

Network-based Enriched Gene Subnetwork Identification: A Game Theoretic Approach

Abolfazl Razi¹, Fatemeh Afghah¹, Salendra Singh², Vinay Varadan²

¹Electrical Engineering & Computer Science, Northern Arizona University,
2112 S Huffer Ln, Flagstaff AZ, USA

²Case Comprehensive Cancer Center, Case Western Reserve University,
11900 Euclid Avenue, Cleveland OH, USA

Correspondence to be addressed to:

Abolfazl Razi (abolfazl.razi@nau.edu) and Vinay Varadan (vxv89@case.edu)

Journal Title: *Biomedical Engineering and Computational Biology*

Relevance to journal: Our manuscript details a novel game theoretic algorithm to identify gene subnetworks associated with clinical and biological phenotypes by leveraging protein-protein interaction networks. While the specific applications focused on in this manuscript are primarily related to cancer biology and prognosis, the principles underlying the game theoretic methodology are broadly applicable to the biomedical field. Thus our manuscript would be of significant interest to researchers spanning both the biomedical engineering and computational biology fields.

Abstract: Identifying subsets of genes that jointly mediate cancer etiology, progression or therapy response remains a challenging problem due to the complexity and heterogeneity in cancer biology, a problem further exacerbated by the relatively small numbers of cancer samples profiled as compared to the sheer number of potential molecular factors involved. Pure data-driven methods that merely rely on multi-omics data have been successful in discovering potentially functional genes but suffer from high false-positive rates and tend to report subsets of genes whose biological interrelationships are unclear. Recently, integrative data-driven models have been developed to integrate multi-omics data with signaling pathway networks in order to identify pathways associated with clinical or biological phenotypes. However, these approaches suffer from an important drawback of being restricted to previously discovered pathway structures and miss novel genomic interactions as well as potential crosstalk among the pathways.

In this article, we propose a novel *Coalition Game*-theoretic approach to overcome the challenge of identifying biologically relevant gene subnetworks associated with disease phenotypes. The algorithm starts from a set of seed genes and traverse a protein-protein interaction (PPI) network to identify modulated subnetworks. The optimal set of modulated subnetworks are identified using *Shapley value* that accounts for both individual and collective utility of the subnetwork of genes. The algorithm is applied to two illustrative applications including the identification of subnetworks associated with i) disease progression risk in response to platinum-based therapy in ovarian cancer and ii) immune infiltration in triple negative breast cancer. The results demonstrate an improved predictive power of the proposed method when compared to state-of-the-art feature selection methods, with the added advantage of identifying novel potentially functional gene subnetworks that may provide insights into the mechanisms underlying cancer progression.

Keywords: Cancer Genomics, Modulated Subnetworks, Coalition Game Theory, Clinical Outcome Prediction, Network Traversal

I. INTRODUCTION

A critical problem in cancer research involves the identification of subset of genes that play crucial roles in different stages of cancer progression from its early stages of carcinogenesis up to the final stage of metastasis. A model that can reliably identify molecular determinants of clinical outcomes would enable the discovery of functional biomarkers predictive of therapy response or disease progression but also provide insights into novel therapeutic targets in this aggressive disease.

Design of data driven models to identify aberrations in genes including differential expression, somatic mutations, copy number alterations, that are associated with clinical outcomes has been center of attention in the past decades [1] [2] [3]. This is accelerated in recent years due to a huge drop in the cost of next generation RNA sequencing and genomics profiling and availability of several public databases such as The Cancer Genome Atlas (TCGA) [4] [5] [6].

Pure data driven models such as GSEA [2] that relate multi-modal genomics measurement to clinical or biological phenotypes demonstrated a great success in discovery of cancer biomarkers and subsets of genes enriched for a desired outcome. However, these methods suffer from a major drawback of high false discovery rate where only a small subset of reported genes play significant roles in the cancer disease and majority of reported genes are false alarms arise from measurement noise, heterogeneity of cancer samples and overfitting issue. Overfitting is unavoidable due to extremely large number of predictors (such as gene expression data) which is in the orders of ten thousands with respect to the samples in the orders of hundreds for most cancers.

The second drawback of these methods are that the reported genes may or may not involve in a shared molecular interactions and hence provide minimal insight for molecular biology scientist to discover the actual underlying biological process that leads to the specific phenotype of interest. To avoid this drawback, recently, integrative models are proposed to integrate multi-omics data with biologically driven pathway networks in order to identify biologically meaningful subnetworks of genes that are enriched for desired biological outcomes. However, these approaches suffer from another important drawback of restricting identified genes to belong to putatively discovered pathway networks and miss the yet undiscovered genomic interactions as well as potential inter-pathway relations. In this work, we take an intermediate approach and use protein-protein interaction (PPI) networks as our basic interaction platform. PPI networks are similar to pathways in that both consist of interacting biomolecular entities effecting specific cellular functions. However, while pathway networks consist of relatively small numbers of deeply-characterized regulatory and signaling events representable as directed graphs, PPI networks capture genome-wide interactions derived from high throughput molecular profiling and large-scale biological screens. PPI networks, which are represented as binary undirected graphs, capture direct, indirect and as yet undiscovered regulatory interactions and can be understood to capture complex cellular logic as simplified connections between pairs of genes [7]. Thus PPI networks provide more flexibility for the discovery of novel biological mechanisms underlying disease phenotypes. Moreover, edge ontology of signaling pathway networks are not well standardized and this difference may cause problem in inference algorithms [8] .

As stated above, to reduce high false reporting rates and enhance generalizability of the developed input-output relation models, sparsity imposing methods are desired. Dimensionality reduction using discriminative component analysis methods such as Canonical Correlation Analysis (CCA), Linear Discriminant Analysis (LDA) and Independent Discriminant Analysis (IDA) are developed to project data into new subspaces, where a few components bear the most discriminative information about data, hence simplifying data storage, prediction and interpretation [9]. Although very efficient in dimensionality reduction, these methods are not ideally suited for the identification of genes driving cancer progression, since the predictors are provided in the transformed subspace [10].

Explicit feature selection methods are divided to wrapper methods and filtering methods [11]. In filtering methods, the predictors are chosen based on their strong connection to the labels with less connection among the features using various geometric or information-theoretic measures, whereas in wrapper methods the features are chosen based on their impact to the classifier. Wrapper methods require exhaustive search and are thereby computationally expensive, while filter methods with geometric distance measures are very fast but incapable of capturing non-linear relations. On the other hand, information-theoretic filtering methods are very powerful but become computationally expensive. Further, they require large numbers of samples in order to obtain reliable empirical information-theoretic measures [12]. In cancer genomics, we are interested in methods that incorporate gene-interaction information such as protein-protein interaction (PPI) or biological pathway network databases in the feature-selection process, in order to identify functionally related sets of genes that jointly discriminate between phenotypes [13], [14].

In this work, we develop a game-theoretic solution that develops pathways emerging from a seed gene set in PPI network by traversing the network to discover the most informative pathways associated with a desired outcome. This algorithm reports a set of compact subnetworks that are collectively associated with the modulation of a specific biological process or clinical outcome, thus facilitating the development of biomarkers using core representative nodes within the identified subnetworks, as opposed to measuring all the genes individually.

We highlight the utility of the proposed algorithm by applying it to two unique and challenging problems that differ both in terms of the nature of the molecular factors involved as well as the phenotype being modeled. In the first application, we focus on identifying gene subnetworks that are jointly associated with platinum-resistance in ovarian cancer, whereas the second application focuses on discovering genomic determinants of immune infiltration in triple negative breast cancer. We show that our algorithm identifies biologically meaningful gene subsets in both applications, while achieving improved statistical association as compared to other feature selection algorithms.

A. Coalition Game Review

In this section, we introduce the concept of coalition game theory and its application in predictive modeling and feature selection considering the synergic predictive power of selected features. Coalition game refers to a class of games, where the players cooperate with one another by forming coalitions [15] as opposed to non-cooperative games in which the

players act individually and compete over a common resource [16]. Coalition games have been recently utilized in feature selection problems to account for the relevance among potentially effective combinations of the features as well as providing a quantitative measure of the impact of each feature on the overall prediction [17] [18] [19] [20] [21]. Coalition game-theoretic based methods can significantly improve the prediction accuracy compared to most existing feature selection techniques that either account only for the impact of individual features on the target labels or consider at most the pairwise correlation. In these conventional approaches, the features that have a low individual impact against the outcome but a considerable contribution when grouped with other features will most likely be filtered out that result in missing actual informative features.

In this work, we propose a novel Network Based Coalition Game (NBCG) algorithm, where the game players are gene subnetworks extracted from the networks. In this algorithm, the game players are subnetworks which are not fixed, but rather developing identities over the game iterations by picking up new genes from the PPI network.

Let N be the number of players, $\mathcal{P} = \{P_1, \dots, P_N\}$ be the set of players and v denote the characteristic function for a transferable utility coalition game (\mathcal{P}, v) . The characteristic function, v is a real-valued function defined on the set of all coalitions, $v: 2^{\mathcal{P}} \rightarrow \mathbb{R}$. If \mathcal{C} denotes a coalition set, $\mathcal{C} \subseteq \mathcal{P}$, the total payoff that can be gained by the members of coalition \mathcal{C} is defined by characteristic function $v(\mathcal{C})$. This function satisfies the two conditions of i) characteristic function of an empty coalition is zero, $v(\emptyset) = 0$, and ii) if \mathcal{C}_i and \mathcal{C}_j ($\mathcal{C}_i, \mathcal{C}_j \subseteq \mathcal{P}$) are two disjoint coalitions, the characteristic function of their union has super-additivity property, as $v(\mathcal{C}_i \cup \mathcal{C}_j) \geq v(\mathcal{C}_i) + v(\mathcal{C}_j)$. The game solutions are determined with possible scenarios that the players can form coalitions and how the total payoff of a coalition is divided amongst the coalition members.

