

Original Article

Development of a novel prognostic model based on clinicopathological and expressional profiles for cutaneous melanoma patients

Yanbin Peng¹, Yunfeng Chu¹, Zhong Chen¹, Wen Zhou¹, Shengxiang Wan¹, Yingfeng Xiao¹, Youlong Zhang², Jialu Li²

¹Department of Microsurgery, Peking University Shenzhen Hospital, Shenzhen, China; ²Department of Biostatistics, Huajia Biomedical Intelligence, Shenzhen, China

Received April 4, 2020; Accepted July 1, 2020; Epub August 15, 2020; Published August 30, 2020

Abstract: Accurate assessment of prognosis is important for the management of heterogeneous survival outcomes of cutaneous melanoma. It is of clinical interest to identify a single set of combinatorial markers that have prognostic value for both recurrence and death events. We enrolled in this study a total of 387 patients from TCGA-SKCM with complete information in failure time and event, clinicopathological variables, expressional and mutational profiles. We performed a composite endpoint-based competing risks analysis to determine the best model in predicting recurrence or death events. We further validated the model performance within pathologically-defined subgroups. The model combined clinicopathological variables and expression markers performed the best among all models in certain time periods. The features selected by this combinatorial model had reasonable prediction performance for both recurrence and death outcomes. The resultant prognostic risk score generated by the model provides a higher-resolution risk stratification within pathologically-defined subgroups. Our study thus provided a new model that can handle both competing events and multiple endpoints. Adding gene-expression information into clinicopathological variables significantly improved the prognostic prediction for specific subgroups.

Keywords: Cutaneous melanoma, prognostic analysis, competing risks survival model

Introduction

Melanoma is an aggressive type of skin cancer, with about 96,480 new cases and 7,230 deaths reported in the USA in 2019 [1-6]. Melanoma has a high overall survival rate if successfully treated with surgery, but the survival rate drops significantly if metastasis occurs [7]. Cutaneous melanoma is a major subtype of melanoma characterized by relatively poor prognosis. Routine prognostic analysis of melanoma uses clinicopathological features including Breslow tumor thickness, age at diagnosis, ulceration, mitotic index and Clark level [8, 9]. The high-throughput screening technology actively developed in recent years has allowed a more thorough investigation of tumor heterogeneity and microenvironment at the molecular level. Many studies thus have worked on finding omics-based prognostic markers for cutaneous melanoma patients.

For univariate signatures, NEK2 upregulation was identified as being correlated with poor recurrence-free survival (RFS) and overall survival (OS) under both univariate and multivariate analyses [10]. The BRAF mutation was associated with a lower risk of recurrence for stage III melanoma patients treated by adjuvant therapy of dabrafenib and trametinib [5]. These single-gene signatures showed an association in prognosis, but had not been specifically evaluated in the prediction of clinical outcomes.

For multiplexed biomarkers, a 31-gene (28 prognostic genes and 3 control genes) signature was identified based on microarray expression data. The signature achieved an AUC of 0.91 in predicting the metastatic status on validation cohort [2]. Another gene expression signature based on RNA sequencing (RNA-seq) data from TCGA-SKCM (The Cancer Genome

