

The class of microarray games and the relevance index for genes

Stefano Moretti · Fioravante Patrone ·
Stefano Bonassi

Received: 21 November 2006 / Accepted: 13 May 2007 / Published online: 27 September 2007
© Sociedad de Estadística e Investigación Operativa 2007

Abstract Nowadays, microarray technology is available to generate a huge amount of information on gene expression. This information must be statistically processed and analyzed, in particular, to identify those genes which are useful for the diagnosis and prognosis of specific diseases. We discuss the possibility of applying game-theoretical tools, like the Shapley value, to the analysis of gene expression data.

Via a “truncation” technique, we build a coalitional game whose aim is to stress the relevance (“sufficiency”) of groups of genes for the specific disease we are interested in. The Shapley value of this game is used to select those genes which deserve further investigation. To justify the use of the Shapley value in this context, we axiomatically characterize it using properties with a genetic interpretation.

Keywords Coalitional game · Shapley value · Power index · Gene expression · Microarray

Mathematics Subject Classification (2000) 91A12 · 91A80 · 92B15 · 92C40

The authors are grateful to two anonymous referees for their extremely helpful comments.

An earlier version of this paper was presented at the VI Spanish Meeting on Game Theory and Practice, July 12–14, 2004, Elche, Spain.

S. Moretti gratefully acknowledges the financial support of the EU project NewGeneris, European Union 6th FP (FOOD-CT-2005-016320).

S. Moretti (✉) · S. Bonassi

Unit of Molecular Epidemiology, National Cancer Research Institute, Largo R. Benzi 10,
16132 Genoa, Italy

e-mail: stefano.moretti@istge.it

S. Bonassi

e-mail: stefano.bonassi@istge.it

F. Patrone

DIPTeM, University of Genova, P.le Kennedy—Pad D, 16129 Genoa, Italy

e-mail: patrone@diptem.unige.it

1 Introduction

Proteins are the structural constituents of cells and tissues and may act as necessary enzymes for biochemical reactions in biological systems. Most genes contain the information for making a specific protein. This information is coded in genes by means of the deoxyribonucleic acid (DNA). *Gene expression* occurs when genetic information contained within DNA is *transcribed* into messenger ribonucleic acid (mRNA) molecules and then *translated* into the proteins.

Nowadays, the microarray technology allows for the quantification of the expression (i.e., the amount of mRNA) for genes under the same biological condition (for instance, a tumor). A microarray works by exploiting the ability of a given mRNA molecule to bind specifically to, or hybridize to, the DNA template from which it originated. By using an array containing many DNA samples, it can be determined, in a single experiment, the expression levels of hundreds or thousands of genes within a cell by measuring the amount of mRNA bound to each site on the array.

There are several different experimental platforms based on microarray technology (see, for instance, Parmigiani et al. 2003). However, a common objective of gene expression microarrays is to consistently generate a matrix of expression data, in which the rows (possibly thousands) index the genes and the columns (usually in the order of tens) index the study samples. Numbers in the matrix represent gene expression values which quantify the level of expression of genes in the samples.

The aim of this work is to address the problem of quantifying the relative relevance of genes in a complex scenario—such as the pathogenesis of a genetic disease—on the basis of the information provided by microarray experiments, taking into account the *level of interaction* among the genes.

Complex experimental artifacts associated with microarray data collection emphasize the need for pre-processing analysis of the data (for instance, the design of the arrays, the quality assessment of the rough data, the normalization procedures), with the goal to reduce systematic errors arising from several experimental procedures. Despite the reduction of experimental bias has been the objective of several works on microarray analysis in the last few years (see, for example, Dudoit et al. 2001; Smith and Speed 2003; Parmigiani et al. 2003), in practice the problem of completely removing the experimental variability is still unsolved and a statistical treatment of the data provided by microarrays is required.

For this reason, in our approach we refer to the *observed average* level of interaction of a group of genes, i.e., the average number of tumor samples in which such a group of genes can be considered responsible, according to a pre-defined causality principle, for the onset of the tumor: the higher the number of samples observed, the lower the probability that chance could affect the inferences provided by the model.

The basic idea of this model comes from the theory of coalitional games. In particular, we consider the framework of simple games, which have been widely applied to the analysis of the power of players in interaction situations as Councils, Parliament, etc. (Shapley and Shubik 1954; Banzhaf 1965; see Owen 1995 for a general introduction to this topic and a summary of these results). We adopt the same formal language of coalitional games for modelling the interaction among genes, considered as players, in connection with a biological condition of interest, e.g., the pathogenesis

of a genetic disease or tumor. The game we consider originates from the comparison of two matrices of gene expression data; one from tumor samples and the other one from normal DNA (referent healthy subjects). We first use a discriminant method on each sample to split the whole set of genes in two sets, i.e., those genes showing an expression ratio largely different from normal samples, and those with expression levels corresponding to normal DNA samples. At this preliminary stage of the model, for each single gene, as in detail explained in Sect. 2, we use the interval boundaries containing most data in the normal distribution of that gene as cut-offs for discrimination (Becquet et al. 2002). We then introduce a causality relation (also called *sufficiency principle*) which directly determines the characteristic function of the game. An interpretation of the biological meaning of a relevance index, used for measuring the “power” of each gene in inducing the tumor, is given and it turns out to coincide with the Shapley value of the game considered.

We start with some preliminary notations in the next section. In Sect. 3 the class of microarray games is introduced starting from the general notion of the *sufficiency principle*, and some basic properties and examples of such games are reported. In Sect. 4 an axiomatic characterization of the Shapley value is given by means of five properties suitable to genetic interpretation of this index. Section 5 concludes with some considerations on related works and future research.

2 Preliminary notations

Let $N = \{1, 2, \dots, n\}$ be a set of n genes, $S_R = \{s_1^R, s_2^R, \dots, s_r^R\}$ be a set of r reference samples, i.e., the set of cells from normal tissues and let $S_D = \{s_1^D, s_2^D, \dots, s_d^D\}$ be the set of cells from tissues with a genetic disease.

The goal of a microarray experiment is to associate to each sample $j \in S_R \cup S_D$ an *expression profile* $A(j) = (A_{ij})_{i \in N}$, where $A_{ij} \in \mathbb{R}$ represents the *expression value* of the gene i in sample j . Globally, such expression values will be indicated as the *data set* of the microarray experiment. In the following we will refer to the data set resulting from the pre-processed method usually called normalization (Dudoit et al. 2001; Smith and Speed 2003), which allows for comparison among expression intensities of genes from different samples. The data set can be expressed in the form of two real valued expression matrices $\mathbf{A}^{S_R} = (A_{ij}^{S_R})_{i \in N, j \in S_R}$ and $\mathbf{A}^{S_D} = (A_{ij}^{S_D})_{i \in N, j \in S_D}$. In summary, we will denote as a *microarray experimental situation* (MES) the tuple $E = \langle N, S_R, S_D, \mathbf{A}^{S_R}, \mathbf{A}^{S_D} \rangle$.

As the first step of our analysis, we are interested in understanding whether genes in samples from the set S_D are abnormally expressed with respect to the expression values showed in samples from the set S_R according to a certain discriminative criterion. More precisely, we refer to the set of abnormally expressed genes as the union of the set of over expressed genes (also called as *up regulated genes*) and the set of under expressed genes (also called *down regulated genes*).

We need to introduce some additional notations to deal with abnormally expressed genes. Note that gene $i \in N$ which results abnormally expressed on a sample $j \in S_D$ can be represented setting to 1 the value of a boolean variable $B_{ij} \in \{0, 1\}$. We call *boolean expression profile* of a sample $j \in S_D$ the vector $\mathbf{B}^j = (B_{ij})_{i \in N}$. A *discriminant method* can be expressed as a map m assigning to each expression profile from

tumor samples a corresponding boolean expression profile. Hence, all the information on the differences of gene expression of sample in S_D from the ones of sample in S_R can be summarized into a *boolean expression matrix* $\mathbf{B}^{E,m} \in \{0, 1\}^{N \times S_D}$.

Since for our purposes the relevant information is contained in the boolean expression matrix $\mathbf{B}^{E,m}$, in the sequel we identify the MES E and the discriminant method m with the matrix $\mathbf{B}^{E,m}$.

