

METHOD

Open Access



I-Boost: an integrative boosting approach for predicting survival time with multiple genomics platforms

Kin Yau Wong¹, Cheng Fan², Maki Tanioka^{2,3}, Joel S. Parker^{2,3}, Andrew B. Nobel^{2,4,5}, Donglin Zeng^{2,5}, Dan-Yu Lin^{2,5*} and Charles M. Perou^{2,3*} 

Abstract

We propose a statistical boosting method, termed I-Boost, to integrate multiple types of high-dimensional genomics data with clinical data for predicting survival time. I-Boost provides substantially higher prediction accuracy than existing methods. By applying I-Boost to The Cancer Genome Atlas, we show that the integration of multiple genomics platforms with clinical variables improves the prediction of survival time over the use of clinical variables alone; gene expression values are typically more prognostic of survival time than other genomics data types; and gene modules/signatures are at least as prognostic as the collection of individual gene expression data.

Keywords: Cancer genomics, Data integration, Gene modules, Variable selection

Background

Prediction of disease outcomes, such as individual patient survival time, is critically important for cancer patients. Traditional prognostic models that rely solely on clinical variables, such as age and tumor stage, fail to account for the molecular heterogeneity of tumors and thus may lead to suboptimal treatment decisions [1]. To remedy this situation, many studies have incorporated gene expression data in survival prediction [2–5].

Large-scale genomics projects such as The Cancer Genome Atlas (TCGA) have generated detailed molecular data on patients with a variety of cancer types. In TCGA, six types of “omics” data have been collected on the same set of patients: DNA copy number variation, somatic mutation, mRNA expression, microRNA expression, DNA methylation, and expression of ~200 proteins/phosphoproteins. The availability of multiple data types has enabled researchers to address a variety of important questions. For example, patients can be more precisely classified into molecular subtypes based on

integrative clustering of multiple genomics data types or platforms [6–8]. In addition, it is possible to identify genes that are related to patient survival time by decomposing the expression of each gene into a component that is explained by the methylation level and a component that is not [9].

One unsolved issue in cancer genomics is the prognostic value of integrated genomics and clinical data versus clinical data only. Yuan et al. [10] compared models with clinical data only versus models with both clinical and genomics data on various cancer types and concluded that genomics data provide only a limited gain in survival prediction accuracy. In their analysis, however, potential differences among data types were not taken into account. For breast cancer, for instance, the combination of genomics and clinical data has been shown to improve outcome predictions [11, 12]. A major goal of the present work is to fully explore the predictive power of integrating clinical and genomics data together.

A second unsolved issue is the prognostic value of individual gene expression values (~25,000) versus a pre-defined set of gene expression signatures or “modules” (~500). Gene modules have been developed for representing distinct cell types (e.g., epithelial, immune, and endothelial), specific biological processes, or activated molecular signaling pathways. They have been shown to

*Correspondence: lin@bios.unc.edu; chuck_perou@med.unc.edu

²Lineberger Comprehensive Cancer Center, University of North Carolina, Chapel Hill, NC 27599, USA

³Department of Genetics, University of North Carolina, Chapel Hill, NC 27599, USA

Full list of author information is available at the end of the article



successfully capture signaling pathway activities or cell type heterogeneity within tumors. We wish to investigate whether individual gene expression data or existing gene modules provide more accurate outcome prediction.

A third unsolved issue is the relative importance of different types of genomics data in outcome prediction. Different data types are collected at different costs and also with widely varying feature spaces. Naturally, not all data types are equally important in outcome prediction. We aim to determine which data types may be omitted from analysis without a significant reduction in prediction accuracy.

An overarching methodological challenge in addressing the aforementioned issues is the identification of genomic variables predictive of survival time when the number of variables is much larger than the sample size. Penalized regression methods, such as least absolute shrinkage and selection operator (LASSO) [13] and elastic net [14], are commonly used to identify important genomic variables. When variables are highly correlated, elastic net tends to have better performance in prediction than LASSO [14]. However, both LASSO and elastic net are generic variable selection procedures that do not distinguish different types of data and thus tend to select more variables from the data types with larger numbers of variables. Because different data types capture different biological structures, both large and small data types may carry important signals. Methods that treat all variables equally may not be able to pick out independent signals from small data types. In addition, LASSO and elastic net impose the same penalty on all regression parameters, which may be overly restrictive because the number of variables and the signal strength vary drastically across data types.

Boosting is an alternative to penalization for model estimation and prediction in high-dimensional settings. It was originally developed for binary classification in machine learning [15, 16]. The idea of boosting is to iteratively reweight the observations, with larger weights given to observations that are misclassified at the previous iteration, and apply simple classifiers on the reweighted data; their results are then combined to produce an aggregated classification procedure. Boosting was later generalized as a forward stagewise additive modeling method for statistical estimation [17, 18], which can be applied to many problems, including regression analysis for survival data [19]. Because of its flexibility in modeling choices and stability in high-dimensional settings, boosting has found applications in genomics studies; see the references in Mayr et al. [20, 21]. As in the case of LASSO and elastic net, however, existing boosting methods, such as component-wise boosting [22], do not distinguish variables of different data types.

To overcome the limitations of LASSO, elastic net, and existing boosting methods, we develop a novel method, termed Integrative Boosting (I-Boost), which combines elastic net with boosting. In I-Boost, the prediction rule is constructed iteratively, where at each iteration, the predictive power of each data type (conditional on the current prediction rule) is evaluated separately and the most predictive data type is selected to update the prediction rule using elastic net. Thus, independent signal from each data type can be incorporated into the prediction rule, and small but predictive data types will not be dominated by data types with large numbers of variables. In addition, the penalties on the regression parameters are learned data-adaptively and separately for different data types. Herein, we demonstrate the advantages of I-Boost using simulation studies and empirical data from the TCGA on patients with eight different cancer types. More importantly, we use I-Boost to address the aforementioned three unsolved issues in cancer genomics.

Results and discussion

Background

Suppose that there are K types of clinical or genomics predictors, with d_k components for the k th type ($k = 1, \dots, K$). For $k = 1, \dots, K$, let $X^{(k)}$ denote the d_k -vector of predictors of the k th type. Write $X = (X^{(1)'}, \dots, X^{(K)'})'$, where A' denotes the transpose of A for any vector or matrix A . Let T denote the survival time of interest. We relate T to X through the proportional hazards model [23], such that the conditional hazard function of T given X takes the form of $h_0(t) \exp(\beta'X)$, where $h_0(t)$ is an arbitrary baseline hazard function, $\beta = (\beta^{(1)'}, \dots, \beta^{(K)'})'$, and $\beta^{(k)}$ is a d_k -vector of regression parameters associated with $X^{(k)}$.

The survival time T is subject to right censoring by C , such that we observe $Y \equiv \min(T, C)$ and $\Delta \equiv I(T \leq C)$, where $I(\cdot)$ is the indicator function. For a study with n patients, the data consist of (Y_i, Δ_i, X_i) ($i = 1, \dots, n$). The partial likelihood [24] for β is

$$L(\beta) = \prod_{i=1}^n \left(\frac{e^{X_i' \beta}}{\sum_{j: Y_j \geq Y_i} e^{X_j' \beta}} \right)^{\Delta_i}.$$

LASSO and elastic net

Because X is high-dimensional, it is not feasible to estimate β by maximizing the partial likelihood. One possible remedy is to impose sparsity assumptions on β and adopt penalization methods, such as LASSO [13] and elastic net [14]. LASSO estimates β by maximizing the L_1 -penalized log-partial likelihood function

$$\log L(\boldsymbol{\beta}) - \lambda \sum_{j=1}^d |\beta_j|,$$

where $d = \sum_{k=1}^K d_k$, and λ is a tuning parameter. Elastic net generalizes LASSO by including an L_2 penalty, such that the objective function becomes

$$\log L(\boldsymbol{\beta}) - \lambda \left\{ \alpha \sum_{j=1}^d |\beta_j| + \frac{1}{2} (1 - \alpha) \sum_{j=1}^d \beta_j^2 \right\},$$

where $\alpha \in [0, 1]$ is a tuning parameter that controls the relative magnitudes of the L_1 and L_2 penalties. (When $\alpha = 1$, elastic net reduces to LASSO.) The implementation of LASSO and elastic net is described in the “Methods” section.

For both LASSO and elastic net, the penalty term dominates under large values of λ , and the parameter estimates tend to be small with some values being exactly zero. Unlike LASSO, elastic net exhibits the grouping effect in that the regression parameters for a group of highly correlated variables tend to be equal, which is desirable in the context of gene selection [14]. Both LASSO and elastic net impose the same penalization on each regression parameter and thus do not distinguish different types of predictors. As a result, these methods may be inefficient when certain data types are much more predictive than others.

I-Boost

To account for the differential predictive power of different data types, we propose a boosting algorithm called I-Boost. Boosting is an iterative optimization algorithm that minimizes a loss function $\ell\{\mathcal{Y}, \mathbf{f}(\mathcal{X})\}$ over a class of functions of predictors $\mathbf{f}(\mathcal{X})$, where $\mathcal{Y} = (Y_1, \dots, Y_n, \Delta_1, \dots, \Delta_n)$, $\mathcal{X} = (X_1, \dots, X_n)$, and $\ell\{\mathcal{Y}, \mathbf{f}(\mathcal{X})\}$ measures the deviation of the prediction $\mathbf{f}(\mathcal{X})$ from the outcome \mathcal{Y} . At each iteration, we update $\mathbf{f}(\mathcal{X})$ additively by the value $\mathbf{b}(\mathcal{X}; \boldsymbol{\beta})$ up to a scaling factor, where \mathbf{b} is a fixed basis function, and $\boldsymbol{\beta}$ is a vector of parameters. Specifically, at the m th iteration, we find $\boldsymbol{\beta}^{(m)}$ that minimizes $\ell\{\mathcal{Y}, \mathbf{f}_{m-1}(\mathcal{X}) + \mathbf{b}(\mathcal{X}; \boldsymbol{\beta}^{(m)})\}$, possibly under some constraints on $\boldsymbol{\beta}^{(m)}$, where \mathbf{f}_{m-1} is the estimate of \mathbf{f} at the $(m-1)$ th iteration. Then, we set $\mathbf{f}_m(\mathcal{X}) = \mathbf{f}_{m-1}(\mathcal{X}) + \nu \mathbf{b}(\mathcal{X}; \boldsymbol{\beta}^{(m)})$ for some fixed step length factor $\nu \in (0, 1]$. We terminate the iterations when some stopping criterion is satisfied.