Marginal importance of player P_i , $\Delta_i(\mathcal{C})$ is defined as its contribution when it joins a coalition \mathcal{C} and is obtained by

$$\Delta_i(\mathcal{C}) = v(\mathcal{C} \cup \{P_i\}) - v(\mathcal{C}). \quad (1)$$

This marginal importance of a player does not reflect a fair share of the player from the characteristic function, since it depends on the order of the players in forming the coalition. To define a fair solution of coalition game (\mathcal{P}, v) , we define the real-valued function, γ that assigns an N -tuple of real numbers, $\gamma(v) = (\gamma_1(v), \gamma_2(v), \dots, \gamma_N(v))$ based on the adopted characteristic function, in which $\gamma_i(v)$ measures the value of player P_i in the game with characteristic function v . The Shapley value can then be defined as a fair unique solution of the game as it assigns a fair quantity for each player based on the average contribution of the player among all possible coalitions with all possible permutations. Formally, the Shapley value of player, $P_i \in \mathcal{P}$ denoted by $\gamma_i(\mathcal{P}, v)$ is defined as the expected marginal importance of player P_i to the set of players who precede this player.

$$\gamma_i(\mathcal{P}, v) = \frac{1}{N!} \sum_{\pi \in \Pi} \Delta_i(\mathcal{C}_i(\pi)), \quad (2)$$

where Π is the set of all $N!$ permutations of \mathcal{P} and $\mathcal{C}_i(\pi)$ is the set of players preceding player P_i in subset \mathcal{C} with permutation π .

In modeling the feature selection problem with a coalition game, the attributes (gene subnetworks extracted from the networks) are defined as the game players, marginal contribution of player, P_i to coalition \mathcal{C} is described as the improvement in the prediction capability of this coalition based on application-dependent evaluation method. The payoff of each coalition \mathcal{C} , $v(\mathcal{C})$, measures the contribution of a coalition to the performance of the predictive model (e.g. classification success rate in supervised learning). In this model, different possible coalitions of genes and pathways are examined to recognize the optimal classification features.

Since in attribute selection problem, the order of features inside a coalition does not change the coalition power, the calculations of Shapley value in (2), can be further simplified by excluding the permutation inside already formed or to be formed coalitions in the average, resulting in the following equation:

$$\gamma_i(\mathcal{P}, v) = \frac{1}{N!} \sum_{\mathcal{C} \subseteq \mathcal{P} \setminus i} \Delta_i(\mathcal{C}) |\mathcal{C}|_i (N - |\mathcal{C}| - 1)!, \quad (3)$$

where $|A|$ denotes the cardinality of set A and $\mathcal{C} \subseteq \mathcal{P} \setminus i$ represents the coalitions to which player P_i does not belong. Moreover, $|\mathcal{C}|_i$ and $|\mathcal{C}|_i (N - |\mathcal{C}| - 1)$ correspond to the permutations of the preceding players and the subsequent players, respectively.

In attribute selection applications with a large number of players, computation of Shapley value for all possible feature coalitions may be computationally intensive. Therefore, we utilize the multi-perturbation Shapley value (MSA) measurement, which is determined using an unbiased estimator based on Shapley value by using sampled permutations of players that form coalitions up to size N' . The idea behind this method is that coalitions of size $N' < N$ are capable of capturing the synergic power of players, hence larger coalitions only present additive power. Therefore, we use

$$\gamma'_i(\mathcal{P}, v) = \frac{1}{|\Pi_{N'}|} \sum_{\pi \in \Pi_{N'}} \Delta_i(\mathcal{C}_i(\pi)), \quad (4)$$

where $\Pi_{N'}$ denotes the sampled permutation on sub-groups of players of size N' . In this work, we use the approximate method of MSA with $N' = 4$. In our proposed algorithm, at each round, the features are randomly divided into groups of size N' . Then, we calculate the corresponding Multi-perturbation Shapely value of feature P_i inside its group, $\gamma'_i(\mathcal{P}, v)$ considering all possible coalitions of size $1 \leq n \leq N'$.

II. METHODS

In this section, we elaborate on the proposed algorithm to find the modulated subnetwork of genes that collectively contribute to a cancer phenotyping. We first note that the genes interact with one another either directly or through other cellular entities such as RNA, proteins, enzymes and protein complexes. These complex interactions are modeled by directed graphs called pathway networks, where each pathway correspond to a biological process active in all or some specific tissue types. Two genes may be indirectly connected to each other through different biological pathways. Each gene is part of a genome that contains a unique sequence of 4 types of nucleobases including Adenine (A), Cytosine (C), Guanine (G), or Thymine (T), and hence encodes for a specific protein through a chain of complex protein synthesis processes including transcription, splicing and translation. The interaction between the protein-products are modeled as undirected graphs called Protein-Protein Interaction (PPI) networks.

A. Proposed Algorithm

In order to identify the subset of genes that play a role in various cancer related subtyping (e.g. success in therapy response for a specific drug), we remind the fact that the interacting genes through some biological processes are more likely to play a collective role in cancer subtyping, since all subtypes are results of an alteration in one or more biological processes. We also note that pure data-driven methods fail in identifying genes which contribute to a specific phenotype due to the well-known large p and small n paradigm, where the number of parameters (genes) overwhelmingly are higher than the number of samples. This causes high false rate and the classical sparsity imposing methods do not perform well due to bias to the utilized dataset and measurement noise. Therefore, a key solution is to use the prior biological knowledge of the interaction among the genes. In this algorithm, we enforce the algorithm output to report the genes from multiple subset of interacting genes as our sparsity imposing method.

In this work, we use PPI networks instead of pathway networks due to the following reasons. Some of the signaling pathways are tissue specific and the regulatory networks may be different from one tissue to another. For instance, the estrogen-receptor signaling pathway may be particularly relevant in breast tissue, while not very active in other tissues. Secondly, these pathway networks differ from one source to another and are subject to constant revisions and modifications. Thirdly, sticking to pathway networks we may miss the interaction between different pathways through as yet undiscovered interactions. PPI networks are also much denser networks than pathway signaling networks with a higher number of connections among nodes. Therefore, using the PPI networks it is less likely to miss any related genes in the final results even though their actual biological interacting mechanism is not yet fully discovered and hence is not present in well curated signaling pathway networks. Further, even if an actual relation among two genes is missing in the PPI network, starting from multiple seed genes prevents losing this missed connection in the final results. The PPI network is provided as an undirected Boolean graph, where the nodes are genes (or their corresponding protein products) and the edges represent biochemical interactions between the connected

nodes [22]. We use the human PPI network and represent it as a $G \times G$ binary matrix A , where $G = 12126$ is the number of genes.

The proposed algorithm is called Network-Based Coalition Game (NBGC) that presents a network traversal algorithm, where the directions of network browsing is determined by a coalition-based game solution. An earlier version of this algorithm is presented in [21] and the modified version of this algorithm along with two illustrative applications are provided next.

The NBGC algorithm as depicted in Figure 1, starts from some initial nodes (seed genes) and then gradually develops links by picking up the nodes from the neighborhood of the subnetwork that present maximal contribution to the desired predictive function. At consecutive iterations each link has the option of extending from left or right end to form a connected link. The algorithm stops if a desired performance criterion is met or the number of nodes reach a predefined limit. The algorithm reports a subset of connected subnetworks that collectively are enriched for an outcome of interest. In order to obtain higher enrichment properties, the algorithm runs in multiple times with different initializations and the round with the best outcome is chosen as the final result. Moreover, we can use averaging methods to report the modulated subnetworks based on their frequency of appearance at multiple runs as will be discussed in the numerical results section. This algorithm is general and can be applied to a wide range of applications. The following is the different blocks of the proposed algorithm for each run.

B. Initializations

To select the subset of genes that are enriched for a phenotype, we start with initial set of $L - 1$ genes, from which the subnetworks emerge. The seed genes may be chosen randomly or using prior biological knowledge based on the application of interest. For instance, one may choose the genes that are most frequently mutated in the cancer being studied that derive the cancer and hence may have crucial impact on the desired cancer related outcome. Another option would be to choose the initial genes by their individual prediction power in the adopted predictive modeling (e.g. the genes with highest association with the target labels). The drawback with both approaches is that they are data-driven methods and may have the issue of biasing to the employed dataset. Therefore, we propose to choose the initialization genes based on their degree distribution in the utilized PPI network. Using hot-spot nodes in PPI network provides more flexibility in browsing the network towards a true functional subnetwork and reduces the chance of missing important subnetworks. This approach also provides higher convergence rate, due to a shorter path from hub nodes to the actual modulated subnetworks.

Another reasonable choice is to choose the initial genes from a set of genes that present high association between the genomics data (e.g. gene expression, somatic mutations and ...) with the target phenotype based on a desired correlation matrix. This is also a reasonable choice, since the highly modulated subnetworks expected to include significantly important individual genes. Therefore starting from these nodes accelerates convergence to the actual modulated subnetworks. The drawback for this approach is the bias to the initial guess and reducing the chance of the discovery of subnetworks that are formed by genes that collectively but not individually impact the desired phenotype. It is notable that the best approach to form the pool

of initialization gene set depends on the application of interest. It is also worth mentioning that regardless of the approach that we use to form the pool of initialization genes, we run the algorithm R times with different seed gene subsets that are randomly chosen from the pool of initialization genes. Therefore, the chance of missing the actual modulated subnetworks is very low. Once the initial genes are chosen, each of them is considered as a subnetwork with only one node to be developed to a higher subnetworks as follows.

C. Subnetwork Expansion

Each internal iteration of algorithm consists of two steps, namely network expansion and contraction. If $\mathcal{S}_i = \{G_{i1}, G_{i2}, \dots, G_{i|\mathcal{S}_i}|\}$, $i \in \{1, 2, \dots, L\}$ is the set of $|\mathcal{S}_i|$ genes forming subnetwork \mathcal{S}_i , then at network expansion step, we add a node to each already developed subnetwork \mathcal{S}_i such that the resulting new network results in a better game theoretic evaluation. Obviously, at the first iteration, each subnetwork consists of a single seed gene (i.e. $|\mathcal{S}_i| = 1, i \in \{1, 2, \dots, L\}$) from which a subnetwork emerges.

