




Measuring players' importance in basketball using the generalized Shapley value

Rodolfo Metulini¹ · Giorgio Gnecco² 

Accepted: 3 March 2022
© The Author(s) 2022

Abstract

Measuring players' importance in team sports to help coaches and staff with the aim of winning the game is gaining relevance, mainly because of the advent of new data and advanced technologies. In this paper we evaluate each player's importance - for the first time in basketball - as his/her average marginal contribution to the utility of an ordered subset of players, through a generalized version of the Shapley value, where the value assumed by the generalized characteristic function of the generalized coalitional game is expressed in terms of the probability a certain lineup has to win the game. In turn, such probability is estimated by applying a logistic regression model in which the response is represented by the game outcome and the Dean's factors are used as explanatory features. Then, we estimate the generalized Shapley values of the players, with associated bootstrap confidence intervals. A novelty, allowed by explicitly considering single lineups, is represented by the possibility of forming best lineups based on players' estimated generalized Shapley values conditional on specific constraints, such as an injury or an "a-priori" coach's decision. A comparison of our proposed approach with industry-standard counterparts shows a strong linear relation. We show the application of our proposed method to seventeen full NBA seasons (from 2004/2005 to 2020/21). We eventually estimate generalized Shapley values for Utah Jazz players and we show how our method is allowed to be used to form best lineups.

Keywords Players' performance · Sports analytics · Logistic regression · Cooperative game theory · National Basketball Association

✉ Giorgio Gnecco
giorgio.gnecco@imtlucca.it
Rodolfo Metulini
rmetulini@unisa.it

¹ Department of Economics and Statistics (DISES), University of Salerno, Via Giovanni Paolo II, 132, 84084 Fisciano, SA, Italy

² Laboratory for the Analysis of Complex Economics Systems (AXES) and Game Science Research Center, IMT School for Advanced Studies, Piazza S. Francesco, 19, 55100 Lucca, Italy

1 Introduction

Data analytics in basketball is nowadays a common practice to help coaches, staff and betting industry about the strategy to adopt (see, e.g., Nikolaidis (2015), Sarlis and Tjortjis (2020)) and it is increasingly used thanks to the large amount of different types of data, that can be classified into two distinct categories: (i) tracking data, generally collected using optical- or device-tracking and processing systems, that capture the movements and trajectories of players or of the ball in the court (Gudmundsson and Horton (2017) providing a quite exhaustive state of the art on this issue)¹; (ii) play-by-play data, which report (summarized in the box-scores) a sequence of relevant events occurring during a game, such as passes and shots as well as technical events, for example fouls and time-outs. The aforementioned data are in use to study the determinants of team performance. For example, by adopting variants of Data Envelopment Analysis (DEA), Yang et al. (2014), Moreno and Lozano (2014) decompose the overall team performance into income and on-court efficiency, while using a Stochastic Frontier approach, Hofler and Payne (2006) investigate how closely NBA teams play up to their potential. Tracking data have been used in Metulini et al. (2018) to measure the role of payers' spacing in team performance.

The recent literature also addressed the topic of the performance of the single player. Players are now being evaluated by using advanced statistical and machine learning techniques that depart from just using basic statistics such as number of points, steals, turnovers. For instance, a large range of works on player performance exists, using the entire box-score (see, e.g., Cooper et al. (2009), Fearnhead and Taylor (2011), Page et al. (2013)), shooting related variables (e.g., Piette et al. (2013), Metulini and Le Carre (2020), Sandri et al. (2020)) or ad-hoc proposed synthetic metrics (Terner and Franks 2021). However, both in offense and in defense, team performance in basketball may be viewed as a network issue, wherein each play represents a *pathway* through which the ball and players cooperatively move from the beginning of the play to the goal (score the basket). It follows that players and team should not be evaluated separately (e.g., by studying players' performance by means of points made or team performance by means of percentage of winnings), but according to comprehensive advanced measures that take into account the team and the players together. In basketball, five players in both home and away teams rotate on the court. The five players of each single team on the court in that moment represent the lineup (or quintet), while the ten players in the court represents the encounters. Ordering players and lineups has been addressed, for example, in Barrientos et al. (2019), via a Bayesian approach, for the analysis of encounters. Kalman and Bosch (2020) examined lineups to group players in order to detect more efficient quintets.

Borrowing the idea of the Shapley value (Shapley 1953) and of the generalized Shapley value (Nowak and Radzik 1994) from Cooperative Game Theory (in the spirit of, e.g., of Hernández-Lamonedá and Sánchez-Sánchez (2010), Hiller (2018)), this paper aims at estimating the importance of single players in basketball in terms of their contributions to the team performance, more precisely by computing their average marginal contribution to the utility of ordered subsets of players of the team (i.e., considering each player playing with ordered subsets of other players in the lineup). For short, this can be also called the average marginal utility (or productivity) of the players. The Shapley value has successfully been used in many political and economic games. The utilization of this value has not massively

¹ The cooperative game theoretical approach for the analysis of the origin of movement addressed in Kolykhalova et al. (2020); Matthiopolou et al. (2020) is worth mentioning here, because it may have applications in the analysis of movement for the case of basketball players.

been percolated to team sports analysis (an exception being the works on Soccer by Auer and Hiller (2015), Hiller (2015)). To the best of our knowledge, the Shapley value has never been used to evaluate players' performance in basketball, except for the conference article by Yan et al. (2020). The Shapley value is the average of the marginal utilities (that can assume both negative and positive values) of a player, one for each combination of players playing together with him. Each marginal utility is computed based on the difference between the values assumed by a function (called characteristic function) that measures cohesion of each combination of players (also called coalition), calculated respectively with and without him in the court. However, it is worth noting that a generalization of the Shapley value to the case of ordered coalitions, which was proposed in the context of generalized coalitional games by Nowak and Radzik (1994), is more suitable than the Shapley value itself for our application to basketball analysis. In the paper, we actually adopt such a generalization of the Shapley value. This is done in order to take into account the fact that only five players for team can play simultaneously, then it is reasonable to assume that every player that *virtually* enters a coalition after the fifth player who enters has a zero marginal utility (instead of a negative marginal utility). However, as it will be explained later in Sect. 2.1, this leads to multiple values for the coalition made of all players in the team (which is called grand coalition). This case is not covered by the classical definition of the Shapley value, but it is addressed by the generalized Shapley value proposed by Nowak and Radzik, which relies on a generalized characteristic function, whose argument is an ordered subset of players (instead of an unordered one).

To correctly determine the generalized characteristic function is a fundamental aspect, as it measures the performance of any specific ordered subset of players when playing together. Yan et al. (2020) themselves, by tackling the problem of estimating the average marginal contributions of the players from a different perspective (i.e., in terms of Shapley values), highlighted the importance of learning from data the characteristic function of the game. Such a remark can be extended to the case of the generalized characteristic function. In the present work, we model the generalized characteristic function in terms of the estimated probability that each lineup has to win the game. In this regard, the paper (according to how our proposed measure is determined) is more closely related to the line of research on estimating players' contribution on winning the game (Deshpande and Jensen 2016), and it deviates from industry-standard measures that have been proposed to evaluate single player's contribution to the team which are based on the difference between the points scored by a player's team and those scored by his/her opponent team during the time that specific player is on the court, such as simple Plus-Minus (PM), regression-based versions of the PM metric to measure player's contribution by accounting for the other players on the court (Adjusted PM, APM; see Rosenbaum 2004), extensions of the APM that include other players' statistics among the explanatories and that control for the team strength (box-score PM, BPM, e.g., Kubatko et al. (2007), Ilardi (2007), Grasseti et al. (2021)) or that try to account for the presence of multicollinearity in APM, e.g., Sill (2010) and Engelmann (2017) based on ridge regression regularization (Regularized APM), or Real Plus Minus (RPM), which normalizes the measure by the number of offensive and defensive possessions. Overall, despite recent PM versions move in the direction of i) not just accounting for scoring factors, and ii) solving for multicollinearity, those issues still deserve more attention (Turner and Franks 2021). Beside PM and its extensions, Win-Shares (WS), calculated using player, team and league statistics, attempts to measure the contribution for team success of its individuals. WS48 (WS per 48 minutes) expresses the WS values in a per-minute basis. Wins Above Replacement (WAR), also referred to WAR Player (WARP), firstly developed for baseball in order to find player's contribution in terms of how many additional wins he/she brings to the team, seeks to evaluate

a player by comparing the performance of a team made up of him/her and four average players with the performance of a team made up of four average players and one replacement-level player. However, despite WS presents the advantage of accounting for the marginal utility of a single player to the win by comparing him/her to an average replacement level, as remarked in Sarlis and Tjortjis (2020), a player WS score is positively influenced by being part of a good team and by the amount of time he/she is on the court. WARP and WS48 outperform WS as they are expressed on a per-minute basis. Value Over Replacement Player (VORP), defined as an estimate of the points per 100 team possessions that a player contributed above a replacement level player, aims at collecting together the advantages of BPM and those of WARP. However, likewise WARP, VORP suffers from issues of multicollinearity (Sarlis and Tjortjis 2020).

The methodological strategy of this paper is composed of three steps, having in mind the aim of choosing a suitable definition for the generalized characteristic function. According to our method, in the first step logistic regression model coefficients based on all NBA data at hand (seventeen seasons) are estimated, where the dependent variable is the dichotomous information about the result of the considered team (called *Outcome*, win=1, defeat=0) and the features used as explanatory variables are represented by appropriate synthetic measures computed based on play-by-play statistics of both the teams in the game (the so called *four Dean's factors*; see Kubatko et al. (2007), Oliver (2004)). In the second step, the logistic regression equation with estimated coefficients (from the first step) is used to derive the probability to win associated with each lineup (i.e., replacing the at-game level explanatory features with those computed at-lineup level). With the probability of winning (used to express the generalized characteristic function) computed for all the lineups, in the third step we compute two different versions (*unweighted* and *weighted*, where the second one accounts for the amount of time players are on the court) of the generalized Shapley value for each player. Lastly, in order to help coaches and staff with lineup management, we propose a greedy approach to find an appropriate lineup. According to this approach, we first choose, among n players in the team, the one with the largest generalized Shapley value (which may be termed the “most important player”). Then, we recompute the generalized Shapley values for all the other players according to the subset of lineups in which the most important player was in (effectively considering a modified game in which the most important player is always present, hence the set of players is successively restricted to the other $n - 1$ players), and we choose the player with the largest generalized Shapley value (“second most important player”), and so on. We repeat the process until the “first five most important players” have been chosen. Variants of this approach can be applied by replacing the “most important player” with a specific player chosen by the coach, or by not considering at all the lineups in which a player was in (under the assumption he/she is unavailable).