In I-Boost, we set the loss function $\ell\{\mathcal{Y}, \mathbf{f}(\mathcal{X})\}$ to be the negative log-partial likelihood function and the basis function to be $\mathbf{b}(\mathcal{X}; \boldsymbol{\beta}^{(m)}) = (X_1^{(k)'} \boldsymbol{\beta}^{(m)}, \dots, X_n^{(k)'} \boldsymbol{\beta}^{(m)})'$, where $X_i^{(k)}$ is the vector of the k th type of predictors for the i th patient, and the data type k is selected data-adaptively. At each iteration, we search over all data types,

select the one that yields the largest decrease in the loss function value at the current iteration, and update (a subset of) the regression parameters corresponding to the selected data type; other parameters are fixed at their current estimated values. To handle high-dimensional data, we impose an elastic net penalty on $\boldsymbol{\beta}^{(m)}$ in the optimization step. Effectively, we perform maximum penalized log-partial likelihood estimation with an offset term $\mathbf{f}_{m-1}(\mathcal{X})$ using a single data type at each iteration. Unlike existing boosting methods, such as component-wise boosting, the basis function in our case is a function of all variables of a data type instead of a single variable. This choice of basis function is motivated by the expectations that some data types are much more predictive than others and that the inclusion of less predictive data types may reduce the prediction accuracy of the model. By considering each data type separately, we perform selection on the data-type level at each iteration.

We propose two versions of I-Boost, namely I-Boost-CV and I-Boost-Permutation, which use cross-validation and permutation, respectively, to choose the tuning parameters of elastic net at each iteration. The permutation procedure randomly permutes the outcome variables in order to remove association between the predictors and the outcome, and the tuning parameters are chosen such that no predictor is selected in half of the permuted data sets. The procedures are described in detail in the “Methods” section.

Simulation studies

We conducted simulation studies to evaluate the performance of LASSO, elastic net, and the two versions of I-Boost. We considered three simulation settings, with different distributions of signals across the data types. In all three settings, a relatively large proportion of the signals is contributed by the clinical variables. The distributions of signals are shown in Fig. 1, and the details of the simulation settings are provided in the “Methods” section.

We assessed the performance of the methods by the quality of prediction and parameter estimation. For prediction, we report the correlation between the estimated risk score $\sum_{k=1}^K X^{(k)'} \hat{\boldsymbol{\beta}}^{(k)}$ and the true risk score $\sum_{k=1}^K X^{(k)'} \boldsymbol{\beta}_0^{(k)}$, where $\hat{\boldsymbol{\beta}}^{(k)}$ and $\boldsymbol{\beta}_0^{(k)}$ are the estimated and true parameter vectors, respectively. A higher correlation represents a greater degree of agreement between the predicted and actual outcomes. We call this measure the risk correlation. For parameter estimation, we report the mean-squared error (MSE), defined as $\sum_{k=1}^K \|\hat{\boldsymbol{\beta}}^{(k)} - \boldsymbol{\beta}_0^{(k)}\|_2^2$.

Figure 1 shows the risk correlation and MSE for elastic net, LASSO, and the two versions of I-Boost based on 1000 replications; the average number of variables

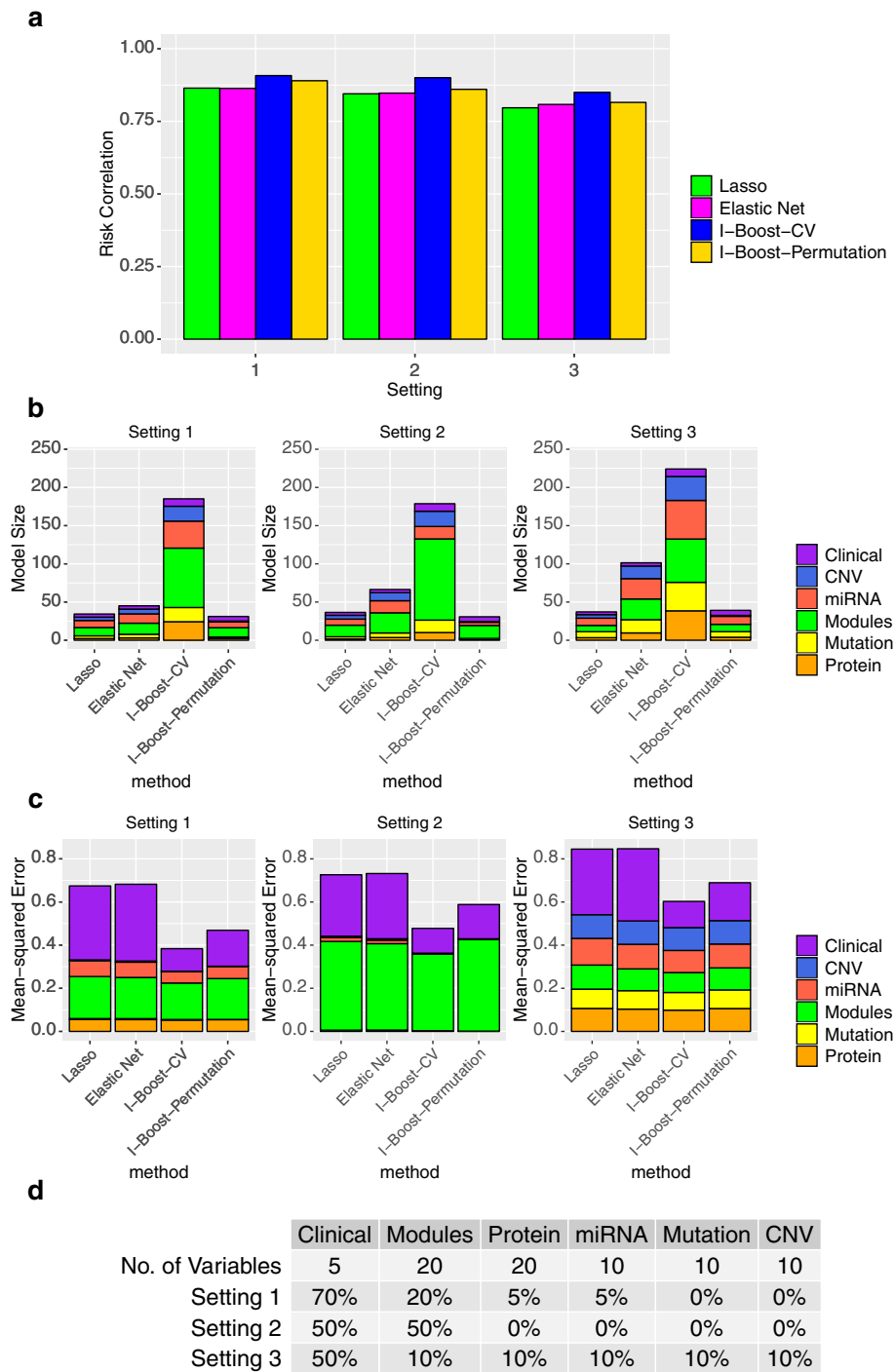


Fig. 1 Simulation settings and results. **a** Prediction accuracy of LASSO, elastic net, I-Boost-CV, and I-Boost-Permutation measured by risk correlation under three different settings. **b** The average number of variables selected by the four methods under three different settings. Different types of the selected variables are represented by different colors. **c** MSE of the four methods under three different settings. The error is decomposed into errors of parameters for different data types, as represented by different colors. **d** Number of signal variables and distribution of signals across different data types for the three simulation settings. The number of signal variables is zero if the proportion of signals of the data type is 0%. Abbreviations are as follows: GeneExp represents individual gene expression, Module represents gene module, Clinical represents clinical variable, CNV represents copy number variant, Mutation represents somatic mutation, miRNA represents microRNA expression, and Protein represents protein expression

selected for each data type is also shown. I-Boost-CV always selects the largest number of variables, followed by elastic net, LASSO, and I-Boost-Permutation. I-Boost-CV selects a large number variables, because it iteratively performs elastic net, and the final model includes selected variables accumulated over all iterations. By contrast, I-Boost-Permutation, though iterative, performs LASSO (which generally selects fewer variables than elastic net) with the tuning parameter selected by the very conservative permutation method [25], so that it selects the least number of variables.

For estimation, the MSE under I-Boost-CV or I-Boost-Permutation is about 20–40% smaller than that under LASSO or elastic net in all settings. Decomposition of the MSE by data types reveals that the MSE for data types with very weak or no signal is small for I-Boost. This result shows that even though I-Boost-CV selects a relatively large number of variables from these data types, the variables generally have very small estimated regression parameters.

For prediction, the two I-Boost methods perform the best overall. In all settings, I-Boost-CV produces more accurate prediction than all other methods. In Settings 1 and 2, where most signals are concentrated on only one or two data types, I-Boost-Permutation produces more accurate prediction than both elastic net and LASSO. In Setting 3, I-Boost-Permutation performs similarly to elastic net, while LASSO performs worse than I-Boost-Permutation. Between the two versions of I-Boost, I-Boost-CV tends to yield better prediction than I-Boost-Permutation, possibly because of the larger number of variables selected by I-Boost-CV. Thus, if the main interest is the selection of relevant variables, then one might consider I-Boost-Permutation for more conservative variable selection, even though this method is somewhat inferior in prediction when compared to I-Boost-CV.

