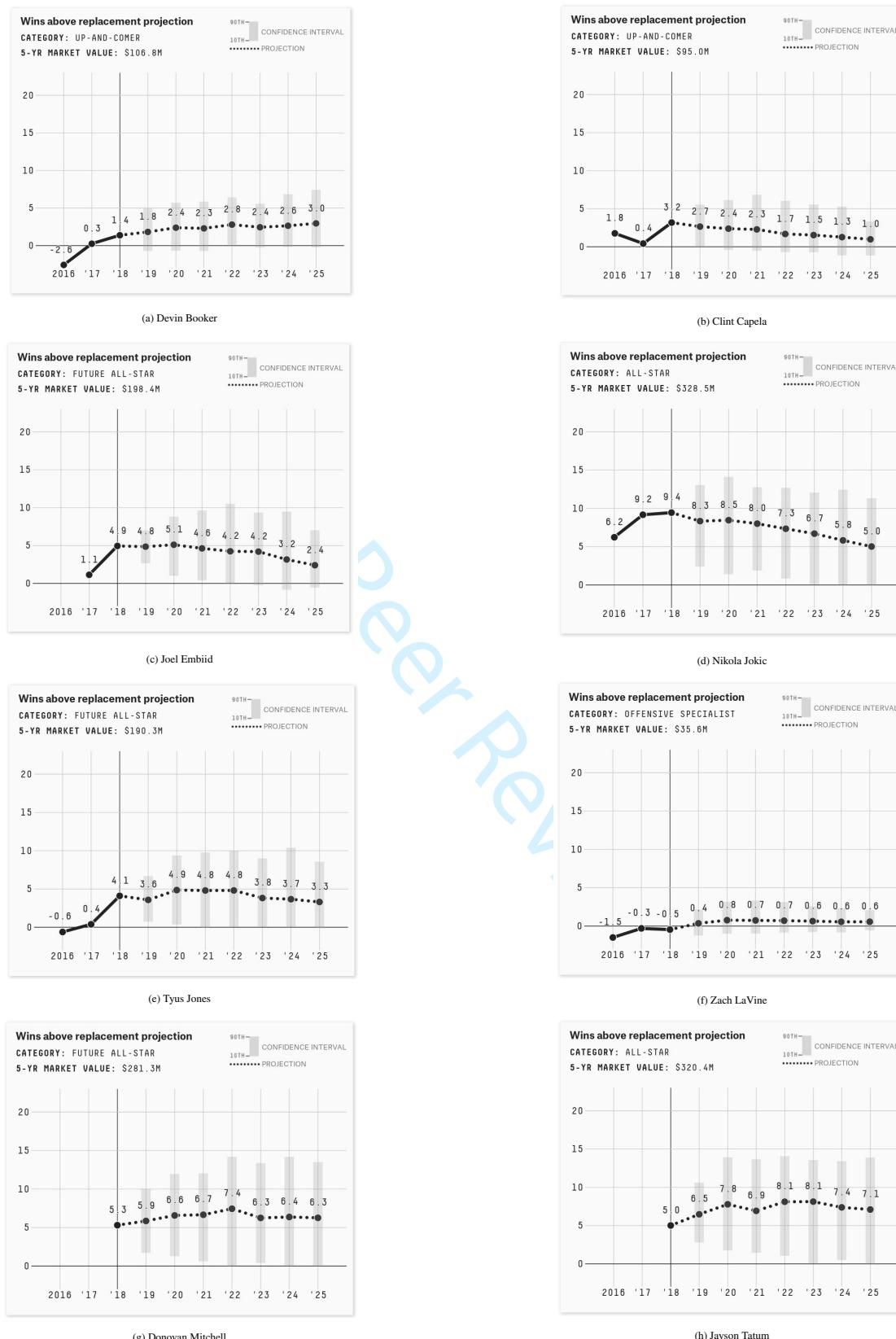
**FIGURE 3** : BPM and WS predictions for Devin Booker, Clint Capela, Joel Embiid, Nikola Jokic, Tyus Jones, Zach LaVine, Donovan Mitchell and Jayson Tatum using only the set of analogue players. Past values are in red and predictions are in green (colors in the online version). The shaded area indicates the limits of the prediction intervals.

to make the 2016-2017 forecasts, but the combination of RPM and BPM is used for the 2018-2019 forecasts (as in 2015-2016 and 2017-2018). According to the developers of CARMELLO, the RPM/BPM blend seems to outperform BPM alone. The RPM statistic quantifies how much a player hurts or helps his/her team when (s)he is on the court. There has been some controversy regarding the validity of RPM, since the computations are not detailed<sup>7</sup>. In fact, the CARMELLO methodology cannot be replicated either. In addition, for seasons before 2000-01, no RPM is available and CARMELLO uses BPM only. The final point that needs to be made is that RPM is not available in our database. Therefore, we would like to draw the reader's attention to the fact that our results are not directly equivalent to those of CARMELLO, since the target variable is not exactly the same. However, both approaches should be complementary.

Fig.3 shows the ROPES forecast obtained for the 8 players selected. In addition, for the sake of a convenient comparison with CARMELLO, Fig. 4 shows the screenshots of the CARMELLO curves for the same 8 players.

In Fig.3 we see that the predictions for Jokic show that his BPM and WS are expected to increase in the next three seasons and will remain quite high for several seasons (though they will go down from age 28). His lower and upper predictions are also high values. The width of the prediction intervals is constant over seasons. In general, the width of the intervals remain quite stable for most players. As a referee rightly mentioned, the intervals show several possible scenarios for some players, going from high to low values, so there is some uncertainty related to the point predictions. The CARMELLO forecast for Jokic indicates some decrease in his performance, but still keeping high numbers (Fig. 4 d). The width of the CARMELLO intervals fluctuates, the uncertainty in the '20 season is bigger than in the '19 and '21 seasons. The BPM-ROPS forecast for Embiid shows that he will improve his BPM in the two coming seasons and then his performance will slowly decline. Regarding WS, it indicates a constant decrease over time. CARMELLO also indicates that Embiid's performance will increase within two seasons and then his values

<sup>7</sup><https://www.boxscoregeeks.com/articles/rpm-and-a-problem-with-advanced-stats>



**FIGURE 4** : CARMELO curves for Devin Booker, Clint Capela, Joel Embiid, Nikola Jokic, Tyus Jones, Zach LaVine, Donovan Mitchell and Jayson Tatum.

12 |

Vinué ET AL

1 will decrease (the intervals move between narrow and wide intervals, Fig. 4 c). The BPM-ROPS forecast for Capela shows a  
2 constant increase in the coming years, keeping good values for many seasons. On the other hand, his WS values will constantly  
3 decrease. In this case, the intervals are widening over the years (especially in the BPM facet). The CARMELO forecast for  
4 Capela is a bit more conservative. Its prediction intervals are a bit wider in the '20 and '21 seasons than in the '19 season and then  
5 narrow a bit again (Fig. 4 b). The BPM-ROPS forecast for Devin Booker and Tyus Jones are somewhat similar to Capela's.  
6 However, their WS-ROPS forecast is a flat arc around 0. The CARMELO predictions for these two players show a certain  
7 increase in their activity (Fig. 4 a and Fig. 4 e). The width of the prediction interval for Booker remains quite constant. For  
8 Jones the width interval fluctuates between some wide and narrow ranges. The BPM-ROPS prediction for Donovan Mitchell  
9 and Jayson Tatum are similar to Jokic's, though their values are not so outstanding. Their WS prediction is a flat arc around  
10 0. The CARMELO predictions for Mitchell and Tatum also show a constant increase in their activity (Fig. 4 g and Fig. 4 h).  
11 In both cases, the width of the prediction interval also fluctuate between some wide and narrow ranges. Finally, the BPM-  
12 ROPES forecast for LaVine shows a constant increase but keeping negative numbers, especially for the next four years. His WS  
13 prediction moves around 0. CARMELO also suggests an ordinary and flat performance in the coming seasons (Fig. 4 f). It is  
14 worth mentioning that some of the WS predictions stop at age 33 because this is the last age for which the set of analogous  
15 players shows values. For Mitchell, Booker and Tatum, WS-ROPS has been too conservative. Higher values would have been  
16 a bit more realistic, since everything seems to indicate that these three players will become very good players in the near future.  
17 Another aspect that demands a careful examination of the results is that the WS prediction for Booker climbs to about 0.5 in the  
18 last two years. These are some pitfalls worth highlighting for end-users of this methodology.

19 As a final point, it is important to remember that statistical models are not completely reliable for long-term forecasts, because  
20 the assumption that the future looks similar to the past slowly breaks down the further we go into the future. So the predictions  
21 should be constantly updated as new data becomes available.

#### 25 26 4.4 | Web application

27 Additionally, an interactive web application available at <https://www.uv.es/vivogui/AppPredPerf.html> allows the user to represent  
28 the BPM and WS forecasting plots for every player in the 2017-2018 season under the age of 24 (154 players in total). A link to  
29 the CARMELO forecast for every player is also provided for easy comparison. The app gives some basic information about the  
30 way it works. It can also be generated from R with these two commands:

31 `library(shiny) ; runUrl('http://www.uv.es/vivogui/softw/AppPredPerf.zip')`

#### 35 36 5 | CONCLUSIONS

37 Basketball, like any other sport, contains a lot of uncertainty. A central issue is to predict future players' performance using past  
38 observations. In spite of the fact that basketball data continues to expand and there is a constant demand for new techniques  
39 that provide objective information to help understand the game, there are not many publicly available projection systems. In  
40 this paper we have presented a methodology to deal with sparse functional data in order to forecast the basketball players'  
41 performance. This has been done by analyzing ROPES and PACE and by including the method of analogues together with  
42 functional archetypoid analysis.

43 ROPES depends on several parameters, so we have carried out a validation study to choose an optimal combination that  
44 provides smooth curves and avoids overfitting (included in the appendix). The combination obtained works well to avoid narrow  
45 intervals and overconfident inferences. A comparison study has also been carried out to compare ROPES with PACE, and with  
46 simple alternatives, such as the average and naïve methods. PACE performed best overall and also in terms of runtime with  
47 respect to ROPES. However, unlike ROPES, it is not possible to obtain prediction intervals with its current computational  
48 implementation. In addition, ROPES also performed better than simple methods. Therefore, we have applied ROPES in the real  
49 case using data between 1973-1974 and 2017-2018 NBA regular seasons.

50 In the sparse case, information from all functions is used to fit each function, so all individuals contribute to a greater or  
51 lesser degree to form the estimations. In order to overcome this problem and to refine the predictions, we have used the so-  
52 called "method of analogues". The idea is to relate a player's curve to one of the possible types of players and then to predict  
53 his performance using only the information about these comparable athletes. In our case, the types of players are given by the  
54 archetypoids of the data set.

Once the computations are finished, an interactive web application shows the plots with the past and future behavior of 2017-2018 NBA players under the age of 24. Two variables have been analyzed: on the one hand, BPM is recognized as the most suitable metric to carry out an analysis involving historical data; on the other hand, WS is another widely-used advanced metric. Adding a second variable allows us to examine differences in career arcs for different aspects of skill. Any other variable can be used. The predictions for 8 players have been presented and a comparison with CARMELO has been done. The implications of the archetypoid coefficients have also been interpreted.