Example 1 Consider an MES $\hat{E} = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ such that \mathbf{A}^{S_R} is reported in the following table

	Sample 1	Sample 2	Sample 3	Sample 4
Gene 1	0.5	0.2	0.3	0.6
Gene 2	12	10	4	5
Gene 3	8	13	20	9
Gene 4	0.8	0.4	1.4	1.1

and \mathbf{A}^{S_D} is reported in the following one

	Sample 1	Sample 2	Sample 3
Gene 1	0.9	0.4	0.7
Gene 2	4.6	15	18
Gene 3	7	21	12
Gene 4	1	0.6	1.6

Now consider a very naive discriminant method \hat{m} for the two classes 1 and 0, where 1 labels abnormally expressed genes and 0 labels normally expressed genes and such that for each $i \in N$ and each $j \in S_D$

$$(\hat{m}(\mathbf{A}^{S_D}, \mathbf{A}^{S_R}))_{ij} = \begin{cases} 1, & \text{if } A_{ij}^{S_D} \geq \max_{h \in S_R} A_{ih}^{S_R} \\ & \text{or } A_{ij}^{S_D} \leq \min_{h \in S_R} A_{ih}^{S_R}, \\ 0, & \text{otherwise.} \end{cases}$$

Then the corresponding boolean expression matrix is the following

$$\mathbf{B}^{\hat{E}, \hat{m}} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Now, let us introduce some basic game theoretical notations. A *coalitional game* or *characteristic-form game* is a pair (N, v) , where N denotes the finite set of *players* and $v : 2^N \rightarrow \mathbb{R}$ the *characteristic function*, with $v(\emptyset) = 0$. If the set N of players is fixed, we identify a coalitional game (N, v) with the corresponding characteristic function v . A group of players $T \subseteq N$ is called a *coalition* and $v(T)$ is called the *value* of this coalition. A coalitional game (N, w) such that $w : 2^N \rightarrow [0, 1]$ is called

a $[0, 1]$ -game. We will denote the class of all $[0, 1]$ -games as \mathcal{W} , with $\mathcal{W} \subsetneq \mathcal{G}$, being \mathcal{G} the class of all coalitional games.¹

Let $\mathcal{C} \subseteq \mathcal{G}$ be a subclass of coalitional games. Given a set of players N , we denote by $\mathcal{C}^N \subseteq \mathcal{G}$ the class of coalitional games in \mathcal{C} with N as set of players.

The *unanimity game* (N, u_R) on $R \subseteq N$ is the game described by $u_R(T) = 1$ if $R \subseteq T$ and $u_R(T) = 0$, otherwise. Every coalitional game (N, v) can be written as a linear combination of unanimity games in a unique way, i.e., $v = \sum_{S \subseteq N, S \neq \emptyset} \lambda_S(v) u_S$ (see, for instance, Owen 1995). The coefficients $\lambda_S(v)$, for each $S \in 2^N \setminus \{\emptyset\}$, are called *unanimity coefficients* of the game (N, v) .

A coalitional game (N, v) is *superadditive* if for all $S, T \subseteq N$ with $S \cap T = \emptyset$, we have $v(S) + v(T) \leq v(S \cup T)$. A coalitional game (N, v) is *monotonic* if for all $S \subseteq T \subseteq N$, we have $v(S) \leq v(T)$. For each $S \subseteq N$, $S \neq \emptyset$, and $i \in S$, the quantity $m_i(v, S) = v(S) - v(S \setminus \{i\})$ is the *marginal contribution* of player i to coalition S . A coalitional game (N, v) is *convex* if the marginal contribution of any player to any coalition is not more than his marginal contribution to a larger coalition, i.e., if it holds that

$$m_i(v, S) \leq m_i(v, T) \tag{1}$$

for all $S \subseteq T \subseteq N$ and all $i \in S$. It is easy to check that convexity implies superadditivity, but not vice versa.

Let $|N|$ be the cardinality of a finite set N . A *payoff vector* or *allocation* (x_1, \dots, x_n) of a coalitional game (N, v) is an $|N|$ -dimensional vector describing the payoffs of the players, such that each player $i \in N$ receives x_i .

A *one-point solution* (or simply a *solution*) for a class \mathcal{C} of coalitional games is a function ψ that assigns a payoff vector $\psi(v)$ to every coalitional game in the class, that is $\psi : \mathcal{C}^N \rightarrow \mathbb{R}^N$.

The most famous solution in the theory of coalitional games is the *Shapley value*, introduced by Shapley (1953). Such a solution can be described in several ways. In view of the analysis of gene relevance in the next section, we introduce the Shapley value ϕ applied to game $(N, v) \in \mathcal{G}^N$ by the general formula

$$\phi_i(v) = \sum_{S \subseteq N: i \in S} \frac{(s-1)!(n-s)!}{n!} m_i(v, S) \tag{2}$$

for each $i \in N$, where $s = |S|$ and $n = |N|$ are the cardinality of coalitions S and N , respectively.

An alternative representation of the Shapley value can be given in terms of the unanimity coefficients $(\lambda_S(v))_{S \in 2^N \setminus \{\emptyset\}}$ of a game (N, v) , that is:

$$\phi_i(v) = \sum_{S \subseteq N: i \in S} \frac{\lambda_S(v)}{|S|} \tag{3}$$

for each $i \in N$.

¹Let T be a set. To denote a subset S of T we use the notation $S \subseteq T$; $S \subsetneq T$ means $S \subseteq T$ and $S \neq T$; $S \not\subseteq T$ means that $S \subseteq T$ is not true.

Another one-point solution for coalitional games is the *Banzhaf value*, introduced by Alon et al. (1999). The Banzhaf value $\beta_i(v)$ of a game $(N, v) \in \mathcal{G}^N$, is defined as follows

$$\beta_i(v) = \sum_{S \subseteq N: i \in S} \frac{1}{2^{n-1}} m_i(v, S) \tag{4}$$

for each $i \in N$. An alternative representation of the Banzhaf value can also be given in terms of unanimity coefficients $(\lambda_S(v))_{S \in 2^N \setminus \{\emptyset\}}$ of a game (N, v) , that is:

$$\beta_i(v) = \sum_{S \subseteq N: i \in S} \frac{\lambda_S(v)}{2^{s-1}} \tag{5}$$

for each $i \in N$.

Finally, a relevant set, possibly empty, of payoff vectors of a coalitional game (N, v) is the *core*, which is defined as follows:

$$\text{core}(v) = \left\{ x \in \mathbb{R}^N \mid \sum_{i \in S} x_i \geq v(S) \ \forall S \in 2^N \setminus \{\emptyset\}; \sum_{i \in N} x_i = v(N) \right\}.$$

3 Interaction among genes

Consider an MES $E = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ and a discriminant method m . In this phase of the analysis we assume that the boolean expression profile $\mathbf{B}^{E,m}(j)$, for each sample $j \in S_D$, is a sufficient condition for the onset of the disease (*sufficiency principle for groups of genes*). Stated differently, a group of genes $U \subseteq N$ which are abnormally expressed in a sample of S_D (according to a discriminant method m applied to the reference expression matrix \mathbf{A}^{S_R}) implies that an individual whose sample has at least all (possibly many more, due to biological and technical bias affecting the data set) the genes in U abnormally expressed (again on the basis of m and \mathbf{A}^{S_R}) should have the disease.

The aim of this work is to give an answer to the following questions: how much relevant for the onset of a tumor are the genes which are abnormally expressed inside the sample S_D ? Is it possible to provide a measure of the power of genes in determining the onset of the tumor in an individual, on the basis of the information collected via samples S_D and S_R and the discriminant method m used?

Consider, for instance, an MES $\tilde{E} = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ and a discriminant method \tilde{m} such that the corresponding boolean expression matrix is

$$\mathbf{B}^{\tilde{E}, \tilde{m}} = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}. \tag{6}$$

On matrix (6) it seems very reasonable to affirm that on the basis of the information collected on the set of samples \mathbf{A}^{S_R} and the discriminant method used \tilde{m} , all the genes

abnormally expressed have the same relevance in causing the tumor, assuming the principle of sufficiency for groups of abnormally expressed genes introduced before.

On the other hand, it could be reasonable to expect experimental situations where there are many boolean expression profiles inside the sample S_D , like in the boolean expression matrix of Examples 1 and 2.

Example 2 Consider again the MES \hat{E} of Example 1 and a more conservative discriminant method \bar{m} such that for each $i \in N$ and each $j \in S_D$

$$(\bar{m}(\mathbf{A}^{S_D}, \mathbf{A}^{S_R}))_{ij} = \begin{cases} 1, & \text{if } A_{ij}^{S_D} \leq p_i^{25\%} \text{ or } A_{ij}^{S_D} \geq p_i^{75\%}, \\ 0, & \text{otherwise,} \end{cases}$$

where $p_i^{25\%}$ and $p_i^{75\%}$ are the 25th and the 75th percentiles of the expression distribution of gene i (i.e., the i th row) in the reference expression matrix \mathbf{A}^{S_R} , for each $i \in N$. The resulting boolean expression matrix is the following

$$\mathbf{B}^{\hat{E}, \bar{m}} = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}.$$

How to deal with these situations?

Given an MES $E = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ and a discriminant method m , first we determine the average number of individuals with the tumor due to the abnormal expression of a given group of genes. For each group $U \subseteq N$, we look at the number of groups of abnormal expressed genes in $\mathbf{B}^{E,m}$ that are included in U . We formalize such a concept via the following definitions (recall that $\mathbf{B}^{E,m}(j)$ denotes the column j , $j \in S_D$, of the boolean expression matrix $\mathbf{B}^{E,m}$).

Definition 1 Let $W \in \{0, 1\}^N$, $n \in \{1, 2, \dots\}$. We define the *support* of W , denoted by $\text{sp}(W)$, by the set

$$\text{sp}(W) = \{i \in \{1, \dots, n\} \mid W_i = 1\}.$$

Example 3 Consider the boolean matrix $\mathbf{B}^{\hat{E}, \hat{m}}$ of Example 1.