We implemented LASSO, elastic net, and the two versions of I-Boost using R-3.2.2 on a 2.93-GHz Xeon Linux computer. On average, performing LASSO, elastic net, I-Boost-Permutation, and I-Boost-CV on one simulated data set (that consists of 500 subjects, 6 data types, and 1294 predictors) takes about 2 min, 14 min, 3 h, and 38 h, respectively. I-Boost-CV is computationally intensive because in each iteration, cross-validation is conducted on a three-dimensional grid. By contrast, in each I-Boost-Permutation iteration, the tuning parameter α is fixed at 1, no cross-validation is involved in the selection of λ , and LASSO is performed only once for each data type. Therefore, I-Boost-Permutation may serve as a computationally efficient alternative to I-Boost-CV.

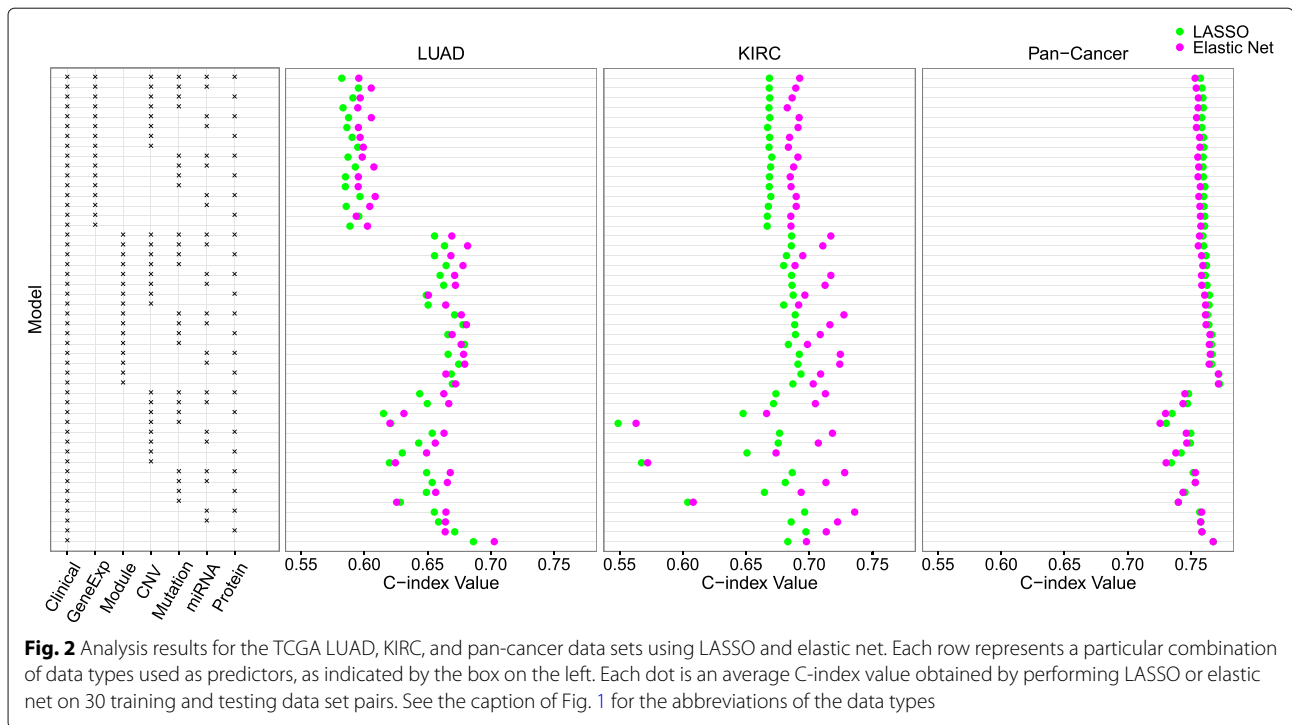
Evaluation of LASSO, elastic net, and I-Boost using TCGA data

We next evaluated the performance of the methods using three TCGA data sets, namely the lung adenocarcinoma (LUAD) data set, the kidney renal clear cell cancer (KIRC) data set, and a pan-cancer data set derived from ~1400 patients that represents eight different tumor types considered by Hoadley et al. [26]; see the “Methods” section for a detailed description of the data sets and the evaluation procedure. For each data set, we first split the data 30 times into training and testing sets. We then performed LASSO, elastic net, and the two versions of I-Boost for various combinations of data types on patients from the training set of each split. For each combination of data types and each split, we calculated the risk scores for patients in the testing set using the estimates from the corresponding training set, and we used the concordance index (C-index) [27] to evaluate the prediction accuracy of the risk scores.

The average C-index values over the splits obtained from LASSO and elastic net are given in Fig. 2. For the KIRC and pan-cancer data sets, the prediction tends to be much better than random (i.e., the C-index values are much larger than 0.5). For the LUAD data set, which has a small sample size, some of the models yield relatively poor prediction (with C-index values smaller than 0.6). For many models, the predictive performance of elastic net is either similar or superior to LASSO.

For LASSO and elastic net, the models containing more data types as predictors do not necessarily perform better than those with fewer data types. One possible explanation is that the extra data types may contain very little relevant information on patient survival, such that adding those data types introduces more noise than signal into the model. In practice, however, it is challenging to decide which data types to consider without prior knowledge of their importance.

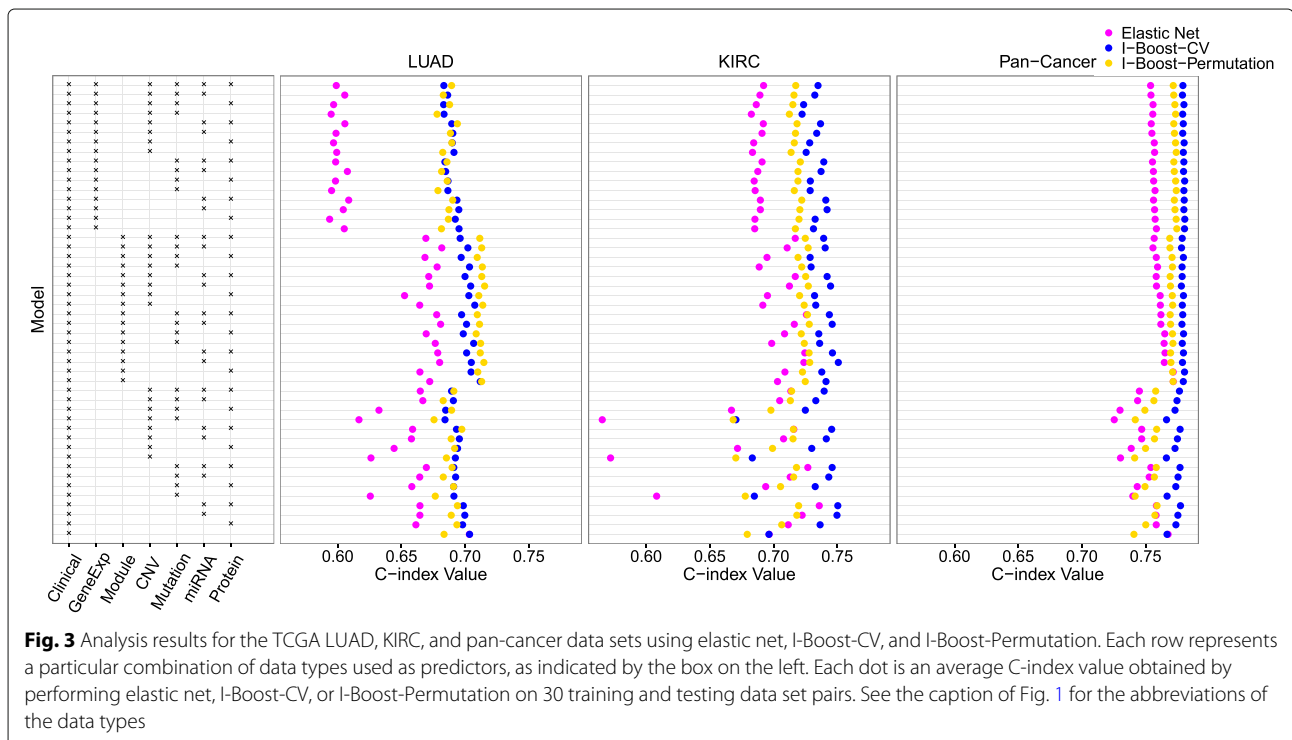
Figure 3 shows the average values of the C-index obtained from elastic net, I-Boost-CV, and I-Boost-Permutation for different models. For the LUAD, KIRC, and pan-cancer data sets, both versions of I-Boost provide better prediction than elastic net in almost all cases. The difference in prediction accuracy between I-Boost and elastic net is particularly large when the sample size is small and the number of predictors is large. The difference is likely due to the fact that I-Boost involves the selection of data types, so that the large and non-predictive data types would not be selected in most iterations, and their presence would not substantially worsen the prediction accuracy. For the KIRC and pan-cancer data sets, I-Boost-CV yields better prediction than I-Boost-Permutation, whereas for LUAD, there are no clear differences between the two methods.