Player forecasting systems are important as a means of summarizing the overall match performance of individual players. Any forecasting method is limited because some aspects such as injury risk or work ethic, which influence future performance, are very difficult to quantify. However, coaches and experts can use these systems to review performances of their own players as well as tracking the performance levels of potential acquisitions. We hope that the approach presented here will provide valuable information about players' overall ability to support decision making. Sparse functional data are very common in sports. Therefore, it is very reasonable to bring methods developed to deal with this kind of data to the field of sports. This methodology can serve as a starting point for further efforts in the same direction. One of the referees suggested us to remark the following two situations that our analysis has not considered: (i) the different amounts of playing time going into each averaged BPM and WS data points. In mathematical terms, this is a case of unequal variances, also called heteroscedasticity; (ii) the pattern of sparsity in the data is not random, since players retiring or leaving the NBA should have low BPM and WS values in these age intervals. Both situations were formulated by the referee. We will consider them in future work. **Following another referee's suggestion, we will also try to compare the players' forecasts using relative rankings in terms of their coefficients from a common archetypoid.** The data and all R code are freely available at <https://www.uv.es/vivogui/software> for reproducibility and further exploration of the results.

## ACKNOWLEDGMENTS

The authors wish to express their gratitude to Alexander Dokumentov and Rob Hyndman for kindly providing the R code to run the ROPES algorithm. The authors are also grateful to the anonymous referees and associated editor for their helpful criticism. G. Vinué worked on the first two revisions of the manuscript as a postdoctoral scholarship holder in international mobility at KU Leuven, Department of Computer Science, Belgium. Funding: This work has been partially supported by the Spanish Ministerio de Ciencia, Innovación y Universidades (AEI/FEDER, UE) Grant DPI2017-87333-R and Universitat Jaume I, UJI-B2017-13.

## Author contributions

G. Vinué planned the design of the research, collected the data, performed the computations, interpreted the results and wrote the manuscript. I. Epifanio planned the design of the research, performed the computations and wrote the manuscript.

## Financial disclosure

None reported.

## Conflict of interest

The authors declare no potential conflict of interests.

## References

- [1] Arndt, C., Brefeld, U., 2016. Predicting the future performance of soccer players. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 9, 373–382, <http://dx.doi.org/10.1002/sam.11321>.
- [2] Aue, A., Dubart Norinho, D., Hörmann, S., 2015. On the Prediction of Stationary Functional Time Series. *Journal of the American Statistical Association* 110 (509), 378–392, <http://dx.doi.org/10.1080/01621459.2014.909317>.

14

Viué ET AL

- [3] Cattelan, M., Varin, C., Firth, D., 2013. Dynamic Bradley-Terry modelling of sports tournaments. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 62 (1), 135–150, <http://dx.doi.org/10.1111/j.1467-9876.2012.01046.x>.
- [4] Cervone, D., D'Amour, A., Bornn, L., Goldsberry, K., 2016. A Multiresolution Stochastic Process Model for Predicting Basketball Possession Outcomes. *Journal of the American Statistical Association* 111 (514), 585–599, <http://dx.doi.org/10.1080/01621459.2016.1141685>.
- [5] Chang, W., Cheng, J., Allaire, J., Xie, Y., McPherson, J., 2015. shiny: Web Application Framework for R. R package version 0.12.2, <https://CRAN.R-project.org/package=shiny>.
- [6] Cutler, A., Breiman, L., 1994. Archetypal Analysis. *Technometrics* 36 (4), 338–347, <http://dx.doi.org/10.2307/1269949>.
- [7] Dai, X., Hadjipantelis, P., Ji, H., Mueller, H.-G., Wang, J.-L., 2016. fdaPACE: Functional Data Analysis and Empirical Dynamics. R package version 0.2.5, <https://CRAN.R-project.org/package=fdaPACE>.
- [8] D'Esposito, M. R., Palumbo, F., Ragozini, G., 2012. Interval Archetypes: A New Tool for Interval Data Analysis. *Statistical Analysis and Data Mining* 5 (4), 322–335, <http://dx.doi.org/10.1002/sam.11140>.
- [9] Di Battista, T., Fortuna, F., 2017. Functional confidence bands for lichen biodiversity profiles: A case study in Tuscany region (central Italy). *Statistical Analysis and Data Mining: The ASA Data Science Journal* 10 (1), 21–28, <https://doi.org/10.1002/sam.11334>.
- [10] Dokumentov, A., 2016. Smoothing, decomposition and forecasting of multidimensional and functional time series using regularisation. Ph.D. thesis, Monash University. Faculty of Business and Economics. Econometrics and Business Statistics, <http://arrow.monash.edu.au/vital/access/manager/Repository/monash:165926>.
- [11] Dokumentov, A., Hyndman, R. J., 2016. Low-dimensional decomposition, smoothing and forecasting of sparse functional data, <http://robjhyndman.com/papers/ROPES.pdf>. Working paper, 1-31.
- [12] Elmore, R., 2018. ballr: Access to Current and Historical Basketball Data. R package version 0.1.1. <https://CRAN.R-project.org/package=ballr>.
- [13] Epifanio, I., 2016. Functional archetype and archetypoid analysis. *Computational Statistics & Data Analysis* 104, 24–34, <http://dx.doi.org/10.1016/j.csda.2016.06.007>.
- [14] Epifanio, I., Ávila, C., Page, Á., Atienza, C., 2008. Analysis of multiple waveforms by means of functional principal component analysis: normal versus pathological patterns in sit-to-stand movement. *Medical & Biological Engineering & Computing* 46 (6), 551–561, <http://dx.doi.org/10.1007/s11517-008-0339-6>.
- [15] Ferraty, F., Vieu, P., 2006. Nonparametric Functional Data Analysis: Theory and Practice. Springer.
- [16] Harrison, A. J., 2014. Applications of functional data analysis in sports biomechanics. In: 32 International Conference of Biomechanics in Sports. pp. 1–9.
- [17] Hollinger, J., 2005. Pro basketball forecast. Potomac Books, Inc., Washington, D.C.
- [18] Hwang, D., 2012. Forecasting NBA player performance using a Weibull-Gamma statistical timing model. In: MIT Sloan Sports Analytics Conference. Boston, MA, USA, pp. 1–10.
- [19] Hyndman, R. J., Athanasopoulos, G., 2013. Forecasting: Principles and Practice. OTexts, <https://www.otexts.org/book/fpp>.
- [20] Hyndman, R. J., Shahid Ullah, M., 2007. Robust forecasting of mortality and fertility rates: A functional data approach. *Computational Statistics & Data Analysis* 51 (10), 4942–4956, <http://dx.doi.org/10.1016/j.csda.2006.07.028>.
- [21] James, G., 2010. The Oxford handbook of functional data analysis. Oxford University Press, Ch. Sparseness and functional data analysis, pp. 298–326.
- [22] Kubatko, J., Oliver, D., Pelton, K., Rosenbaum, D. T., 2007. A Starting Point for Analyzing Basketball Statistics. *Journal of Quantitative Analysis in Sports* 3 (3), 1–10, <http://dx.doi.org/10.2202/1559-0410.1070>.

- [23] Nguyen, H. D., Ullmann, J. F. P., McLachlan, G. J., Voleti, V., Li, W., Hillman, E. M. C., Reutens, D. C., Janke, A. L., 2018. Whole-volume clustering of time series data from zebrafish brain calcium images via mixture modeling. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 11 (1), 5–16, <https://doi.org/10.1002/sam.11366>.
- [24] Oliver, D., 2004. Basketball on paper: Rules and tools for performance analysis. Potomac Books, Inc., Washington, D.C.
- [25] R Core Team, 2019. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, <https://www.R-project.org/>.
- [26] Ragozini, G., Palumbo, F., D'Esposito, M. R., 2017. Archetypal analysis for data-driven prototype identification. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 10 (1), 6–20, <http://dx.doi.org/10.1002/sam.11325>.
- [27] Ramsay, J. O., Silverman, B., 2002. Applied Functional Data Analysis. Methods and Case Studies. Springer.
- [28] Ramsay, J. O., Silverman, B., 2005. Functional Data Analysis, 2nd Edition. Springer.
- [29] Ramsay, J. O., Hooker, G., Graves, S. 2009. Functional Data Analysis with R and MATLAB. Springer.
- [30] Salador, K., 2011. Forecasting Performance of International Players in the NBA. In: MIT Sloan Sports Analytics Conference. Boston, MA, USA, pp. 1–18, <http://www.sloansportsconference.com/wp-content/uploads/2011/08/Forecasting-Performance-of-International-Players-in-the-NBA.pdf>.
- [31] Shang, H. L., Hyndman, R. J., 2017. Grouped Functional Time Series Forecasting: An Application to Age-Specific Mortality Rates. *Journal of Computational and Graphical Statistics* 26 (2), 330–343, <http://dx.doi.org/10.1080/10618600.2016.1237877>.
- [32] Shea, S. M., 2014. Basketball analytics: Spatial tracking. Createspace, Lake St. Louis, MO.
- [33] Shea, S. M., Baker, C. E., 2013. <http://www.basketballanalyticsbook.com/>.
- [34] Shea, S. M., Baker, C. E., 2013. Basketball analytics: Objective and efficient strategies for understanding how teams win. Advanced Metrics, LLC, Lake St. Louis, MO.
- [35] Silver, N., 2018. CARMELO NBA player projections. <https://fivethirtyeight.com/features/our-nba-player-projections-are-ready-for-2018-19/>, <https://fivethirtyeight.com/features/whats-new-in-our-nba-player-projections-for-2017-18/>, <https://fivethirtyeight.com/features/whats-new-in-our-nba-projections-for-2016-17/>, <https://fivethirtyeight.com/features/how-were-predicting-nba-player-career/>, <https://projects.fivethirtyeight.com/carmelo/>.
- [36] Simpson, G., Oksanen, J., 2016. analogue: Analogue and Weighted Averaging Methods for Palaeoecology. R package version 0.17-0, <https://CRAN.R-project.org/package=analogue>.
- [37] Stekler, H. O., Vaughan Williams, L. (Editors), 2010. Sports forecasting (Special issue). *International Journal of Forecasting* 26 (3), 1–3, <http://dx.doi.org/10.1016/j.ijforecast.2009.12.005>.
- [38] Strumbelj, E., Vračar, P., 2012. Simulating a basketball match with a homogeneous Markov model and forecasting the outcome. *International Journal of Forecasting* 28 (2), 532–542, <http://dx.doi.org/10.1016/j.ijforecast.2011.01.004>.
- [39] Viboud, C., Boelle, P.-Y., Carrat, F., Valleron, A.-J., Flahault, A., 2003. Prediction of the Spread of Influenza Epidemics by the Method of Analogues. *American Journal of Epidemiology* 158 (10), 996–1006, <https://doi.org/10.1093/aje/kwg239>.
- [40] Vinué, G., 2017. Anthropometry: An R Package for Analysis of Anthropometric Data. *Journal of Statistical Software* 77 (6), 1–39, <https://doi.org/10.18637/jss.v077.i06>.
- [41] Vinué, G., Epifanio, I., 2017. Archetypoid analysis for sports analytics. *Data Mining and Knowledge Discovery*, 1–35, <https://doi.org/10.1007/s10618-017-0514-1>.
- [42] Vinué, G., Epifanio, I., Alemany, S., 2015. Archetypoids: A new approach to define representative archetypal data. *Computational Statistics and Data Analysis* 87, 102–115, <http://dx.doi.org/10.1016/j.csda.2015.01.018>.