Then $\text{sp}(\mathbf{B}^{\hat{E}, \hat{m}}(1)) = \{1, 3\}$, $\text{sp}(\mathbf{B}^{\hat{E}, \hat{m}}(2)) = \{2, 3\}$, and $\text{sp}(\mathbf{B}^{\hat{E}, \hat{m}}(3)) = \{1, 2, 4\}$.

Definition 2 Let $E = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ be an MES and let m be a discriminant method. We define the corresponding *microarray game* as the coalitional game (N, v) , where

- The set of genes N is the set of players.
- The characteristic function v assigns to each coalition $T \in 2^N \setminus \{\emptyset\}$ the average number of samples with tumor determined by T according to the sufficiency principle for groups of genes.

More precisely, we define $v(T)$, for each $T \in 2^N \setminus \{\emptyset\}$, as the value

$$v(T) = \frac{|\Theta(T)|}{|S_D|}, \tag{7}$$

where $|\Theta(T)|$ is the cardinality of the set

$$\Theta(T) = \{k \in S_D \mid \text{sp}(\mathbf{B}^{E,m}(k)) \subseteq T, \text{sp}(\mathbf{B}^{E,m}(k)) \neq \emptyset\} \tag{8}$$

and $v(\emptyset) = 0$.

The class of microarray games will be denoted with the symbol \mathcal{M} . Note that the class of microarray games is a proper subclass of the class of $[0, 1]$ -games. More precisely, if $|S_D|$ is fixed, the characteristic function of a microarray game may take values only in the set $\{\frac{k}{|S_D|} \mid k = 0, \dots, |S_D|\}$.

Remark 1 Condition $\text{sp}(\mathbf{B}^{E,m}(k)) \neq \emptyset$ in relation (8) is due to practical considerations concerning the interpretation of the sufficiency principle for groups of genes on samples where genes do not show any abnormal expression properties. We are assuming that such a sample contributes to decrease the level of association between the abnormal expression of genes and the disease in all coalitions $S \subseteq N, S \neq \emptyset$. Consider for instance an MES $\dot{E} = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ and a discriminant method \dot{m} such that the corresponding boolean expression matrix is

$$\mathbf{B}^{\dot{E},\dot{m}} = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}.$$

Since the sample represented on the third column of the boolean matrix $\mathbf{B}^{\dot{E},\dot{m}}$ is a vector of zeros, it is easy to check in relation (7) that the numerator is always smaller than $|S_D| = 3$ and, consequently, $v(T) < 1$, for each $T \subseteq N$.

Let $E = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ be an MES and let m be a discriminant method. According to equality (7), an equivalent way to calculate the corresponding microarray game v is as a sum of unanimity games as follows

$$v = \frac{1}{|S_D|} \sum_{j \in S_D: \text{sp}(\mathbf{B}^{E,m}(j)) \neq \emptyset} u_{\text{sp}(\mathbf{B}^{E,m}(j))}, \tag{9}$$

where $u_{\text{sp}(\mathbf{B}^{E,m}(j))}$ is the unanimity game on $\text{sp}(\mathbf{B}^{E,m}(j)) \subseteq N$, for each $j \in S_D$.

Alternatively, it is possible to rewrite (9) in terms of the unanimity coefficients of a microarray game v . In formula

$$v = \sum_{S \subseteq N: S \neq \emptyset} \lambda_S u_S, \tag{10}$$

where $\lambda_S = \frac{\bar{\lambda}_S}{|S_D|}$ and $\bar{\lambda}_S = |\{k \in S_D | \text{sp}(\mathbf{B}^{E,m}(k)) = S\}|$ is the number of occurrences of the coalition S as support in the boolean expression matrix $\mathbf{B}^{E,m}$.

Example 4 Consider again the boolean matrix $\mathbf{B}^{\hat{E},\hat{m}}$ of Example 1. By (9) the corresponding microarray game $(\{1, 2, 3, 4\}, v)$ is such that

$$v = \frac{1}{3}(u_{\{1,3\}} + u_{\{2,3\}} + u_{\{1,2,4\}}).$$

It follows that $v(1) = v(2) = v(3) = v(4) = 0$, $v(1, 2) = v(1, 4) = v(2, 4) = v(3, 4) = 0$, $v(1, 3) = v(2, 3) = v(1, 3, 4) = v(2, 3, 4) = v(1, 2, 4) = \frac{1}{3}$, $v(1, 2, 3) = \frac{2}{3}$, $v(1, 2, 3, 4) = 1$.

Proposition 1 introduces some basic properties of microarray games.

Proposition 1 *Let $(N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R})$ and m be an MES and a discriminant method, respectively, and let v be the corresponding microarray game in \mathcal{M}^N . Then v is a monotonic and convex coalitional game.*

Proof First, note that by (10) microarray games are positive linear combination of unanimity games. It is a well known fact that unanimity games are monotonic and convex games. The proof of Proposition 1 immediately follows from the fact that a linear combination of monotonic and convex games with nonnegative coefficients preserves monotonicity and convexity, respectively. □

At this point, the main question addressed by our work can be reformulated in the following terms: is it possible to employ the standard theory of coalitional games to measure the relevance of each gene in determining the onset of a tumor on the basis of the microarray experimental situation observed and the discriminant method used?

In the last fifty years, many studies have addressed the goal of evaluating the power of players (e.g., members of councils, voters in an electoral system, parties in a parliament, individual components of a complex system, etc.) in *simple games*, which are coalitional games whose characteristic function can only assume values 1 (for winning coalitions, e.g., coalitions which are able to force the endorsement of a motion) or 0 (for losing coalitions) (see, for instance, Owen 1995 for a general introduction to simple games and political applications; Ramamurthy 1990 for an application in the framework of reliability theory, to investigate the relationship between the operating state of a complex system and the operating state of its individual components). In such contexts, the idea was to evaluate the amount of power of players according to the role covered by each of them in supporting the goal of each possible coalition. In practice, the Shapley value (Shapley 1953; Shapley and Shubik 1954), as well as other solutions for coalitional games, has been used as a power index, i.e., a function assigning to each simple game (N, w) , where N is a finite set of players, a vector of $|N|$ real numbers indicating the power of each player in (N, w) (Shapley and Shubik 1954; Banzhaf 1965).

Taking into consideration such applications, it is not counterintuitive to look at the Shapley value of a microarray game as a possible measure of the relevance of the

involved genes. Section 4 is devoted to support this idea by means of an axiomatic characterization of the Shapley value on the class of microarray games, satisfying properties which have a sound interpretation in the genomic scenario.

Note that if an MES is given, the computation of the Shapley value $\phi(v)$ of the corresponding microarray game $v \in \mathcal{M}^N$, in virtue of (3) and (10), is very easy, independently from the number of genes involved. More precisely, we have that

$$\phi_i(v) = \frac{1}{|S_D|} \sum_{S \subseteq N: i \in S} \frac{\bar{\lambda}_S}{|S|} \quad (11)$$

for each $i \in N$. If the original MES is not given, and only a microarray game $v \in \mathcal{M}^N$ is given, formula (11) does not apply and the computation of the Shapley value of v may be very hard if the number of genes is high. Since we are interested to practical situations, where the original MES is always given as the output of a microarray experiment, we do not further investigate the computational aspects related to the Shapley value calculation when only the characteristic function (and not the original MES) is known.

We conclude this section introducing some examples of application of the Shapley value to microarray games. In particular, Example 8 deals with an application on a microarray game arising from an MES with real microarray data provided by the literature.

Example 5 The Shapley value of the microarray game in Example 4 is $(\frac{5}{18}, \frac{5}{18}, \frac{1}{3}, \frac{1}{9})$. This means that on the basis of the corresponding MES \hat{E} and the discriminant method \hat{m} the Shapley value of the microarray game states that the most important attribute in determining the tumor onset—on the average—is gene 3, followed by genes 1 and 2 with the same score and gene 4.

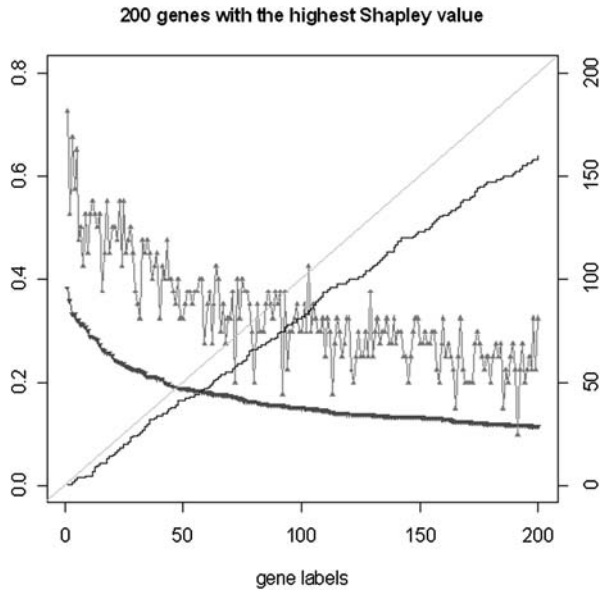
Example 6 The Shapley value of the microarray game corresponding to the boolean matrix in Example 2 is $(\frac{2}{9}, \frac{1}{3}, \frac{2}{9}, \frac{2}{9})$. On the basis of the considerations detailed in Example 5, we obtain that the most important gene in determining the tumor onset, on the average, is gene 2, followed by gene 1, 3 and 4 with the same score.