Overall, in this work we propose a new method to measure players' contributions to the team that can be adopted to help coaches and staff in retrieving insights on which players and lineup to choose and that gathers most of the advantages (and avoids disadvantages) of industry-standard measures. In fact, similarly to BPM, our measure presents the advantage of accounting for both scoring and non scoring, offensive and defensive factors. Moreover, as we will show in Sect. 4, in estimating the winning probabilities we adopt a set of box-score synthetic measures (the 4 Dean's factors) which presents an extremely high goodness of fit. It is also worth noting in this regard that the adopted parameters – associated to box-score synthetic measures – are estimated on a very large dataset spanning seventeen seasons and may be seen as the weights to assign to a synthetic measure for modeling a mechanism for obtaining a win. Moreover, similarly to what WARP and VORP do with replacement level player, with our approach we consider marginal utilities of players' considering lineups.

But we do that by explicitly considering all the lineups in which he/she has played with (so, adopting a more “holistic” approach), without the need to consider a level for the replacement player and avoiding multicollinearity issues², see Mishra (2016). It is also worth noting that our metrics, which are proportional to suitable averages of winning probabilities of lineups to which each player belong³, are expressed in terms of a solution concept from cooperative game theory (i.e., in terms of the generalized Shapley value), whereas other industry-standard measures lack such a game-theoretical interpretation. Moreover, explicitly taking into account single lineups, permits us to evaluate players’ average marginal contributions conditional on the presence of specific players on the court. In such a way, a lineup management is permitted. To make an example, let us suppose the Utah Jazz coach is sure to include Rudy Gobert in the lineup, and he wants to compose the lineup conditional on the presence of him. With our method we can compute the generalized Shapley values of the other players according to the lineups in which Gobert was on the court. Another advantage of our method is that the generalized Shapley value, on which it is based, has an axiomatic characterization expressed in terms of simple properties (Michalak et al. 2014) that can be easily transferred to team sports and to basketball in particular. Moreover, in case one wants to add more features to increase the goodness of fit of the model to predict the winning probabilities, he/she just has to change only the specific definition of the generalized characteristic function considered in our method⁴, letting everything else unaltered. Moving to disadvantages, our method presents a limitation: the generalized Shapley value of a player, to be estimated, needs a large number of different lineups containing that player, due to the fact that the variance of its estimate is inversely proportional to the number of such lineups (a similar result holds for the Shapley value; see Castro et al. (2009), of which it constitutes a generalization). This fact does not currently allow us to compute players’ average marginal contributions with the estimated generalized Shapley values at the single game level (however, they may be estimated, e.g., by considering the lineups occurring in half a season, instead of the whole season). This aspect deserves to be addressed in future developments.

The paper is structured as follows: in Sect. 2 we define the adopted methods, Sect. 3 introduces the data, and Sect. 4 presents the application to NBA data. Section 5 concludes the work with a discussion.

2 Methods

2.1 The game-theoretical approach

In cooperative game theory, a generalized coalitional game (see, e.g., Nowak and Radzik (1994)) is defined as a pair (N, v) ; where $N = \{1, 2, \dots, n\}$ is the player set (whose cardinality is denoted either by $|N|$ or n), and v is the generalized characteristic (or utility) function,

² For instance, in the context of regression-based PM measures, in the case of multicollinearity the estimated coefficients of the regression that are associated with the individual players (hence, the values themselves assumed by such regression PM measures) can be quite sensitive to changes in the data. In our logistic regression model, instead, the coefficients of the regression are not directly associated with the individual players, so the sensitivity of the estimate of the vector of coefficients of the regression to changes in the data is less important (i.e., one can limit to focus the attention only on the prediction capability of the model, which turns out to be quite high).

³ Lineups to which he/she does not belong are not considered in the computation, but are implicitly taken into account when one evaluates the generalized Shapley value of another player.

⁴ See Sect. 2.1 for details about its precise definition.

which assigns to every ordered list (or ordered coalition) T extracted from the set N a certain worth $v(T)$ reflecting the abilities of such an ordered coalition; i.e., denoting by \mathcal{T} the set of such ordered coalitions, one has $v : \mathcal{T} \rightarrow \mathbb{R}$ such that $v(\emptyset) = 0$. This definition differs from the related and more commonly known definition of a coalitional game, whose characteristic function is defined on the set of (unordered) coalitions (Chapter 17 in Maschler et al. 2013). To clarify the concepts above, the set $\{1, 2, 3\}$ is an unordered coalition, whereas the sequences $(1, 2, 3)$, $(1, 3, 2)$, $(2, 1, 3)$, $(2, 3, 1)$, $(3, 1, 2)$, $(3, 2, 1)$ are all the 6 corresponding ordered coalitions.

In contrast to a classic game, a sport game is characterized by first attributing nonzero worth to coalitions with a given cardinality $m < n$, then extending the definition of the (generalized) utility function to the other coalitions, following a suitable rule (see some examples in the remaining of this section). Such coalitions are either ordered or unordered, depending on the model adopted. In the case of basketball, we have $m = 5$, whereas n is the total number of players rotating in the court in a game. Suitable specifications of the generalized utility function in the context of basketball data analysis are reported later in this section: see Eqs. (2, 3, 4, 5, 6) in the following.

Given the framework above, a solution to a coalitional game can be interpreted as a way to distribute the utility of the coalition made of all the n players (which is also called “grand coalition” in the game-theoretical literature) among its members. This is achieved by following, e.g., suitably formalized fairness and efficiency principles (Chapter 17 in Maschler et al. 2013). In the case of a generalized coalitional game (Nowak and Radzik 1994), the utility of the grand coalition is replaced by the average of the utilities of all possible ordered coalitions of n players (in total, using the factorial notation, there are $n!$ such ordered coalitions). It is well-known that a solution to a (generalized) coalition game can be interpreted as a way to rank all the players of a team, by attributing a numerical value to each of them, which represents a measure of his/her importance in the team. For instance, in the case of a coalitional game, the Shapley value of a player is a suitable average of the marginal utility provided by that player when he/she joins a properly-generated random coalition of players (Chapter 17 in Maschler et al. 2013). In the case of a generalized coalition game, a similar interpretation holds for the generalized Shapley value (introduced in the next paragraph). By providing suitable ways of ranking players in a team, these interpretations justify the application of these concepts to basketball data analysis⁵. As argued later in this article, however, in the context of such analysis, the generalized Shapley value is a better metric than the Shapley value. This motivates the use of the former concept in the present work.

The generalized Shapley value (Nowak-Radzik value, see Nowak and Radzik (1994)) of player $i = 1, \dots, n$ in a generalized coalition game is the average of the marginal contribution of that player when he enters an ordered subcoalition of a random ordered coalition T with cardinality $|T| = n$. The average is taken over all such ordered coalitions T of the players, giving the same weight to each T , according to the following formula:

$$\phi_i^{NR}(N, v) = \frac{1}{n!} \sum_{T \in \mathcal{T} \text{ with } |T|=n} (v((T(i), i)) - v(T(i))) , \quad (1)$$

where $T(i)$ denotes the ordered (sub)coalition made by the predecessors of i in the permutation T , and $(T(i), i)$ denotes the ordered (sub)coalition made by $T(i)$ followed by i (more

⁵ An additional motivation is that the Shapley value can be proved to be the unique solution to a coalitional game – intended as an allocation of the total worth of the team – that satisfies four quite natural axioms (properties) called *symmetry*, *null player*, *efficiency*, and *additivity* (Chapter 17 in Maschler et al. 2013). An analogous axiomatic characterization holds for the generalized Shapley value for the case of a generalized coalitional game, see Michalak et al. (2014).

details about the notation and some examples to better understand some terms in Eq. (1) are reported in footnote 6⁶. Because of Eq. (1), the generalized Shapley value of player i can be also called the average marginal utility of that player⁷.

In our application to basketball analysis, we consider two possible choices for the generalized characteristic function $v(\cdot)$ in Eq. (1), which are denoted respectively by $v_1(\cdot)$ and $v_2(\cdot)$. In the first case, we model the generalized characteristic function $v_1(\cdot)$ in terms of the probability $P(Win)$ of winning the game for any specific quintet of players. In the second case, we model the generalized characteristic function $v_2(\cdot)$ in terms of both the probability $P(Win)$ of winning the game for any specific quintet and the probability of occurrence $P(Occ)$ of that quintet on the court. In this way, the resulting generalized Shapley values $\phi_i^{NR}(N, v_1)$ and $\phi_i^{NR}(N, v_2)$ provide different measures of the importance of each player in a basketball team, because the first one takes into account only the probability of winning of each lineup he is part of, whereas the second one depends also on the probability of playing for each such lineup. For this reason, in Sect. 4, the two generalized Shapley values $\phi_i^{NR}(N, v_1)$ and $\phi_i^{NR}(N, v_2)$ are denoted, respectively, as “unweighted generalized Shapley value” and “weighted generalized Shapley value”.

To construct the generalized characteristic functions $v_1(\cdot)$ and $v_2(\cdot)$, we follow the next steps. First, we consider the case in which the arguments of the generalized characteristic functions $v_1(\cdot)$ and $v_2(\cdot)$ have cardinality $m = 5$. Then, when $|(T(i), i)| = 5$, we let

$$v_1((T(i), i)) = P(Win)_{(T(i), i)} \quad (2)$$

be the probability of winning the game for the ordered (sub)coalition of players $(T(i), i)$, and similarly, when $|T(i)| = 5$, we let

$$v_1(T) = P(Win)_{T(i)} \quad (3)$$

be the probability of winning the game for the ordered (sub)coalition of players $T(i)$, which does not contain player i . Then, after suitably extending the definition of $v_1(\cdot)$ to ordered coalitions having a number of players different from 5 (see the next Eq. (6)), we compute $\phi_i^{NR}(N, v_1)$ according to Eq. (1).