16

Vinu  ET AL

- 1 [43] Vra ar, P., Strumbelj, E., Kononenko, I., 2016. Modeling basketball play-by-play data. *Expert Systems with Applications* 44, 58 66, <http://dx.doi.org/10.1016/j.eswa.2015.09.004>.
- 2 [44] Wakim, A., Jin, J., 2014. Functional Data Analysis of Aging Curves in Sports, <http://arxiv.org/abs/1403.7548>, 1-25.
- 3 [45] Yao, F., M ller, H.-G., Wang, J.-L., 2005. Functional Data Analysis for Sparse Longitudinal Data. *Journal of the American*  
4 Statistical Association 100 (470), 577 590, <http://dx.doi.org/10.1198/016214504000001745>.
- 5 [46] Zorita, E., Von Storch, H., 1999. The Analog Method as a Simple Statistical Downscaling Technique: Comparison  
6 with More Complicated Methods. *Journal of Climate* 12, 2474 2489, [http://dx.doi.org/10.1175/1520-0442\(1999\)012<2474:TAMAAS>2.0.CO;2](http://dx.doi.org/10.1175/1520-0442(1999)012<2474:TAMAAS>2.0.CO;2).
- 7 [47] Zimmermann, A., 2016. Basketball predictions in the NCAAB and NBA: Similarities and differences. *Statistical Analysis*  
8 and Data Mining: The ASA Data Science Journal 9, 350 364, <http://dx.doi.org/10.1002/sam.11319>.

## AUTHOR BIOGRAPHY

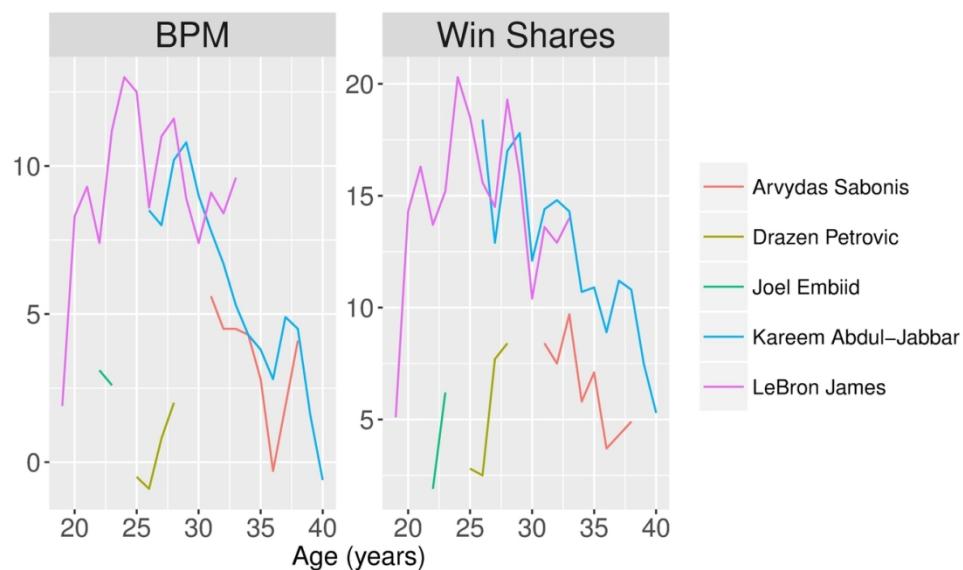


**G. Vinu .** Guillermo Vinu  was born in Valencia (Spain) in 1985. He graduated in Mathematics (2008) and completed his PhD in Statistics and Optimization (2014), both at the University of Valencia. His research interests are concerned on computational statistics with applications in industry and sports analytics (<http://orcid.org/0000-0003-2083-8276>).

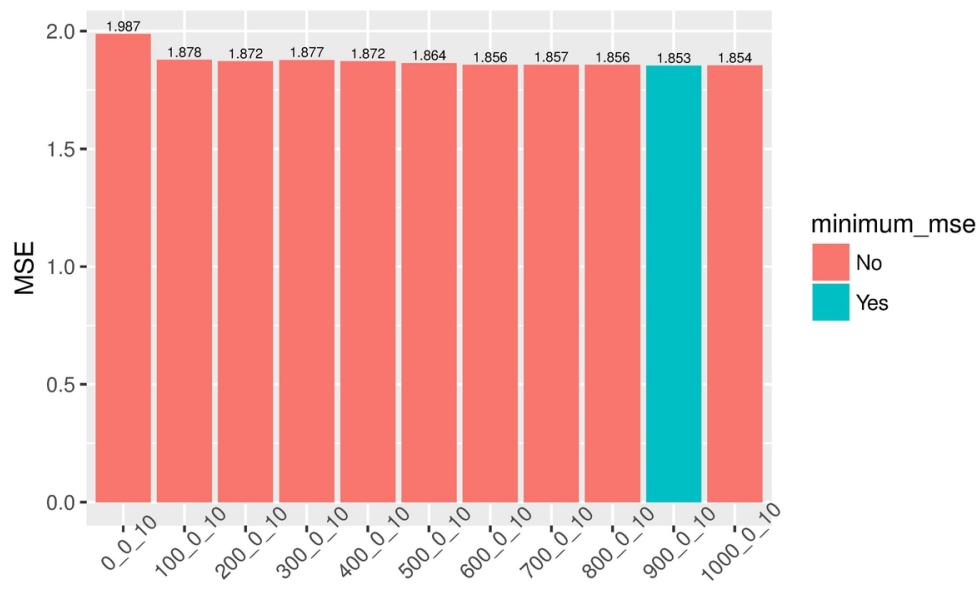


**I. Epifanio.** Irene Epifanio was born in Valencia (Spain) in 1975. She graduated in Mathematics (1997) and received his PhD in Statistics (2002), both at the University of Valencia. She is currently a Titular Professor at the Department of Mathematics at the Universitat Jaume I (Spain). Her current research interests include statistical learning, image analysis and functional data analysis (<http://orcid.org/0000-0002-6973-311X>).

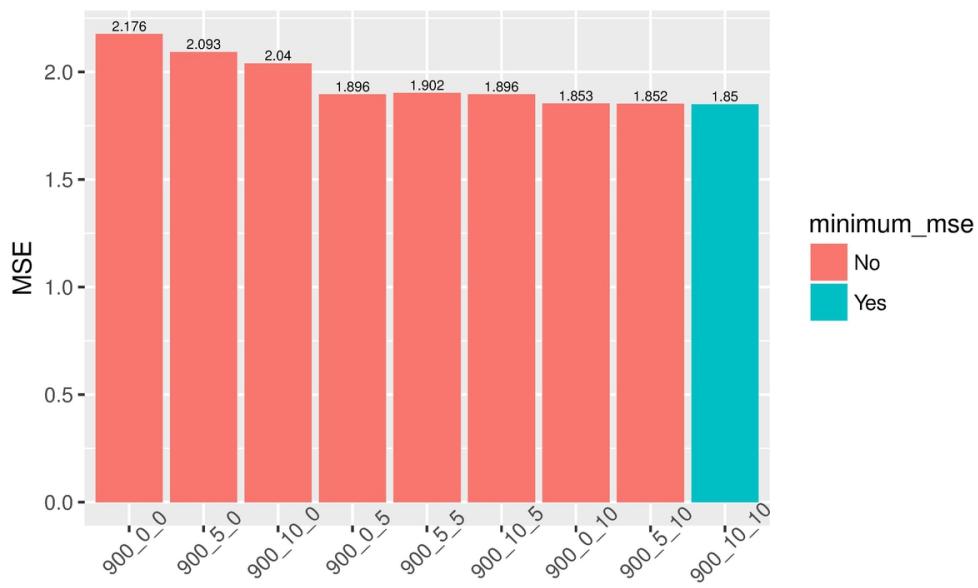
**How to cite this article:** Vinu  G., and I. Epifanio (2019), Forecasting basketball players' performance using sparse functional data, *Statistical Analysis and Data Mining*, 2019;:-.



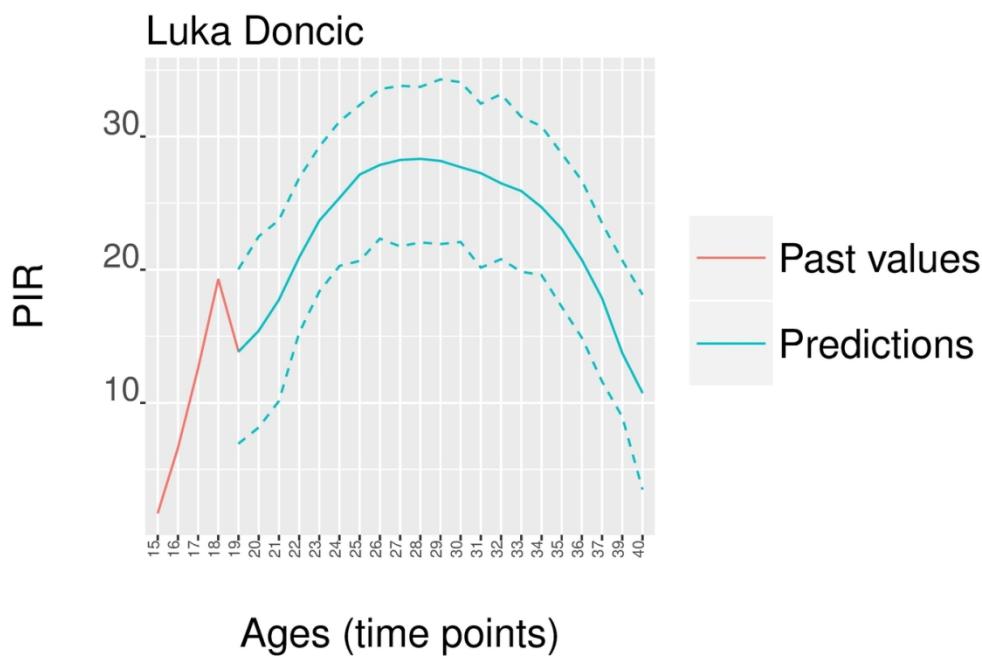
116x67mm (300 x 300 DPI)

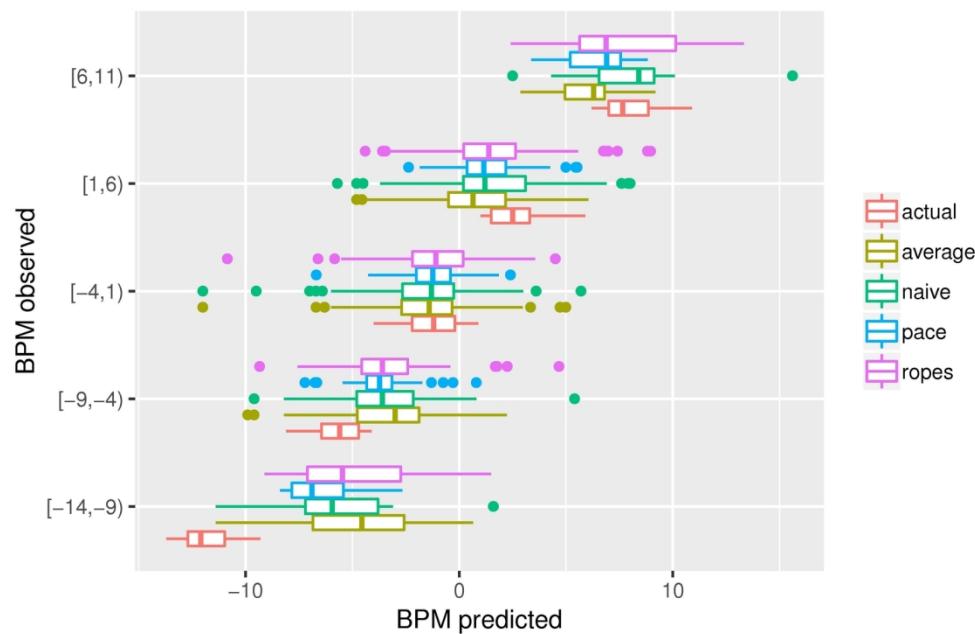


101x67mm (300 x 300 DPI)