Example 7 Consider again the boolean expression matrix (6). The Shapley value of the corresponding microarray game is $(\frac{1}{2}, 0, \frac{1}{2}, 0)$.

Example 8 We introduce here a preliminary application of our model to a real MES $E_c = \langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ where \mathbf{A}^{S_D} and \mathbf{A}^{S_R} represent the tumor/normal data set (freely obtainable from the web²) containing expression levels of a set N of 2000 genes measured using *Affymatrix oligonucleotide* microarrays for a set S_D of 40 tumor samples and a set S_R of 22 normal samples, in total 62 samples from colon tissues (Alon et al. 1999). After the preprocessing stage performed by the Bioconductor specific software for microarray analysis (Gentleman et al. 2004), the discriminant method \hat{m} introduced in Example 1 is applied in order to provide the

²<http://microarray.princeton.edu/oncology/affydata/index.html>.

Fig. 1 On the x -axis, top 200 genes are labelled in decreasing order according to the Shapley value. For each gene $i \in \{1, \dots, 200\}$, the y_i -coordinate on the left y -axis of each triangle point-down represents the Shapley value $\phi_i(v_c)$ (values must be divided by 100); the y_i -coordinate on the left y -axis of each triangle point-up equals the ratio $\omega_i(v_c)$; the y_i -coordinate on the right y -axis of the increasing stair steps line equals the number of genes with Shapley value not smaller than $\phi_i(v_c)$ that at the same time have a ratio $\omega(v_c)$ not smaller than $\omega_k(v_c)$, where k is the i th gene when genes are ranked in decreasing order according to $\omega(v_c)$



boolean expression matrix $B^{E_c, \hat{m}}$, which finally produces the corresponding microarray game (N, v_c) .

The Shapley value $\phi(v_c)$ of the microarray game (N, v_c) is computed by means of the procedure suggested by (11), implemented in the programming language R (R Development Core Team 2004). Also the discriminant methods and other procedures for the management of data sets used in this application is implemented using the language and environment R . For graphical reasons, only the Shapley value distribution of the 200 genes with highest Shapley value is depicted in Fig. 1.

For each gene $i \in N$, the ratio of samples such that gene i takes value 1 in the Boolean matrix $B^{E_c, \hat{m}}$ is also computed. In formula,

$$\omega_i(v_c) = \frac{|\{j \in S_D | B_{ij}^{E_c, \hat{m}} = 1\}|}{|S_D|} \tag{12}$$

for each gene $i \in N$.

In Fig. 1 the ratio $\omega_i(v_c)$ is plotted (triangle point-up) for the 200 genes with highest Shapley value. The two graphs show that if relevant genes are selected as the first m genes with highest Shapley value, these genes usually do not coincide with the first m genes with highest $\omega(v_c)$ ratio, for each $m \in \{1, \dots, 200\}$. The number of genes among the first m with highest Shapley value which are also among the first m with highest $\omega(v_c)$ is computed, for each $m \in \{1, \dots, 200\}$, and plotted as a stair steps line in Fig. 1. The overlap number is found to vary between 100% and 45% (median 81%) of the number m of selected genes, with $m \in \{1, \dots, 200\}$.

In Table 1, the first ten genes with highest Shapley value on the microarray game (N, v_c) are indicated.

Some of the genes selected were previously observed in association with the colon cancer (Fujarewicz and Wiench 2003): the vasoactive intestinal peptide (VIP) has

Table 1 Top ten genes ranked according to the Shapley value. In bold we indicate those genes which are also ranked among the top ten genes according to the Banzhaf value calculated according to relation (5)

Gene number	Gene name	Shapley $\times (10^{-3})$
Z50753	H. sapiens mRNA for GCAP-II/ uroguanylin precursor	3.83
H17434	NUCLEOLIN (HUMAN)	3.56
H06524	GELSOLIN PRECURSOR, PLASMA (HUMAN)	3.34
H72234	DNA-(APURINIC OR APYRIMIDINIC SITE) LYASE (HUMAN)	3.33
M36634	Human vasoactive intestinal peptide (VIP) mRNA, complete cds.	3.23
U06698	Human neuronal kinesin heavy chain mRNA, complete cds.	3.21
H61410	PLATELET GLYCOPROTEIN IV (H. sapiens)	3.14
R39209	HUMAN IMMUNODEFICIENCY VIRUS TYPE I ENHANCER-BINDING PROTEIN 2 (H. sapiens)	3.13
M58050	Human membrane cofactor protein (MCP) mRNA, complete cds.	3.09
H08393	COLLAGEN ALPHA 2(XI) CHAIN (H. sapiens)	3.01

been suggested to promote the growth and proliferation of tumor cells; the membrane cofactor protein (MCP) represents a possible mechanism of the ability of the tumor to evade destruction by the immune system (tumor *escape*); gelsolin is a protein which acts as both a regulator and an effector of *apoptosis*, i.e., the mechanism responsible for the physiological deletion of cells. DNA-apurinic or apyrimidinic site lyase protein plays an important role in DNA repair and in resistance of cancer cells to radiotherapy (Moler et al. 2000).

4 An axiomatic characterization of the Shapley value with genetic interpretation

A *gene regulatory pathway (GRP)* is a collection of genes in a cell which interact with each other, dynamically orchestrating the level of expression of the genes in the collection. A simple GRP would consist of one or more input signaling pathways, regulatory proteins that integrate the input signals, several target genes and the RNA and proteins produced from those target genes (Bower and Bolouri 2001). Inferring GRPs from gene expression data is a crucial step to understand the function of genes in the onset of a genetic disorder. Yet the mechanisms that control gene expression in many, if not most, GRPs, are only beginning to be elucidated. For example, *Thymidylate synthase (TS)* plays an important role in chemotherapy for colon cancer and is thought to be one of the target genes that the *E2F1 transcription factor* binds to and regulates. However, the GRP governing the expressions of genes TS and E2F1 in primary colon cancer specimens remains unclear (Kasahara et al. 2000).

One of the main difficulties to understand the main mechanisms governing gene regulatory pathways is the high number of genes involved in a microarray study. A strategy to reduce the complexity related with the high number of genes involved in GRP is to filter out noisy, irrelevant and redundant genes (Jager 2006). The expression of noisy genes is strongly affected by noise in the measurements, which stems from experimental variation but can also reflect sampling effects. Irrelevant genes are those that are not related to the disease, for instance, genes that are constantly expressed at the same level in normal and tumor samples. Redundant genes are those that are highly correlated to other genes and, in fact, are regulated by others which can be considered the biologically relevant genes, i.e., those genes which are primarily responsible for the onset of the genetic disorder and could be used for a diagnostic microarray design.

In this section we claim that the Shapley value can be used for the selection of biologically relevant genes. We support this claim by means of the so-called ‘property driven’ (or ‘axiomatic’) approach, that is: we shall be able to characterize the Shapley value of a microarray game using basic properties. The aim of these properties is to state how a relevance index should behave in very simple situations of genes interaction.

In the direction of characterizing the Shapley value by means of properties with a genetic interpretation, first we need to embed the notion of GRP into the context of microarray games using a game theoretical terminology. In this respect, the definition of *partnership of genes* plays a key role.

Definition 3 Let $v \in \mathcal{M}^N$. A coalition $S \in 2^N \setminus \{\emptyset\}$ such that for each $T \subsetneq S$ and each $R \subseteq N \setminus S$

$$v(R \cup T) = v(R) \quad (13)$$

is a *partnership of genes* in the microarray game v .

Note that the concept of partnership in coalitional games has been introduced in Kalai and Samet (1988) in a general context not involving genes in order to represent those coalitions in a coalitional game v that behaves like one individual, since all its sub-coalitions are powerless. There are at least two important reasons supporting the decision to adopt the definition of partnership of genes as a good representation of a GRP in the microarray game context.

First, the definition of partnership does not require any a priori information on the exact regulatory mechanisms among genes inside the GRP. Due to the high complexity of a GRP, this kind of information is not yet available for many genes. For example, as we already mentioned, it is not yet clarified which genes in the GRP including TS and E2F1 regulates each other in primary colon cancer cells.

Second, the definition of partnership requires that it is not possible to recognize a proper subgroup of genes which directly interact with an external gene or group of genes in provoking the onset of the tumor. This is a necessary condition for a collection of genes to behave as a GRP, seen as network of genes able to specify the identity and level of expression of groups of target genes. On the other hand, it is possible that the joint value of a coalition made by two disjoint partnerships would be

greater than the mere sum of their single values, keeping into account the possibility that distinct GRPs may interact inside a cell.