Similarly, to define the other generalized characteristic function $v_2(\cdot)$, we replace Equations (2) and (3) respectively with

$$v_2((T(i), i)) = P(Occ)_{(T(i), i)} P(Win)_{(T(i), i)} \quad (4)$$

⁶ Here, we follow a similar notation as the one used in Michalak et al. (2014). First, the elements of each ordered coalition $T \in \mathcal{T}$ are denoted by $T_1, \dots, T_{|T|}$, where the index refers to the order according to which a player enters that ordered coalition, in a “virtual” process of its construction. The ordered coalition made by the single element i is denoted by i itself. For any two disjoint ordered coalitions $T^{(1)}, T^{(2)} \in \mathcal{T}$, one denotes by $(T^{(1)}, T^{(2)})$ the ordered coalition obtained by concatenating $T^{(1)}$ and $T^{(2)}$, i.e., the ordered coalition in which the elements of $T^{(1)}$ (which occur in the order associated with $T^{(1)}$) precede those of $T^{(2)}$ (which occur in the order associated with $T^{(2)}$). For instance, if $T^{(1)} = (4, 1)$ and $T^{(2)} = (3, 6, 5)$, then $(T^{(1)}, T^{(2)}) = (4, 1, 3, 6, 5)$. For any ordered coalition T and any player i , $T(i)$ denotes the ordered (sub)coalition formed by the players that precede i in T (this coincides with T if i is not present in T). As an example, if $T = (2, 1, 3, 5, 4)$, then one has $T(5) = (2, 1, 3)$, and $(T(5), 5) = (2, 1, 3, 5)$.

⁷ The generalized Shapley value (1) is formally quite similar to the classical Shapley value (Maschler et al. 2013, Chapter 17), the main difference being that the generalized characteristic function $v(\cdot)$ takes into account the order according to which the players form an ordered subcoalition of T , whereas the analogous characteristic function used in the definition of the Shapley value is defined on the set of unordered subcoalitions of the grand coalition made by all the players (i.e., the order in which the players occur when forming the grand coalition is not taken into account to define the characteristic function, but only the marginal utility of each player).

and

$$v_2(T) = P(Occ)_{T(i)} P(Win)_{T(i)}, \quad (5)$$

where $P(Occ)_{(T(i),i)}$ and $P(Occ)_{T(i)}$ are the probabilities of occurrence on the court of the ordered (sub)coalitions of players $(T(i), i)$ and $T(i)$, respectively.

Details about how to model the probabilities introduced above are reported later at the end of this section and in Sect. 2.2. In particular, we assume that they do not depend on the order in which the players of the two specific ordered (sub)coalitions appear, respectively, in Eqs. (2), (3), (4), and (5)⁸.

As discussed in the next paragraph, the following way of extending the definitions of the generalized characteristic functions $v_k(\cdot)$ (for $k = 1, 2$) also to the other ordered coalitions with cardinality different from 5 justifies why the generalized Shapley value is used instead of the Shapley value for our application to basketball analysis:

$$v_k(T) = \begin{cases} 0 & \text{if } |T| < m = 5 \\ v_k(\{T_1, T_2, T_3, T_4, T_5\}) & \text{if } |T| \geq m = 5. \end{cases} \quad (6)$$

It follows by construction that the marginal contribution of any player that (virtually) enters in sixth position is always 0⁹. Since $v_k(\{T_1, T_2, T_3, T_4, T_5\})$ (the worth of a specific quintet of players) depends on T_1, T_2, T_3, T_4 , and T_5 , this means that the grand coalition has different worth, depending on the order of appearance of its players in it (only the first 5 players - the ones on the court - being actually considered, without taking into account their internal order). This prevents the application of the Shapley value to the specific case, since such application requires a unique worth for the grand coalition, which does not depend on the order in which the players enter it. In contrast, it justifies the application of the generalized Shapley value, which has not such a constraint. This issue is further discussed in the Appendix.

From a computational point of view, for each player i , the evaluation of the generalized Shapley value (1) with each of the two specifications $v_k(\cdot)$ for the generalized characteristic function requires considering the worth of $C(n-1, 4) = \frac{(n-1)!}{4!(n-1-4)!}$ different (unordered) quartets of players in the court together with player i . Since, by varying player i , not all the resulting $C(n, 5) = \frac{n!}{5!(n-5)!}$ quintets are observed, we adopt the following kind of approximation for the generalized Shapley value (1)¹⁰, which follows from its interpretation as the average marginal contribution of player i when he enters a random subcoalition: the quintets observed are interpreted as i.i.d. realizations of quintets $\{T_1, T_2, T_3, T_4, T_5\}$ obtained from a random permutation T (all permutations being equally likely), and player i is assumed to have the same probability of forming a subcoalition with the other players of the specific quintet when he enters first, second, third, fourth, or fifth. In this way, the average marginal contribution in (1) is replaced by an empirical average marginal contribution, based on the observed quintets, and taking into account that each player has probability $\frac{5}{n}$ of entering in

⁸ In this article, When the value of $v(T) = v((T_1, T_2, \dots, T_{|T|}))$ depends only on the elements $T_1, T_2, \dots, T_{|T|}$, but not on their order, it is denoted by $v(\{T_1, T_2, \dots, T_{|T|}\})$.

⁹ It is worth noting that, beside the generalization of the Shapley value developed in Nowak and Radzik (1994) and adopted in this work, another similar generalization exists in the literature, and is due to Sanchez and Bergantiños, see Sanchez and Bergantiños (1997) (refer also Michalak et al. (2014) for a comparison of the two generalizations). However, such a generalization cannot be applied (at least not in a straightforward manner) to the present setting because it can be interpreted as the Shapley value of an “average” coalitional game obtained by averaging the worths of all permutations of each ordered coalition, making it difficult to get a 0 marginal contribution for every player that (virtually) enters in sixth position.

¹⁰ A similar method is often used also for an approximate evaluation of the Shapley value itself (Castro et al. 2009).

one of the first 5 positions. Moreover, the values assumed by the two generalized characteristic value functions $v_k(\cdot)$ in correspondence of the observed quintets are estimated based on suitable features (see the next section). In the following, such estimates are denoted as $\hat{v}_k(\cdot)$. In summary, denoting by \mathcal{L}_i the set of observed (unordered) lineups (quintets) in which player i appears, one gets the following estimate of his generalized Shapley value:

$$\hat{\phi}_i^{NR}(N, v_k) = \frac{5}{n} \frac{1}{5|\mathcal{L}_i|} \sum_{L \in \mathcal{L}_i} (\hat{v}_k(L) - 0) = \frac{1}{n|\mathcal{L}_i|} \sum_{L \in \mathcal{L}_i} \hat{v}_k(L). \quad (7)$$

Equation (7) expresses the average value of a quintet in which player i occurs, multiplied by the factor $\frac{1}{5}$. The presence of such a factor is motivated by the fact that, for any specific quintet, each player has the same probability of being the last player to join all the other members of that quintet (in the “virtual” process of construction of the lineups, starting from the ordered coalitions of size n). As already mentioned, the other factor $\frac{5}{n}$ expresses the probability that player i enters in one of the first 5 positions (otherwise, his marginal contribution is 0). Equation (7) also clarifies that, when the generalized characteristic function $v_1(\cdot)$ is considered, the generalized Shapley value of player i is proportional to the average winning probability of a quintet in which he occurs. Instead, when the generalized characteristic function $v_2(\cdot)$ is considered, the value of each quintet depends also on its probability of being on the court. In a sense, $v_1(\cdot)$ represents the “instantaneous” utility (winning probability) of each quintet, whereas $v_2(\cdot)$ could represent its “integrated” or “weighted” utility, taking into account its probability of occurrence on the court.

As a simple illustrative example of application of Eq. (7), let $i = 3$, $n = 13$ and $\mathcal{L}_3 = \{\{1, 2, 3, 4, 5\}, \{1, 2, 3, 6, 8\}, \{3, 5, 7, 9, 12\}\}$. Moreover, suppose $\hat{v}_1(\{1, 2, 3, 4, 5\}) = 0.6$, $\hat{v}_1(\{1, 2, 3, 6, 8\}) = 0.7$, and $\hat{v}_1(\{3, 5, 7, 9, 12\}) = 0.8$. In this case, from Eq. (7), one gets

$$\hat{\phi}_3^{NR}(N, v_1) = \frac{1}{13 \cdot 3} (0.6 + 0.7 + 0.8) \simeq 0.054. \quad (8)$$

The example reported above is only illustrative because, to obtain a reliable estimate $\hat{\phi}_3^{NR}(N, v_1)$ of the generalized Shapley value $\phi_3^{NR}(N, v_1)$, a quite large number $|\mathcal{L}_3|$ of lineups containing player 3 is needed. Indeed, assuming as an example that all the lineups containing player i are equally likely and that they are sampled independently, neglecting the issue of possible repetitions in lineup sampling, and supposing for simplicity $\hat{v}_k(\cdot) = v_k(\cdot)$, Equation (7) can be interpreted as a Monte Carlo estimate of the generalized Shapley value $\phi_i^{NR}(N, v_k)$. It easily follows from the theory of Monte Carlo sampling that such estimate is unbiased, and its variance is proportional to $1/|\mathcal{L}_i|$. The proof of a similar result is well-known for the case in which the generalized Shapley value is replaced by the Shapley value, and the latter is approximated by a standard Monte Carlo estimate: see, e.g., (Proposition 3.1 in Castro et al. 2009) and its proof reported therein¹¹.

In Sect. 2.2, estimates of the winning probabilities $P(Win)$ of the various quintets (which enter the definitions of both $v_1(\cdot)$ and $v_2(\cdot)$) are provided. Moreover, the occurrence probability $P(Occ)$ of the each quintet (which enters the definition of $v_2(\cdot)$ only) is estimated as the ratio between the number of minutes played by that quintet in the considered period and the number of minutes played by all the quintets in the same period (the latter, regardless of the presence or absence of any specific player i in the quintet).

¹¹ The results of the analysis change only slightly in the case of sampling with possible lineup repetitions (under which samples become dependent). Indeed, in this variation, the estimate is still unbiased, and its variance is even slightly smaller than the one obtained in the case of sampling without replacement, see Rice (2005).

2.2 Estimates of lineups' winning probabilities using the logistic regression model

There is a plethora of studies on estimating the probability to win a basketball game using different sources of data and various statistics and machine learning techniques.