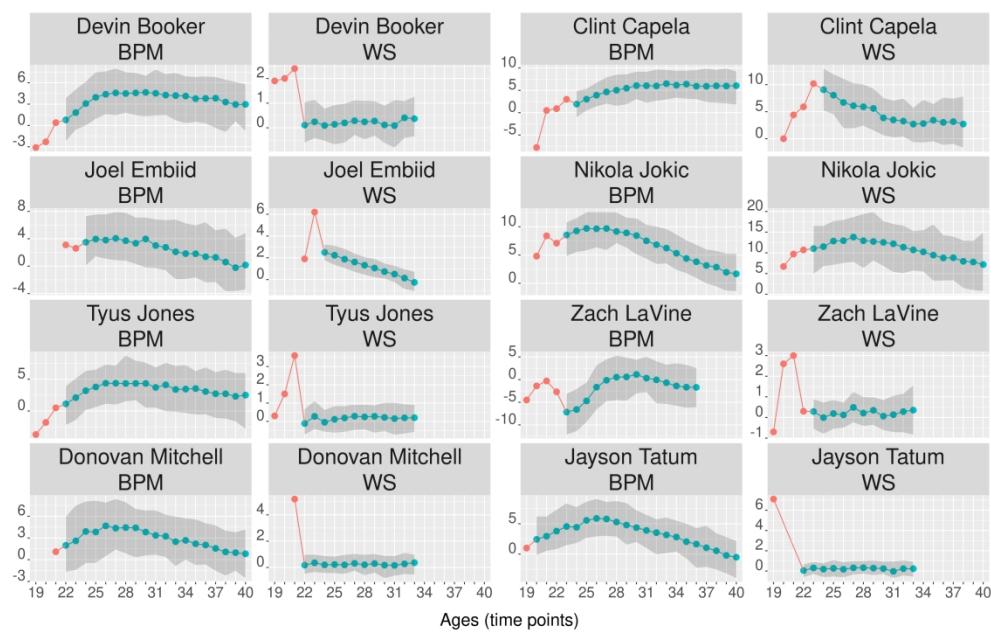


101x67mm (300 x 300 DPI)





152x98mm (300 x 300 DPI)



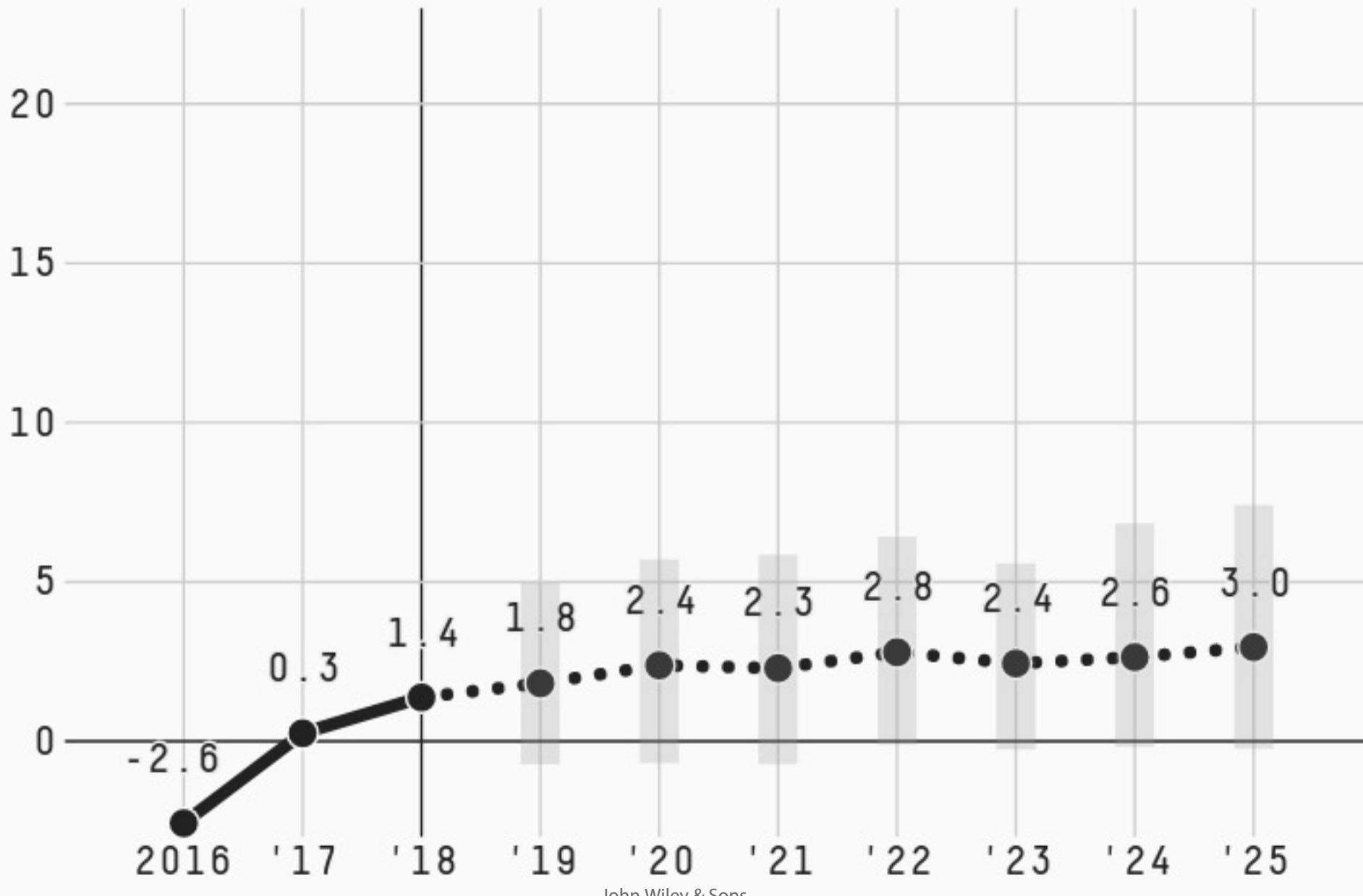
406x254mm (300 x 300 DPI)