First, we form a candidate gene set for each subnetwork. If $\mathcal{S}_i = \{G_{i1}, G_{i2}, \dots, G_{i|\mathcal{S}_i}|\}$ is the set of $|\mathcal{S}_i|$ genes forming subnetwork \mathcal{S}_i , then the inclusion candidate gene set Ω_i for subnetwork \mathcal{S}_i is defined as the set of genes in the direct neighborhood of the subnetwork or equivalently the nodes with direct links to any of the subnetwork nodes in the PPI network, i.e. $\Omega_i = \{G_{ik} : G_{ij} \in \mathcal{S}_i, A(G_{ij}, G_{ik}) = 1\}$ for all $k \in \{1, 2, \dots, G\}$ and $j \in \mathcal{S}_i$.

Then, we run a coalition game for the players including the current subnetworks $\mathcal{S}_i, i \in \{1, 2, \dots, L\}$ and the set of the candidate genes. For each subnetwork \mathcal{S}_i , we calculate Shapley value for each candidate gene $G_{ik} \in \Omega_i$ by forming a coalition game with player set $\mathcal{P} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_{N_L}, G_{ik}\}$, where $|\mathcal{P}| = L + 1$. The Shapley value for the candidate gene quantifies the expected marginal importance of the candidate gene when forming coalitions with different combinations of the current subnetworks as elaborated in section II. Evaluation of the characteristic function for each coalition $v(\mathcal{C}_i)$ depends on the desired objective based on the application as described in section IV for two proposed applications. The candidate gene with the maximum Shapley value, $G_i^{(\text{add})} = \underset{G_{ik} \in \Omega_i}{\operatorname{argmax}} \gamma_{G_{ik}}(v)$ is chosen to join the subnetwork \mathcal{S}_i , provided that its Shapley value exceeds a predefined threshold value T_{add} (i.e. $\gamma_{G_i^{(\text{add})}}(v) \geq T_{\text{add}}$). This condition is required since there might be cases that none of candidate genes contribute significantly to the desired objective or even contribute negatively. In such situation a stop flag F_i is set for the subnetwork meaning that the expansion of the corresponding subnetwork \mathcal{S}_i is terminated and hence is excluded from the expansion steps in the subsequent iterations. The termination flag F_i is also set if the number of nodes in each subnetworks exceeds a predefined value S_{max} . This procedure is repeated for all active subnetworks.

Note on Computation Complexity: It is noteworthy that in contrast to conventional games, game players are not static entities and rather they are developing entities over the consecutive iterations of the coalition game. This implementation significantly reduces the complexity of game with respect to the conventional method of considering each individual gene as a game player.

We note that the computational complexity of the algorithm is exponentially proportional to the number of players. Since we repeat the game for each candidate gene, the complexity of the proposed game is $O(|\Omega_1 \cup \dots \cup \Omega_L| \times 2^{1+L})$. If we use MSA method with maximum coalition size $L' < 1 + L$, then complexity reduces to $O(|\Omega_1 \cup \dots \cup \Omega_L| \times 2^{L'})$. Therefore, the complexity of the game exponentially grows with the number of subnetworks and hence does not increase by the game evolution. This complexity is much less than the conventional method of treating each gene as a game player that provides the complexity of $O(2^{|\Omega_1 \cup \dots \cup \Omega_L| + |\mathcal{S}_1 \cup \dots \cup \mathcal{S}_L|})$, which grows exponentially with the number of genes forming the subnetworks and hence increases over the consecutive iterations.

D. Subnetwork Contraction

In the early version of this algorithm that is presented in [21], each iteration of algorithm includes only network expansion step. One drawback to this approach is that each subnetwork evolves from a randomly selected seed gene and expands towards a modulated subnetwork. Therefore, the subnetworks always include the seed genes and in order to find the best modulated subnetwork, we need to start from a large number of sets of different initialization genes in order to determine the best set of modulated subnetworks enriched for a desired outcome. In this modified version, we include an additional step of network contraction by removing the nodes whose contribution to the desired outcome in terms of Shapley value is below a predefined limit. Therefore, the subnetworks have the flexibility of moving away from the seed genes along the PPI network to the modulated subnetworks. This is due to the fact that the contribution of each gene to the desired outcome may change as the subnetworks evolve over time.

In order to identify the genes with low (or negative) Shapley values, we first define a set of exclusion candidate genes Ψ_i for each subnetwork \mathcal{S}_i . The algorithm chooses the subnetwork edges including the nodes that are connected only to one node in the developed subnetwork, hence their removal does not break down the subnetwork into two disjoint subnetworks. For instance, if $\theta_j = \{G_k: G_k \in \mathcal{S}_i, A(G_k, G_{ij}) = 1\}$ is the set of neighbor nodes for gene $G_{ij} \in \mathcal{S}_i$, then $G_{ij} \in \Psi_i$ if and only if we have $|\theta_j \cap \mathcal{S}_i| = 1$.

To evaluate each candidate gene G_{ik} , we exclude it from the corresponding subnetwork and then run a coalition game for the set of players including the current subnetworks and the candidate genes of interest; i.e. $\mathcal{P} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_i \setminus G_{ik}, \dots, \mathcal{S}_{N_L}, G_{ik}\}$, where $|\mathcal{P}| = L + 1$. The Shapley value for the candidate gene G_{ik} quantifies the expected marginal importance of the candidate gene in coalition with different combinations of the current subnetworks. The gene with minimal Shapley value $G_i^{(\text{remove})} = \underset{G_{ik} \in \Omega_i}{\operatorname{argmin}} \gamma_{G_{ik}}(v)$ has the lowest contribution to the game outcome and hence can be removed from the host subnetwork if its contribution is below a predefined threshold value T_{removal} (i.e. $G_i^{(\text{remove})} \leq T_{\text{removal}}$). The removal threshold T_{removal} should be chosen conservatively and below the expansion threshold T_{add} to avoid removing informative nodes ($T_{\text{removal}} < T_{\text{add}}$). A reasonable value is $T_{\text{removal}} = 0$ to remove only the nodes with negative contribution to the game outcome.

Repeating this game for all removal candidate genes for all subnetworks, provides the resulting subnetworks at the end of the current iteration.

E. Subnetwork evaluation and termination criteria check

At each algorithm iteration, the game-based expansion and contraction steps provide a set of at most L subnetworks. As mentioned earlier, subnetworks are allowed to join each other and form a new larger subnetwork. At the end of each iteration, the performance of the obtained subnetworks denoted by $Acc^{(t)}$ is evaluated based on the appropriate performance metric for the application of interest; i.e. $Acc^{(t)} = v(\mathcal{S}_1^{(t)} \cup \mathcal{S}_2^{(t)} \cup \dots \cup \mathcal{S}_L^{(t)})$. The iterative algorithm stops if one of the following terminations conditions are met:

- i) If the termination flags for all of the subnetworks are set either due to lack of an informative node in their neighborhood (i.e. $\gamma_{G_{ik}}(v) < T_{add}$ for $\forall G_{ik} \in \Omega_i$) or due to reaching the maximum number of nodes for a subnetwork (i.e. $|\mathcal{S}_i^{(t)}| \geq S_{max}$, where postscript (t) denotes iteration t).
- ii) If the number of total collected genes exceeds a predefined value (i.e. $\sum_{i=1}^N |\mathcal{S}_i^{(t)}| \geq N_{max}$).
- iii) If the number of iterations reach its maximum value N_T (i.e. $t > N_T$).
- iv) If all subnetworks remain unchanged during the current iteration and no gene is added or removed from a subnetwork. (i.e. $\mathcal{S}_i^{(t)} = \mathcal{S}_i^{(t-1)}$ for $i = 1, 2, \dots, L$).
- v) If the desired performance at the end of iteration t denoted by $Acc^{(t)}$ based on the collected genes forming the developed subnetworks reaches a predefined desired value (i.e. $Acc^{(t)} \geq Acc_{max}$).

We note that the last two conditions are the most important conditions controlling the performance of the algorithm by timely terminations. Other termination criteria are defined to avoid non-termination of the algorithm when it does not converge to a meaningful set of subnetworks, and thus these conditions should be chosen loose enough to avoid early terminations.

If the termination criterion is satisfied, the algorithm terminates and the subnetworks along with the obtained performance are reported. Otherwise, the next iteration is executed.

As mentioned earlier, since the algorithm starts with randomly chosen seed genes, the algorithm is executed multiple times (denoted by R) and the resulting subnetworks for the algorithm execution that yield the best performance are reported. The summary of each run of the algorithm is presented in Figure 2.

III. RESULTS

In this section, the proposed algorithm is utilized to solve two important cancer-related problems including i) identification of gene subnetworks whose expressions mediate platinum-based therapy response in ovarian cancer, ii) identification of gene subnetworks whose somatic disorders significantly impact immune system scores in breast cancer. The main difference between utilizing the proposed algorithm for the two applications is coalition evaluation method based on the characteristic of training dataset. There are also some other minor differences that will be elaborated by detailing utilization of the proposed NBGC algorithm for both applications in the following sections.

A. *Therapy response prediction in ovarian cancer*

In this application, we are interested in finding subset of biologically relevant gene sets that significantly impact therapy response in ovarian cancer. Ovarian cancer is an extremely aggressive disease with poor overall outcomes due to late diagnosis and lack of targeted therapies [4]. Furthermore, a majority of ovarian cancer patients progress or suffer cancer recurrence within 5 years of frontline platinum-based therapy. Identification of gene subnetworks associated with resistance to platinum-based therapy resistance [23] [24] would enable the discovery of functional biomarkers and novel therapeutic targets in this aggressive disease. Accordingly, we here apply our algorithm to identify gene subnetworks within the human PPI network whose joint expression levels are associated with recurrence free survival on platinum-based therapy in ovarian cancer. The molecular data is obtained from The Cancer Genome Atlas (TCGA) dataset [25] and includes 201 cancer samples with their gene expression levels and clinical responses. The clinical response data includes survival information (death or cancer progression) after platinum-based chemotherapy. We first divide the samples into two cohorts of poor and good survival rates. The poor survival cohort includes samples with events during the first 6 months of receiving platinum therapy, excluding patients who left the study (censored samples). Patients who survive at least 6 months without cancer progression are included in the good survival cohort. Therefore, therapy response classes are represented by a binary vector $y_{[n_S \times 1]}$, where $n_S = 201$ is the number of samples.