Loeffelholz et al. (Loeffelholz et al. 2009), by applying different variants of Artificial Neural Networks (ANNs, Zhang (2000)) on box-scores of 620 NBA games, predicted match outcome with a correct winner prediction percentage of 74.33%. Miljkovic et al. (Miljković et al. 2010) used Naive Bayes Classifier (Langley et al. 1992) to predict both the outcome and the spread for 778 2009–2010 NBA games using 141 box-score features as covariates, and obtained an accuracy of 67%. Beckler et al. (Beckler et al. 2013), using box-score data for seasons from 1991–1992 to 1996–1997, were able to achieve up to 73% accuracy for the NBA outcome prediction, using four standard binary classification algorithms: (i) Linear Regression (with binary outcome), (ii) Support Vector Machines (Cortes and Vapnik 1995), (iii) Logistic Regression (Hosmer and Lemeshow 2013) and (iv) ANNs. Cheng et al. (Cheng et al. 2016) used a Maximum Entropy approach (Jaynes 1957) to predict outcomes, with an accuracy of 74.4%, in NBA playoff, using box-score data from seasons from 2007–08 to 2014–15. Thabtah et al. (2019) obtained an accuracy in prediction of outcomes up to more than 80% applying Naive Bayes, ANNs and Logistic Model Trees (Landwehr et al. 2005), based on box scores from NBA finals from 1980 to 2017.

We adopt a logistic regression model strategy to estimate the probability of winning the match.

Traditionally, some features from the play-by-play of the game are used as a set of explanatories for winning prediction. Play-by-play information are usually extracted from the box-scores which are, generally, freely available online. Among them, we can cite:

- Shooting features, such as 2 points and 3 points field goals made (and missed) by the two teams, free throws made (and missed),
- Offensive features, such as offensive rebounds grabbed and assists served by the two teams,
- Defensive features, such as defensive rebounds, fouls made, steals and blocks made by the two teams.

Recently, syntheses of these features have been preferred to the use of simpler features. For example, the measurement of the number of points made (drawn) per 100 offensive (defensive) possessions by the team has been exploited. In particular, the relevance of the four Dean's factors (see Kubatko et al. (2007), Oliver (2004)) to predict the probability to win the game is well-known and agreed in the literature. The four factors are the following:

- Shooting: effective field goal percentage (eFG%): $\frac{(FG+0.5*3P)}{FGA}$,
- Turnovers: turnover percentage (TOV%): $\frac{TOV}{(FGA+0.44*FTA+TOV)}$,
- Offensive rebound percentage (ORB%): $\frac{ORB}{(ORB+OppDRB)}$,
- Free throws percentage (FT%): $\frac{FT}{FGA}$,

where FG is the number of field goals made by the considered team, $3P$ is the number of 3 points field goals made, TOV is the number of turnovers, FGA is the total number of attempted shots, FTA is the number of free throws attempted, ORB and DRB are, respectively, the number of offensive and defensive rebounds, Opp stays for opponent team (so, $OppDRB$ is the number of defensive rebounds grabbed by the opponent team), and FT is the number of made free throws.

Oliver (2004) argued that shooting counts for 40%, turnovers for 25%, rebounding for 20%, and free throws for 15%. Moreover, the effectiveness of these measures was demonstrated as

the coefficient of determination obtained by fitting a linear regression model to real data is about 0.9. Because basketball is played by two teams, we have to take the same four factors for the opponent team, in order to account for both teams' features. This means that the four Dean's factors become actually eight in our analysis (herein after, we make use of the notation *off* for the Dean's factors associated with the considered team, and of the notation *def* for the Dean's factors associated with the opponent team).

The logistic regression model reads as in the next equation:

$$\log \frac{P(Y_i = 1 | X)}{P(Y_i = 0 | X)} = X_i \beta, \quad (9)$$

where the left part of the equation represents the log-odd of Y_i conditional on X . Y is the response binary variable representing the outcome of the games, $Y_i \in \{0, 1\}$, $i = 1, \dots, g$, where g is the number of games. X_i is the i -th row of the design matrix X with g rows and p columns ($p = 8$, the eight Dean's factors used as explanatory variables, $eFG\%_{off}$, $eFG\%_{def}$, $TOV\%_{off}$, $TOV\%_{def}$, $ORB\%_{off}$, $ORB\%_{def}$, $FT\%_{off}$, $FT\%_{def}$, computed at the game level). β is a vector containing the p regression parameters associated with the explanatory variables. These parameters have to be estimated from the data.

A logistic regression model cannot be applied in a straightforward manner to estimate the probability to win for each lineup, because a single lineup does not play the full match, thus making it impossible to determine the *Outcome* variable (win = 1, defeat = 0) for that quintet. To deal with this issue, we adopt the following strategy: we use, as training set, the matrix of features X including Dean's factors computed at the single game level (i.e., each row of the dataset corresponds to a single game), in order to estimate, through a vector $\hat{\beta}$, the coefficients of the logistic regression model (i.e., the ones contained in the vector β). Then, using the dataset \tilde{X} where the Dean's factors are computed at the single lineup level (i.e., each row of the dataset corresponds now to a lineup), we predict the probability to win the game $P(Win)_{L_j}$ on each lineup L_j by using the vector $\hat{\beta}$ estimated from the training set and the Dean's factors associated with that lineup and its adversary lineups (an average of the Dean's factors is taken, considering that the adversary lineups are not fixed). More in detail, let \tilde{X}_j the j -th row of the matrix \tilde{X} with l rows (where l is the number of lineups considered) and $p = 8$ columns (expressing the eight Dean's factors computed at the lineup level), and let $\hat{\beta}$ be the vector of estimated coefficients from the previous training step. The above-described strategy makes it possible to express the probability to win the game for the lineup L_j as in the next equation:

$$P(Win)_{L_j} = \frac{\exp(\tilde{X}_j \hat{\beta})}{1 + \exp(\tilde{X}_j \hat{\beta})}, j = 1, \dots, l. \quad (10)$$

In the third step of our methodological approach, we apply the estimates of the winning probabilities provided by Eq. (10) to approximate the generalized characteristic function $v_1(\cdot)$ used as input for the computation of the generalized Shapley value $\phi_i^{NR}(N, v_1)$, as in Eqs. (2) and (3). Similarly, we combine such estimates with the ones of the occurrence probabilities (see the end of Sect. 2.1) to approximate the generalized characteristic function $v_2(\cdot)$ used as input for the computation of the generalized Shapley value $\phi_i^{NR}(N, v_2)$, as in Eqs. (4) and (5). Finally, the generalized Shapley values are approximated by using Eq. (7).

3 Data

The data used in the application have been extracted from the play-by-play of all NBA games (both regular seasons and play-offs) for the seasons from 2004–2005 to 2020–2021, for a total of 17 seasons. For each game and for both home and away teams, we have the information about the event type (such as start/end of the period, made/missed 2 points shot, made/missed 3 points shot, made/missed free throw, offensive/defensive rebound, assists, steal, block, foul) associated with the exact moment in which that event happens and also associated with the lineups of both the two teams. When the event is a shot (made or missed) we also have available the position on the court in terms of x -axis and y -axis coordinates, respectively related to court length and court width. Features related to the dependent variable (*Outcome*) and to Dean's factors for both X and \tilde{X} have been generated from those of the play-by-play dataset, which has been made available thanks to a friendly agreement with BigDataBall Company (UK) (www.bigdataball.com), which collected the play-by-play of all the NBA regular season and play-off games starting from 2004. We have retrieved computed WS, WS48, BPM and VORP values from Basketball Reference website (www.basketball-reference.com/).

4 Application

We apply the proposed strategy to the above-described NBA dataset. Firstly, we estimate a logistic regression model as in Eq. (9), based on the training set represented by the full set of games, which refers to the seasons from 2004–2005 to 2020–2021 (both regular seasons and play-offs) in order to increase as much as possible the sample dimension. Such training set counts for $g = 21$, 735 games. Starting from the dataset described in Sect. 3, we retrieve, at the single game level, the values of the dichotomous variable (*Outcome*), which assumes value 1 if the considered team won the game, 0 otherwise¹². We also retrieve the values of the Dean's factors, where, as already mentioned, with the term *off* we refer to the factors computed for the considered team, and, with the term *def*, to the factors computed for the opponent team. In order to let the effect on the Outcome of Dean's factors comparable, we have normalized all the features using a z-score transformation. We also compute, at the single game level, other box-score statistics, such as assists (*AST*), blocks (*BLK*), and fouls (*FLS*), for both teams. Those variables have also been normalized to z-scores (Table 1).

Despite Kubatko et al. (2007), for a similar analysis, estimated a linear regression model along with the Ordinary Least Squares (OLS) method, we can assert that logistic regression is preferable to OLS linear regression when the dependent variable presents a binary outcome. In fact, using OLS, the predicted values for the probability of success may be smaller than zero or greater than one (Wooldridge 2010). Logistic regression results, together with the McFadden pseudo R^2 (McFadden 1979), are reported in Table 1.

The high value for the McFadden pseudo R^2 (0.581) for the model with just Dean's factors as features (first column of Table 1) indicates a very good model fit. In fact, McFadden R^2 values tend to be considerably lower than those of the traditional R^2 index and values between 0.2 to 0.4 represent an excellent fit (McFadden 1979). The addition to the model of assists, blocks and fouls as explanatory variables increases the McFadden R^2 at a very limited extent, despite the coefficients associated with both assists' (*AST*) and fouls' (*FLS*) features turn out to be statistically significant.

¹² It has not been possible to distinguish between home and away team. For this reason, the considered team is selected randomly between the two teams of the game.