# Wins above replacement projection

CATEGORY: UP-AND-COMER

5-YR MARKET VALUE: \$106.8M

90TH-  
CONFIDENCE INTERVAL  
10TH-  
PROJECTION

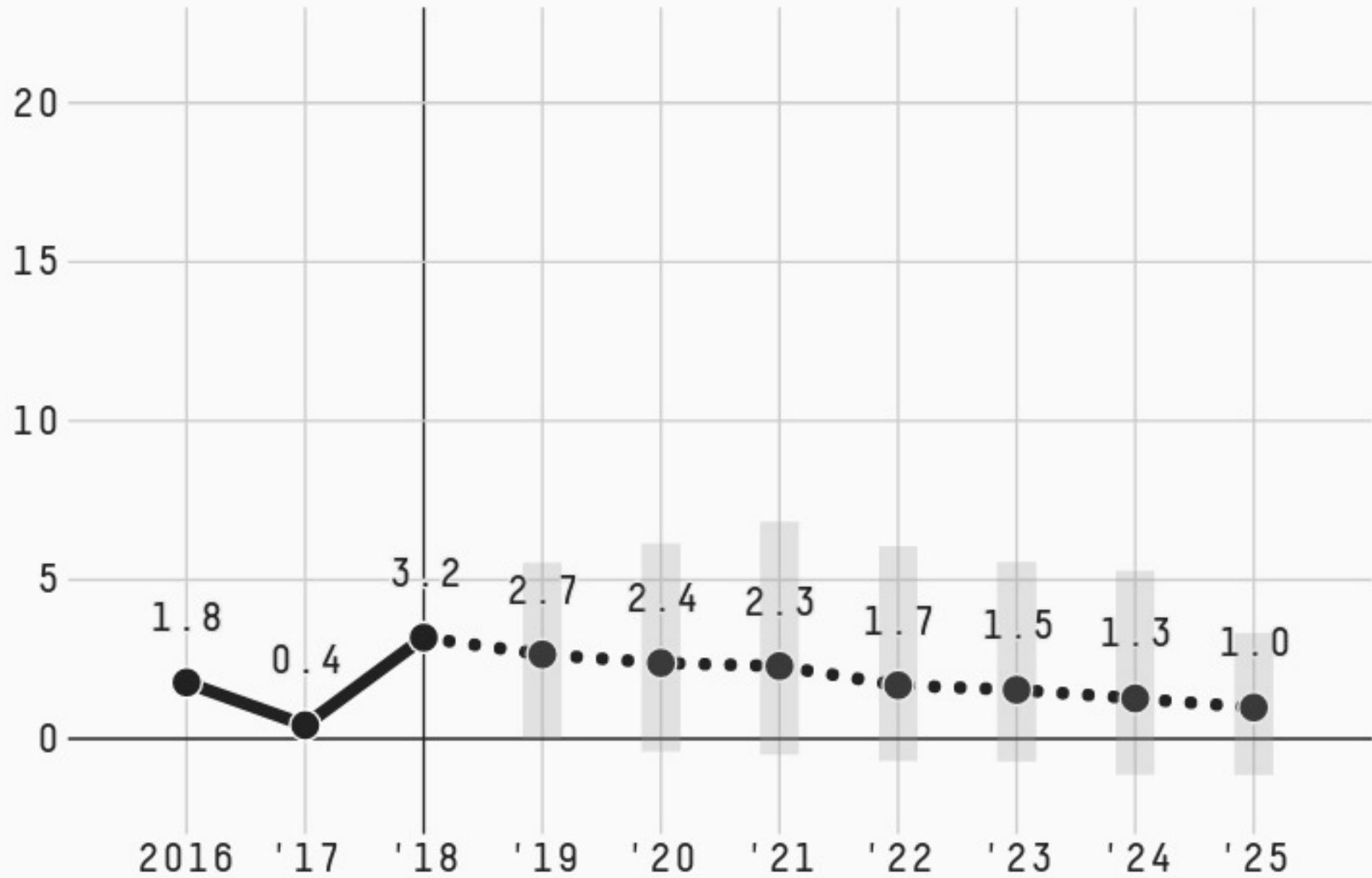


# Wins above replacement projection

CATEGORY: UP-AND-COMER

5-YR MARKET VALUE: \$95.0M

90TH—  
10TH—  
CONFIDENCE INTERVAL  
PROJECTION

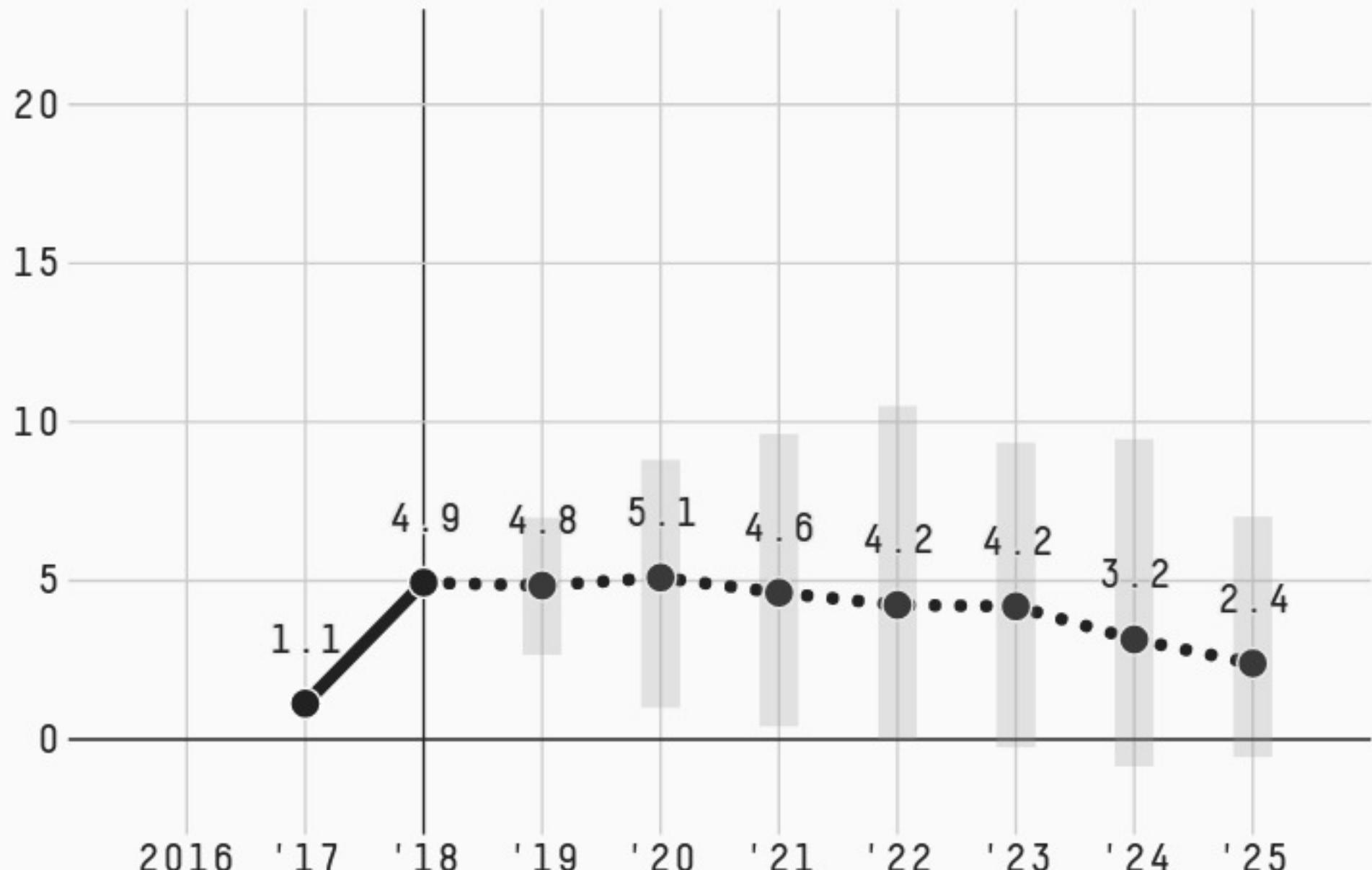


# Wins above replacement projection

CATEGORY: FUTURE ALL-STAR

5-YR MARKET VALUE: \$198.4M

90TH-  
10TH-  
CONFIDENCE INTERVAL  
\*\*\*\*\* PROJECTION

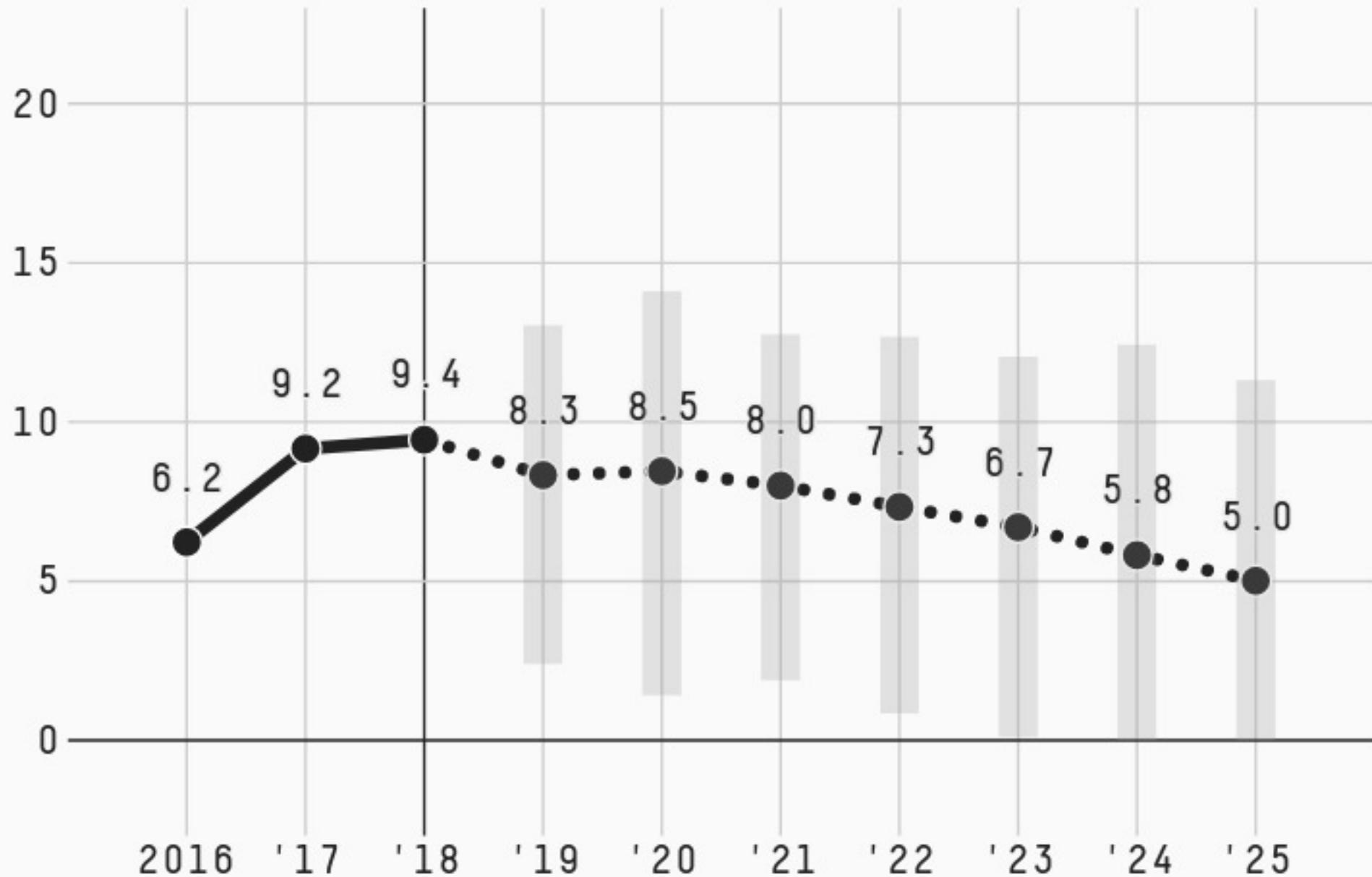


# Wins above replacement projection

CATEGORY: ALL-STAR

5-YR MARKET VALUE: \$328.5M

90TH-  
CONFIDENCE INTERVAL  
10TH-  
\*\*\*\*\* PROJECTION

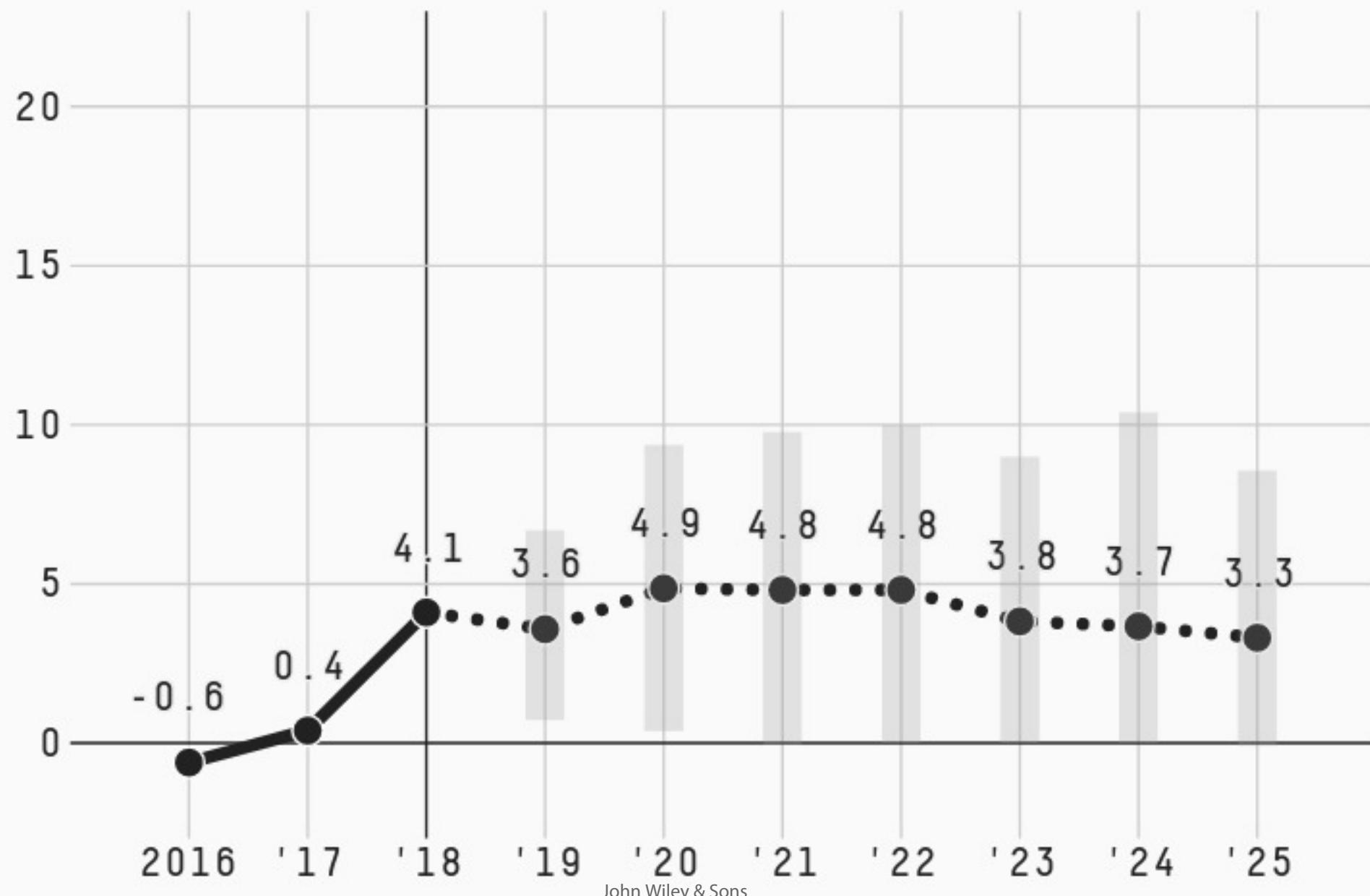


# Wins above replacement projection

CATEGORY: FUTURE ALL-STAR