The dataset $X_{[n_S \times n_G]}$ includes continuous-valued gene expression data for $n_G = 9544$ genes and $n_S = 201$ samples. The genes are the intersection of genes with available expression data and the genes whose protein product are included in the PPI networks. Noting the datatype, we choose the classification rate (based on 5-fold cross-validation) as the characteristic function $v(\cdot)$ in the proposed game theoretic NBGC algorithm. We choose to use the binary class labels as therapy response identifiers in order to accelerate algorithm runs. Given these binary class labels for phenotyping, a natural selection for characteristic function is classification accuracy. Further, we use binary classification during training and test phases of the proposed algorithm to obtain the results, whereas we use the method of Kaplan-Meier estimation after K-means clustering in order to compare the reported genes by our algorithm and other competitor feature selection methods to avoid bias to a specific test method.

We use the top-100 genes with highest degree in the PPI network (hot-spot nodes) as initialization gene pool. We only use network expansion step in the implementations for this application. The rest of the game parameters are set as follows:

The number of subnetworks are set to $L = 5$, which results in $N=L+I=6$ players including the candidate gene, maximum group size $N' = 3$, maximum number of collected genes at each subnetwork $S_{\max} = 20$, maximum number of total collected genes $N_{\max} = 100$, maximum number of internal iterations $N_T = 100$. The threshold on the Shapley value for including a new gene into the subnetwork is $T_{add}=0$ meaning that the gene with maximal Shapley value is accepted if it improves the classification accuracy (in average) when joins the previously formed coalitions). The desired classification rate to stop the algorithm is set to $Acc_{\max} = 0.98$. Large value of Acc_{\max} reduces false report rate by sticking to the algorithm runs that provide very impactful subset of gene subnetworks.

The algorithm runs for $R = 100$ rounds using randomly selected genes among the initialization gene pool. We first rank the genes based on their degree in PPI network and select top-100 genes as initialization pool, and then we randomly draw a subset of 4 genes for each algorithm execution. At each round of execution, the algorithm reports a collection of subnetworks that are highly associated with the survival outcomes. The proposed solution can be integrated with any binary classification method. In this work, we arbitrarily use the SVM classification with RBF kernel. However, the obtained results is not sensitive to the choice of classifier and the numerical results show negligible change when using other classifiers (such as Random Forest, Naive Bayes, Bayes Net, and KNN).

We compare the results with the same number of genes obtained using two benchmark solution categories. We apply state-of-the-art feature selection methods including Correlation based subset evaluation (CFS), Chi-square test based subset evaluation (Chi-Square) and mutual-information based subset evaluation method (Gain-Ratio). Additionally, two representative wrapper methods including Best First Search (BFS) method with Naive-Bayes and ranker method with SVM classifier were also applied. These methods report the most informative genes that may or may not belong to connected subnetworks. We also compare the proposed method with a network-based traversal method, where the subnetworks are initiated from the same initial gene seeds as in our proposed method. Instead of using the Shapely value, genes from the connected subnetwork in proximity of the seed genes are selected using a random walk until it collects the same number of genes as the proposed method. This whole procedure is repeated for $R=100$ times and the set of gene-subnetworks which provides the best prediction accuracy is selected for comparison with the proposed methodology.

In order to compare the relevance of the obtained gene sets across methods, in addition to the classification accuracy based on phenotype survival rates, we also compare the discriminative power of the gene sets in terms of continuous-valued survival probabilities. Therefore, patients are clustered using K-means clustering based on the gene expression data for the selected genes reported by the different methods. Then, we estimate the survival probability for each cluster using the standard method of Kaplan-Meier estimation followed by survival difference estimation using the log-rank test method that provides the probability of obtaining such a

difference purely by chance (p-value). The results of these comparisons are provided in Table. 1 and Figure 3.

Table 1 compares the proposed solution with the aforementioned methods. For competitor methods, where the sorted list of genes are provided based on various correlation or information based measures, we choose the same number of the top-genes that are reported by our method, which is 18. We note that 18 is not a predefined value, but rather it is the number of genes that is reported by the algorithm, when it meets the stop criterion that is a set of rules based on various parameters as detailed in section II-E. It is seen that the proposed method outperforms the other competitor methods. The main cause is that the proposed coalition-based solution considers the collective power of gene sets based on Shapley value concept. This is in particular interesting, since the competitor methods are not restricted to choose the genes from connected subnetworks. The proposed solution provides more insightful and clinically relevant gene subnetworks.

The set of subnetworks identified by our proposed algorithm for $L = 5$ is depicted in Figure 3a. Subnetworks 1, 2, 4 correspond to i) vascular endothelial regulation, ii) TGFb signaling and iii) cell cycle progression and apoptosis pathways, respectively. These pathways belong to well-known hallmarks of cancer, thus suggesting that our proposed methodology is able to identify potentially functional pathways mediating therapy resistance. The subnetwork 5 joins subnetwork 4 at iteration 3 and subnetwork 3 stops extending at iteration 3, since no new informative neighbor genes were available. Indeed, while the TGF-beta signaling has previously been implicated as playing an oncogenic role in epithelial cancers, its role in mediating platinum resistance has not been widely explored. In addition, Subnetwork 3, which includes the gene *TXNDC9*, is also a novel finding in this disease context. Notably, while the exact function of *TXNDC9* is as yet unknown, it is likely involved in cell differentiation, and has been shown to be associated with increased risk of progression in colon cancer [26]. Thus, our findings using the NBCG algorithm showing *TXNDC9* being associated with platinum resistance in ovarian cancer makes it further likely that this gene could be a potential therapeutic target across disease contexts. These results therefore reveal the ability of NBCG to identify novel mechanisms and likely therapeutic targets across cancers.

Figure 3b presents the survival curves for patients clustered into two groups using K-means clustering based on the expression level of the genes obtained from the proposed game-theoretic method (Figure 3a). The result demonstrates that the proposed solution can identify gene subnetworks with higher survival discriminatory power as compared to the estimates from the best feature selection method [in this case, CFS] (Figure 4a), and the Optimal Network-based Random-walk solution (Figure 4b).

B. Impact of Somatic Aberrations on Immune system scores in triple-negative breast cancer (TNBC)

Multiple studies have demonstrated that evaluating the extent of tumor infiltrating lymphocytes within triple negative breast tumors can provide significant prognostic information [27] [28] [29] [30] in this highly aggressive subtype of breast cancer. Indeed, we [31] and others [32] [33] [34] have shown that gene expression based immune signatures can provide both information on prognosis and response to therapy in specific subsets of breast

cancers. However, it is as yet unclear why certain subsets of triple negative breast cancers exhibit high levels of immune surveillance, while others escape it. Evidence in colon cancer suggests that genomic aberrations such as microsatellite instability are likely to contribute to increased immune surveillance [35], thus potentiating the use of novel immune therapies [36] [37] in this cancer subtype. We therefore hypothesize that genomic aberrations affecting key pathways involved in maintaining genomic fidelity are more likely associated with high lymphocytic infiltration in triple negative breast cancer. Accordingly, we here apply our algorithm to identify mutated gene subnetworks that are associated with differential immune infiltration in TNBC.

We leveraged gene expression data from the TCGA and estimated the level of immune infiltration in the tumor samples using a previously published algorithm, ESTIMATE [38] that is based on the expression levels of a 140-gene immune signature. Briefly, a single-sample gene set enrichment methodology was used to derive the index of immune activity (Immune Index) by comparing the expression levels of the 140 signature-genes against the background expression of all genes on the array [38]. The Immune Index values for all samples in this study were estimated using the R-package associated with the published algorithm. Somatic mutations were determined using whole exome sequencing data while somatic copy-number alterations were obtained using SNP-arrays. Subsequently, we applied our algorithm to identify gene subnetworks whose mutations or copy-number alterations were significantly associated with high Immune Indexes.

Accordingly, we first obtained gene-level somatic copy-number alterations by mapping the sCNA loci on to genes, with log-ratios > 0.1 being considered amplification events while log-ratios < -0.1 being considered as deletion events. Both amplification and deletion events mapped to a copy number alteration event represented by '1' in the corresponding binary matrix. Similarly, we incorporated gene-level mutation information, with a gene considered mutated if it harbors either missense, nonsense single nucleotide changes, Insertion/Deletions, and frame-shift indels. Silent mutations were excluded for the purposes of this analysis.

Therefore, both Mutation and sCNA data are $n_S \times n_G$ binary matrices denoted by $X_{[n_S \times n_G]}$, where n_S is the number of samples and n_G is the number of genes satisfying two conditions: i) profiling information is available and ii) they are included in the employed PPI network. We tried three different scenarios by incorporating mutation data only, copy number alteration data only and the combination of both. In the later scenario, '1' in the binary data matrix refers to copy number alteration event, a mutation event or both. The data matrix includes binary information for $n_G = 9976$ genes and $n_S = 109$ TNBC samples. The triple-negative breast cancers were identified using immunohistochemistry data available from the TCGA for the tumor samples [39] [40].

The Immune scores are incorporated as a $n_S \times n_G$ column vector of continuous-valued immune score denoted by $y_{[n_S \times 1]}$. In evaluating each coalition of genes, we are interested in finding how somatic aberration are associated with the immune scores. In this regard, for each coalition, we divide the training samples into two neutral and impacted cohorts. The first cohort includes the sample for which at least one of the genes experience an aberration event, whereas the second cohort includes the samples with all genes in their normal situations. Therefore, if X is the $n_S \times n_G$ input data matrix (either sCNA, Mut or combination) for n_S

samples and total n_G genes, and $\mathcal{S}_i = \{G_{i1}, G_{i2}, \dots, G_{i|\mathcal{S}_i|}\}$ is the coalition of genes in subnetwork \mathcal{S}_i , then the neutral and impacted cohorts are:

$$X_0(\mathcal{S}_i) = \{s: \sum_{g \in \mathcal{S}_i} X(s, g) = 0\}, \quad X_1(\mathcal{S}_i) = \{s: \sum_{g \in \mathcal{S}_i} X(s, g) > 0\}. \quad (5)$$

In order to evaluate the separation between the two clusters in terms of immune scores y , we use fisher index :

$$F_i(\mathcal{S}_i) = \frac{|\mu_{x_0} - \mu_{x_1}|}{\sigma_{X_0}^2 + \sigma_{X_1}^2}, \quad (6)$$

where μ_{x_i} and $\sigma_{x_i}^2$ are respectively the mean and variance of two cohorts. For this application, we use three different input types including i) sCNA, ii) Mut and iii) sCNA +Mut, where in the last scenario ‘1’ means a sCNA event, a mutation event or both. In this application, we used both network expansion and contraction steps.