We call *relevance index* for genes a solution $F : \mathcal{M}^N \rightarrow \mathbb{R}^N$ on the class of microarray games with the set of genes N as the set of players. Some interesting properties for relevance indexes, related to the concept of partnership of genes, are the following.

Property 1 *Let $v \in \mathcal{M}^N$. The solution F has the Partnership Rationality (PR) property, if*

$$\sum_{i \in S} F_i(v) \geq v(S)$$

for each $S \in 2^N \setminus \{\emptyset\}$ such that S is a partnership of genes in the game v .

The PR property determines a lower bound of the power of a partnership, i.e., the total relevance of a partnership of genes in determining the onset of the tumor in the individuals should not be lower than the average number of cases of tumor enforced by the partnership itself.

Property 2 *Let $v \in \mathcal{M}^N$. The solution F has the Partnership Feasibility (PF) property, if*

$$\sum_{i \in S} F_i(v) \leq v(N)$$

for each $S \in 2^N \setminus \{\emptyset\}$ such that S is a partnership of genes in the game v .

On the contrary of PR, the PF properties determines an upper bound of the power of a partnership, i.e., the total relevance of a partnership of genes in determining the tumor onset in the individuals should not be greater than the average number of cases of tumor enforced by the grand coalition.

Together, PR and PF identify a nonnegative scale to quantify the relevance of genes in provoking a disease, assigning value 1 to a partnership of genes that alone, according to the sufficiency principle, is responsible for the disease in all the tumor samples.

Basically, the criterium to compare the relevance of different partnerships is their values in the microarray game, but some extra relevance can be attributed to target genes in the partnerships according to their role in all possible coalitions.

Property 3 *Let $v \in \mathcal{M}^N$. The solution F has the Partnership Monotonicity (PM) property, if*

$$F_i(v) \geq F_j(v)$$

for each $i \in S$ and each $j \in T$, where $S, T \in 2^N \setminus \{\emptyset\}$ are partnerships of genes in v such that $S \cap T = \emptyset$, $v(S) = v(T)$, $v(S \cup T) = v(N)$, $|S| \leq |T|$.

The PM property is very intuitive: consider two disjoint partnerships of genes enforcing the same average number of cases of tumor in the set of samples. If the genes outside the union of those two partnerships are irrelevant—that is, they do not contribute in increasing the average number of tumors—then genes in the smaller partnership should receive a higher relevance index than genes in the bigger one, where the likelihood that some genes are redundant is higher.

The next two properties do not involve the concept of partnership of genes.

Property 4³ *Let $v_1, \dots, v_k \in \mathcal{M}^N, k > 1$. The solution F has the Equal Splitting (ES) property, if*

$$F\left(\frac{\sum_{i=1}^k v_i}{k}\right) = \frac{\sum_{i=1}^k F(v_i)}{k}.$$

Remark 2 We want to reassure the reader that $\frac{\sum_{i=1}^k v_i}{k} \in \mathcal{M}^N$.

First, consider a sequence of k MESs $E_1, \dots, E_k, k > 1$, with the same set of genes N and such that the cardinality of the set of samples with the disease in each MES E_i is the same, for each $i \in \{1, \dots, k\}$. For each $i \in \{1, \dots, k\}$, let m_i be a discriminant method that can be applied on the MES E_i . Consequently, the k boolean expression matrices $\mathbf{B}_1^{E_1, m_1}, \dots, \mathbf{B}_k^{E_k, m_k}$, have the same number of columns and the same number of rows. Let $v_1, \dots, v_k \in \mathcal{M}^N$ be the microarray games defined on the boolean matrices $\mathbf{B}_1^{E_1, m_1}, \dots, \mathbf{B}_k^{E_k, m_k}$, respectively. Then $\frac{\sum_{i=1}^k v_i}{k}$ coincides with the microarray game corresponding to the boolean expression matrix obtained juxtaposing the matrices $\mathbf{B}_1^{E_1, m_1}, \dots, \mathbf{B}_k^{E_k, m_k}$.

On the other hand, we want to prove that given k microarray games $v_1, \dots, v_k \in \mathcal{M}^N, k > 1$, we have $\frac{\sum_{i=1}^k v_i}{k} \in \mathcal{M}^N$, independently from the fact that the boolean expression matrices corresponding to such games have the same number of columns. In order to prove this fact for $k = 2$ a cumbersome notation is needed. Let $\mathbf{B}^1 \in \{0, 1\}^{N \times S_{D1}}$ and $\mathbf{B}^2 \in \{0, 1\}^{N \times S_{D2}}$ be two boolean expression matrices with the same set of genes N and where $S_{D1} = \{s_1^{D1}, s_2^{D1}, \dots, s_l^{D1}\}$ and $S_{D2} = \{s_1^{D2}, s_2^{D2}, \dots, s_p^{D2}\}$ are two sets of tumor samples. Let $\oplus : \mathbb{R}^{N \times S_{D1}} \times \mathbb{R}^{N \times S_{D2}} \rightarrow \mathbb{R}^{N \times \{s_1, s_2, \dots, s_{l+p}\}}$ be a matrix operator such that if $\mathbf{C} = \mathbf{B}^1 \oplus \mathbf{B}^2$ then $\mathbf{C}(s_j) = \mathbf{B}^1(s_j^{D1})$ for each $j \in \{1, \dots, l\}$ and $\mathbf{C}(s_{j+l}) = \mathbf{B}^2(s_j^{D2})$ for each $j \in \{1, \dots, p\}$. Let $v_1, v_2 \in \mathcal{M}^N$ be two microarray games, obtained by Definition 2 on \mathbf{B}^1 and \mathbf{B}^2 , respectively. It is easy to check that the game $\frac{v_1+v_2}{2}$ is the microarray game corresponding to the boolean expression matrix $(\bigoplus_{i=1}^l \mathbf{B}^1) \oplus (\bigoplus_{i=1}^p \mathbf{B}^2)$. For $k > 2$ similar arguments hold, too.

The ES property requires that the average relevance index of genes in two or more different microarray games $v_1, \dots, v_r \in \mathcal{M}^N$ with the same set of genes, even arising

³Assuming the continuity of F , it can be proved, using functional equation theory, that the ES property is equivalent to the simpler property of requiring that F satisfies $F(\frac{v+w}{2}) = \frac{F(v)+F(w)}{2}$ for each pair $v, w \in \mathcal{M}^N$.

from MES provided by different laboratories, must be equal to the relevance index of genes in the average game $\frac{\sum_{i=1}^r v_i}{r}$. There is no reference in the definition of the ES property neither to the accuracy of the relevance index in each different microarray game v_i , $i \in \{1, \dots, r\}$, nor to the accuracy of the relevance index in the average microarray game $\frac{\sum_{i=1}^r v_i}{r}$. For a discussion on how to estimate the accuracy of the relevance index in a microarray game see Moretti (2006a).

The ES property simply underlies a principle of equivalence of reliability levels for microarray games arising from different MES. If we have reasons to assume that different microarray games are equally reliable, then the ES property advises to behave in a natural way to summarize the relevance indexes: making the average. For example, it could be the case of microarray games arising from the equal splitting of the same MES. Let $\langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ be an MES and let S_{D_1}, \dots, S_{D_m} form a partition of the set of samples S_D such that $|S_{D_1}| = |S_{D_2}| = \dots = |S_{D_m}|$. If the ES property holds, then the relevance index computed on the microarray game corresponding to $\langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ equals the average of the relevance indices computed on the microarray games arising from the microarray experimental situations $\langle N, S_{D_1}, S_R, \mathbf{A}^{S_{D_1}}, \mathbf{A}^{S_R} \rangle, \dots, \langle N, S_{D_m}, S_R, \mathbf{A}^{S_{D_m}}, \mathbf{A}^{S_R} \rangle$, respectively. In fact, there is no reason to assume that an MES $\langle N, S_{D_q}, S_R, \mathbf{A}^{S_{D_q}}, \mathbf{A}^{S_R} \rangle$ should be more reliable than $\langle N, S_{D_t}, S_R, \mathbf{A}^{S_{D_t}}, \mathbf{A}^{S_R} \rangle$, for some $q, t \in \{1, \dots, m\}$; moreover, the relevance index of genes in the MES $\langle N, S_D, S_R, \mathbf{A}^{S_D}, \mathbf{A}^{S_R} \rangle$ is independent from the equal splitting partition $\{S_{D_1}, \dots, S_{D_m}\}$ chosen.

The fifth and last property involves the definition of *null gene* of a game (N, v) , that is a player $i \in N$ such that $v(S \cup i) = v(S)$ for each $S \subseteq N \setminus \{i\}$.

Property 5 *Let $v \in \mathcal{M}^N$. The solution F has the Null Gene (NG) property, if for each null gene $i \in N$*

$$F_i(v) = 0.$$

The interpretation of the NG property is straightforward: if a player does not contribute anything to each coalition $S \in 2^N$ then he gets null relevance.