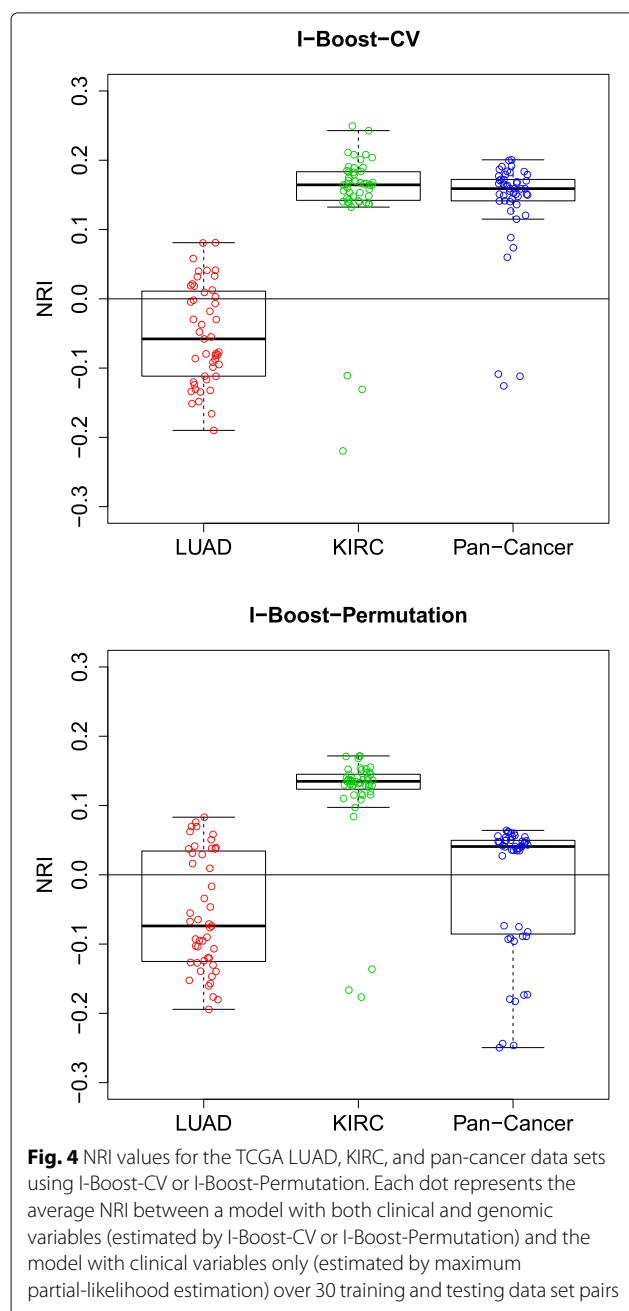
Prognostic value of integrated clinical and genomics data

To assess whether the genomic variables provide extra predictive power in the presence of the clinical variables, we computed the net reclassification improvement (NRI) [28, 29] values between the models with both clinical and genomic variables (estimated by I-Boost-CV

or I-Boost-Permutation) and the model with clinical variables only (estimated by maximum partial likelihood estimation). The NRI compares a model of interest with a baseline model and measures how much a subject's predicted risk under the model of interest, relative to that under the baseline model, aligns with the subject's survival



time. For instance, an NRI of 0.2 means that by switching from the baseline model to the model of interest, the proportion of high-risk subjects being reassigned a larger predicted risk is on average larger, by a value of 0.2, than the proportion of low-risk subjects being so reassigned; here, high-risk or low-risk subjects refer, respectively, to those with survival times shorter or longer than a fixed threshold, which we set to be 3 years throughout the paper. (See the “Methods” section for a theoretical definition of the NRI.) The average NRI values over data splits are shown in Fig. 4.



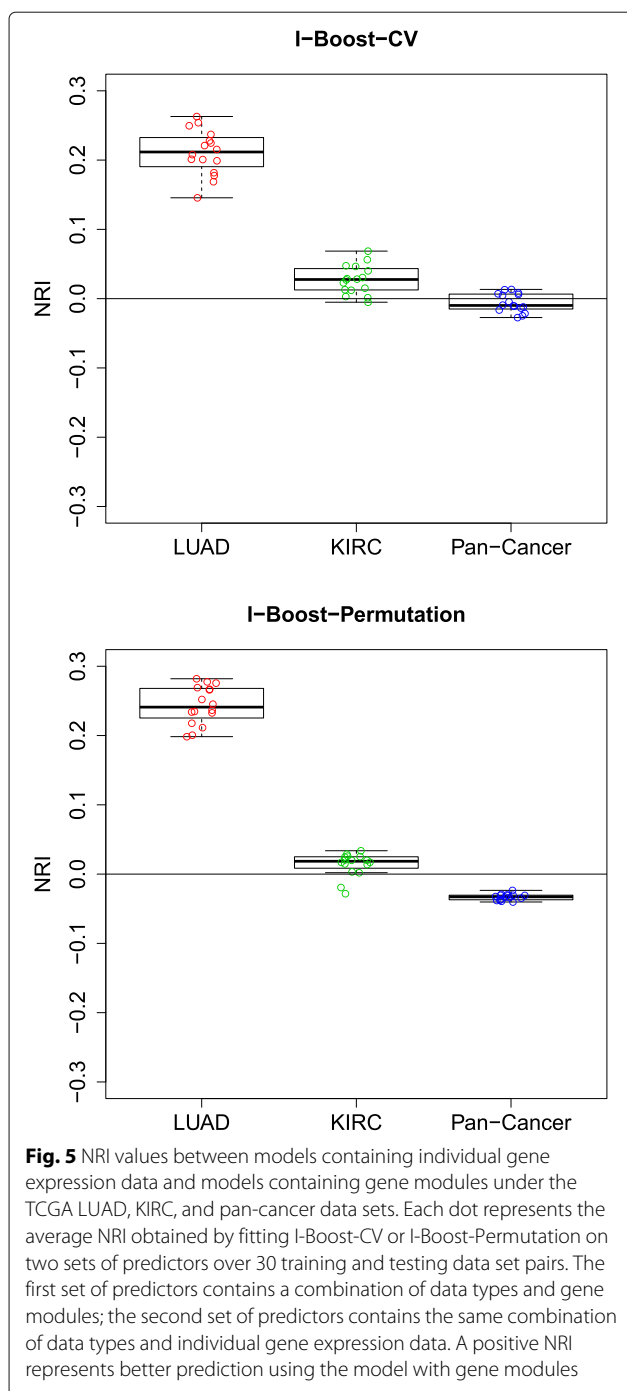
The patterns of the results from I-Boost-CV and I-Boost-Permutation are similar. For the KIRC and pan-cancer data sets, the majority of the models that contain both clinical and genomic variables yield positive NRI, which implies that they provide better prediction than the model with clinical variables only. Most NRI values under I-Boost-CV are close to 0.2; in biomarker studies, an NRI of 0.2 is considered an intermediate-level improvement [30]. For the LUAD data set, only a few models that contain both clinical and genomic variables provide better prediction than the model with clinical variables only. These results indicate that in certain cancer types, genomic variables contribute to survival prediction in the presence of clinical variables, and the magnitude of the contribution can be large. When the same comparisons are made using LASSO or elastic net, however, the inclusion of genomic variables in the models does not appreciably improve prediction.

Evaluation of gene expression modules

To compare the performance of gene modules versus individual gene expression data, we calculated the NRI values between models with each type of gene expression data separately. Specifically, for each combination of data types other than individual gene expression data and gene modules, we computed the NRI between the model with those data types and gene modules (estimated by I-Boost-CV or I-Boost-Permutation) and that with those data types and individual gene expression data. The NRI values are shown in Fig. 5. Under both methods, the use of gene modules leads to substantially better prediction than the use of expression data of all individual genes for the LUAD data set. For the KIRC and pan-cancer data sets, the performance of the two types of gene expression data is similar, and there is no strong evidence favoring gene modules or individual gene expression data on the basis of prediction accuracy. Nevertheless, because gene modules are smaller in number and much easier to interpret, we generally recommend the use of gene modules over individual gene expression data.

Comparison among genomics data types

To evaluate the relative prognostic value of each genomics data type, we formed a series of nested models as follows. We began by setting the model with clinical variables only as the first member of the series of models. At each later step, we computed the NRI between each model that contains all currently included data types and an extra genomics data type and the model included at the previous step. The model that yielded the largest NRI was set to be the next member of the series of models. The process was repeated until all data types were included. (Individual gene expression data were not considered in this analysis.) At each step of the process, the data type



that yielded the largest improvement in predictive power (over the data types already included) was selected, so that more predictive data types tend to be included earlier, and the order in which the data types entered the models reflects their relative importance. We performed this procedure for elastic net and the two versions of I-Boost. For the LUAD, KIRC, and pan-cancer data sets, the NRI values and their 95% confidence intervals for the series of models are plotted in Fig. 6, and the data type selected at

each step is shown. We also plotted the C-index against the number of variables selected for each model.

Because different methods vary in their abilities to extract useful information from given data types, the orders of data types determined by the methods are generally different. For the LUAD, KIRC, and pan-cancer data sets, the NRI under I-Boost-CV or I-Boost-Permutation tends to be positive or around zero with the inclusion of each new data type. This indicates that I-Boost extracts useful information from each additional data type and that its performance tends not to be worsened by the inclusion of additional variables.

I-Boost-Permutation always selects the smallest number of variables, followed by elastic net and I-Boost-CV. This finding is consistent with the conclusions from the simulation studies. Because the C-index obtained by I-Boost-Permutation is higher in most cases than that obtained by elastic net, we conclude that I-Boost-Permutation provides the same or better prediction using fewer variables than elastic net.

For the LUAD and pan-cancer data sets, gene modules are the first genomics data type selected under both versions of I-Boost, and the inclusion of gene modules leads to considerable improvement in prediction accuracy. For the KIRC data set, miRNA expression data are first selected by I-Boost-CV, while gene modules are first selected by I-Boost-Permutation. For I-Boost-CV, however, the model with clinical variables and gene modules yields an NRI of 0.19, which represents a substantial improvement over the model with clinical variables only. The confidence intervals of the NRI include zero due to the small sample sizes of the testing data sets. Nevertheless, the pattern of consistent positive NRI values shown in Fig. 4 and the fact that the NRI values are averages over 30 data splits suggest that the improvement in prediction accuracy is robust. For both versions of I-Boost, after the inclusion of the first genomics data type, the improvement in prediction accuracy with the inclusion of additional data types is marginal. We conclude that gene modules are overall the most predictive genomics data type, and the remaining genomics data types tend not to provide extra predictive power beyond clinical variables and gene modules.