For each scenario, we execute the algorithm for $R = 1000$ rounds using randomly chosen seed genes from the pool of top-100 genes based on their degrees in the PPI network. We also tried other initializations methods including: i) top 100 most frequently mutated genes in BRCA and top top-100 highly correlated genes whose individual mutations which divide the samples between two cohorts with highest fisher index for their immune scores. The numerical results suggest that due to the flexibility and mobility of the formed subnetworks by the algorithm which is provided by network expansion and contraction steps, the result are not sensitive to the choice of initialization method and perform almost equally. However, the proposed method provides the desired performance with lower number of genes due to the existence of shorter paths from randomly selected seed genes to the modulated subnetworks.

The rest of game-theoretic algorithm parameters are set to: maximum number of subnetworks $L = 4$ resulting in $N = 5$ players including the candidate gene, maximum number of players in a coalition $N' = 4$, maximum number of total collected genes $N_{max} = 100$, , maximum number of collected genes in each subnetwork $S_{max} = 50$, and maximum number of internal iterations $N_T = 100$. The browsing parameters such as minimum improvement that controls adding a node to a subnetwork T_{add} , minimum performance degradation by removing a node from a subnetwork $T_{removal}$, and desired performance Acc_{max} that controls the network browsing and stop criterion can be defined in terms of the distance between the resulting cohort centers $|\mu_{x_1} - \mu_{x_0}|$ for the given coalition as well as fisher indexes $F_i(\mathcal{S}_i)$. Since the phenotypes are continuous-valued immune scores, a natural selection for characteristic function is their fisher index since it reflects the separability of clusters obtained based on the mutations in the examined gene subnetworks. Considering the range of immune scores [-110: 8765], we set these parameters based on the cohort center distances to the follows: $Acc_{max} = 5000$, $T_{add} = 0$, $T_{removal} = -200$. Acc_{max} is chosen heuristically considering the range of immune scores. Before elaborating on the resulting subnetworks, we first evaluate the consistency of the proposed algorithm.

C. Convergence analysis

In order to evaluate the robustness of the proposed method, we study the consistency of results when initialized with different seed genes. In our case, the initialization gene pool includes 100 genes, where at each run we choose $L = 4$ seed genes to develop L subnetworks. The total number of different seed genes is $\binom{100}{L} = 3.9 \times 10^6$, therefore the with high probability each round of algorithm starts with a different seed genes. We are interested in evaluating how consistent are the results, when starting with different seed genes.

Each round of the proposed NBGC algorithm provides up to $L(r) \leq L$ modulated subnetwork with the resulting gene set $\mathcal{G}(r)$ with total of $N(r) = |\mathcal{G}(r)| \leq \min(L S_{\max}, N_{\max}) = N_{\max}$ genes. Therefore, the total number of genes reported by the algorithm at R rounds are limited to:

$$k_{\max} = \cup_{r=1}^R \mathcal{G}(r) \leq \sum_{r=1}^R N(r) \leq RN_{\max} \quad (7)$$

Since, the ground truth is not known for this application, we first rank the set of reported genes based on their frequency of appearance in multiple rounds of the algorithms and then consider the top- k most frequently reported genes $\mathcal{G}_k = \{g: \sum_{r=1}^R 1(g \in \mathcal{G}(r)) \geq k\}$ as the optimal result, where $1(\cdot)$ is the indicator function.

For each round of algorithm, $n_k(r)$ denotes the number of reported genes that are in top- k genes, i.e. $n_k(r) = |\mathcal{G}(r) \cap \mathcal{G}_k|$. These genes are considered consistent. The rest of $N - n_k(r)$ genes are considered false reports which are inconsistent with the average results obtained by multiple rounds of the algorithm. This is depicted in Figure 5. For instance, the round $r = 3$ corresponds to an algorithm round with poor consistency.

The consistency score of round r denoted by $c_k(r)$ is defined as the rate of the reported genes that belong to \mathcal{G}_k :

$$c_k(r) = \frac{|\mathcal{G}(r) \cap \mathcal{G}_k|}{|\mathcal{G}(r)|} = \frac{|\mathcal{G}(r) \cap \mathcal{G}_k|}{N(r)} \quad (8)$$

The average of $c_k(r)$ over all R algorithm rounds (i.e. $\mathbf{c}_k = \frac{1}{R} \sum_{r=1}^R c_k(r)$) provides the consistency curves that is depicted in Figure 6(a), where gene acceptance rate is defined as $|\mathcal{G}_k|/|\mathcal{G}_{k_{\max}}|$. The frequency of appearance of top k genes in multiple runs of the proposed algorithm is depicted in Figure 6(b). In this Figure, the frequency of appearance for any of the top-100 genes are depicted. For instance, the first top-gene appear at about 50% of the rounds of algorithm (approximately 500 out of 1000). Also, the top-20 genes also appear in about 20% of the algorithm runs. Therefore, using only 5 rounds of the proposed algorithm with different initialization ensures capturing of the top-20 genes with high probability. The consistency results in Figures 6(a) and 6(b) are valid for the three different scenarios using sCNA, Mut and combined data types.

Association between Mutated Subnetworks and Immune Scores

Figure 7 represent the obtained results for the mutation data. Each execution of algorithm provides top modulated subnetworks along with the resulting fisher index for immune score between two cohorts, with and without modulated subnetworks. Then, we report the modulated subnetworks for the algorithm round that provides the highest fisher index.

Given the evidence for convergence of the proposed algorithm, we evaluated the biological significance of the identified subnetworks. We specifically focused on identifying gene subnetworks harboring somatic mutations that are associated with high immune infiltration (Figure 7), given the clinical significance of high immune infiltration in TNBCs. As detailed in Figure 7, we identified a densely connected subnetwork of genes that jointly exhibited a very strong association with high immune index (Fisher's Index = 3.25), with TNBC samples harboring mutations in these genes exhibiting significantly higher Immune Index (Differential Immune Index = 4948), as compared to the rest of the cohort. We performed pathway enrichment analysis on this gene set using the National Cancer Institute's Pathway Interaction Database (NCI-PID) [41], a curated collection of known biomolecular interactions and key signaling pathways associated with cancer, to evaluate if genes belonging to a specific cancer-related pathways were enriched within the subnetwork, followed by assessment of false discovery rate using the Benjamini-Hochberg False Discovery Rate (FDR) [42] methodology. Our analysis revealed significant enrichment of ATR signaling ($P \ll 10^{-3}$; $FDR \ll 10^{-2}$) and p53 pathways ($P \ll 10^{-3}$; $FDR \ll 10^{-2}$), pointing to the likelihood that TNBC tumors deficient in DNA damage repair mechanisms are more likely to trigger enhanced immune surveillance. Indeed, mutations in MDM2, CHEK2, and BRCA2, all of whom are key players in DNA damage repair were frequently identified as significantly associated with differential immune index (Figure 7) by our algorithm. These findings strongly suggest that immune surveillance in TNBC tumors is associated with disruptions in DNA repair pathways, a finding consistent with previous reports of increased immune cell infiltrates in high-grade serous ovarian cancers harboring BRCA1 or BRCA2 mutations [43].

However, the immunogenic role of these pathways in TNBC has not been previously reported. Additionally, the NBCG algorithm also identified aberrations in *AKT1* and *PTEN*, both of which belong to the oncogenic PI3K/AKT pathway to be associated with increased immune infiltration (Figure 6). Taken together, these novel findings, if confirmed in additional studies, would enable the identification of tumor-specific neoantigens, while also identifying targets for vaccines and adoptive immune cell therapies. Furthermore, given that multiple trials are ongoing to evaluate the benefit of PD-L1/PD-1 blockade in TNBC, our findings could be leveraged to identify a mutational signature of benefit from immune checkpoint inhibitor therapies.

Furthermore, given that multiple trials are ongoing to evaluate the benefit of PD-L1/PD-1 blockade in TNBC, our findings, if validated, could be leveraged to identify a mutational signature of benefit from immune checkpoint inhibitor therapies.

IV. CONCLUSIONS

A novel Coalition game theory based algorithm is developed using PPI networks to identify gene subnetworks that are associated with clinical or biological phenotypes. This algorithm implements a novel subnetwork selection mechanism by utilizing two network expansion and contraction steps. The proposed Network Based Coalition Game (NBCG) algorithm initiates subnetworks from randomly selected or predefined seed nodes and expands the subnetworks by collecting nodes from their neighborhoods based on the concept of Shapley value, where the subnetworks and the candidate genes are the players of the developed coalition game. This approach considers the collective power of the subnetworks using Shapley value based on an application-independent characteristic function. To avoid the bias to seed nodes, an additional step of network contraction is developed by which the subnetworks remove uninformative nodes and thereby are capable of moving away from the seed nodes toward more informative network sections, in case the originally selected seed nodes are not significantly useful. This algorithm is general in the sense that it is capable of identifying subnetworks of connected nodes that maximizes a desired objective based on the collective power of the set of selected network nodes.

The proposed NBCG algorithm is utilized to find subnetworks of genes whose profiling information is associated with a desired clinical outcome or a biological mechanism. As an illustrative example, we employed this algorithm to determine gene subnetworks whose expression levels are associated with progression-free survival in patients with ovarian cancer treated with platinum-based chemotherapy. Additionally, we employed the NBCG algorithm to identify gene networks harboring genomic aberrations that are associated with and immune cell infiltration in triple negative breast cancers. The proposed method improves upon the state of the art feature selection methods in identifying the most informative gene sets with the added advantage of using protein-protein interaction networks to identify functionally related gene sets that jointly discriminate between phenotypes. Additionally, this approach takes into account the collective power of subnetworks using the concept of Shapley value, as opposed to techniques that grow each subnetwork individually.