Remark 3 It is well known in literature that the Shapley value satisfies the NG property on each class of coalitional games $\mathcal{C}^N \subseteq \mathcal{G}^N$. The ES property directly follows from Remark 2 together with additivity and homogeneity of the Shapley value ϕ on \mathcal{G}^N , that is $\phi(\alpha v + \beta w) = \alpha \phi(v) + \beta \phi(w)$ for each $v, w \in \mathcal{G}^N$ and each $\alpha, \beta \in \mathbb{R}$.

Remark 4 Another fact known in the literature is that the Banzhaf value satisfies the NG property on each class of coalitional games $\mathcal{C}^N \subseteq \mathcal{G}^N$. From considerations similar to ones expressed in Remark 3, it follows that the Banzhaf value satisfies the ES property as well.

The following Lemmas 1, 2 and Proposition 2 play a role in the axiomatic characterization of the Shapley value on the class of microarray games introduced in Theorem 1. First, we need to introduce some notions. Let $v \in \mathcal{M}^N$ and let $S \in 2^N \setminus \{\emptyset\}$ be a partnership in v . Then it is trivial to prove that each $T \subseteq S$ is a partnership itself.

A maximal partnership $S \in 2^N \setminus \{\emptyset\}$ in v is a maximal subset of N with the property to be a partnership in v .

Note that, by Definition 3, it immediately follows that all one player coalitions are partnerships in v . One easily obtains that the collection of maximal partnerships in v forms a partition of N . For instance, in the microarray game $(\{1, 2, 3, 4\}, v)$ of Example 4, the set of all of the maximal partnerships is $\{\{1\}, \{2\}, \{3\}, \{4\}\}$ and coincides with set of all the partnerships in v ; whereas in the microarray game of Example 7, the set of all of the maximal partnerships is $\{\{1, 3\}, \{2, 4\}\}$.

Lemma 1 *Let $v \in \mathcal{M}^N$ and let $S \in 2^N \setminus \{\emptyset\}$ be a partnership in v . Then the Shapley value attributes the same relevance index to players in S .*

Proof Let $\phi(v)$ be the Shapley value on the game v . For each $U \subseteq N$ such that $i \in U$ the marginal contribution of player $i \in S$ is the following

$$\begin{aligned} &v(U) - v(U \setminus \{i\}) \\ &= v([U \cap S] \cup [U \setminus S]) - v([(U \cap S) \setminus \{i\}] \cup [U \setminus S]) \\ &= \begin{cases} v(U \setminus S) - v(U \setminus S) & \text{if } U \cap S \neq S \\ v(U) - v(U \setminus S) & \text{if } U \cap S = S \end{cases} \\ &= \begin{cases} 0 & \text{if } U \cap S \neq S \\ v(U) - v(U \setminus S) & \text{if } U \cap S = S, \end{cases} \end{aligned} \tag{14}$$

where the second equality follows by Definition 3 on partnership S .

Then, the marginal contribution of each player $i \in S$ to coalition U is different from zero only if S is a subset of U , which means that by (2) the Shapley value of each player $i \in S$ is

$$\begin{aligned} \phi_i(v) &= \sum_{U \subseteq N: i \in U} \frac{(u-1)!(n-u)!}{n!} (v(U) - v(U \setminus \{i\})) \\ &= \sum_{U \subseteq N: S \subseteq U} \frac{(u-1)!(n-u)!}{n!} (v(U) - v(U \setminus S)), \end{aligned}$$

proving that the Shapley value is the same for each player $i \in S$. □

Remark 5 Let $v \in \mathcal{M}^N$ and let $S \in 2^N \setminus \{\emptyset\}$ be a partnership in v . Note that by formula (4) and relation 14, following arguments similar to ones used in proof of Lemma 1, we have that the Banzhaf value of each player $i \in S$ is

$$\begin{aligned} \beta_i(v) &= \sum_{U \subseteq N: i \in U} \frac{1}{2^{n-1}} (v(U) - v(U \setminus \{i\})) \\ &= \sum_{U \subseteq N: S \subseteq U} \frac{1}{2^{n-1}} (v(U) - v(U \setminus S)), \end{aligned}$$

proving that also the Banzhaf value is the same for each player $i \in S$.

Lemma 2 *Let $v \in \mathcal{M}^N$ and let $S \in 2^N \setminus \{\emptyset\}$ be a partnership in v . Then*

$$v(U) = 0$$

for each $U \subsetneq S$.

Proof Suppose on the contrary $v(U) \neq 0$. Then, by Definition 2, $v(R \cup U) > v(R)$ for each $R \subseteq N \setminus U$, which yields a contradiction by Definition 3. \square

Proposition 2 *The Shapley value satisfies the properties PM, PR, PF.*

Proof Let $v \in \mathcal{M}^N$ and let $\phi(v)$ be the Shapley value on the game v .

- (i) Let S and T be two disjoint partnerships such that $v(S) = v(T)$ and $v(S \cup T) = v(N)$. If $v(N) = 0$, then $\phi_i(v) = 0$ for each $i \in S \cup T$ and the PM property is satisfied.

Consider now the case $v(N) > 0$. First note that if S and T are subsets of two different maximal partnerships U and V , respectively, then $S = U$ and $T = V$. In fact, suppose on the contrary that $S \subsetneq U$ or $T \subsetneq V$. By condition $v(S) = v(T)$ and Lemma 2 we have $v(S) = v(T) = 0$, and then, by Definition 3, it follows $v(S \cup T) = 0 \neq v(N)$, which yields a contradiction. Consequently, we have just to prove that the PM property holds under the following two cases:

- (i.i) S and T are two different maximal partnerships. By condition $v(S \cup T) = v(N)$ and Definition 2, it turns out that $v(U) = v(U \cap (S \cup T))$ for each $U \subseteq N$. By Lemma 2 and Definition 3 $v(R) = 0$ for each $R \subseteq S \cup T$, with $S, T \not\subseteq R$. Hence, by relation (10) it is possible to write the game v in terms of unanimity games in the following way

$$v = \lambda_S u_S + \lambda_T u_T + \lambda_{S \cup T} u_{S \cup T}.$$

By relation (11), we have that

$$\phi_i(v) = \frac{1}{|S_D|} \left(\frac{\bar{\lambda}_S}{|S|} + \frac{\bar{\lambda}_{S \cup T}}{|S| + |T|} \right) \tag{15}$$

for each $i \in S$ and

$$\phi_j(v) = \frac{1}{|S_D|} \left(\frac{\bar{\lambda}_T}{|T|} + \frac{\bar{\lambda}_{S \cup T}}{|S| + |T|} \right) \tag{16}$$

for each $j \in T$. Moreover, by Lemma 2, we have that $\lambda_S = \frac{\bar{\lambda}_S}{|S_D|} = \frac{|\Theta(S)|}{|S_D|} = v(S)$ and $\lambda_T = \frac{\bar{\lambda}_T}{|S_D|} = \frac{|\Theta(T)|}{|S_D|} = v(T)$. Since $v(S) = v(T)$, we can rewrite relations (15) and (16), respectively, in the following way

$$\phi_i(v) = \frac{v(S)}{|S|} + \frac{1}{|S_D|} \frac{\bar{\lambda}_{S \cup T}}{|S| + |T|} \tag{17}$$

for each $i \in S$ and

$$\phi_j(v) = \frac{v(S)}{|T|} + \frac{1}{|S_D|} \frac{\bar{\lambda}_{S \cup T}}{|S| + |T|} \tag{18}$$

for each $j \in T$. By relation (17) and (18) we immediately have that for each $i \in S$ and each $j \in T$

$$\phi_i(v) \geq \phi_j(v) \Leftrightarrow |S| \leq |T|,$$

proving that the PM property holds.

- (i.ii) S and T are subsets of the same maximal partnership. Then, by Lemma 1, the Shapley value is the same for each $i \in S \cup T$, proving that the PM property holds. This concludes the proof that the Shapley value satisfies the PM property.
- (ii) The convexity of microarray games by Proposition 1 guarantees that the Shapley value $\phi(v)$ is in the core of the microarray game v . The PR property follows directly from coalitional rationality of core allocations.
- (iii) For each $S \in 2^N \setminus \{\emptyset\}$ such that S is a partnership in v , by monotonicity of v and the fact that $\phi(v)$ is in the core of the microarray game v we have $\sum_{i \in N \setminus S} \phi_i(v) \geq v(N \setminus S) \geq 0$. On the other hand, by efficiency of the Shapley value, $\sum_{i \in N} \phi_i(v) = v(N)$ and then $\sum_{i \in S} \phi_i(v) \leq v(N)$, which proves that the Shapley value satisfies the PF property. □

Remark 6 Let a finite set N be given. Note that the PR, PF, PM, ES, and NG properties previously introduced are defined for relevance indexes, i.e., solutions for games in \mathcal{M}^N . This is due to the fact that in this paper we focus on the interpretation of such properties in the biological context of microarray experiments. However, from a mathematical point of view, these five properties are meaningful also for solutions on the class of coalitional games \mathcal{G}^N . In fact, the definition of partnership has been introduced in Kalai and Samet (1988) for games in \mathcal{G}^N .