5-YR MARKET VALUE: \$190.3M

90TH-  
10TH-  
CONFIDENCE INTERVAL  
PROJECTION

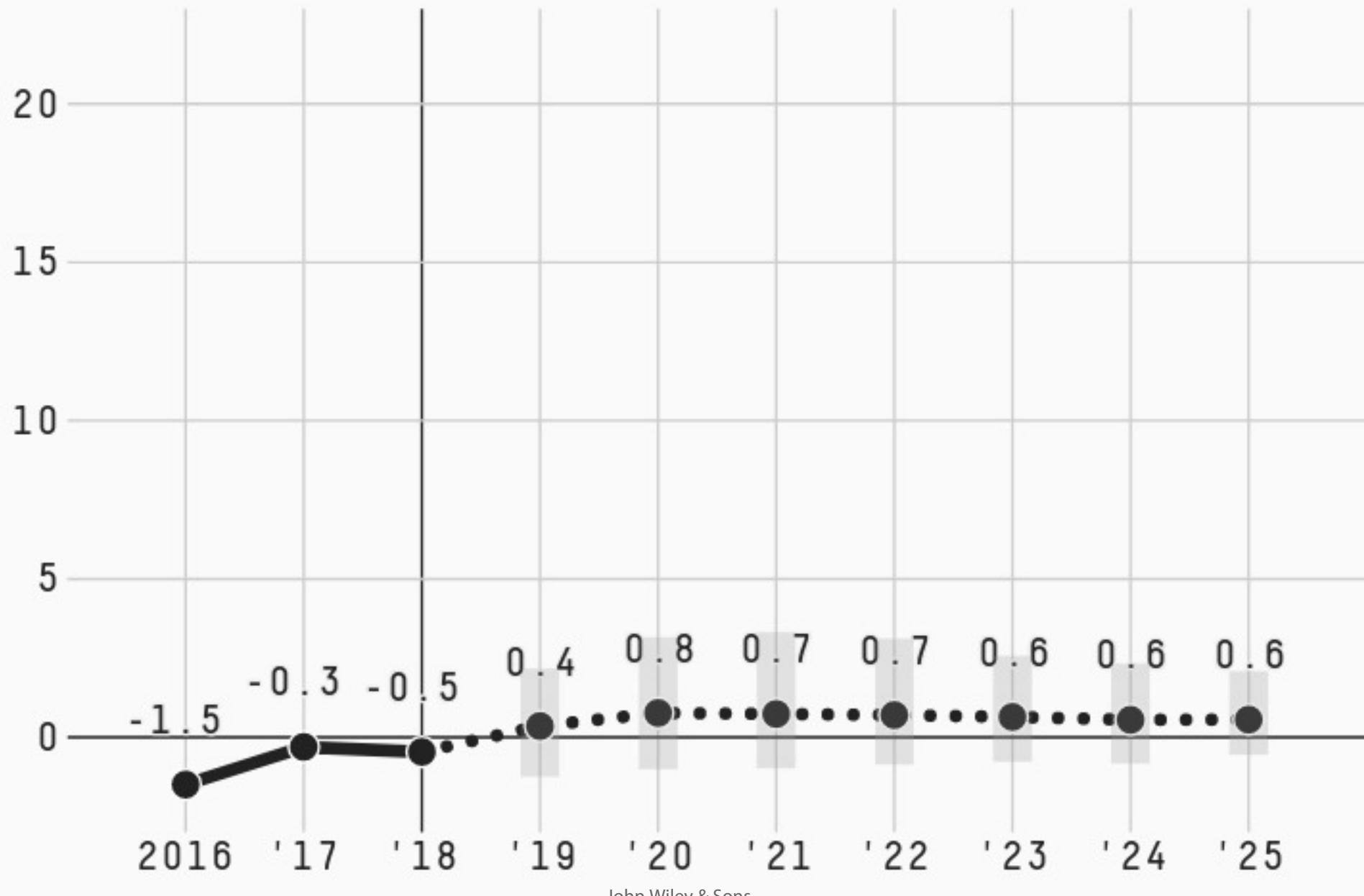


# Wins above replacement projection

CATEGORY: OFFENSIVE SPECIALIST

5-YR MARKET VALUE: \$35.6M

90TH-  
10TH-  
CONFIDENCE INTERVAL  
\*\*\*\*\* PROJECTION

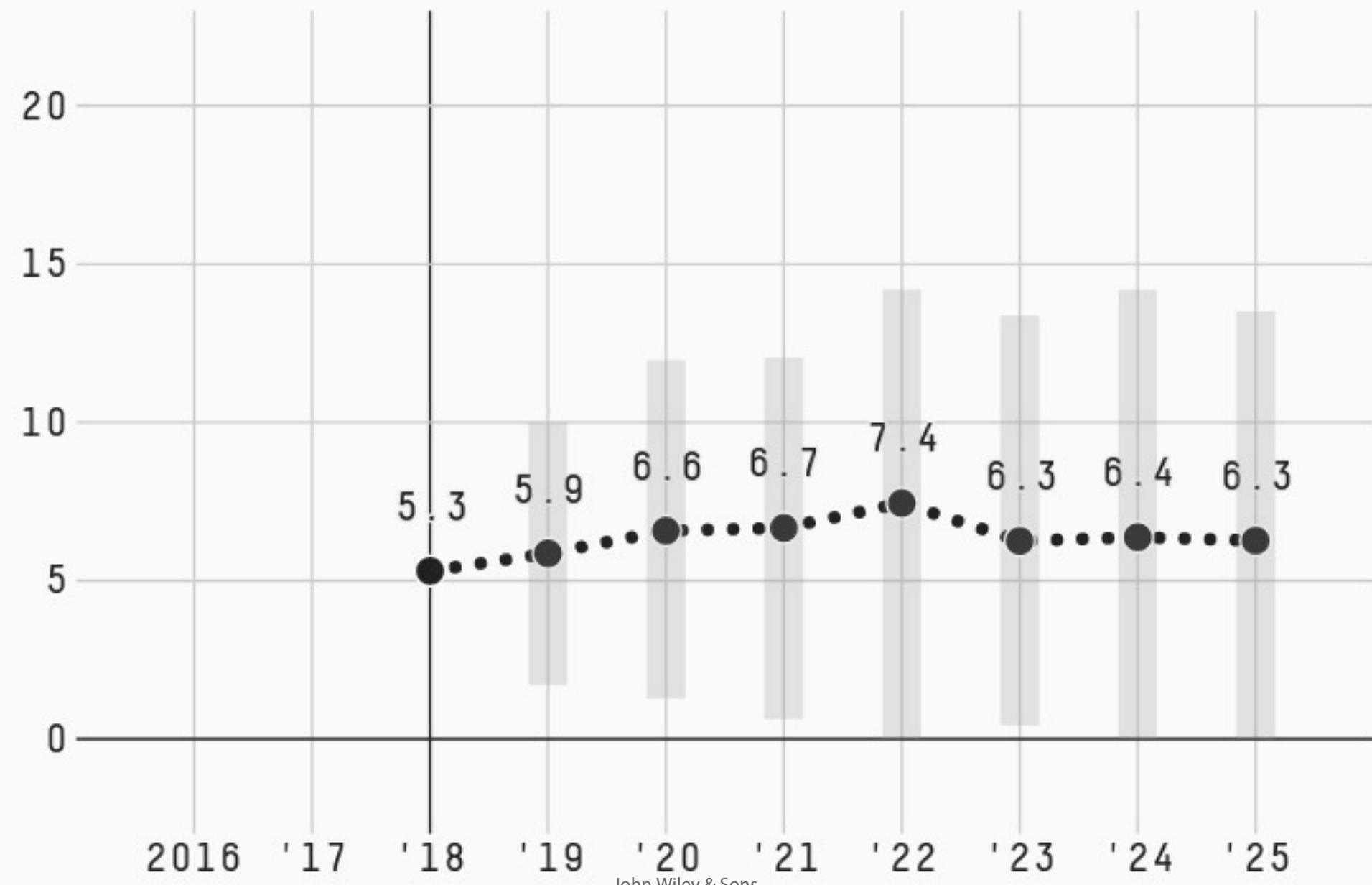


# Wins above replacement projection

CATEGORY: FUTURE ALL-STAR

5-YR MARKET VALUE: \$281.3M

90TH - CONFIDENCE INTERVAL  
10TH - PROJECTION

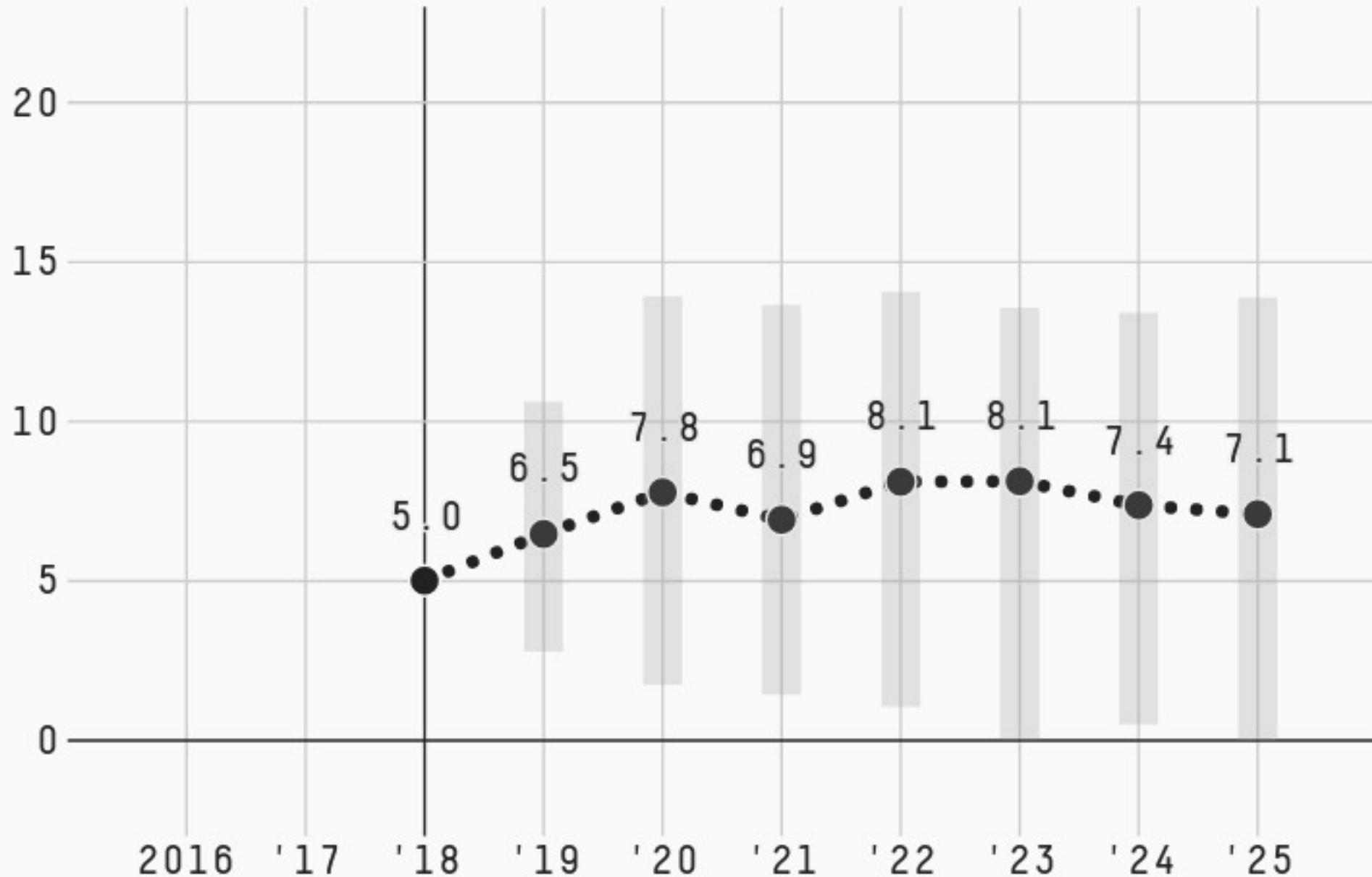


# Wins above replacement projection

CATEGORY: ALL-STAR

5-YR MARKET VALUE: \$320.4M

90TH-  
10TH-  
CONFIDENCE INTERVAL  
PROJECTION



1

2

**ARTICLE TYPE****Forecasting basketball players' performance using sparse functional data. APPENDIX<sup>†</sup>**G. Vinué<sup>\*1</sup> | I. Epifanio<sup>2</sup>