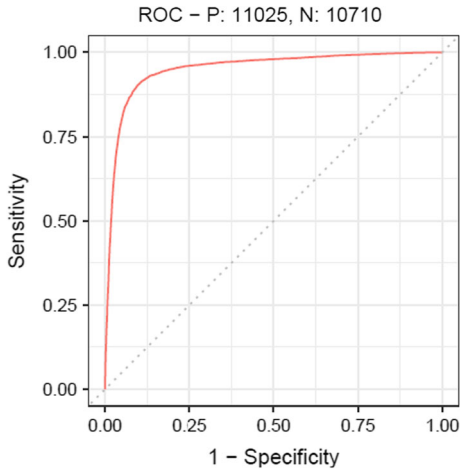
Table 1 Logistic regression with Maximum Likelihood (ML) - Results on the training set. $g=21,735$ games

	Dependent variable: outcome	
	(1)	(2)
eFG%_Off	10.255*** (0.147)	9.530*** (0.157)
eFG%_Def	-10.225*** (0.146)	-9.546*** (0.156)
TOV%_Off	-1.850*** (0.038)	-1.668*** (0.040)
TOV%_Def	1.749*** (0.037)	1.600*** (0.038)
ORB%_Off	0.998*** (0.026)	0.964*** (0.028)
ORB%_Def	-0.998*** (0.026)	-0.970*** (0.028)
FT%_Off	0.780*** (0.029)	0.620*** (0.041)
FT%_Def	-0.810*** (0.029)	-0.610*** (0.042)
AST_Off		0.681*** (0.032)
AST_Def		-0.598*** (0.031)
BLK_Off		-0.017 (0.032)
BLK_Def		0.058* (0.033)
FLS_Off		-0.553*** (0.035)
FLS_Def		0.526*** (0.035)
Constant	0.063*** (0.023)	0.071*** (0.024)
Observations	21,735	21,735
Log Likelihood	-6,304.493	-5,742.271
Akaike Inf. Crit.	12,626.990	11,514.540
McFadden pseudo R^2	0.581	0.619

Note: * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

It is worth discussing the effect of each feature on the probability of winning the game. Note that, being all the explanatory features expressed as z-scores, the coefficients can be compared in terms of the effect that each feature has on the log-odd of the outcome (winning the game). Offensive effective field goal percentage (eFG%_Off), being the sign of its estimated coefficient positive, plays a positive role on winning the game. More specifically, we

Fig. 1 Receiving Operation
Characteristic curve computed
from the full sample of 21,735
games: 11,025 positives
(outcome = 1) and 10,710
negatives (outcome = 0)



can say that, a unitary increase on the normalized offensive effective field goal percentage increases the log-odd of winning the game by about 10 (10.255). Its defensive counterpart (eFG%_Def) plays a negative role (a unitary increase on opponents' eFG% decreases the log-odd by 10.255). Offensive turnover percentage (TOV%_Off) decreases the probability to win (the log-odd decreases of 1.850 for an unitary increase of TOV%_Off) and TOV%_Def increases the same probability (the log-odd increases by 1.749). Offensive rebounds percentage (ORB%_Off) and defensive rebounds percentage (ORB%_Def) play, respectively, a positive (log-odd increases by 0.998) and a negative (log-odd decreases by 0.998) role on the winning probability. Offensive free throws percentage (FT%_Off) positively influences the probability to win (log-odd increase by 0.780), while defensive free throws percentage (FT%_Def) negatively influences the same probability (log-odd decreases by 0.810). All the coefficients are statistically significant for a level of $\alpha = 0.01$. Moreover, overall, the signs of the estimated coefficients are all consistent to those expected.

To quantify the robustness of our logistic regression model in terms of classification performance, its Receiving Operation Characteristic (ROC) curve (Krzanowski 2009) is depicted in Fig. 1. The Area Under the Curve (AUC) associated with this ROC, using as validation set the training set itself (which is justified by its large dimension ($g = 21,735$) compared with the number ($p = 8$) of explanatory variables), stands to 0.951, that is a very high value. A k -fold cross validation with $k = 10$ (McLachlan et al. 2005) is also performed as a robustness check. An average AUC of 0.946 is obtained. The Hit-rate accuracy measure (Bensic et al. 2005), which is computed as the ratio between the number of correctly classified games¹³ and the total number of games in the sample, stands to 0.903.

In light of the extremely good logistic regression model fit, we do not think it is required to include further explanatory features in the model. We use the following values (as in the column 2 of Table 1) to define the vector $\hat{\beta}$ to be used to determine $P(\text{Win})$ for each lineup L_j in the second step of our analysis:

$$[10.255, -10.255, -1.850, 1.749, 0.998, -0.998, 0.780, -0.810]'$$

¹³ I.e., the sum of the number of games where a win has been predicted for the considered team when it actually won, and the number of games where a defeat has been predicted for the considered team when it actually lost.

We choose to estimate the winning probabilities for the lineups of the Utah Jazz team¹⁴ in the 2020–21 regular season because they eventually finished the regular season with the best NBA record (52-20). As testified by online articles¹⁵, a heated public debate is underway on social networks on which Jazz player turns to be the most important one for the team, with a general preference for Rudy Gobert.

For the analysis, only the 107 different lineups that were on the court more than a cut-off value of a total of 4 minutes in the 2020/21 regular season have been considered. The duration these lineups were on the court covers about 86% of the total time of play of Utah Jazz during the full 2020/21 regular season. In those 107 lineups, a total of 11 different players rotating on the court has been counted. This constitutes a quite representative sample of the total number (426) of different possible lineups extracted from this set of 11 players.

Applying Eq. (10) to the considered lineups, we compute the value assumed by the winning probability for each lineup, then we determine the (estimates of the) generalized Shapley values, as in Eq. (7), for the following 11 players: Donovan Mitchell (*guard*), Bojan Bogdanović (*forward*), Joe Ingles (*forward*), Rudy Gobert (*center*), Mike Conley (*guard*), Jordan Clarkson (*guard*), Royce O'Neale (*forward*), Georges Niang (*forward*), Derrick Favors (*center*), Miye Oni (*guard*), Trent Forrest (*guard*).

In details, specializing Eq. (7) to the cases $k = 1, 2$ and $n = 11$, we compute the two following versions of the generalized Shapley value:

$$UWGS_i = \frac{1}{11|\mathcal{L}_i|} \sum_{L \in \mathcal{L}_i} \hat{v}_1(L), \quad (11)$$

$$WGS_i = \frac{1}{11|\mathcal{L}_i|} \sum_{L \in \mathcal{L}_i} \hat{v}_2(L). \quad (12)$$

Equation (11) represents the unweighted version of the generalized Shapley value. According to this version, the winning probability of each lineup contributes with the same weight, regardless of percentage of time spent by the lineup on the court. Equation (12), instead, represents the weighted version of the generalized Shapley value. This takes also into account the percentage of time spent by each lineup on the court (so, the winning probabilities of different lineups have typically different weights).

Table 2 reports the resulting generalized Shapley values for the Utah Jazz players.

It is interesting to note that the players who were on the court more often, such as Gobert (74.2% of the total time), O'Neale (77.6%) and Bogdanović (74.0%), present a higher UWGS (0.0487, 0.0446 and 0.0439, respectively) compared with bench players, such as Niang and Favors (0.0413 and 0.0360, respectively). An exception is Conley, who just played for the 54.1% of the time (due to an injury) but presents an UWGS of 0.0504 (best in its team). Another interesting result is that of Clarkson, who won the title of 6th man of the league (best player coming off the bench): he obtains a small UWGS (0.0360, just 3 team mates in its team make it worse).

¹⁴ The Utah Jazz is an NBA basketball team based in Salt Lake City. The Jazz competes in the Western Conference, Northwest Division. The team has been playing its home games at Vivint Arena since 1991. This franchise began playing as an expansion team in 1974 with the name of New Orleans Jazz as a tribute to New Orleans' history of jazz music. The Jazz moved to Salt Lake City (state of Utah) in 1979. From the late 1980s to the beginning of 2000s, the Jazz had its successful period thanks to a famous duo formed by John Stockton and Karl Malone, who helped the team to reach two consecutive NBA finals, in 1997 and 1998. When both Stockton and Malone moved on in 2003, the team went through an about 10 years dark period. With the development of Rudy Gobert and Donovan Mitchell into All-Stars, the Jazz launched itself back into title contention.

¹⁵ E.g., <https://thejnotes.com/2021/11/06/important-player-utah-jazz/2/>.

Table 2 Generalized Shapley values for the 11 selected players of the Utah Jazz team in the 2020/21 regular season (with rank in brackets). n lineups is the number of different lineups where that player was in; %time is the percentage of time that player was on the court, with respect to the time played by all the 107 considered lineups; the expressions of the two generalized Shapley values are detailed in Eqs. (11) and (12)

Player (i)	n lineups	%time	UWGS (rank)	WGS*100 (rank)
Royce O'Neale	78	77.6	0.0446 (4)	0.0380 (5)
Bojan Bogdanović	73	74.0	0.0439 (6)	0.0382 (4)
Rudy Gobert	71	74.2	0.0487 (2)	0.0454 (2)
Donovan Mitchell	67	66.5	0.0445 (5)	0.0389 (3)
Joe Ingles	65	54.3	0.0452 (3)	0.0359 (7)
Jordan Clarkson	61	44.1	0.0360 (8)	0.0242 (8)
Mike Conley	49	54.1	0.0504 (1)	0.0510 (1)
Derrick Favors	36	25.7	0.0324 (9)	0.0200 (9)
George Niang	30	24.7	0.0413 (7)	0.0368 (6)
Miye Oni	4	3.5	0.0005 (10)	0.0004 (10)
Trent Forrest	1	1.0	0.0002 (11)	0.0002 (11)

The second version of the generalized Shapley (WGS) is weighted by the amount of time each player i was on the court, for each lineup he took part of. This represents a different measure since it can happen that a certain player obtains a high value of marginal utility when playing in the lineup where it is mostly employed, but a low value when playing in a lineup where it is employed only for a limited amount of time. The WGS associates a higher weight with the marginal utility value obtained in the first lineup, and a lower weight to the marginal utility value obtained in the second lineup (i.e., it gives more importance to lineups in which player i plays more often). It is interesting to note that, despite, in general, players obtain smaller WGS compared to UWGS, Conley's weighted value is higher than the unweighted (0.0510 vs. 0.0504), while Clarkson, who already has a small UWGS, makes it even worse according to WGS (0.0242). Similar is the case of Favors, whose weighted version is way worse than the unweighted one (0.0200 vs. 0.0324).

Bootstrap (Efron 1992) confidence intervals are provided for both the weighted and unweighted versions of the generalized Shapley values, by resampling with replacement the lineups $n_r = 200$ times¹⁶. Box plots with the 99% bootstrap confidence intervals are displayed in Fig. 2. Moreover, Fig. 3 demonstrates the robustness of our measures to changes in the sample of lineups considered (i.e., according to different bootstrap samples). Globally, it looks that Conley ranks first most of the times, Gobert ranks second most of the times. Niang, Bogdanović, Mitchell, O'Neale and Ingles compete for the positions 3rd to 7th, Clarkson ranks 8th most of the times, Favors ranks 9th most of the times and Oni and Forrest compete for the last two positions.