We also evaluated the prognostic value of genomics data in the absence of clinical data. The average C-index values for combinations of genomics data types over 30 training and testing data splits for the LUAD, KIRC, and pan-cancer data sets are given in Additional file 1: Fig. S1. The maximum C-index values obtained using genomics data types alone are 0.64, 0.72, and 0.74 in the LUAD, KIRC, and pan-cancer data sets, respectively; they are substantially smaller than the corresponding maximum values obtained using both clinical and genomics data. For the LUAD data set, miRNA expression data alone yield the

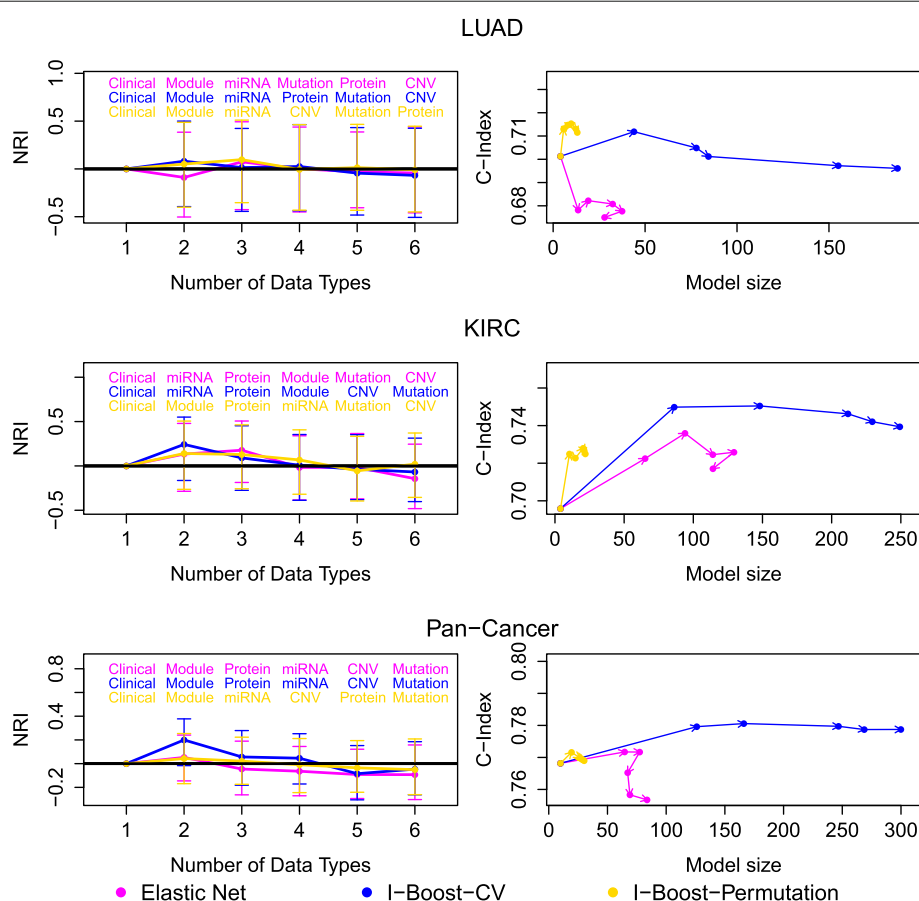


Fig. 6 Analysis results for the TCGA LUAD, KIRC, and pan-cancer data sets, using elastic net, I-Boost-CV, and I-Boost-Permutation on nested models. In the left panel, the leftmost dots are fixed at zero, and each remaining dot represents the average NRI obtained by fitting elastic net, I-Boost-CV, or I-Boost-Permutation over 30 training and testing data set pairs. Each dot except the leftmost dots represents the maximum NRI between a model that contains one more data type than the model corresponding to the dot on the left and the model corresponding to the dot on the left. Above each dot, the name of the additional data type is included. In the right panel, the average C-index values and the average numbers of selected variables for the models shown in the left panel are plotted. The arrows indicate the orders of models with respect to the number of data types they contain. See the caption of Fig. 1 for the abbreviations of the data types

largest C-index, whereas for the KIRC data set, the combination of miRNA expression and protein expression data yield the largest C-index. For the pan-cancer data set, the C-index values for combinations of genomics data types with individual gene expression data are almost identical and are larger than those obtained without individual gene expression data.

Important predictors for the LUAD, KIRC, and pan-cancer data sets

To obtain the final models of important predictors, we performed I-Boost-Permutation on the LUAD, KIRC, and pan-cancer data sets. The final models are shown in Tables 1, 2, and 3 for the LUAD, KIRC, and pan-cancer data sets, respectively. The predictors that are also selected by LASSO, elastic net, and I-Boost-CV are marked.

Age and pathological nodal status are negatively associated with survival time in the LUAD, KIRC, and pan-cancer data sets. Age has been reported to be prognostic for many cancer types [31–33]. In the analysis of the pan-cancer data set, cancer types were selected, which is logical, since the survival time is known to depend on cancer types [26]. Thus, the tissue of origin remains an important prognostic factor. Among the gene modules, Glycolysis_signature and MUnknown_24 are negatively associated with survival time in the LUAD and pan-cancer data sets; these two modules are correlated with Hypoxia signatures among a set of 1198 TCGA breast cancer patients. Likewise, Pcorr_IGS_Correlation and Activate.Endothelium, which are negatively associated with survival time in the pan-cancer data set, are correlated with proliferation signatures; the latter are known to be negatively associated with survival time.

Table 1 Analysis results from I-Boost-Permutation for the TCGA LUAD data set

Predictor	Estimate
Module_UNC_MPYMT_NEU_Cluster_Median_BMC.Med.Genomics.2011_PMIID.21214954*	− 0.2587
Mutation_HMCN1*	− 0.0498
Mutation_FAT3*	− 0.0363
Clinical_gender_female_0	− 0.0082
miRNA_hsa-miR-181c-5p*	− 0.0040
Mutation_AHNAK2*	− 0.0031
Mutation_LAMA2*	− 0.0001
CNV_BeroukhimS2.19p12-66*	0.0217
Module_UNC_Glycolysis_Signature_Median_BMC.Med.2009_PMIID.19291283*	0.0368
Module_IMMUNE_Bindea_Cell_Th2_cells_Median_Immunity.2013_PMIID.24138885*	0.0547
miRNA_hsa-miR-582-3p*	0.1438
Clinical_age*	0.1586
Clinical_pathologic_N*	0.4700

Note: "Estimate" is the estimate of the log hazard ratio under the Cox proportional hazards model, where a positive value represents an increase of the hazard. The predictors are standardized to have unit standard deviation. Gender is coded as female = 0 and male = 1; pathologic stage T is dichotomized into T1 (0) and T2–T4 (1); pathologic stage N is dichotomized into N0 (0) and N1–N3 (1). Predictors that are also selected by LASSO, elastic net, and I-Boost-CV are marked with an asterisk (*)

In contrast, signatures of CD8 T cells, non-inflammatory breast cancer (nIBC and MM_Red2), and luminal features (Mature_LuminalUp, GP7_estrogen signaling, HS_Green1, HS_Green8, LUMINAL_Cluster, Duke_Module06_er, Pcorr_Dasatinib_L_Correlation, and HS_Green18) are positively associated with survival time in the KIRC or pan-cancer data sets. The NEU_cluster module is positively associated with survival time in the LUAD data set, which is biologically significant because this module represents epithelial luminal cell differentiation and thus tracks more differentiated and lower grade lung cancers. The selected features, many of which are also selected by other variable selection methods, have significant biological implications and demonstrate the robustness of the I-Boost methodology.

Conclusions

In this paper, we present a novel method, termed I-Boost, for variable selection and outcome prediction that is especially powerful when one wishes to simultaneously consider multiple genomics and/or proteomics data types. We used simulation studies and real data to demonstrate that in the presence of multiple data types with diverse signal strength, I-Boost produces better outcome prediction than LASSO and elastic net. We proposed two versions of I-Boost, namely I-Boost-CV and I-Boost-Permutation. I-Boost-CV yields

Table 2 Analysis results from I-Boost-Permutation for the TCGA KIRC data set

Predictor	Estimate
Protein_AR*	− 0.1056
Module_IMMUNE_Bindea_Cell_CD8 T cells_Median_Immunity.2013_PMIID.24138885*	− 0.0696
Module_Mature_LuminaUp_Median_Nat.Med.2009_PMIID.19648928*	− 0.0676
Module_UNC_MM_Red2_Median_BMC.Med.Genomics.2011_PMIID.21214954*	− 0.0671
Module_GP7_Estrogen signaling: r=0.97	− 0.0580
miRNA_hsa-miR-10b-3p*	− 0.0369
Module_UNC_HS_Green1_Median_BMC.Med.Genomics.2011_PMIID.21214954	− 0.0369
miRNA_hsa-miR-192-5p	− 0.0291
Protein_Src_pY416*	− 0.0279
miRNA_hsa-miR-425-3p*	− 0.0166
Module_UNC_LUMINAL_Cluster_Median_BMC.Med.Genomics.2011_PMIID.21214954	− 0.0159
Module_UNC_HS_Green8_Median_BMC.Med.Genomics.2011_PMIID.21214954*	− 0.0129
Protein_PRAS40_pT246*	− 0.0107
Module_UNC_Duke_Module06_er_Median_Mike_PMIID.20335537*	− 0.0027
Module_Pcorr_squamoid_PLOS.2012_PMIID.22590557	0.0052
Clinical_pathologic_N	0.0060
miRNA_hsa-miR-21-5p	0.0068
Module_UNC_MM_p53null.Basal_Median_Genome.Biol.2013_PMIID.24220145*	0.0069
Protein_Caveolin-1	0.0124
miRNA_hsa-miR-21-3p	0.0129
Protein_TIGAR	0.0259
miRNA_hsa-miR-92b-3p*	0.0274
miRNA_hsa-miR-223-3p*	0.0313
miRNA_hsa-miR-130a-3p*	0.0572
miRNA_hsa-miR-222-3p*	0.0583
Protein_IGFBP2*	0.0631
miRNA_hsa-let-7a-3p*	0.0686
Clinical_age	0.1101
Module_UNC_Scorr_Basal_Correlation_JCO.2009_PMIID.19204204*	0.1370
Clinical_pathologic_T*	0.2470

Note: See Note of Table 1

more accurate prediction than I-Boost-Permutation, but it generally selects many more variables and is computationally more intensive. By contrast, I-Boost-Permutation is computationally efficient and selects much fewer variables, which may be preferable for follow-up experiments.