The above mentioned illustrative applications confirm the utility of the NBCG algorithm in identifying subnetworks of functionally related genes and proteins that are associated with cancer phenotypes, thus enabling the discovery of novel biomarkers and therapeutic targets in cancer.

V. REFERENCES

- [1] J. Lamb, E. Crawford, D. Peck, J. Modell, I. Blat, M. Wrobel and e. al., "The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease," *Science*, vol. 313, pp. 1929-1935, 2006.
- [2] A. Subramanian, P. Tamayo, V. Mootha, S. Mukherjee, B. Ebert, M. Gillette and e. al., "Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles," *PNAS*, vol. 102, pp. 15545-15550, 2005.

- [3] R. Tibshirani, T. Hastie, B. Narasimhan and G. Chu, "Diagnosis of multiple cancer types by shrunken centroids of gene expression," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 99, pp. 6567-6572, 2002.
- [4] D. Bell and e. al., "Integrated genomic analyses of ovarian carcinoma," *Nature*, vol. 474, pp. 609-15, 2011.
- [5] C. Perou and e. al., "Comprehensive molecular portraits of human breast tumours," *Nature*, vol. 490, pp. 61-70, 2012.
- [6] R. Kucherlapati and e. al, "Comprehensive molecular characterization of human colon and rectal cancer," *Nature*, vol. 487, pp. 330-337, 2012.
- [7] P. Creixell and e. al, "Pathway and network analysis of cancer genomes," *Nature Methods*, vol. 12, pp. 615-621, 2015.
- [8] L. J. Lu and e. al., "Comparing Classical Pathways and Modern Networks: Towards the Development of an Edge Ontology," 2015.
- [9] A. Sarveniazi, "An Actual Survey of Dimensionality Reduction," *American Journal of Computational Mathematics*, vol. 4, no. 2, pp. 55-72, 2014.
- [10] L. van der Maaten, E. Postma and H. van den Herik, "Dimensionality Reduction: A Comparative Review," *Machine Learning*, vol. 10, pp. 66-71, 2009.
- [11] G. Ghandrashekar and F. Sahin, "A survey on feature selection methods," *Elsevier Computers & Electrical Engineering*, vol. 40, no. 1, pp. 16-28, 2014.
- [12] Y. Saeys, I. Inza and P. and Larranage, "A review of feature selection techniques in bioinformatics," *Bioinformatics*, vol. 23, pp. 2507-2517, 2007.
- [13] J. Soul, T. Hardingham, R. Boot-Handford and J. Schwartz, "PhenomeExpress: A refined network analysis of expression datasets by inclusion of known disease phenotypes," *Scientific Reports*, 2015.
- [14] H. Chuang, E. Lee, Y. Liu, D. Lee and I. T, "Network-based classification of breast cancer metastasis," *Molecular Systems Biology*, vol. 3, no. 140, 2007.
- [15] D. Ray, *A Game-Theoretic Perspective on Coalition Formation*, New York: Oxford University Press, 2007.
- [16] J. Nash, "Non-Cooperative Games," *The Annals of Mathematics, Second Series*, vol. 54, no. 2, pp. 286-295, 1951.
- [17] A. Razi, F. Afghah, A. Belle and K. N. K. Ward, "Blood Loss Severity Prediction using Game Theoretic Based Feature Selection," in *2014 IEEE-EMBS International Conference on Biomedical and Health Informatics (BHI)*, pp. 776-780, 2014.
- [18] X. Sun, Y. Liu, J. Li, J. Zhu, H. Chen and X. and Liu, "Feature Evaluation and Selection with Cooperative Game Theory," *Pattern Recognition*, vol. 45, no. 8, p. 2992-3002, 2012.
- [19] F. Afghah, A. Razi and K. Najarian, "A Shapley Value Solution to Game Theoretic-based Feature Reduction in False Alarm Detection," in *Neural Information Processing Systems (NIPS), Workshop on Machine Learning in Healthcare*, arXiv:1512.01680, 2015.
- [20] F. Afghah, A. Razi, S. Soroushmehr, S. Molaei, H. Ghanbari and K. and Najarian, "A Game Theoretic Predictive Modeling Approach to Reduction of False Alarm," in *2015 International Conference for Smart Health (ICSH'15)*, Mayo Clinic, 2015.
- [21] A. Razi, F. Afghah and V. Varadan, "Identifying Gene Subnetworks Associated with Clinical Outcome in Ovarian Cancer using Network Based Coalition Game," in *37th Annual International Conference of the IEEE Engineering in Medicine and Biology Conference (EMBC'15)*, Aug. 2015.
- [22] S. Razick, G. Magklaras and I. Donaldson, "iRefIndex: a consolidated protein interaction database with provenance," *BMC Bioinformatics*, vol. 9, 2008.
- [23] R. Tothill, A. Tinker, J. George, R. Brown, S. Fox, S. Lade and e. al., "Novel molecular subtypes of serous and endometrioid ovarian cancer linked to clinical outcome," *Clin Cancer Research*, vol. 14, pp. 5198-5208, 2008.

- [24] K. Wrzeszczynski, V. Varadan, S. Kamalakaran, D. Levine, N. Dimitrova and R. Lucito, "Integrative prediction of gene function and platinum-free survival from genomic and epigenetic features in ovarian cancer," *Methods in Molecular Biology*, vol. 1049, pp. 35-51, 2013.
- [25] "The Cancer Genome Atlas," [Online]. Available: <http://cancergenome.nih.gov/>.
- [26] A. Lu, X. Wangpu, D. Han, H. Feng, J. Zhao, J. Ma and a. et., "TXNDC9 expression in colorectal cancer cells and its influence on colorectal cancer prognosis," *Cancer investigation*, vol. 30, pp. 721-726, 2012.
- [27] R. Salgado, C. Denkert, S. Demaria, N. Sirtaine, F. Klauschen, G. Pruneri and e. al., "The evaluation of tumor-infiltrating lymphocytes (TILs) in breast cancer," *Annals of oncology: official journal of the European Society for Medical Oncology / ESMO*, vol. 26, no. recommendations by an International TILs Working Group 2014., pp. 259-271, 2015.
- [28] S. Loi, S. Michiels, R. Salgado, N. Sirtaine, V. Jose, D. Fumagalli and e. al., "Tumor infiltrating lymphocytes are prognostic in triple negative breast cancer and predictive for trastuzumab benefit in early breast cancer: results from the FinHER trial," *Annals of oncology : official journal of the European Society for Medical Oncology / ESMO.*, vol. 25, pp. 1544-1550, 2014.
- [29] S. Loi, "Tumor-infiltrating lymphocytes, breast cancer subtypes and therapeutic efficacy," *Oncimmunology*, vol. 2, p. e24720, 2013.
- [30] L. S., N. Sirtaine, F. Piette, R. Salgado, G. Viale, F. Van Eenoo and e. al., "Prognostic and predictive value of tumor-infiltrating lymphocytes in a phase III randomized adjuvant breast cancer trial in node-positive breast cancer comparing the addition of docetaxel to doxorubicin with doxorubicin-based chemotherapy: BIG 02-98," *Journal of Clinical Oncology*, vol. 31, pp. 860-867, 2013.
- [31] V. Varadan, H. Gilmore, K. Miskimen, D. Tuck, S. Parsai, A. Awadallah, I. Krop, E. Winer, V. Bossuyt, G. Somlo, M. Abu-Khalaf, M. Fenton, W. Sikov and L. Harris, "Immune Signatures Following Single Dose Trastuzumab Predict Pathologic Response to Preoperative Trastuzumab and Chemotherapy in HER2-Positive Early Breast Cancer," *Cancer Research*, p. in press, 2016.
- [32] E. Perez, E. Thompson, K. Ballman, S. Anderson, Y. Asmann, K. Kalari and e. al., "Genomic Analysis Reveals That Immune Function Genes Are Strongly Linked to Clinical Outcome in the North Central Cancer Treatment Group N9831 Adjuvant Trastuzumab Trial.," *Journal of Clinical Oncology*, vol. 33, pp. 701-708, 2015.
- [33] M. Iglesia, B. Vincent, J. Parker, K. Hoadley, L. Carey, C. Perou and e. al., "Prognostic B-cell signatures using mRNA-seq in patients with subtype-specific breast and ovarian cancer," *Clinical Cancer Research*, vol. 20, pp. 3818-3129, 2014.
- [34] C. Gu-Trantien, S. Loi, S. Garaud, C. Equeter, M. Libin, A. de Wind and e. al., "CD4(+) follicular helper T cell infiltration predicts breast cancer survival," *The Journal of clinical investigation*, vol. 123, pp. 2873-2892, 2013.
- [35] D. Tougeron, E. Fauquembergue, A. Rouquette, F. Le Pessot, R. Sesboue, M. Laurent and e. al., "Tumor-infiltrating lymphocytes in colorectal cancers with microsatellite instability are correlated with the number and spectrum of frameshift mutations," *Modern pathology: an official journal of the United States and Canadian Academy of Pathology*, vol. 22, pp. 1186-1195, 2009.
- [36] N. Llosa, M. Cruise, A. Tam, E. Wicks, E. Hechenbleikner, J. Taube and e. al., "The vigorous immune microenvironment of microsatellite instable colon cancer is balanced by multiple counter-inhibitory checkpoints," *Cancer discovery*, vol. 5, pp. 43-51, 2015.
- [37] D. Le, J. Uram, H. Wang, B. Bartlett, H. Kemberling, A. Eyring and e. al., "PD-1 Blockade in Tumors with Mismatch-Repair Deficiency," *The New England Journal of Medicine*, vol. 372, pp. 2509-2520, 2015.
- [38] K. Yoshihara, M. Shahmoradgoli, E. Martinez, R. Vegesna, H. Kim, W. Torres-Garcia and e. al., "Inferring tumour purity and stromal and immune cell admixture from expression data," *Nature communications*, vol. 4, p. 2612, 2013.
- [39] R. R. Bastien and e. al, "PAM50 Breast Cancer Subtyping by RT-qPCR and Concordance with Standard Clinical Molecular Markers," *BMC Medical Genomics*, vol. 5, no. 44, pp. 1-12, 2012.