<sup>1</sup>Department of Statistics and O.R., University of Valencia, 46100 Burjassot, Spain

<sup>2</sup>Dept. Matemàtiques and Institut de Matemàtiques i Aplicacions de Castelló, Campus del Riu Sec. Universitat Jaume I, 71 Castelló, Spain

**Correspondence**

\*G. Vinué. Email: guillermo.vinue@uv.es

**Summary**

Statistics and analytic methods are becoming increasingly important in basketball. In particular, predicting players' performance using past observations is a considerable challenge. The purpose of this study is to forecast the future behavior of basketball players. The available data are sparse functional data, which are very common in sports. So far, however, no forecasting method designed for sparse functional data has been used in sports. A methodology based on two methods to handle sparse and irregular data, together with the analogous method and functional archetypoid analysis is proposed. Results in comparison with traditional methods show that our approach is competitive and additionally provides prediction intervals. The methodology can also be used in other sports when sparse longitudinal data are available.

**KEYWORDS:**

Forecasting, Functional data analysis, Archetypal analysis, Functional sparse data, Basketball

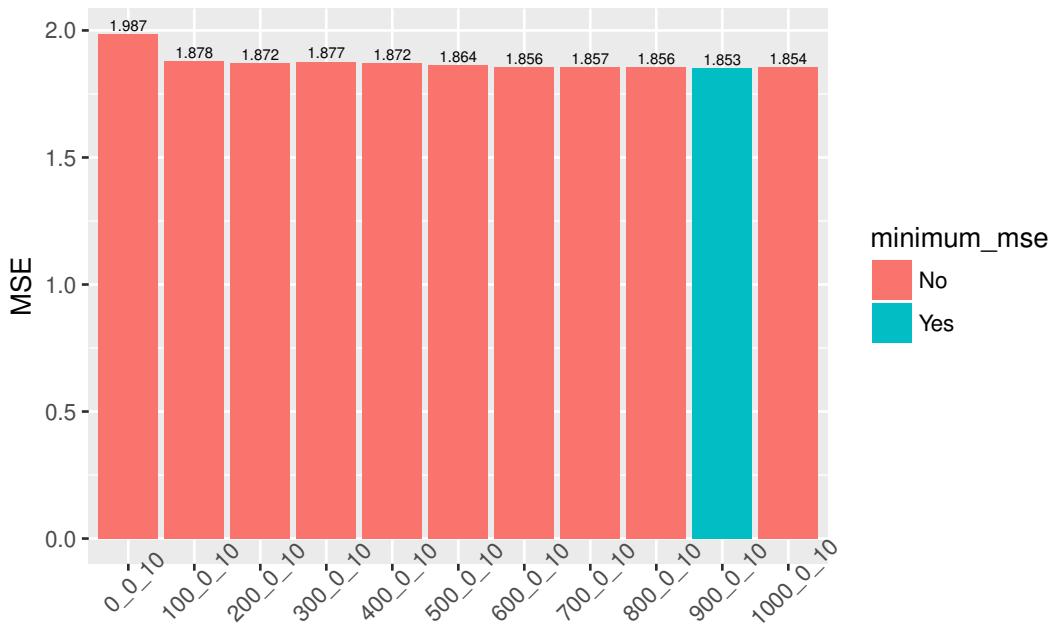
**1 | SELECTION OF PARAMETERS**

ROPEs depends on three tuning parameters ( $\lambda_2$ ,  $\lambda_1$ ,  $\lambda_0$ ), which have to be chosen to guarantee that the model itself returns predictions with enough accuracy. We evaluate the precision of the model's prediction in terms of the mean squared error (MSE). MSE measures the average of the squares of the differences between the predicted values  $\hat{y}$  and the true values  $y$  across all individual estimates  $i$ , as shown in Eq. (1).

$$MSE = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \frac{1}{n} \sum_{i=1}^n e_i^2 \quad (1)$$

We adopt MSE since ROPEs uses it to measure the error term. In order to select the parameters, we proceed as follows: our goal will be to predict the BPM in the 2017-2018 season, for the players who played at least in one season before the 2017-2018 season and who also played in the 2017-2018 season itself. The justification for doing this is related to sporting reasons. In sports, when coaches and managers are building their rosters, it is highly important for them to have a basic idea about how players will perform during the following season. Of course, they would also like to know the players' performance in the long term, but most rosters are built according to the most immediate season. This would allow them to decide whether the current roster should remain the same for the next season or whether some players should be replaced. This procedure makes sense because we will consider the previous performance of all the players selected, but we are only interested in predicting their BPM for the next season, by taking into account each player's data and the information about the other players. This procedure is more computationally efficient than the leave-one-out approach.

<sup>†</sup>The data and software associated with this paper are available at <https://www.uv.es/vivigui/software>.



**FIGURE 1** Averaged MSE across folds for every combination of lambdas when  $\lambda_2$  is moving and  $\lambda_0, \lambda_1$  are fixed. The combination for the smallest MSE is highlighted in green (colors in the online version).

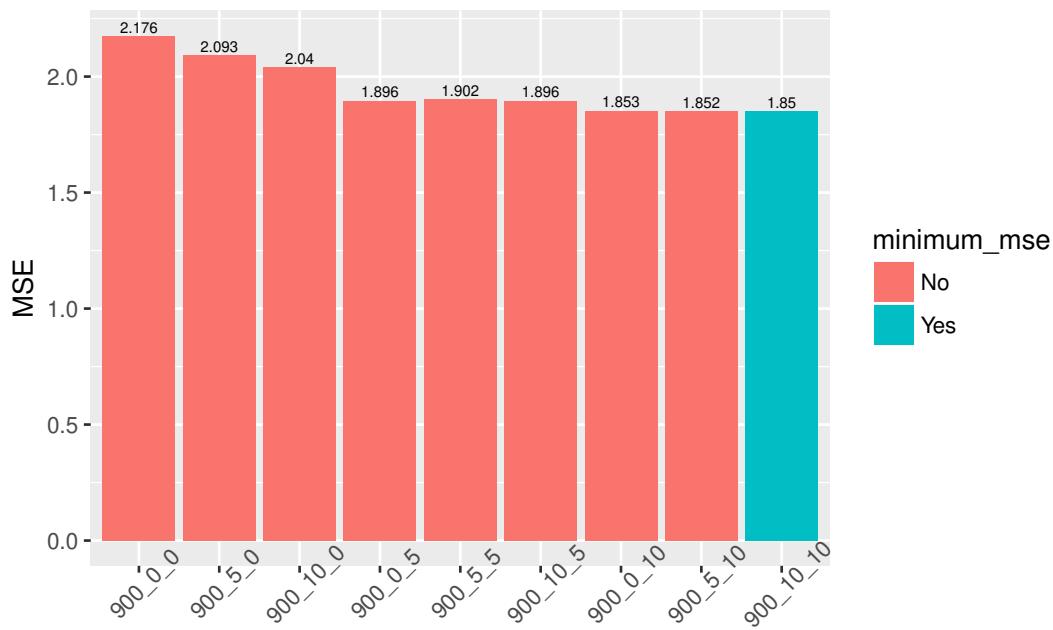
From the 3075 players, there are 385 who played in the 2017-2018 season and at least another season before that. Firstly, we split our data into a training+validation set (*TrVaSet*) with 2690 (3075 – 385) players and a test set with the 385 aforementioned players. No test set player belongs to *TrVaSet*. We will select the optimal combination of  $\lambda$ 's with *TrVaSet*. To do this, we have carried out an exhaustive 10-fold cross-validation for each parameter combination  $j$ . Inside each fold  $i$ , we use 60% observations as the training set and the remaining observations as the validation set. The mean squared error,  $MSE_{ij}$ , is then computed on the observations of the validation set. Finally, we obtain the  $MSE_j$  for every  $j$  by averaging across folds. For the validation players, their BPM value is replaced by NA in the  $Y$  matrix. In the  $W$  masking matrix the 1 value is then replaced by 0.

Let's define the parameter combinations. The first step is to optimize the parameter  $\lambda_2$ , setting  $\lambda_1 = 0$  and  $\lambda_0 = 10$ . The parameter  $\lambda_2$  takes values in a sequence from 0 to 1000 in increments of 100. In this way, the first blend is  $(\lambda_2 = 0, \lambda_1 = 0, \lambda_0 = 10)$ , the second is  $(\lambda_2 = 100, \lambda_1 = 0, \lambda_0 = 10)$  and so on. We are looking for smooth curves, so we place more emphasis on  $\lambda_2$  because it is related to the second derivative and this derivative is strongly related to the smoothness of the curve. This is justified because if the second derivative is a smooth curve, both the first derivative and the original function will also be smooth. From the definition of derivative, it directly follows that if a function has a first derivative at any point, then it does not have a sharp bend (v-shape) at that point (the same can be said of the second derivative with respect to the first derivative). See for example [2, Section 5.2.2] for further insights. The opposite is not always true.

Fig. 1 shows the averaged MSE across folds for every combination when only  $\lambda_2$  is moving. The smallest MSE was for the combination with  $\lambda_2 = 900$ .

Once the optimal  $\lambda_2$  has been found, we then adjust  $\lambda_1$  and  $\lambda_0$  as well. Both  $\lambda_0$  and  $\lambda_1$  take values in a sequence from 0 to 10 in increments of 5. Fig. 2 shows the averaged MSE across folds for every combination when both  $\lambda_0$  and  $\lambda_1$  are moving. The smallest MSE was for the combination  $(\lambda_2 = 900, \lambda_1 = 10, \lambda_0 = 10)$ .

The grid search procedure is chosen since it is a traditional way of performing hyperparameter optimization and can return results in a reasonable amount of time. This second point was particularly important because ROPES is computationally expensive. However, we would like to emphasize that we are not arguing that the grid search is the most efficient approach. There might be other more efficient alternatives than the grid search. We aim to investigate them as part of our future work.



**FIGURE 2** Averaged MSE across folds for every combination of lambdas when  $\lambda_2$  is fixed and  $\lambda_0, \lambda_1$  are moving. The combination for the smallest MSE is highlighted in green (colors in the online version).