To the sake of comparison with existing industry standard measures for players' contribution, we analyse the Pearson correlation of weighted and unweighted generalized Shapley values with WS, WS48, BPM and VORP, and the Kendall's Tau correlations between the

¹⁶ Although the variance of the Monte Carlo estimate of each generalized Shapley value has been addressed at the end of Sect. 2.1, such variance is difficult to compute exactly (only loose upper bounds on it are expected to be obtained quite easily, proceeding, e.g., in a similar way as in Gnecco et al. (2021), where the case of the variance of the Monte Carlo estimate of each Shapley value was investigated, for a different application of coalitional games). So, this justifies the generation of bootstrap confidence intervals in the present context.

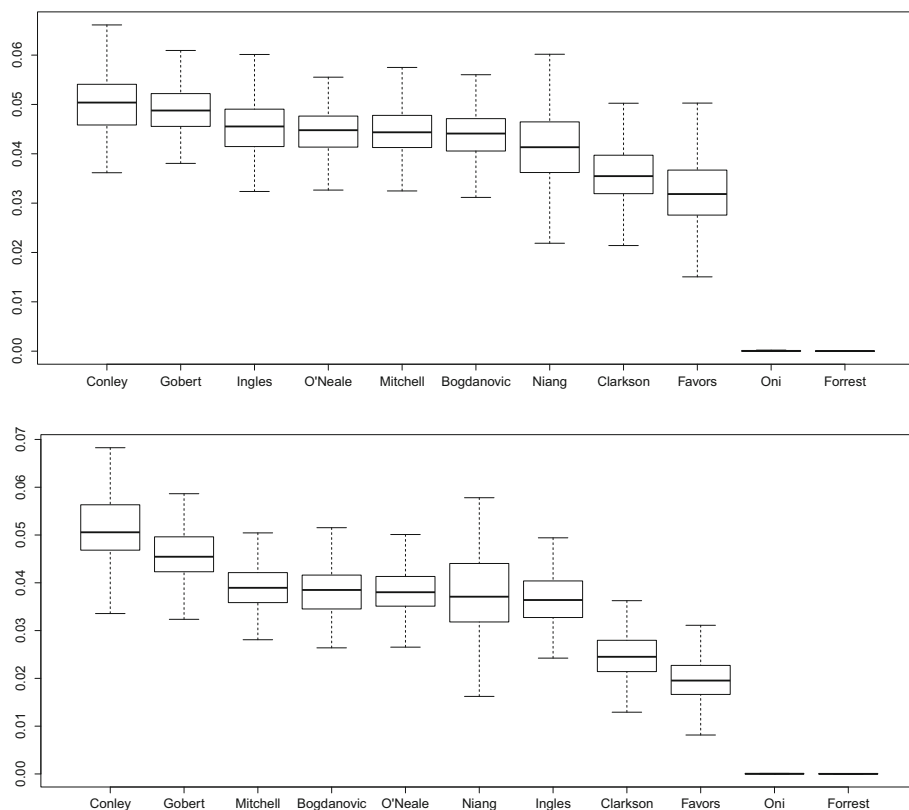


Fig. 2 Box plots with 1st, median and 3rd quartiles (boxes) and the 99% confidence intervals (whiskers) from $n_r = 200$ bootstrap samples, for the 11 Utah Jazz players during season 2020/21. UWGS (top chart), WGS (bottom chart)

rankings of players according to the same measures, as shown in Table 3 for the Utah Jazz players during the season 2020/21.

We can also use the generalized Shapley values to suggest best lineups. According to the UWGS, the five players with the largest average marginal utility are Conley, Gobert, Ingles, O'Neale and Mitchell (2 guards, 2 forwards and 1 center), while according to the WGS, the first five players in the ranking are Conley, Gobert, Mitchell, Bogdanović and O'Neale (still, 2 guards, 2 forwards and 1 center). Both lineups can be employed during the game. So, according to the Utah Jazz case study, the adoption of a more sophisticated (either classical or generalized) Shapley value definition with constraints, which takes into account that certain coalitions are excluded a-priori (Hiller 2018), seems to be not necessary.

However, it may be that, for each of the two cases above, the lineup made with those 5 players does not turn out to be the best one in terms of winning probability (finding it would involve a combinatorial optimization problem). Moreover, it is interesting to study the players' average marginal contributions conditional on the presence of a certain teammate (or certain teammates) on the court, or conditional on his/their absence. In the present context, the determination of average marginal conditional contributions is possible as we explicitly account for single lineups. More precisely, if the average is made conditional on the presence of $l_p \leq 4$ teammates, it is enough to modify Eq. (7) by reducing the number n of players of

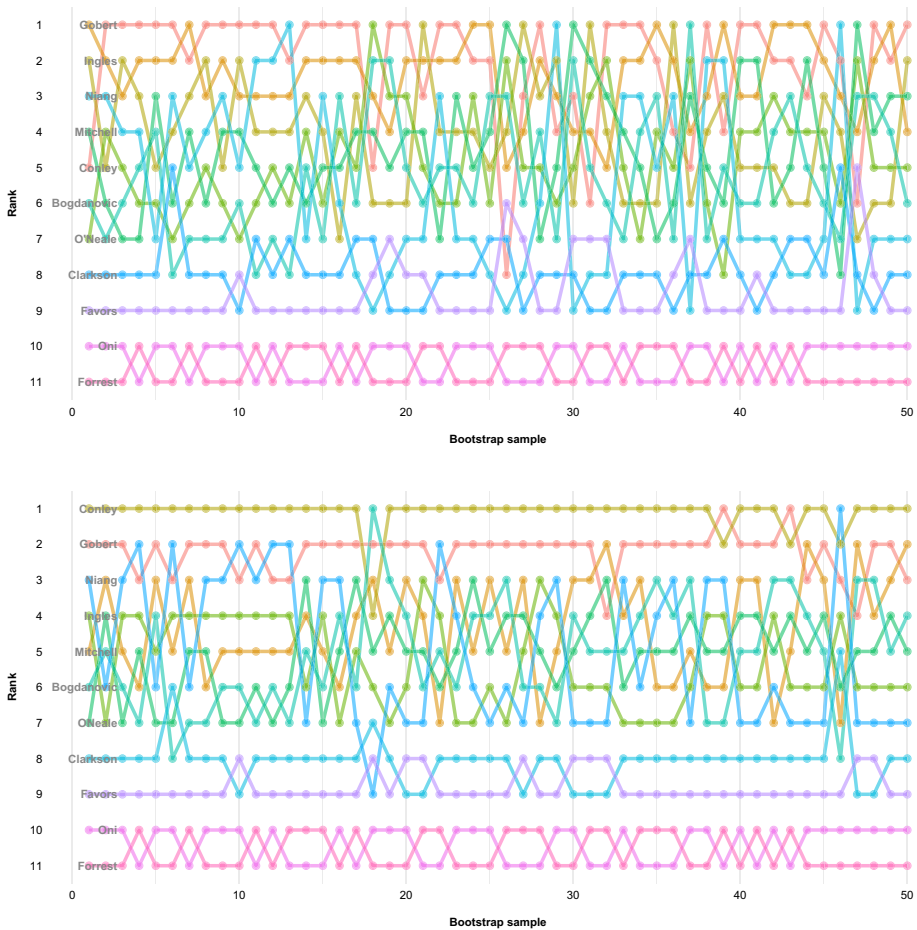


Fig. 3 Bump chart reporting the rank of Utah Jazz players according to UWGS (top) and to WGS (bottom) and the bootstrap samples (first 50, for clarity). Season 2020/21

the team by l_p (as l_p players are now fixed in the lineup, whereas only $n - l_p$ are “rotating” to produce each average), and restricting the lineups \mathcal{L}_i to those containing both the specific player i and the l_p fixed players. Similarly, if the average is made conditional on the absence of $l_a < n - 5$ teammates, it is enough to modify Eq. (7) by reducing the number n of players of the team by l_a (as l_a players are now absent, whereas only $n - l_a$ are “rotating” to produce each average), and restricting the lineups \mathcal{L}_i to those containing the specific player i and excluding each of the l_a absent players.

These variations may be interesting because a player might perform better, e.g., when a specific teammate (or subset of teammates) is on the court compared to the case in which that teammate (or subset of teammates) is not in. Furthermore, in specific circumstances, it might be necessary to form lineups by considering specific constraints, e.g., the coach may want to build the lineup around a specific player (so, we propose to pick the 4 teammates that better perform conditional on the presence of him/her on the court). Similarly, the presence of an

Table 3 Pearson correlation (top) and Tau-Kendall rank correlation among the two generalized Shapley values and the 4 industry-standard measures adopted for players' contributions. Utah Jazz players. Season 2020/21

	WS	WS48	BPM	VORP	UWGS	WGS
<i>Pearson</i>						
WS	1.00	.786	.851	.929	.822	.850
WS48		1.00	.841	.750	.627	.710
BPM			1.00	.943	.794	.833
VORP				1.00	.751	.784
UWGS					1.00	.968
WGS						1.00
<i>Kendall's Tau</i>						
WS	1.00	.709	.709	.836	.855	.708
WS48		1.00	.709	.655	.709	.636
BPM			1.00	.764	.709	.636
VORP				1.00	.691	.582
UWGS					1.00	.782
WGS						1.00

injured player would modify the number of players effectively available (so, we propose to form the best lineup conditional on his/her absence).

We present two specific applications. The first one is related to a case in which the coach wants to make Rudy Gobert part of the lineup. To do so, we propose to apply a greedy algorithm, that, by re-computing the generalized Shapley values conditional on specific constraints, permits, using a multi-step approach, to find approximately the best lineup. The first step of the algorithm aims at finding the second player for the lineup. We find him to be the player with the highest average marginal contribution conditional on the presence of Gobert (i.e., just considering the lineups in which he is on the court). In the second step we find the third player, i.e., the one whose average marginal contribution, conditional on the presence on the court of both Gobert and the other chosen player, is the highest. The algorithm continues until five players have been chosen.