- [40] C. Sweeney, "Intrinsic Subtypes from PAM50 Gene Expression Assay in a Population-Based Breast Cancer Cohort: Differences by Age, Race, and Tumor Characteristics," *AACR Cancer Epidemiol Biomarkers Prev*, vol. 23, no. 5, pp. 714-724, 2014.
- [41] C. Chaefer, K. Anthony, S. Krupa, J. Buchoff, M. Day, T. Hannay and e. al., "PID: the Pathway Interaction Database," *Nucleic acids research*, vol. 37, pp. D674-D679, 2009.
- [42] Y. Hochberg and Y. Benjamini, "More powerful procedures for multiple significance testing," *Statistics in medicine*, vol. 9, pp. 811-818, 1990.
- [43] J. McAlpine, H. Porter, M. Kobel, B. Nelson, L. Prentice, S. Kalloger and e. al., "BRCA1 and BRCA2 mutations correlate with TP53 abnormalities and presence of immune cell infiltrates in ovarian high-grade serous carcinoma.," *Modern pathology*, vol. 25, no. 5, pp. 740-750, 2012.
- [44] J. Fan, R. Samworth and Y. and Wu, "Ultrahigh dimensional feature selection: Beyond the linear model," *Journal of Machine Learning Research*, vol. 10, pp. 2013-2038, 2009.
- [45] L. S. Shapley, "A value for n-person games," in *Contributions to the Theory of Games, volume II*, H.W. Kuhn and A.W. Tucker (eds.), Princeton University Press, 1953, pp. 307-317.
- [46] A. Razi, F. Afghah and V. Varadan, "Identifying Gene Subnetworks Associated with Clinical Outcome in Ovarian Cancer using Network Based Coalition Game," in *International Conference of the IEEE Engineering in Medicine and Biology Conference (EMBC)*, Milan, 2015.
- [47] A. Subramaniana, P. Tamayoa, V. Moothaa and e. al., "Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles," *PNAS*, vol. 102, no. 43, p. 15545–15550, 2005.

TABLES

Table I. Comparison of genes selected using the proposed method and other state-of-the-art feature selection methods based on its prediction accuracy and survival outcome separation. The results are corresponding to the first 18 genes reported by each method.

Method	Log-rank Test P-value	Prediction Success Rate
CFS	0.01814	0.6488
Chi-Square	0.25505	0.6667
Gain-Ratio	0.47773	0.5179
Best First Search	0.07646	0.5714
SVM: Ranker	0.09190	0.5714
Optimal Random-Walk	0.08060	0.6190
Proposed NBCG	0.00004	0.7262

FIGURES

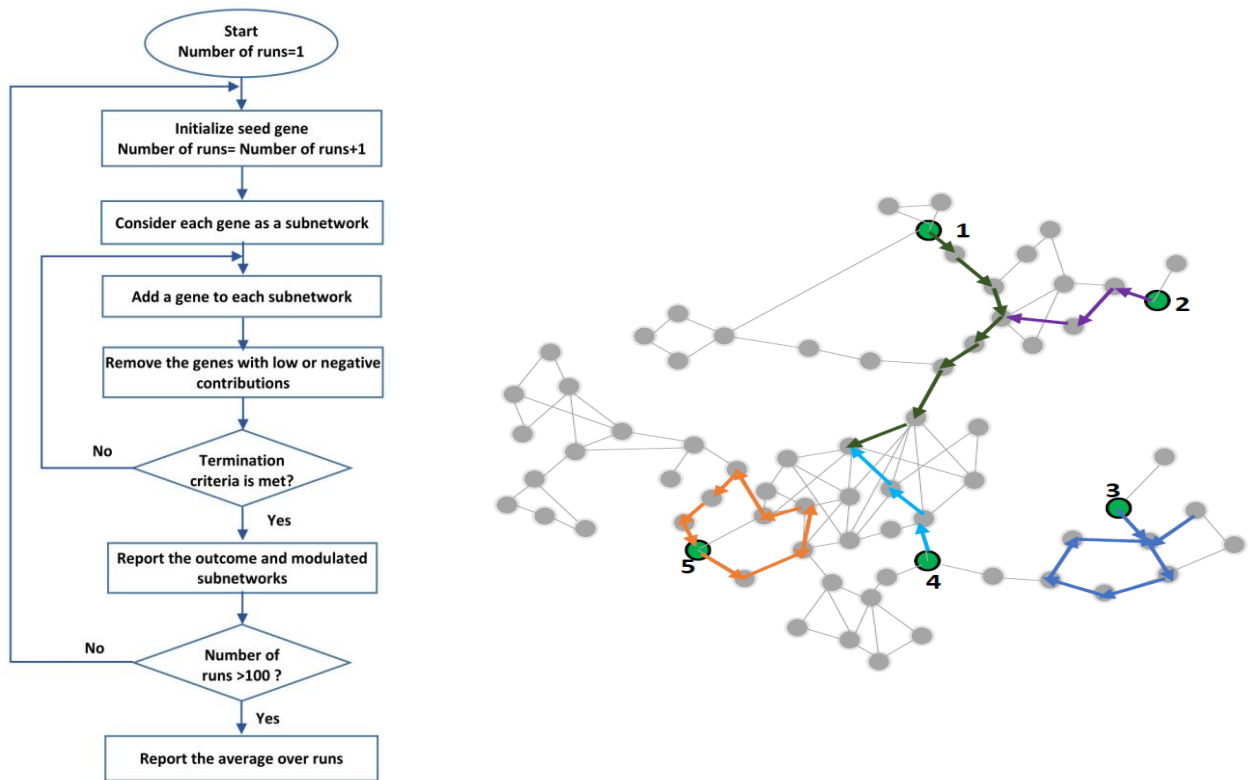


Figure 1. An example of developing 5 modulated subnetworks over a PPI network. The initialization seed genes are marked with green color. The paths may loop back to themselves or join each other to form chain, star and loop configurations.

Algorithm: Network-Based Coalition Game (NBCG) to identify modulated gene subnetworks**Inputs:**

- Training and test datasets $(X^{train}, y^{train}, X^{test}, y^{test})$
- Set of 100 initial genes \mathcal{G}_s
- PPI network in terms of binary matrix A
- Evaluation method: characteristic function $v(\cdot)$.
- Algorithm parameters $(L, L', Acc^{max}, T_{add}, T_{removal}, S_{max}, N_{max}, N_T)$

Outputs: Set of subnetworks $(\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_L)$ and the obtained performance metric $Acc = v(\mathcal{S}_1 \cup \mathcal{S}_2 \cup \dots \cup \mathcal{S}_L)$

1) Perform initializations:

- a) Randomly choose L seed genes $G_{i1}, i = 1, 2, \dots, L$ from the set of initial genes \mathcal{G}_s
- b) For $i=1$ to L : Initialize Subnetwork i : $\mathcal{S}_i^{(t=1)} = \{G_{i1}\}$, Reset stop flag $F_i = 0$
- c) Reset number of iteration: $t = 0$

2) Loop

Increment number of iteration: $t = t + 1$

- a) Perform network expansion: For subnetworks $i=1$ to L
 - Create set of inclusion candidate nodes: set of direct neighbor nodes Ω_i
 - For all $G_{ik} \in \Omega_i$, calculate Shapley value $\gamma_{G_{ik}}(v)$
 - Find the best inclusion candidate gene: $G_i^{(add)} = \underset{G_{ik} \in \Omega_i}{argmax} \gamma_{G_{ik}}(v)$
 - If $\gamma_{G_i^{(add)}}(v) \geq T_{add}$ then $\mathcal{S}_i^{(t)} = \mathcal{S}_i^{(t-1)} \cup G_i^{(add)}$, else $F_i = 1$
 - If $|\mathcal{S}_i| > S_{max}$ then $F_i = 1$

End for

- b) Perform network contraction: For subnetworks $i=1$ to L
 - Create set of exclusion candidate nodes: set of end-point nodes Ψ_i
 - For all $G_{ik} \in \Omega_i$, calculate Shapley value $\gamma_{G_{ik}}(v)$
 - Find the best removal candidate gene: $G_i^{(remove)} = \underset{G_{ik} \in \Omega_i}{argmin} \gamma_{G_{ik}}(v)$
 - If $\gamma_{G_i^{(remove)}}(v) \leq T_{remove}$ then $\mathcal{S}_i^{(t)} = \mathcal{S}_i^{(t-1)} \setminus G_i^{(remove)}$

End for

- c) Subnetwork evaluation and termination criteria check:

- Calculate performance of the obtained subnetworks: $Acc^{(t)} = v(\mathcal{S}_1^{(t)} \cup \mathcal{S}_2^{(t)} \cup \dots \cup \mathcal{S}_L^{(t)})$
- If $\sum_{i=1}^N |\mathcal{S}_i^{(t)}| \geq N_{max}$ or $\forall F_i = 1$ or $\forall \mathcal{S}_i^{(t)} = \mathcal{S}_i^{(t-1)}$ or $Acc^{(t)} \geq Acc_{max}$ or $t > N_T$, exit loop

End Loop

Figure 2. Network-based coalition game (NBGC) algorithm to identify modulated gene subnetworks using PPI networks