## 2 | PROJECTIONS FOR INTERNATIONAL PLAYERS

Due to the increasing number of foreign players entering the NBA, the identification of foreign rising stars has become very important. However, there have been very few attempts in this regard. In [1], international players in the NBA were examined, but only in terms of the salaries they earned. A paper that is devoted more to forecasting is [3], where the author investigates which international statistics and features can be related to success in the NBA and makes some predictions. CARMELO also includes projections for European players, but only based on their biographical data, not on their statistics<sup>1</sup>.

We focus on the ACB league, which is the top professional basketball division in Spain and the strongest domestic league in the world besides the NBA. Data are collected using the **BAwiR** R package [4]. Regular season average statistics from the 1985-1986 to 2017-2018 seasons are analyzed. The target variable considered is the Performance Index Rating (PIR), which is the most widely used metric in European basketball leagues to measure the players' performance. Until more advanced metrics are developed for European leagues, PIR remains as the generally accepted way to rank players. Its formula is:

$$\frac{(\text{PTS} + \text{TRB} + \text{AST} + \text{STL} + \text{BLKfV} + \text{PFrv}) - (\text{FGmissed} + \text{FTmissed} + \text{TOV} + \text{BLKag} + \text{PFcm})}{\text{FGmissed} + \text{FTmissed} + \text{TOV} + \text{BLKag} + \text{PFcm}} \quad (2)$$

We propose to use the combination of archetypoid analysis and ROPES to provide some insights into the potential of young ACB-grown players of becoming a star in the NBA, by considering their PIR values. In particular, we will analyze Luka Doncic, who has been selected with the number 3 pick in the 2018 NBA draft. This means that there are high expectations of him. Doncic, born in 1999 in Slovenia, was playing for Real Madrid until the end of the 2017-2018 season.

After applying archetypoid analysis to the sparse functional database, four archetypoids were obtained (career PIR shown in brackets): Paco Vázquez (3.54), Arvydas Sabonis (28.4), Troy Bell (-7.5) and Rudy Fernández (12.8). Arvydas Sabonis is the representative of super star players, in line with the expected results. He is a basketball legend, member of the FIBA Hall of

<sup>1</sup><https://fivethirtyeight.com/features/whats-new-in-our-nba-player-projections-for-2017-18/>

1

2 | Vinué ET AL

3 **TABLE 1** Similarity of current NBA players who played in their early years in the ACB, to the four archetypoids according to  
4 the  $\alpha$  coefficients.

5 Player	6 Draft (Pick)	7 Archetypoids			
		P. Vázquez	A. Sabonis	T. Bell	R. Fernández
P. Gasol	2001 (3)	0.13	0.25	0.1	0.52
R. Rubio	2009 (5)	0	0.6	0	0.4
B. Biyombo	2011 (7)	0.01	0.32	0.31	0.36
K. Porzingis	2015 (4)	0.04	0.25	0.26	0.45
M. Hezonja	2015 (5)	0.12	0.2	0.4	0.28
L. Doncic	2018 (3)	0	0.54	0	0.46

14  
15 Fame and of the Naismith Memorial Basketball Hall of Fame. He has the PIR record in one ACB game with 66. He played  
16 seven seasons in the NBA, where he was a member of the 1996 NBA All-Rookie First Team.  
17

18 Rudy Fernández is the representative of great players, but a step below the super star level. He has been one of the most  
19 complete European players during the 21st century, and is very skillful both offensively and defensively. He played an important  
20 role in all the major achievements of the Spanish national team over the last 15 years. He played four seasons in the NBA, where  
21 he was a member of the 2009 NBA All-Rookie Second Team.

22 Paco Vázquez represents the mid-class players with a remarkable career. He played 515 ACB games and won the ACB  
23 championship in the 1997-1998 season. He also played in the Spanish national team in one European championship.

24 Troy Bell represents short-term or negligible performance players. He only appeared in 2 games with Real Madrid, where he  
25 signed a two-month contract, playing a total of 44 minutes. Bell was drafted 16th overall in the first round of the 2003 NBA Draft.

26 Over the past 20 years, more and more players have entered the NBA from the ACB. From all these ACB players, the player  
27 who has had the greatest impact is Pau Gasol, who was the 2002 NBA Rookie of the Year, won two NBA championships with  
28 the Lakers in 2009 and 2010 and played six NBA All-Star Games. Another player with an already long career in the NBA is  
29 Ricky Rubio, who has been playing there for the last seven seasons. Rubio became the youngest player ever to play in the ACB  
30 at age 14. The last three ACB players who were selected in the top 10 of the NBA Draft before Doncic were Bismack Biyombo  
31 (in 2011), Kristaps Porzingis (in 2015) and Mario Hezonja (in 2015).

32 It is interesting then to compare the  $\alpha$  values of Gasol, Rubio, Biyombo, Porzingis and Hezonja with those of Doncic. Table  
33 shows their  $\alpha$  values and the details of their draft picks.

34 According to Doncic's  $\alpha$  values, we see that his activity in terms of his PIR can be compared with a combination of Sabonis'  
35 and Rudy's performance, which indeed indicates his outstanding performance. This is in line with high expectations for the NBA  
36 and the fact that he is considered one of the most promising European prospects ever. We see that Rubio's ACB performance was  
37 also great. Biyombo is a mixture between Rudy and Sabonis, but also with Bell, so there could be some uncertainty regarding  
38 his performance in the NBA. Both Rubio and Biyombo are now consolidated players in the NBA.  
39

40 On the contrary, Porzingis and Hezonja do not look so good, especially Hezonja, whose performance is most closely related  
41 to Bell's. Hezonja's impact on the NBA has been very limited so far. It would probably have been better for Hezonja to stay in  
42 the ACB for a couple more years to further develop his skills, before moving on to play in the NBA. In the case of Gasol, we  
43 would like to point out that he played extremely well in just one ACB season before moving on to the NBA. This is a plausible  
44 reason why his alpha value is not higher than that of Sabonis, as we could expect. However, his greatest alpha is for Rudy, which  
45 is the other representative of great players.  
46

47 Fig. 3 shows the aging curve prediction for Luka Doncic, which can be used as a proxy of his future NBA potential. We see  
48 that his performance will increase in the coming years, achieving high PIR values. We can also compare our forecast with that  
49 of CARMELO<sup>2</sup>. CARMELO also predicts increasingly good performance in the coming years. Both methodologies agree in  
50 that Doncic will have a good future. Our analysis could be used as a guide to anticipate future performance in the NBA.  
51  
52  
53  
54  
55  
56

57 <sup>2</sup><https://projects.fivethirtyeight.com/carmelo/luka-doncic/>

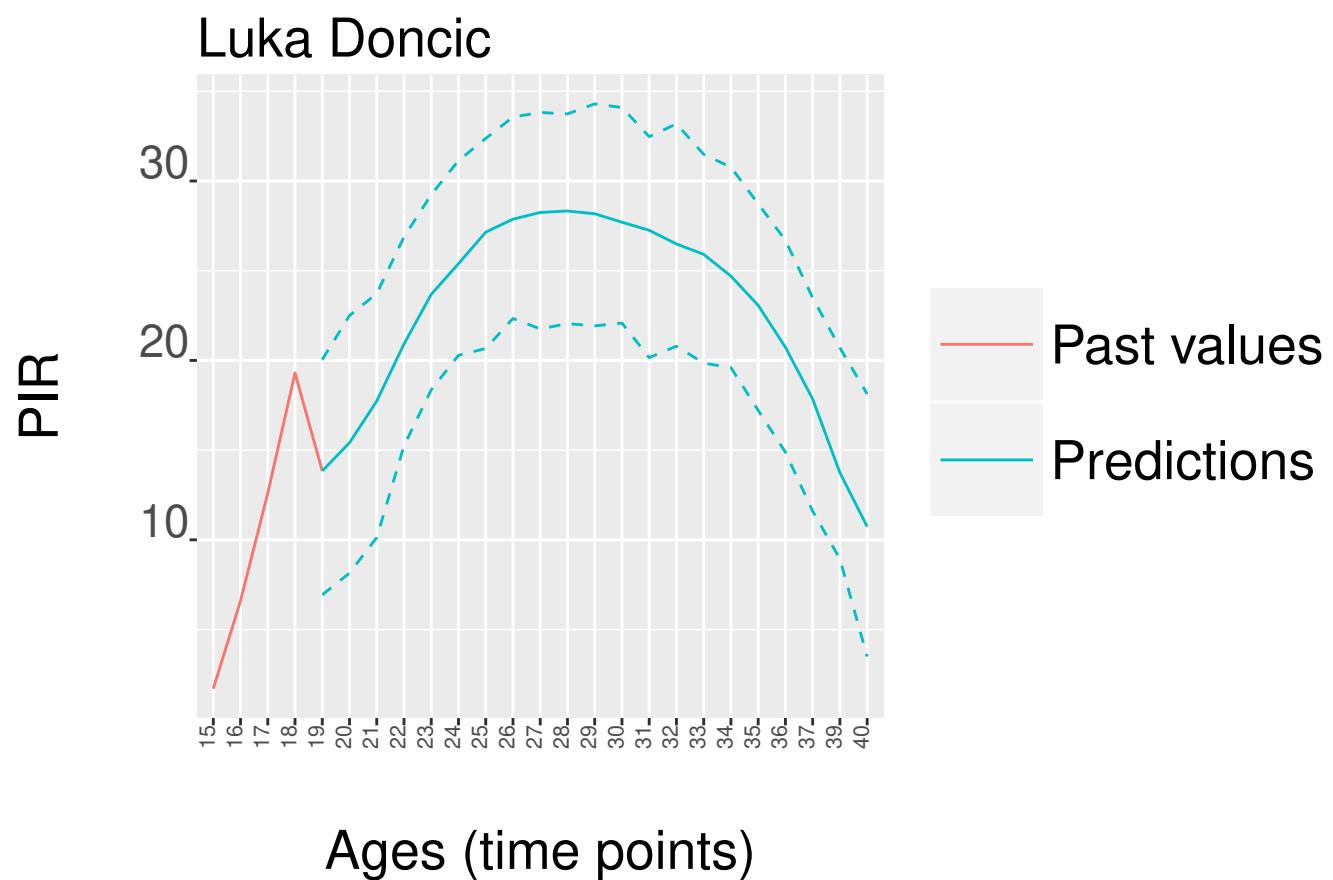


FIGURE 3 PIR prediction for Luka Doncic (colors in the online version).

## References

- [1] Eschker, E., Perez, S. J., Siegler, M. V., 2004. The NBA and the influx of international basketball players. *Applied Economics* 36(10), 1009–1020, <https://doi.org/10.1080/0003684042000246713>.
- [2] Ramsay, J. O., Silverman, B., 2005. *Functional Data Analysis*, 2nd Edition. Springer.
- [3] Salador, K., 2011. Forecasting Performance of International Players in the NBA. In: MIT Sloan Sports Analytics Conference. Boston, MA, USA, pp. 1–18, <http://www.sloansportsconference.com/wp-content/uploads/2011/08/Forecasting-Performance-of-International-Players-in-the-NBA.pdf>.
- [4] Vinué, G., 2019. *BAwiR: Analysis of Basketball Data*. R package version 1.2, <https://CRAN.R-project.org/package=BAwiR>.

**How to cite this article:** Vinué G., and I. Epifanio (2019), Forecasting basketball players' performance using sparse functional data, *Statistical Analysis and Data Mining*, 2019;:-.