As we already observed in Remark 3, the Shapley value, as a solution for coalitional games in \mathcal{G}^N , satisfies the NG and ES properties even on the class \mathcal{G}^N . On the other hand, we are not allowed to claim the same statement with respect to the other properties, PR, PF, and PM. Consider, for instance, the coalitional game $(\{1, 2\}, v)$ such that $v(\{1, 2\}) = 1$, $v(\{1\}) = -2$ and $v(\{2\}) = 0$. Obviously, v is not a microarray game (coalition $\{1\}$ gets a negative value) and it is easy to check that $\{1\}$ and $\{2\}$ are two partnerships in the game v . But, in this game v , the Shapley value allocates $-\frac{1}{2}$ to player 1 and $\frac{3}{2}$ to player 2, and, consequently, the PF property is not satisfied by the Shapley value on partnership $\{2\}$. We skip the easy proof that neither the PM property nor the PR property is satisfied by the Shapley value on the class of coalitional games \mathcal{G}^N .

Remark 7 Let $v \in \mathcal{M}^N$ and let $\beta(v)$ be the Banzhaf value on the game v . Following the same line (i) in the proof of Proposition 2 it is possible to prove that the Banzhaf value satisfies the PM property. Specifically, by formula (5), the use of relations (17)

and (18) in (i.i) of the proof would be replaced, respectively, by the use of the following relations

$$\beta_i(v) = \frac{v(S)}{2^{|S|-1}} + \frac{1}{|S_D|} \frac{\bar{\lambda}_{S \cup T}}{2^{|S|+|T|-1}} \tag{19}$$

for each $i \in S$ and

$$\beta_j(v) = \frac{v(S)}{2^{|T|-1}} + \frac{1}{|S_D|} \frac{\bar{\lambda}_{S \cup T}}{2^{|S|+|T|-1}} \tag{20}$$

for each $j \in T$; the use of Lemma 1 in (i.ii) of the same proof would be replaced by the use of Remark 5.

The Banzhaf value satisfies also the PF property. In fact, for each $S \in 2^N \setminus \{\emptyset\}$ we have

$$\begin{aligned} \sum_{i \in S} \beta_i(v) &= \sum_{i \in S} \frac{1}{|S_D|} \sum_{S \subseteq N: i \in S} \frac{\bar{\lambda}_S}{2^{|S|-1}} \\ &\leq \sum_{i \in S} \frac{1}{|S_D|} \sum_{S \subseteq N: i \in S} \frac{\bar{\lambda}_S}{|S|} = \sum_{i \in S} \phi_i(v) \leq 0, \end{aligned} \tag{21}$$

where the first equality follows by relations (5) and (10), the second equality by relation (11) and the second inequality by Proposition 2 (PF property of the Shapley value ϕ).

In conclusion, by Remark 4, it follows that the Banzhaf value satisfies the properties NG, ES, PM, and PF on the class of microarray games \mathcal{M}^N .

It is easy to find a microarray game where the Banzhaf value does not satisfy the PR property. Consider, for example, the unanimity game $(\{1, 2, 3, 4\}, u_{\{1,2,3\}})$. Then, $\beta_i(u_{\{1,2,3\}}) = \frac{1}{2^{3-1}} = \frac{1}{4}$ for each i in the partnership $\{1, 2, 3\}$ and $u_{\{1,2,3\}}(1, 2, 3) = 1 > \frac{3}{4} = \sum_{i \in \{1,2,3\}} \beta_i(u_{\{1,2,3\}})$.

We conclude this remark noticing that the normalized Banzhaf value does not satisfy ES.

Theorem 1 *Let a finite set N be given. The Shapley value on the class \mathcal{M}^N of microarray games is the unique relevance index which satisfies the properties PR, PF, PM, ES, and NG.*

Proof We already know by Proposition 2 and Remark 3 that the Shapley value satisfies the five properties PR, PF, PM, ES, and NG. To prove the uniqueness consider a map $\psi : \mathcal{M}^N \rightarrow \mathbb{R}^N$ satisfying PR, PF, PM, ES, and NG.

Consider the unanimity game $(N, u_S) \in \mathcal{M}^N$, where $S \in 2^N \setminus \{\emptyset\}$. First note that players $j \in N \setminus S$ are null genes. Then, by NG property, $\psi_j(u_S) = 0$ for each $j \in N \setminus S$.

Moreover, it is easy to see that S is a maximal partnership in u_S . Then, by Lemma 2, for each pair of nonempty sets $U, W \subseteq S$ such that $U \cap W = \emptyset$ and $U \cup W = S$, $u_S(U) = u_S(W) = 0$ and $u_S(U \cup W) = u_S(S) = u_S(N)$. Since PM property holds for ψ , then $\psi_i(u_S) = \psi_j(u_S)$ for each $i, j \in S$.

It follows that $\sum_{i \in S} \psi_i(u_S) = |S| \psi_k(u_S)$, with $k \in S$. By PR, $|S| \psi_k(u_S) \geq 1$ and, by PF, $|S| \psi_k(u_S) \leq 1$. Hence, for each $S \in 2^N \setminus \{\emptyset\}$ and each $i \in N$

$$\psi_i(u_S) = \begin{cases} \frac{1}{|S|}, & \text{if } i \in S, \\ 0, & \text{otherwise.} \end{cases} \tag{22}$$

Finally, for each $i \in N$ we have

$$\begin{aligned} \psi_i(v) &= \psi \left(\frac{\sum_{S \subseteq N: S \neq \emptyset} \bar{\lambda}_S u_S}{|S_D|} \right) = \frac{1}{|S_D|} \sum_{S \subseteq N: S \neq \emptyset} \bar{\lambda}_S \psi_i(u_S) \\ &= \frac{1}{|S_D|} \sum_{S \subseteq N: i \in S} \frac{\bar{\lambda}_S}{|S|}, \end{aligned} \tag{23}$$

where the first equality follows by (10), the second equality by the ES property, and the third equality by relation (22).

According to (11), it has been proved that $\psi(v) = \phi(v)$, where $\phi(v)$ is precisely the Shapley value on the microarray game v . □

Let N be a finite set of players. Note that, following the same steps of the proof of Theorem 1, it is also possible to characterize the Shapley value on the class \mathcal{M}^N using the properties NG and ES together with the efficiency⁴ and symmetry⁵ axioms in Shapley (1953). The NG property and the efficiency and symmetry axioms together yield a value that is uniquely determined on unanimity games. Combined with the ES property and (23), this yields the uniqueness result on the class of microarray games \mathcal{M}^N .

5 Conclusions

In this paper we introduce an application of coalitional games to gene expression analysis related with disease onset. We also present and discuss an axiomatic characterization of the Shapley value aimed at identifying a relevance index for genes.

Many models for data analysis have been presented in the literature for inferring, from a matrix of gene expression data, the role of genes, their interaction and their behavior when changes in condition of the biological system occur (Moler et al. 2000; Su et al. 2003). So far, mathematical techniques used for extracting information from gene expression microarrays can be classified into three main groups: *statistical methods* used for identifying genes that are regulated by different conditions of interest, e.g., to find single genes or groups of genes which show a statistically significant

⁴Let $v \in \mathcal{C}^N \subseteq \mathcal{G}^N$. The solution F satisfies the *efficiency* axiom, if $\sum_{i \in N} F_i(v) = v(N)$.

⁵Let $v \in \mathcal{C}^N \subseteq \mathcal{G}^N$. Two players $i, j \in N$ are *symmetric* in v if for each $S \subseteq N \setminus \{i, j\}$, $v(S \cup \{i\}) = v(S \cup \{j\})$. The solution F satisfies the *symmetry* axiom, if for each $i, j \in N$ which are symmetric in v we have $F_i(v) = F_j(v)$.

difference in the expression levels in tumor samples with respect to normal samples (Fujarewicz and Wiench 2003; Storey and Tibshirani 2003); *unsupervised analysis techniques*, used as a method to identify groups of genes or samples with similar behavior (Golub et al. 1999; Alon et al. 1999); *class prediction tools*, used to classify samples into known categories of morphology, known biological features, clinical outcomes, etc., according to gene expression patterns (Dudoit and Fridlyand 2003; Golub et al. 1999).

The novelty of the approach with respect to the classical methods is essentially twofold. First, the class of coalitional games used, called the class of microarray games, provides the effective opportunity to describe the association between the global expression of each coalition of genes and a genetic disease or another biological condition of interest and, as a consequence, to incorporate in the successive analysis all possible genes interaction ties related with the biological condition. For example, it is possible to describe the association between the over-expression or the under-expression properties of genes in each coalition and the tumor or the effect of a treatment in samples.

Even considering all possible subsets of genes, which means increasing a lot the level of complexity of the analysis, no strong assumptions on the expression probability distributions have been done. In fact, the characteristic function of a microarray game relies completely on the observed experimental gene expression matrix. The very relevant assumption in this context, is the definition of the causality relation (also called *sufficiency principle*) which incorporates the criterium used to establish whether the expression levels of genes in a coalition are associated or not with the biological condition of interest. All the information on gene associations stored in the characteristic function of a microarray game can be successively exploited to quantitatively resume the role of each gene in each possible coalition by means of the application of solution concepts for coalitional games.