Table 3 Analysis results from I-Boost-Permutation for the TCGA pan-cancer data set

Predictor	Estimate
Module_Pcorr_Dasatinib_L_Correlation_Cancer.Res.2007_PMID.17332353*	− 0.1517
Module_UNC_MS_CD44_DOWN_Median_PNAS.2009_PMID.19666588*	− 0.0447
Module_UNC_HS_Green18_Median_BMC.Med.Genomics.2011_PMID.21214954*	− 0.0396
Module_UNC_MPYMT_NEU_Cluster_Median_BMC.Med.Genomics.2011_PMID.21214954*	− 0.0351
Module_UNC_MN0tch4_Median_BMC.Med.Genomics.2011_PMID.21214954*	− 0.0290
miRNA_hsa-miR-101-3p*	− 0.0282
Module_IMMUNE_Bindea_Cell_CD8 T cells_Median_Immunity.2013_PMID.24138885*	− 0.0224
Module_Shipitsin_CD44_B_Median_Cancer.Cell.2007_PMID.17349583*	− 0.0184
Protein_p38_pT180_Y182*	− 0.0182
CNV_wa.9.p*	− 0.0148
Module_Inflammatory_Breast_Cancer_491_nIBC_CCR.2013_PMID.23396049*	− 0.0033
miRNA_hsa-miR-34a-5p*	0.0001
Protein_Dvl3	0.0002
Protein_PA1-1	0.0008
Module_UNC_Glycolysis_Signature_Median_BMC.Med.2009_PMID.19291283	0.0012
CNV_Basal.13q34-86*	0.0084
Module_UNC_ADM_S100A10_A110NDGR1_Cluster_Median_BMC.Med.Genomics.2011_PMID.21214954	0.0147
Module_Pcorr_IGS_Correlation_NJEM.2007_PMID.17229949*	0.0162
Module_Extensive_Residual_Disease_ER54_Median_JAMA.2011_PMID.21558518*	0.0170
Clinical_gender_female_0	0.0300
Module_UNC_Activate_Endothelium_Median_Clin.Exp.Metastasis.2014_PMID.23975155*	0.0317
Module_UNC_MUnknown_24_Median_BMC.Med.Genomics.2011_PMID.21214954*	0.0444
miRNA_hsa-miR-582-3p*	0.0575
Clinical_LUAD	0.0771
Clinical_HNSC	0.0787
Clinical_BLCA	0.0812
Clinical_KIRC	0.0915
Module_UNC_Duke_Module20_stat3_Median_Mike_PMID.20335537*	0.1088
Clinical_pathologic_N*	0.1725
Clinical_pathologic_T*	0.1943
Clinical_age*	0.3288

NOTE: For cancer type, BRCA is the reference group. For the interpretations of other variables and parameters, see Note of Table 1

Consistent with the current literature, we found that clinical variables are strong predictors of survival time. With I-Boost, we were able to build upon the clinical variables and extract additional useful information from genomic variables in order to improve the prediction; the improvement that we obtained with I-Boost was considerably larger than that obtained by either LASSO or elastic net. We also compared the use of individual gene expression data versus gene modules and found that the use of gene modules leads to improvement in prediction accuracy and more interpretable results. When we considered the selected I-Boost models, clinical variables (e.g., age, tumor size, and pathological nodal status) were strong predictors of survival. The I-Boost methods also selected several gene modules that were previously identified as prognostic of outcomes, whether positive or negative.

Our study has limitations. The main limitation is that the LUAD and KIRC data sets pertain to a relatively small number of patients, with an even smaller number of observed events. This limitation motivated us to combine eight solid epithelial tumor types to form a large pan-cancer data set. The analyses on the pan-cancer data might not properly account for heterogeneity across different cancer types. Another limitation of our study is that the quality of the clinical data varies across different cancer types; for example, the follow-up time for some cancer types was quite short.

In summary, we demonstrated that the performance of I-Boost is superior to that of elastic net and LASSO and that the performance of gene modules is superior to that of the totality of individual genes. The I-Boost methodology is applicable to any disease states where multiple types of genomics and/or proteomics data are available and thus has potential applications beyond cancer studies.

Methods

Data description

TCGA provides a large open-access database that includes clinical and genomics data for patients with 33 cancer types or subtypes. Herein, we focused on eight cancer types or subtypes, namely, LUAD, KIRC, colon adenocarcinoma (COAD), rectal adenocarcinoma (READ), lung squamous cell carcinoma (LUSC), bladder urothelial carcinoma (BLCA), breast invasive carcinoma (BRCA), and head and neck squamous cell carcinoma (HNSC). For clinical variables, somatic mutation, copy number variation, mRNA expression, and miRNA expression, data on 2272 patients were obtained from the December 22, 2012, Pan-Cancer-12 data freeze from the Sage Bionetworks Repository Synapse [34]; the data were previously processed and described by Hoadley

et al. [26]. Protein expression data were downloaded from Broad GDAC Firehose [35] for a subset of 1779 patients included in the data set of Hoadley et al. [26].

Clinical variables included gender, age, pathological stages T and N, and cancer type. In all analyses, COAD and READ were considered as one cancer type. For mRNA expression data, we used RNA-seq by Expectation-Maximization (RSEM) [36] to quantify the transcript abundances measured by RNA sequencing and used the log₂-transformed up-quantile-normalized RSEM values of 12,434 genes. The RNA sequencing was performed at the University of North Carolina at Chapel Hill [37–39]. Gene level expression data are also available on the Broad GDAC Firehose [35]. For mutation data, we used the single nucleotide variant calls, which were de-duplicated and re-annotated using the Ensembl version 69 transcript database. A total of 130 genes with non-synonymous mutations in more than 10% of the whole sample were included for the analyses. The combined mutation annotation format file is available from the Synapse resource. For miRNA expression data, we used the read count data for 305 normalized expressions, which were compiled into an abundance matrix for 5p and 3p mature miRBase strands [37]. For reverse-phase protein arrays, we used the level 3 normalized data for 136 proteins or phospho-proteins. For copy number data, SNP6.0 array-based gene-level somatic copy number alteration data were generated from the GISTIC analysis [40]. The input data matrix is available in Synapse at syn1710678. We used the copy number values for 216 cancer-specific segments, which are frequently altered in cancer of various types including breast cancer, and segments for all chromosome arms (a total of 41 segments) [41, 42].

We defined gene modules as sets of co-expressed genes that are considered to be functional units in breast cancer. We built a collection of 497 gene modules. The modules were constructed on the basis of 73 publications or results from the Gene Set Enrichment Analysis [43]. A partial list of the modules appears in Fan et al. [12]. Among the modules, 461 are median expression values for homogeneously expressed genes, 33 are correlations of expression values with predetermined gene centroids, and 3 are built from previously published gene expression prognostic models.

After removing patients with missing data, the total sample size was 1420, including 202 LUAD patients and 195 KIRC patients. All survival times were censored at 5 years if the patients were still in the study at that time point. For the pan-cancer data set, the median follow-up time was 16.8 months, and the censoring rate was 77.6%. For the subset of LUAD patients, the median follow-up time was 13.9 months, and the censoring rate

was 71.3%. For the subset of KIRC patients, the median follow-up time was 28.9 months, and the censoring rate was 63.6%.

LASSO and elastic net

We implemented LASSO and elastic net using the R-package “glmnet” [44] and used fivefold cross-validation to select the tuning parameters. For elastic net, cross-validation was performed over a two-dimensional grid of (α, λ) , while for LASSO, α was set to be 1. For elastic net, the grid for α was chosen to be (0.05, 0.1, 0.2, ..., 1.0), and a grid for λ was chosen separately for each α using the default settings of glmnet. (A minimum value of 0.05 was considered for α , because α too close to 0 may result in too many variables being selected; in particular, no sparsity is imposed if $\alpha = 0$.) To make the selection procedure more stable, we repeated the split and evaluation procedure five times, and the cross-validation errors were averaged over the five repetitions.

I-Boost

The I-Boost algorithm is given as follows:

- 1 Set $f_{0,i} = 0$ for $i = 1, \dots, n$, and let $\mathbf{f}_0 = (f_{0,1}, \dots, f_{0,n})'$.
- 2 Consider $m = 1, 2, \dots$:

- (a) For a given $k_m \in \{1, \dots, K\}$, calculate

$$\boldsymbol{\beta}^{(m)} = \operatorname{argmax}_{\boldsymbol{\beta}} \left\{ \log L^{(k_m)}(\mathbf{f}_{m-1}; \boldsymbol{\beta}) - p^{(k_m)}(\boldsymbol{\beta}; \alpha_m, \lambda_m) \right\}$$

using the coordinate-descent algorithm [44], where

$$L^{(k)}(\mathbf{f}; \boldsymbol{\beta}) \equiv \prod_{i=1}^n \left(\frac{e^{f_i + \mathbf{X}_i^{(k)'} \boldsymbol{\beta}}}{\sum_{j: Y_j \geq Y_i} e^{f_j + \mathbf{X}_j^{(k)'} \boldsymbol{\beta}}} \right)^{\Delta_i}$$

is the partial likelihood with offset term \mathbf{f} and covariates $\mathbf{X}^{(k)}$, α_m and λ_m are tuning parameters, $\mathbf{f} = (f_1, \dots, f_n)'$, and

$$p^{(k)}(\boldsymbol{\beta}; \alpha, \lambda) \equiv \lambda \left\{ \alpha \sum_{j=1}^{d_k} |\beta_j| + \frac{1}{2} (1 - \alpha) \sum_{j=1}^{d_k} \beta_j^2 \right\}$$

is the elastic net penalty. The selection of k_m , α_m , and λ_m is described below.