Results of this application based on WGS are reported in Table 4. In the first step, Mitchell is chosen because he reports the highest WGS conditional on the presence of Gobert. In this step Conley and O'Neale present similar WGS to that of Mitchell. In the second step, Conley is chosen because he reports the highest WGS conditional on the presence of Gobert and Mitchell on the court. This time, Conley's WGS is way larger than those of the other players. In the third step, Ingles is chosen because he reports the highest WGS conditional on the presence of Gobert, Mitchell and Conley on the court. However, O'Neale and Bogdanović report very similar WGS. In the fourth and last step, O'Neale is chosen because he reports the highest WGS conditional on the presence of Gobert, Mitchell, Conley and Ingles on the court.

The second application considers the case in which Mike Conley is injured and it is still based on WGS. The greedy algorithm starts by finding the player with the highest value conditional on the absence of Conley and proceeds similarly as in the first application. In the first step (Table 5), Ingles is chosen because he reports the highest WGS conditional on the absence of Conley. In the second step, Bogdanović is chosen because he reports the highest WGS conditional on the absence of Conley and the presence of Mitchell on the court. In the third step, Niang is chosen because he reports the highest WGS conditional to the absence of

Table 4 Greedy algorithm results for the case of Rudy Gobert chosen by the coach

Player	(1st step)	(2nd step)	(3rd step)	(4th step)
Donovan Mitchell	0.0551	—	—	—
Mike Conley	0.0535	0.0821	—	—
Joe Ingles	0.0383	0.0547	0.0941	—
Royce O’Neale	0.0512	0.0707	0.0924	0.1076
Bojan Bogdanović	0.0485	0.0664	0.0922	n.a.
Jordan Clarkson	0.0337	0.0231	n.a.	n.a.
George Niang	0.0398	n.a.	n.a.	n.a.
Miye Oni	0.0000	n.a.	n.a.	n.a.
Trent Forrest	n.a.	n.a.	n.a.	n.a.
Derrick Favors	n.a.	n.a.	n.a.	n.a.

Conditional WGS of the algorithm steps for each player are reported. In the first step, Mitchell is chosen because he reports the highest WGS conditional on the presence of Gobert. In the second step, Conley is chosen because he reports the highest WGS conditional on the presence of Gobert and Mitchell on the court. In the third step, Ingles is chosen because he reports the highest WGS conditional on the presence of Gobert, Mitchell and Conley on the court. In the fourth and last step, O’Neale is chosen because he reports the highest WGS conditional on the presence of Gobert, Mitchell, Conley and Ingles on the court. “n.a.” means that the related player never played conditional on the presence on the court of chosen (until that step) players

Conley and the presence of Ingles and Bogdanović on the court. In the fourth step O’Neale and Gobert are chosen because they report the highest WGS conditional to the absence of Conley and the presence of Ingles, Bogdanović and Niang on the court. It worth noting that the lineups chosen in the two applications differ. The most surprisingly evidence is the absence of Mitchell in a lineup where Conley is injured. This evidence makes our proposal based on the computation of average conditional marginal contributions relevant, as players’ average marginal contributions actually appear different when evaluated conditional on the presence/absence of specific teammates on the court.

5 Conclusions

Data analytics in basketball, likewise in other professional team sports, is increasingly adopted, and in recent years has percolated to many academic fields, such as Computer Science, Applied Mathematics, Statistics, Management Science, and Economics. In the era of *Big Data*, where online platforms make available live streams of large amounts of data, more and more often managers and staff of professional teams face the need for extracting useful information for the monitoring of the performance of their team, as well as their single players.

Aware of the fact that a team may be viewed as a network of players who cooperate with the same purpose, in this manuscript we have been dedicated to the development of a methodological strategy aimed at measuring the average marginal contributions of the players to achieve the goal of their team (i.e., to win the match).

Each player’s average marginal utility has been computed here by means of a generalized (both unweighted and weighted) version of the Shapley value, where the generalized characteristic function has been expressed in terms of the probability to win the game, in turn estimated using a logistic regression strategy.

Table 5 Greedy algorithm results for the case of Mike Conley injured. Conditional WGS of the algorithm steps for each player are reported

Player	(1ststep)	(2ndstep)	(3rdstep)	(4thstep)
Joe Ingles	0.0341	–	–	–
Bojan Bogdanović	0.0316	0.0487	–	–
George Niang	0.0204	0.0206	0.0809	n.a.
Rudy Gobert	0.0337	0.0385	0.0529	0.0925
Royce O’Neale	0.0307	0.0459	0.0579	0.0925
Donovan Mitchell	0.0292	0.0440	0.0596	n.a.
Jordan Clarxson	0.0139	0.0131	0.0258	n.a.
Derrick Favors	0.0225	0.0370	0.0580	n.a.
Miye Oni	0.0004	0.0005	n.a.	n.a.
Trent Forrest	0.0002	n.a.	n.a.	n.a.

In the first step, Ingles is chosen because he reports the highest WGS conditional on the absence of Conley. In the second step, Bogdanović is chosen because he reports the highest WGS conditional on the absence of Conley and the presence of Ingles on the court. In the third step, Niang is chosen because he reports the highest WGS conditional on the absence of Conley and the presence of Ingles and Bogdanović on the court. In the fourth step O’Neale and Gobert are chosen because they reports the highest WGS conditional on the absence of Conley and the presence of Ingles, Bogdanović and Niang on the court. “n.a.” means that the related player never played conditional on the presence on the court of chosen (until that step) players

With this work we place ourselves in the literature aimed at finding a measure for the player’s contribution by proposing an approach that gathers most of the advantages (and avoids disadvantages) of the industry-standard ones, and that permits to do lineup management. Moreover, we do so by using (for the first time, to the best of our knowledge) a game-theoretical approach based on generalized Shapley values, that has never been applied to basketball before. In summary, the generalized Shapley value evaluates the importance of each player by considering him/her as an individual who is a member of a larger team, and achieves this goal in a way that is more structured than other industry-standard measures. Indeed, in order to evaluate a player-specific quantity (his/her importance in the team), our proposed approach uses features at the team level as a starting point, following a multi-step approach. First, an underlying machine-learning model (logistic regression, in our case) is trained to predict the winning probability based on several features associated to a lineup. In this way we are addressing the performance of the whole lineup in this first step, not of the individual (which is addressed in a successive step). It is worth remarking that these features turn out to have a high predictive capability. In the final step, by averaging on several lineups, we obtain an estimate of the importance of each player.

Overall, this work targets to help managers, coaches and the staff for planning their strategies about the team and the players, by providing them a robust measure of player’s marginal utility along with a strategy for the management of the lineup. The proposed approach of analysis could be used by coaches, e.g., in the following ways:

- by monitoring the features associated with a lineup, the coach could opt for replacing a player, in case the current estimate of the winning probability turned out to be too low. Of course, this would require estimating the coefficients of the logistic regression model using the data available before the current match (this is not an issue, since also data from past seasons could be used to that aim);

- the choice of a specific lineup could be guided by players' ranks based on generalized Shapley values, possibly conditional on other constraints (e.g., the presence/absence of a specific player in the lineup). This would require estimating the generalized Shapley value using the data available before the current match, but coming from the same season, because of the possibly different composition of the team in consecutive seasons. In order to increase the amount of observed lineups available to get this estimate, one could take into account also data coming from the preseason and (when available) from practice games, especially for the first part of the season, for which no data coming from other official matches in the season are available.

We conclude discussing possible future developments. First, it may be interesting to employ a version of the generalized Shapley value that i) excludes a-priori some coalitions, in such a way to account for impossible lineups, ii) it is allowed to be used for the analysis of player's contribution at single game level. Second, it may be worth assessing the impact of including additional features (for example, those coming from the position of the ball, which may be retrieved by the joint use of computer vision and machine learning techniques, see Giuffrida et al. (2019)) to the model of the generalized characteristic function.

Acknowledgements The authors would like to thank Serhat Ugur (CEO, BigDataBall) for having kindly given access to his data, Sergio Pietro Destefanis (University of Salerno), an anonymous reviewer and the associate editor for their precious feedback. This work was made in the framework of the project FARB-ORSA 198493 (for R.M.) funded by University of Salerno, and the Galileo 2021 project "Automatic Movement Analysis Techniques for Applications in Cognitive/Motor Rehabilitation" (for G.G.). The authors contributed equally to the work.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Appendix: alternative definition of the generalized characteristic function

As an alternative to Eq. (6), another way of defining $v_k(T)$ for $k = 1, 2$ could be the following:

$$v_k(T) = \begin{cases} 0 & \text{if } |T| < m = 5 \\ v_k(\{T_1, T_2, T_3, T_4, T_5\}) & \text{if } |T| = m = 5 \\ 0 & \text{if } |T| > m = 5. \end{cases} \quad (13)$$

In this case, the Shapley value and the generalized Shapley value (1) would be the same, not depending the value assumed by the function $v_k(\cdot)$ on T on the order of the elements belonging to T that are successive to the fifth one. Nevertheless, assuming positive values for $v_k(\{T_1, T_2, T_3, T_4, T_5\})$, the marginal utility of each player $i \neq T_1, T_2, T_3, T_4, T_5$ when entering in the sixth position would be negative and would not depend on the particular choice of the player $i \neq T_1, T_2, T_3, T_4, T_5$ (i.e., it would be independent of the ability of that player). Concluding, by replacing Eq. (6) with Eq. (13), one would achieve a unique worth (i.e., 0) for the grand coalition (which is necessary for the computation of the Shapley value), at the cost of making negative the marginal contribution of each player that enters in sixth position

(regardless of the actual worth of that player). For this reason, we have preferred to use the Definition (6), which better agrees with intuition, and motivates the use of the generalized Shapley value (1) in the present context, instead of simply the Shapley value.