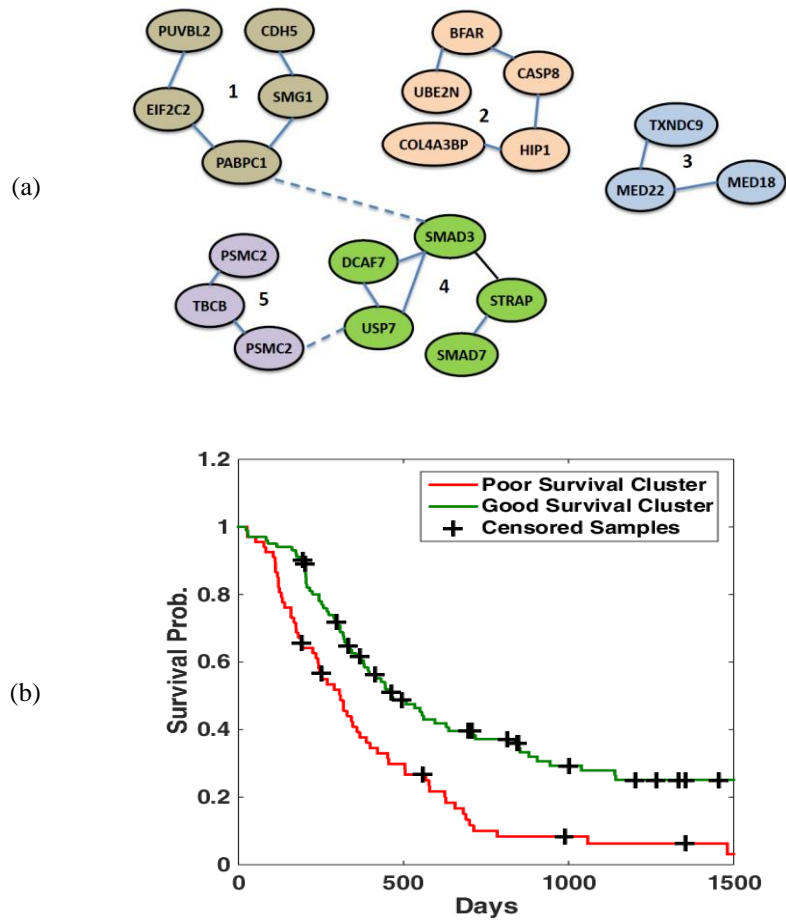


Figure 3. (a) Sample subnetworks reported by the proposed algorithm. Start genes are marked dark. The numbers show the sequence of subnetwork forming in consecutive algorithm iterations. (b) Platinum-free survival, p -value = 4×10^{-5} .

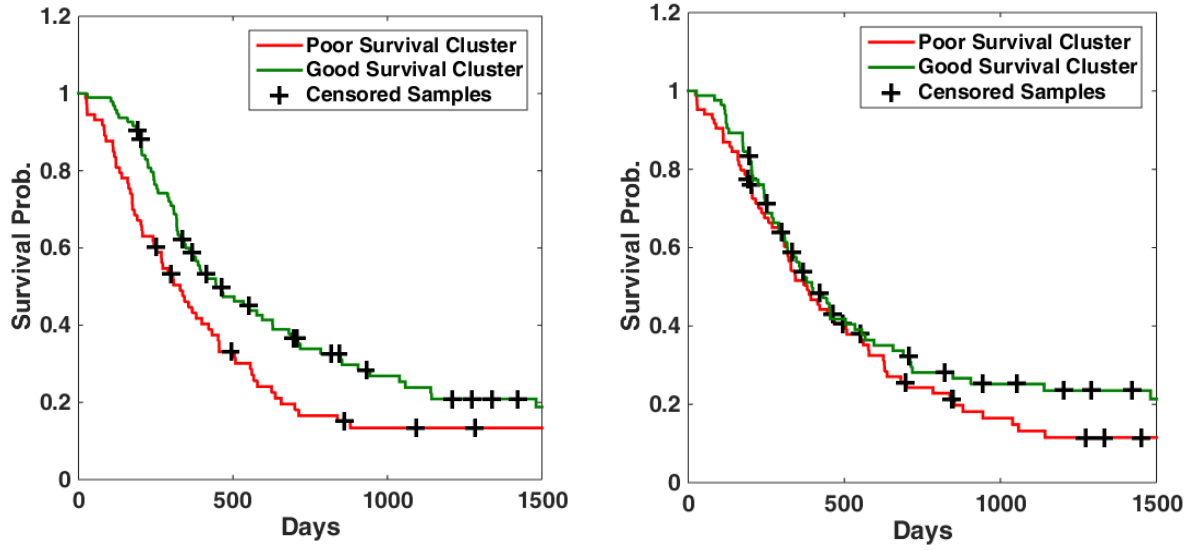


Figure 4: Survival probability obtained by Kaplan-Meier Estimate for the cancer samples clustered using the genes that are selected by (a) best classification method (CFS), p -value=0.018 and (b) optimal network based random walk method, p -value = 0.08.

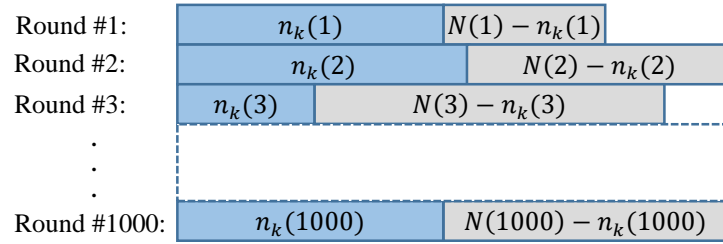


Figure 5. Each round of algorithm provides up to L subnetworks and a total set of gene shown by $\mathcal{G}(\mathbf{r})$ with cardinality $N(\mathbf{r}) = |\mathcal{G}(\mathbf{r})|$. The resulting gene set includes $n_k(\mathbf{r}) = |\mathcal{G}(\mathbf{r}) \cap \mathcal{G}_k|$ genes that belong to the top- k gene set and considered consistent.

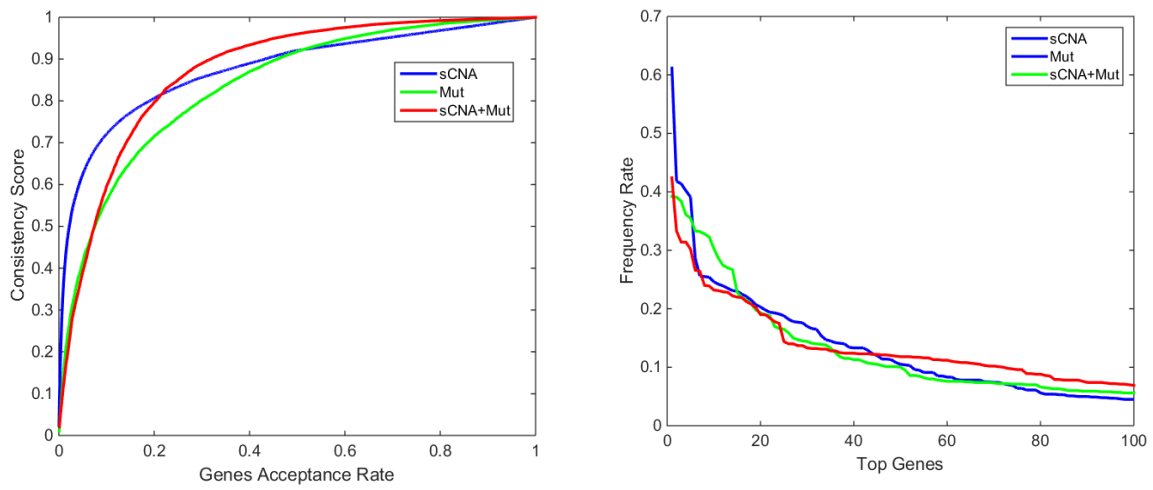


Figure 6. (a) Average consistency scores of the proposed algorithm using different data types. (b) frequency of appearance of top- k genes in mutiple runs of the algorithm.

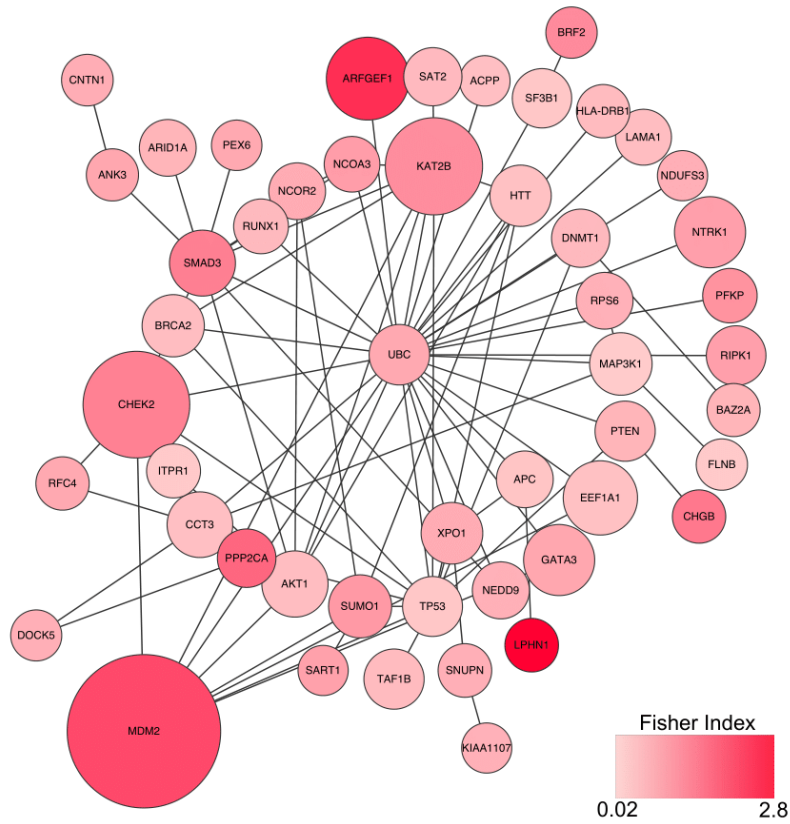


Figure 7. Developed gene subnetworks using the NBCG algorithm including genes, when mutated, are associated with increased immune cell infiltration within tumors. The network starts with 4 randomly selected hot-spot genes including *UBC*, *SMAD1*, *MEPCE*, *TP53* and develop 4 subnetworks that ultimately join each other and form the following connected subnetworks. The size of the nodes represents the frequency of their appearance in the reported list for multiple algorithm execution (e.g. *MCM2* is the most frequently reported gene by different algorithm executions). The intensity of node colors refer to their individual associations with the immune scores in terms of fisher index between the immune scores of two cohorts with and without mutations in this specific gene (e.g. the gene *LHHN1* demonstrates a highest association with immune scores).