- (b) Set $f_{m,i} = f_{m-1,i} + v \mathbf{X}_i^{(k_m)'} \boldsymbol{\beta}^{(m)}$ for $i = 1, \dots, n$ with $v = 0.1$ and $\mathbf{f}_m = (f_{m,1}, \dots, f_{m,n})'$.

At the m th iteration, only the regression parameters corresponding to the k_m th data type are updated. We refer

to the d -vector with value $\beta^{(m)}$ at the positions corresponding to the k_m th data type and zero elsewhere as the current estimate at the m th iteration. The current estimate at each iteration contributes to the final parameter estimate additively, and the final parameter estimate is simply the sum of the current estimates obtained from all steps multiplied by ν .

I-Boost-CV and I-Boost-Permutation use cross-validation and permutation, respectively, to choose $(k_m, \alpha_m, \lambda_m)$ at step 2(a). For I-Boost-CV, we adopt fivefold cross-validation separately at each iteration over a three-dimensional grid on $\{1, \dots, K\} \times [0.05, 1] \times (0, \lambda_m^{(\max)})$ for $(k_m, \alpha_m, \lambda_m)$, where $\lambda_m^{(\max)}$ is a value large enough to shrink the current estimate to zero.

For I-Boost-Permutation, we first perform LASSO separately for each data type $X^{(k)}$ ($k = 1, \dots, K$) with tuning parameter $\lambda_m^{(k)}$, where $\lambda_m^{(k)}$ is selected using the permutation method proposed by Sabourin et al. [25]; the permutation method is only applicable to LASSO. The procedure is motivated by the principle that in a null model, i.e., in the absence of any relevant predictors, the tuning parameters should be chosen such that no variable is selected. The permutation selection procedure first generates hypothetical null models by randomly permuting $(Y_i, \Delta_i, f_{m-1,i})$ B times at each iteration, so that in each permuted data set, the association between the predictors and the outcome (and the offset term) is removed. The procedure then finds the smallest λ such that no variable is selected for each permuted data set and selects the median of the B values of λ . For the k th data type ($k = 1, \dots, K$), let $\lambda_m^{(k)}$ be the selected tuning parameter and $\beta_m^{(k)}$ be the corresponding LASSO estimate. We select k_m based on the partial-likelihood value at $\beta_m^{(k)}$, i.e., $k_m = \operatorname{argmax}_k L^{(k)}(f_{m-1}; \beta_m^{(k)})$, and set $\alpha_m = 1$ and $\lambda_m = \lambda_m^{(k_m)}$.

Empirical studies suggested that a small value of the step length factor ν often improves and almost never worsens the performance of boosting [45]. Therefore, it is recommended that ν is chosen to be as small as possible while the algorithm remains computationally feasible. In the settings we have considered, the performance of I-Boost is not sensitive to ν within the range of $\nu \in (0.05, 0.5)$. Therefore, we set ν to a moderately small value of 0.1.

Conventional boosting methods require a stopping criterion to avoid over-fitting. In our experience, however, because the tuning parameters are selected separately at each iteration for I-Boost, they eventually lead to shrinkage of all (current) parameter estimates. Therefore, we do not adopt a separate procedure to determine the stopping time of the iteration. We terminate the iteration when f_m remains constant for five consecutive iterations.

Simulation studies

In the simulation studies, we considered all data types except individual gene expression data. For each simulation data set, we generated the predictors by sampling without replacement whole vectors of predictors from the TCGA pan-cancer data set. We generated the survival time from a proportional hazards model with the baseline hazard function $h_0(t) = t$ and generated the censoring time from an exponential distribution with a mean chosen to result in censoring proportion of about 50%. We set the sample size n to 500 in all settings.

The regression parameters were chosen to produce a different proportion of signals across data types, where the signal of data type k is defined to be $\operatorname{Var}(X^{(k)'} \beta_0^{(k)})$, and the predictors were standardized. The variables with non-zero regression parameters, hereafter referred to as signal variables, were chosen to be weakly correlated. We considered three settings, with the distributions of signals and number of signal variables shown at the bottom of Fig. 1. In all settings, the signals of all data types sum to 1.2, and the regression parameters of signal variables of the same data type are equal; based on simulation studies not presented, the relative performance of different methods is very similar under different values of total signal. In Setting 1, the clinical variables contain much stronger signals than the other data types. Mutation and copy number variation data do not contain any signal. In Setting 2, all signals are concentrated on the clinical variables and gene modules, and the two data types equally share the signals. In Setting 3, the clinical variables contain the most signals, and the remaining signals are evenly distributed across the other data types.

Because we considered a total of six data types, I-Boost-CV is computationally demanding. To lessen the computational burden, we set $\nu = 0.2$ instead of the value 0.1 used in real data analysis.

Assessment of prediction

To assess an analysis method, we split the data into 30 training and testing sets with a 3:2 ratio of sample sizes. We used the R-package “sampling” [46] to perform the data split, such that the distributions of the clinical variables in the training and testing sets are approximately equal. We performed the analysis on the training sets, and the results were assessed on the corresponding testing sets using the C-index. For each split of the data, we repeated this estimation-validation procedure on different combinations of data types as predictors. We only consider combinations of data types that include clinical variables, because clinical variables are almost always considered in practice, and one of the main objectives of this paper is to evaluate the prognostic value of the combination of genomics and clinical data. The analyses were conducted

on the 30 splits of the data and on the 48 combinations of data types for the LUAD, KIRC, and pan-cancer data sets.

To quantify the prediction accuracy, we used the C-index. Let T_i be the survival time and X_i be a vector of predictors for the i th subject, and let β be a vector of regression parameters. The risk score is defined as $X_i'\beta$. If T_i and $X_i'\beta$ are continuous, then the C-index is defined as $P(X_i'\beta > X_j'\beta \mid T_i < T_j)$. The C-index is the probability that for a random pair of subjects in which the first subject has a shorter survival time, the risk score for the first subject is higher. Thus, C-index measures how well the risk score aligns with the actual survival time. For each pair of training and testing sets, we set β to be the parameter estimate obtained from the training set and estimated the C-index for the testing set using the method of Pencina and D'Agostino [27]. If no variable was selected, then a C-index value of 0.5 was assigned.

We used the NRI to compare the prediction accuracy of a model of interest and a baseline model. Let T be the survival time, X and \tilde{X} be vectors of predictors for the model of interest and the baseline model, respectively, and β and $\tilde{\beta}$ be the corresponding vectors of regression parameters. Let q and \tilde{q} be the (estimated) survival probabilities at a fixed time point t_0 given univariate covariates $X'\beta$ and $\tilde{X}'\tilde{\beta}$, respectively, where t_0 is a survival-time threshold, such that subjects with $T < t_0$ are considered high risk. The NRI is defined as $P(q < \tilde{q} \mid T < t_0) - P(q < \tilde{q} \mid T > t_0)$. A large NRI means that by switching from the baseline model to the model of interest, the direction of change of the predicted risk aligns with the actual survival time for a large proportion of subjects. To compute the NRI between two models using a pair of training and testing sets, we set $(\beta, \tilde{\beta})$ to be the parameter estimates obtained from the training set and calculated (q, \tilde{q}) on the testing set. We estimated the NRI and its confidence interval on the testing set using the method of Uno et al. [29]. The reported NRI values and confidence limits are the average values over 30 training and testing data splits. Note that this NRI is one half the value of the NRI (> 0) defined in Pencina et al. [28, 30].

Additional file

Additional file 1: Fig. S1. The prognostic value of genomics data types. (PDF 38 kb)

Acknowledgements

Not applicable.

Funding

This work was supported with funds from the NCI Breast SPOR program (P50-CA58223-09A1), the Breast Cancer Research Foundation, the Susan G. Komen, The V Foundation for Cancer Research, and by National Institutes of Health grants R01CA148761, R01GM047845, R01HG009974, and P01CA142538.

Availability of data and materials

For clinical variables, somatic mutation, copy number variation, mRNA expression, and miRNA expression, data were obtained from Synapse [34] (<https://www.synapse.org/#!Synapse:syn2468297>). Protein expression data were downloaded from Broad GDAC Firehose [35] (<https://doi.org/10.7908/C11G0KM9>). The processed data set used in this paper is deposited on Zenodo [47] (<https://doi.org/10.5281/zenodo.2530387>).

The I-Boost software is distributed under the MIT license and can be downloaded from Github [48] (<https://github.com/alexwky/I-Boost>) or Zenodo [49] (<https://doi.org/10.5281/zenodo.2529986>). The simulated data sets and the codes to reproduce all the analyses presented in this paper are available on GitHub [50] (<https://github.com/alexwky/I-Boost-Paper2019>) and deposited on Zenodo [51] (<https://doi.org/10.5281/zenodo.2532847>).

Authors' contributions

DYL and CMP coordinated the overall studies. KYW performed the statistical analyses. KYW, ABN, DZ, and DYL developed the statistical methodologies. CF and JSP coordinated the processing of the data. KYW, MT, DYL, and CMP wrote the paper, which all authors reviewed. All authors read and approved the final manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

CMP is an equity stock holder, consultant, and Board of Director Member of BioClassifier LLC and GeneCentric Diagnostics. CMP is also listed as an inventor on patents on the Breast PAM50 and Lung Cancer Subtyping assays. The other authors declare that they have no competing interests.

Publisher's Note

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

Author details

¹Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hom, Hong Kong. ²Lineberger Comprehensive Cancer Center, University of North Carolina, Chapel Hill, NC 27599, USA. ³Department of Genetics, University of North Carolina, Chapel Hill, NC 27599, USA. ⁴Department of Statistics and Operations Research, University of North Carolina, Chapel Hill, NC 27599, USA. ⁵Department of Biostatistics, University of North Carolina, Chapel Hill, NC 27599, USA.