References

- Auer, B. R., & Hiller, T. (2015). On the evaluation of soccer players: A comparison of a new game-theoretical approach to classic performance measures. *Applied Economics Letters*, 22(14), 1100–1107.
- Barrientos, A. F., Sen, D., Page, G. L., & Dunson, D. B. (2019). Bayesian inferences on uncertain ranks and orderings. arXiv preprint [arXiv:1907.04842](https://arxiv.org/abs/1907.04842).
- Beckler, M., Wang, H., & Papamichael, M. (2013). NBA oracle. Zulettz besucht am, 17(20082009.9).
- Bensic, M., Sarlija, N., & Zekic-Susac, M. (2005). Modelling small-business credit scoring by using logistic regression, neural networks and decision trees. *Intelligent Systems in Accounting, Finance & Management*, 13(3), 133–150.
- Castro, J., Gómez, D., & Tejada, J. (2009). Polynomial calculation of the Shapley value based on sampling. *Computers & Operations Research*, 36(5), 1726–1730.
- Cheng, G., Zhang, Z., Kyebambe, M. N., & Kimbugwe, N. (2016). Predicting the outcome of NBA playoffs based on the maximum entropy principle. *Entropy*, 18(12), 450.
- Cooper, W. W., Ruiz, J. L., & Sirvent, I. (2009). Selecting non-zero weights to evaluate effectiveness of basketball players with DEA. *European Journal of Operational Research*, 195, 563–574.
- Cortes, C., & Vapnik, V. (1995). Support-vector networks. *Machine Learning*, 20(3), 273–297.
- Deshpande, S. K., & Jensen, S. T. (2016). Estimating an NBA player's impact on his team's chances of winning. *Journal of Quantitative Analysis in Sports*, 12(2), 51–72.
- Efron, B. (1992). Bootstrap methods: Another look at the jackknife. *Breakthroughs in Statistics* (pp. 569–593). New York: Springer.
- Engelmann, J. (2017). Possession-based player performance analysis in basketball (adjusted+/-and related concepts). In *Handbook of statistical methods and analyses in sports* (pp. 231–244, 1st edn). New York: Chapman and Hall/CRC.
- Fearnhead, P., & Taylor, B. M. (2011). On estimating the ability of NBA players. *Journal of Quantitative Analysis in Sports*, 7(3). <https://doi.org/10.2202/1559-0410.1298>.
- Giuffrida, D., Benetti, G., De Martini, D., & Facchinetti, T. (2019). Fall detection with supervised machine learning using wearable sensors. In *2019 IEEE 17th International Conference on Industrial Informatics (INDIN)* (Vol. 1, pp. 253–259). IEEE. Helsinki, Finland.
- Gnecco, G., Hadads, Y., & Sanguineti, M. (2021). Public transport transfers assessment via transferable utility games and Shapley value approximation. *Transportmetrica A: Transport Science*, 17(4), 540–565.
- Grassetti, L., Bellio, R., Di Gaspero, L., Fonseca, G., & Vidoni, P. (2021). An extended regularized adjusted plus-minus analysis for lineup management in basketball using play-by-play data. *IMA Journal of Management Mathematics*, 32(4), 385–409.
- Gudmundsson, J., & Horton, M. (2017). Spatio-temporal analysis of team sports. *ACM Computing Surveys (CSUR)*, 50(2), 1–34.
- Hernández-Lamonedá, L., & Sánchez-Sánchez, F. (2010). Rankings and values for team games. *International Journal of Game Theory*, 39(3), 319–350.
- Hiller, T. (2018). The effects of excluding coalitions. *Games*, 9(1). <https://doi.org/10.3390/g9010001>.
- Hiller, T. (2015). The importance of players in teams of the German Bundesliga in the season 2012/2013—a cooperative game theory approach. *Applied Economics Letters*, 22(4), 324–329.
- Hiller, T. (2018). *On the stability of couples*. *Games*, 9(3), 48.
- Hofler, R. A., & Payne, J. E. (2006). Efficiency in the National Basketball Association: A stochastic frontier approach with panel data. *Managerial and Decision Economics*, 27(4), 279–285.
- Hosmer, D. W., Jr., & Lemeshow, S. (2013). *Applied logistic regression & sturdivant*. Hoboken: Wiley.
- Ilardi, S. (2007). Adjusted plus-minus: An idea whose time has come. Retrieved from 82games.com (<http://www.82games.com/ilardi1.htm>).
- Jaynes, E. T. (1957). Information theory and statistical mechanics. *Physical Review*, 106(4), 620.
- Kalman, S., & Bosch, J. (2020) NBA lineup analysis on clustered player tendencies: A new approach to the positions of basketball & modeling lineup efficiency of soft lineup aggregates. 42 Analytics.
- Kolykhalova, K., Gnecco, G., Sanguineti, M., Volpe, G., & Camurri, A. (2020). Automated analysis of the origin of movement: An approach based on cooperative games on graphs. *IEEE Transactions on Human-Machine Systems*, 50(6), 550–560.

- Krzanowski, W. J. (2009). *ROC curves for continuous data & Hand* (1st edn) Boca Raton: CRC Press, New York: Chapman and Hall/CRC.
- 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). <https://doi.org/10.2202/1559-0410.1070>.
- Landwehr, N., Hall, M., & Frank, E. (2005). Logistic model trees. *Machine Learning*, 59(1–2), 161–205.
- Langley, P., Iba, W., & Thompson, K. (1992). An analysis of Bayesian classifiers. In AAAI '92: *Proceedings of the tenth national conference on Artificial intelligence*, pp. 223–228. San Jose, CA: AAAI Press.
- Loeffelholz, B., Bednar, E., & Bauer, K. W. (2009). Predicting NBA games using neural networks. *Journal of Quantitative Analysis in Sports*, 5(1). <https://doi.org/10.2202/1559-0410.1156>.
- Maschler, M., Solan, E., & Zamir, S. (2013). *Game Theory*. Cambridge: Cambridge University Press.
- Matthiopoloulou, O., Bardy, B., Gnecco, G., Motter, D., Sanguineti, M., & Camurri, A. (2020). A computational method to automatically detect the perceived origin of full-body human movement and its propagation. *ICMI '20 Companion: Companion Publication of the 2020 International Conference on Multimodal Interaction*, pp. 449–453.
- McFadden, D. (1979). Quantitative methods for analysing travel behavior of individuals: Some recent developments. In D. Hensher & P. Stopher (Eds.), *Behavioral travel modeling* (pp. 279–318). London: Croom-Heim.
- McLachlan, G. J., Do, K. A., & Ambrose, C. (2005). *Analyzing microarray gene expression data*. Hoboken: Wiley.
- Metulini, R., & Le Carre, M. (2020). Measuring sport performances under pressure by classification trees with application to basketball shooting. *Journal of Applied Statistics*, 47(12), 2120–2135.
- Metulini, R., Manisera, M., & Zuccolotto, P. (2018). Modelling the dynamic pattern of surface area in basketball and its effects on team performance. *Journal of Quantitative Analysis in Sports*, 14(3), 117–130.
- Michalak, T. P., Szczepański, P. L., Rahwan, T., Chrobak, A., Brânzei, S., Wooldridge, M., & Jennings, N. R. (2014). Implementation and computation of a value for generalized characteristic function games. *ACM Transactions on Economics and Computation*, 2(4), 1–35. <https://doi.org/10.1145/2665007>.
- Miljković, D., Gajić, L., Kovačević, A., & Konjović, Z. (2010). The use of data mining for basketball matches outcomes prediction. In *IEEE 8th International Symposium on Intelligent Systems and Informatics*, pp. 309–312.
- Mishra, S. K. (2016). Shapley value regression and the resolution of multicollinearity. Available at SSRN, 2797224. <https://doi.org/10.2139/ssrn.2797224>.
- Moreno, P., & Lozano, S. (2014). A network DEA assessment of team efficiency in the NBA. *Annals of Operations Research*, 214(1), 99–124.
- Nikolaïdis, Y. (2015). Building a basketball game strategy through statistical analysis of data. *Annals of Operations Research*, 227(1), 137–159.
- Nowak, A., & Radzik, T. (1994). The Shapley Value for n-person games in generalized characteristic function form. *Games and Economic Behavior*, 6(1), 150–161.
- Oliver, D. (2004). *Basketball on paper: Rules and tools for performance analysis* (1st ed.). Sterling: Potomac Books, Inc.
- Oliver, D. (2004). Roboscout and the four factors of basketball success. *Journal of Basketball studies* (blog). Retrieved from http://www.rawbw.com/~deano/articles/20040601_roboscout.htm.
- Page, G. L., Barney, B. J., & McGuire, A. T. (2013). Effect of position, usage rate, and per game minutes played on NBA player production curves. *Journal of Quantitative Analysis in Sports*, 9(4), 337–345.
- Piette, J., Anand, S., & Zhang, K. (2013). Scoring and shooting abilities of NBA players. *Journal of Quantitative Analysis in Sports*, 6(1). <https://doi.org/10.2202/1559-0410.1194>.
- Rice, J. A. (2005). *Mathematical statistics and data analysis* (2nd ed.). Wadsworth: Belmont.
- Rosenbaum, D. (2004). Measuring how NBA players help their teams win. Retrieved from 82Games.com (<http://www.82games.com/comm30.htm>).
- Sanchez, E., & Bergantiños, G. (1997). On values for generalized characteristic functions. *OR Spectrum*, 19, 229–234.
- Sandri, M., Zuccolotto, P., & Manisera, M. (2020). Markov switching modelling of shooting performance variability and teammate interactions in basketball. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 69(5), 1337–1356.
- Sarlis, V., & Tjortjis, C. (2020). Sports analytics-evaluation of basketball players and team performance. *Information Systems*, 93, 101562.
- Shapley, L. S. (1953). A value for n-person games. In H. W. Kuhn & A. W. Tucker (Eds.), *Contributions to the theory of games* (Vol. 2, pp. 307–17). Princeton, NJ: Princeton University Press.
- Sill, J. (2010). Improved NBA adjusted+/-using regularization and out-of-sample testing. In *Proceedings of the 2010 MIT Sloan Sports Analytics Conference*.

- Terner, Z., & Franks, A. (2021). Modeling player and team performance in basketball. *Annual Review of Statistics and Its Application*, 8, 1–23.
- Thabtah, F., Zhang, L., & Abdelhamid, N. (2019). NBA game result prediction using feature analysis and machine learning. *Annals of Data Science*, 6(1), 103–116.
- Wooldridge, J. M. (2010). *Econometric analysis of cross section and panel data*. Cambridge: MIT press.
- Yang, C. H., Lin, H. Y., & Chen, C. P. (2014). Measuring the efficiency of NBA teams: Additive efficiency decomposition in two-stage DEA. *Annals of Operations Research*, 217(1), 565–589.
- Yan, T., Kroer, C., & Peysakhovich, A. (2020). Evaluating and rewarding teamwork using cooperative game abstractions. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, & H. Lin (Eds.), *Advances in Neural Information Processing Systems* (Vol. 33, pp. 6925–6935).
- Zhang, G. P. (2000). Neural networks for classification: A survey. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 30(4), 451–462.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.