The second novelty of the approach presented in this paper is based on the idea of application of solution concepts to microarray games, and on the strong connection between game theory and the property driven (known also as ‘axiomatic’ in the game theoretical literature) approach commonly used for studying the properties of solution concepts. Usually, the interpretation of the results obtained by classical statistical procedures are strongly dependent from the theoretical model used for the analysis or from strong assumptions about the reference population from which the samples are collected. The axiomatic method in the game theoretical approach offers the possibility to overturn this view: only weak assumptions on the population are needed and what is strongly outlined a priori are the boundaries for a plausible interpretations of the results. In the game theoretical approach, the result is the outcome of a solution concept applied to a microarray game built on a gene expression matrix. Its interpretation is contextualized ex-ante by means of sound basic properties, that have to be satisfied by a numerical representation of the relevance of each gene in associating the expressions of coalitions with the genetic disease of interest. This view is particular valuable in the genomic field, which is still a relatively young research topic, and the evidences to support strong hypothesis on the reference populations or the application of sophisticated mathematical models are still far from to be clear.

Recently, coalitional games have been used in gene expression analysis in Moretti (2006b) where a method based on the framework of *minimum cost spanning trees* (mcost) has been introduced to first represent the similarity between pairs of genes and, second, to implement the notion of association for coalitions of genes by means of mcost games. A specific solution for mcost games, the P -value (Branzei et al. 2004), has been used as a relevance index for genes instead of the Shapley value. For such an approach, a discriminant method to dichotomize the expression matrix is not required but some arbitrariness is introduced in the model for the selection of an appropriate notion of similarity between pairs of genes and, successively, for the definition of the level of similarity for coalitions of genes.

In Fragnelli and Moretti (2007) *classification games* with genes in the role of players have been studied to analyze the power of groups of genes to classify samples into the right classes (for instance, the class of normal tissues or the class of tumor tissues). Classification games turn out to be closely related to microarray games and, on some numerical examples, the Shapley value and the Interaction index (Grabisch and Roubens 1999) have been studied as methods for selection of genes with high performance in sample classification.

Coalitional games have been previously used in gene analysis also in a work by Kaufman et al. (2004) as an application of the Multi-perturbation Shapley value Analysis (MSA) (Keinan et al. 2004). The aim of this work was to identify the importance in terms of causal responsibility of some genes in performing a certain function in yeast cells. In their approach, Kaufman et al. (2004) evaluate the value of each coalition as a measure of the biological system's performance for a certain function (e.g., the ability of the system to survive the UV irradiation). In order to obtain such a value for each coalition, they carried out a series of experiments where genes of each different subset of n genes were perturbed concomitantly; on each experiment the performance score was also measured and the score assigned to the corresponding subset of perturbed genes, finally obtaining a coalitional game. Since 2^n experiments were needed to obtain a coalitional game, implying the impossibility to deal with the complete structure of the game both for practical and computational reasons, the authors suggested two complementary approaches: (a) the use of mathematical predictors on the available data set to predict the missing performance scores (Dudoit and Fridlyand 2003; Golub et al. 1999); (b) limiting the focus to one and two dimensional interactions (Grabisch and Roubens 1999; Keinan et al. 2004). In our application setting, where samples of tumoral individuals are involved, of course, we cannot imagine to perform such perturbation experiments. On the other hand, the interpretation of the Shapley value as a measure of the functional causal contribution of genes in a biological system, as provided by Kaufman et al. (2004) seems to corroborate our interpretation of the Shapley value as indicator of the relevance of genes in tumor onset.

References

- Alon U, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, Levine AJ (1999) Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissue probed by oligonucleotide arrays. *Proc Natl Acad Sci USA* 96:6745–6750

- Banzhaf JF III (1965) Weighted voting doesn't work: a game theoretic approach. *Rutgers Law Rev* 19:317–343
- Bequet C, Blachon S, Jeudy B, Boulicaut JG, Gandrillon O (2002) Strong-association-rule mining for large-scale gene-expression data analysis: a case study on human SAGE data. *Genome Biol* 3(12)
- Bower JM, Bolouri H (eds) (2001) Computational modelling of genetic and biochemical networks computational molecular biology series. MIT Press, Cambridge
- Branzei R, Moretti S, Norde H, Tijs S (2004) The P -value for cost sharing in minimum cost spanning tree situations. *Theory Decis* 56:47–61
- Dudoit S, Fridlyand J (2003) Classification in microarray experiments. In: Speed TP (ed) Statistical analysis of gene expression microarray data. Chapman & Hall/CRC, London/Boca Raton, pp 93–158
- Dudoit S, Yang YH, Luu P, Speed TP (2001) Normalization for cDNA microarray data. In: Bittner ML, Chen Y, Dorsel AN, Dougherty ER (eds) Microarrays: optical technologies and informatics. Proceedings of SPIE, vol 4266, pp 141–152
- Fragnelli V, Moretti S (2007) A game theoretical approach to the classification problem in gene expression data analysis. *Comput Math Appl*, doi:[10.1016/j.camwa.2006.12.088](https://doi.org/10.1016/j.camwa.2006.12.088)
- Fujarewicz K, Wiench M (2003) Selecting differentially expressed genes for colon tumor classification. *Int J Appl Math Comput Sci* 13(3):327–335
- Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, Ellis B, Gautier L, Ge Y, Gentry J, Hornik K, Hothorn T, Huber W, Iacus S, Irizarry R, Leisch F, Li C, Maechler M, Rossini AJ, Sawitzki G, Smith C, Smyth G, Tierney L, Yang JYH, Zhang J (2004) Bioconductor: open software development for computational biology and bioinformatics. *Genome Biol* 5:80, <http://genomebiology.com/2004/5/10/R80>
- Golub T, Slonim D, Tamayo P, Huard C, Gaasenbeek M, Mesirov J, Coller H, Loh M, Downing J, Caligiuri M, Bloomfield C, Lander E (1999) Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science* 286:531–537
- Grabisch M, Roubens M (1999) An axiomatic approach to the concept of interaction among players in cooperative games. *Int J Game Theory* 28:547–565
- Jager J (2006) Deriving small diagnostic biomarker panels from genome wide. PhD Dissertation, Max Planck Institute for Molecular Genetics
- Kalai E, Samet D (1988) Weighted Shapley values. In: Roth A (ed) The Shapley value, essays in honor of Lloyd S. Shapley. Cambridge University Press, Cambridge, pp 83–100
- Kasahara M, Takahashi Y, Nagata T, Asai S, Eguchi T, Ishii Y, Fujii M, Ishikawa K (2000) Thymidylate synthase expression correlates closely with E2F1 expression in colon cancer. *Clin Cancer Res* 6:2707–2711
- Kaufman A, Kupiec M, Ruppin E (2004) Multi-knockout genetic network analysis: the Rad6 example. In: Proceedings of the 2004 IEEE computational systems bioinformatics conference (CSB'04), August 16–19, 2004, Stanford, California
- Keinan A, Sandbank B, Hilgetag CC, Meilijson I, Ruppin E (2004) Fair attribution of functional contribution in artificial and biological networks. *Neural Comput* 16(9):1887–1915
- Moler EJ, Chow ML, Mian IS (2000) Analysis of molecular profile data using generative and discriminative methods. *Physiol Genomics* 4:109–126
- Moretti S (2006a) Game Theory applied to gene expression analysis. PhD Dissertation, University of Genoa, Italy
- Moretti S (2006b) Minimum cost spanning tree games and gene expression data analysis. In: ACM international conference proceeding series, p 199
- Owen G (1995) Game theory, 3rd edn. Academic, San Diego
- Parmigiani G, Garret ES, Irizarry RA, Scott SL (2003) The analysis of gene expression data: an overview of methods and software. In: Parmigiani G, Garret ES, Irizarry RA, Zeger SL (eds) The analysis of gene expression data: methods and software. Springer, New York
- R Development Core Team (2004) R: a language and environment for statistical. R foundation for statistical computing, Vienna, Austria, 2004, ISBN 3-900051-00-3, <http://www.R-project.org>
- Ramamurthy KG (1990) Coherent structures and simple games. Kluwer Academic, Dordrecht
- Shapley LS (1953) A value for n -person games. In: Kuhn HW, Tucker AW (eds) Contributions to the theory of games II. Annals of mathematics studies, vol 28. Princeton University Press, Princeton, pp 307–317
- Shapley LS, Shubik M (1954) A method for evaluating the distribution of power in a committee system. *Am Political Sci Rev* 48:787–792

- Smith K, Speed T (2003) Normalization of cDNA microarray data. *Methods* 31:265–273
- Storey JD, Tibshirani R (2003) SAM thresholding and false discovery rates for detecting differential gene expression in DNA microarrays. In: Parmigiani G, Garret ES, Irizarry RA, Zeger SL (eds) *The analysis of gene expression data: methods and software*. Springer, New York
- Su Y, Murali TM, Pavlovic V, Schaffer M, Kasif S (2003) RankGene: identification of diagnostic genes based on expression data. *Bioinformatics* 19(12):1578–1579