Received: 15 May 2018 Accepted: 23 January 2019

Published online: 07 March 2019

References

- Shedden K, Taylor JM, Enkemann SA, Tsao MS, Yeatman TJ, Gerald WL, et al. Gene expression-based survival prediction in lung adenocarcinoma: a multi-site, blinded validation study. *Nat Med*. 2008;14:822–7.
- West M, Blanchette C, Dressman H, Huang E, Ishida S, Spang R, et al. Predicting the clinical status of human breast cancer by using gene expression profiles. *Proc Natl Acad Sci USA*. 2001;98:11462–7.
- Beer DG, Kardia SL, Huang CC, Giordano TJ, Levin AM, Misek DE, et al. Gene-expression profiles predict survival of patients with lung adenocarcinoma. *Nat Med*. 2002;8:816–24.
- Shipp MA, Ross KN, Tamayo P, Weng AP, Kutok JL, Aguiar RC, et al. Diffuse large B-cell lymphoma outcome prediction by gene-expression profiling and supervised machine learning. *Nat Med*. 2002;8:68–74.
- van't Veer LJ, Dai H, van de Vijver MJ, He YD, Hart AA, Mao M, et al. Gene expression profiling predicts clinical outcome of breast cancer. *Nature*. 2002;415:530–6.
- Shen R, Olshen AB, Ladanyi M. Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics*. 2009;25:2906–12.
- Mo Q, Wang S, Seshan VE, Olshen AB, Schultz N, Sander C, et al. Pattern discovery and cancer gene identification in integrated cancer genomic data. *Proc Natl Acad Sci USA*. 2013;110:4245–50.

8. Lock EF, Hoadley KA, Marron JS, Nobel AB. Joint and individual variation explained (JIVE) for integrated analysis of multiple data types. *Ann Appl Stat*. 2013;7:523–42.
9. Wang W, Baladandayuthapani V, Morris JS, Broom BM, Manyam G, Do KA. iBAG: integrative Bayesian analysis of high-dimensional multiplatform genomics data. *Bioinformatics*. 2013;29:149–59.
10. Yuan Y, Van Allen EM, Omberg L, Wagle N, Amin-Mansour A, Sokolov A, et al. Assessing the clinical utility of cancer genomic and proteomic data across tumor types. *Nat Biotechnol*. 2014;32:644–52.
11. Parker JS, Mullins M, Cheang MC, Leung S, Voduc D, Vickery T, et al. Supervised risk predictor of breast cancer based on intrinsic subtypes. *J Clin Oncol*. 2009;27:1160–7.
12. Fan C, Prat A, Parker JS, Liu Y, Carey LA, Troester MA, et al. Building prognostic models for breast cancer patients using clinical variables and hundreds of gene expression signatures. *BMC Med Genomics*. 2011;4:3.
13. Tibshirani R. Regression shrinkage and selection via the lasso. *J R Stat Soc Series B Stat Methodol*. 1996;58:267–88.
14. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J R Stat Soc Series B Stat Methodol*. 2005;67:301–20.
15. Schapire RE. The strength of weak learnability. *Mach Learn*. 1990;5:197–227.
16. Freund Y, Schapire RE. A decision-theoretic generalization of on-line learning and an application to boosting. *J Comput Syst Sci*. 1997;55:119–39.
17. Breiman L. Arcing classifier (with discussion). *Ann Stat*. 1998;26:801–49.
18. Friedman J, Hastie T, Tibshirani R. Additive logistic regression: a statistical view of boosting (with discussion). *Ann Stat*. 2000;28:337–407.
19. Hothorn T, Bühlmann P, Dudoit S, Molinaro A, Van Der Laan MJ. Survival ensembles. *Biostatistics*. 2005;7:355–73.
20. Mayr A, Binder H, Gefeller O, Schmid M. The evolution of boosting algorithms. *Methods Inf Med*. 2014;53:419–27.
21. Mayr A, Binder H, Gefeller O, Schmid M. Extending statistical boosting. *Methods Inf Med*. 2014;53:428–35.
22. Bühlmann P, Yu B. Boosting with the L2 loss: regression and classification. *J Am Stat Assoc*. 2003;98:324–39.
23. Cox DR. Regression models and life-tables. *J R Stat Soc Series B Stat Methodol*. 1972;34:187–220.
24. Cox DR. Partial likelihood. *Biometrika*. 1975;62:269–76.
25. Sabourin JA, Valdar W, Nobel AB. A permutation approach for selecting the penalty parameter in penalized model selection. *Biometrics*. 2015;71:1185–94.
26. Hoadley KA, Yau C, Wolf DM, Cherniack AD, Tamborero D, Ng S, et al. Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin. *Cell*. 2014;158:929–44.
27. Pencina MJ, D'Agostino RB. Overall C as a measure of discrimination in survival analysis: model specific population value and confidence interval estimation. *Stat Med*. 2004;23:2109–23.
28. Pencina MJ, D'Agostino Sr RB, Steyerberg EW. Extensions of net reclassification improvement calculations to measure usefulness of new biomarkers. *Stat Med*. 2011;30:11–21.
29. Uno H, Tian L, Cai T, Kohane IS, Wei L. A unified inference procedure for a class of measures to assess improvement in risk prediction systems with survival data. *Stat Med*. 2013;32:2430–42.
30. Pencina MJ, D'Agostino RB, Pencina KM, Janssens ACJ, Greenland P. Interpreting incremental value of markers added to risk prediction models. *Am J Epidemiol*. 2012;176:473–81.
31. Lieu CH, Renfro LA, De Gramont A, Meyers JP, Maughan TS, Seymour MT, et al. Association of age with survival in patients with metastatic colorectal cancer: analysis from the ARCAD Clinical Trials Program. *J Clin Oncol*. 2014;32:2975–82.
32. de la Rochefordière A, Campana F, Fenton J, Vilcoq J, Fourquet A, Asselain B, et al. Age as prognostic factor in premenopausal breast carcinoma. *Lancet*. 1993;341:1039–43.
33. Asmis TR, Ding K, Seymour L, Shepherd FA, Leighl NB, Winton TL, et al. Age and comorbidity as independent prognostic factors in the treatment of non-small-cell lung cancer: a review of National Cancer Institute of Canada Clinical Trials Group trials. *J Clin Oncol*. 2008;26:54–9.
34. Sage Bionetworks Repository Synapse. Multiplatform analysis of 12 cancer types to identify integrative subtypes; <https://www.synapse.org/#Synapse:syn2468297>. Accessed 12 Oct 2015.
35. Broad Institute TCGA Genome Data Analysis Center. Analysis-ready standardized TCGA data from Broad GDAC Firehose 2016_01_28 run. Broad Institute of MIT and Harvard; <https://doi.org/10.7908/C11G0KM9>. Accessed 26 Jun 2017.
36. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics*. 2011;12:323.
37. The Cancer Genome Atlas Research Network. Comprehensive molecular portraits of human breast tumours. *Nature*. 2012;490:61–70.
38. The Cancer Genome Atlas Research Network. Comprehensive molecular characterization of human colon and rectal cancer. *Nature*. 2012;487:330–7.
39. The Cancer Genome Atlas Research Network. Comprehensive genomic characterization of squamous cell lung cancers. *Nature*. 2012;489:519–25.
40. Zack TI, Schumacher SE, Carter SL, Cherniack AD, Saksena G, Tabak B, et al. Pan-cancer patterns of somatic copy number alteration. *Nat Genet*. 2013;45:1134–40.
41. Beroukhi R, Mermel CH, Porter D, Wei G, Raychaudhuri S, Donovan J, et al. The landscape of somatic copy-number alteration across human cancers. *Nature*. 2010;463:899–905.
42. Chao HH, He X, Parker JS, Zhao W, Perou CM. Micro-scale genomic DNA copy number aberrations as another means of mutagenesis in breast cancer. *PLoS ONE*. 2012;7:e51719.
43. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci USA*. 2005;102:15545–50.
44. Simon N, Friedman J, Hastie T, Tibshirani R. Regularization paths for Cox's proportional hazards model via coordinate descent. *J Stat Softw*. 2011;39:1–13.
45. Friedman JH. Greedy function approximation: a gradient boosting machine. *Ann Stat*. 2001;29:1189–232.
46. Tillé Y, Matei A. Sampling: survey sampling. R package version 2.8. 2016.
47. Wong KY, Fan C, Maki T, Parker JS, Nobel AB, Zeng D, et al. I-Boost: an integrative boosting approach for predicting survival time with multiple genomics platforms. Processed data; 2019. <https://doi.org/10.5281/zenodo.2530387>. Accessed 4 Jan 2019.
48. Wong KY, Fan C, Maki T, Parker JS, Nobel AB, Zeng D, et al. I-Boost: an integrative boosting approach for predicting survival time with multiple genomics platforms. Source code Github repository; 2019. <https://github.com/alexwky/I-Boost>. Accessed 4 Jan 2019.
49. Wong KY, Fan C, Maki T, Parker JS, Nobel AB, Zeng D, et al. I-Boost: an integrative boosting approach for predicting survival time with multiple genomics platforms. Source code; 2019. <https://doi.org/10.5281/zenodo.2529986>. Accessed 4 Jan 2019.
50. Wong KY, Fan C, Maki T, Parker JS, Nobel AB, Zeng D, et al. I-Boost: an integrative boosting approach for predicting survival time with multiple genomics platforms. Code Github repository; 2019. <https://github.com/alexwky/I-Boost-Paper2019>. Accessed 7 Jan 2019.
51. Wong KY, Fan C, Maki T, Parker JS, Nobel AB, Zeng D, et al. I-Boost: an integrative boosting approach for predicting survival time with multiple genomics platforms. Code; 2019. <https://doi.org/10.5281/zenodo.2532847>. Accessed 7 Jan 2019.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