Prognosis of cutaneous melanoma

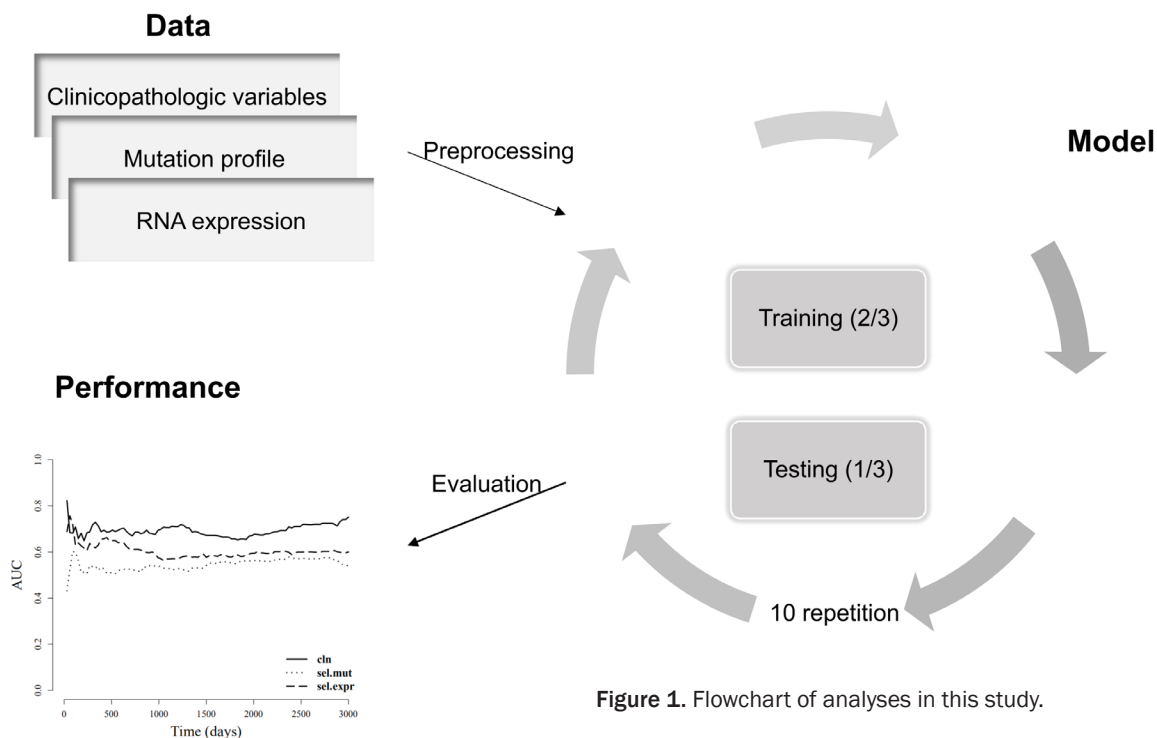


Figure 1. Flowchart of analyses in this study.

Atlas-Skin Cutaneous Melanoma) dataset [11] was used to calculate sample-specific Leukocyte Infiltration Scores (LIS), which were found to be associated with survival in independent datasets [12]. There were also studies focused on developing prognostic models by combining multiple types of data. Jayawardana et al [13] developed models to classify 1-year and 4-year survival status based on clinicopathological, mutations, mRNA, microRNA, and protein information as well as different combinations of each. They found that models based on the combination of clinicopathological variables and mRNA profiles performed the best under a cross-validation framework. Jiang et al [14] used sparse PCA and partial least squares methods to develop models using multi-dimensional omics data. Their methods adequately considered all profile information and had achieved the highest C-index values of OS prediction so far. Although these analyses provided effective prognostic signatures, few studies focused on developing a single combinatorial signature for both OS and RFS analysis.

Our work presented a composite endpoint-based competing risks model based on clinicopathological variables and gene expression data. This model can not only correct the bias

caused by competing risk events, but also selects a single set of features that are associated with both OS and RFS. We found that the model performed significantly better than models based only on clinicopathological variables or models with other types of data integration.

Methods

Data sources

Our study was performed on the TCGA-SKCM dataset which contain 470 cutaneous melanoma patients. The study work flowchart was shown in **Figure 1**. We removed clinicopathological variables with a high proportion of missing values and filtered out extremely unbalanced categorical variables. This resulted in six clinicopathological variables for further analysis: age at diagnosis, gender, tumor status, primary location, AJCC (the American Joint Committee on Cancer) pathologic tumor stage and AJCC nodes pathologic stage.

The sketch of the disease process after initial treatment was shown in **Figure S1**. The initial treatment date was also defined as the initial pathological diagnosis date for this dataset. In this data, one may progress to recurrence after

initial treatment (1 to 2) or die directly without experiencing recurrence (1 to 3), the latter of which would prevent one from having recurrence. Such a death was then the competing risks event for the recurrence event. Routine survival analysis could induce bias here because the competing risks event should not be treated as the censoring event. We therefore analyzed recurrence-related survival by competing risks modeling. The failure time was the time interval from initial treatment to the corresponding failure event of interest. To include as much data as possible, we set the composite endpoint time as the time of recurrence if one has the record, or as the time of death if one has the record of recurrence but no detail recurrence time information. For those patients who had recurrence but didn't have the time of recurrence or death, the composite endpoint time was set as the last follow-up time. This composite endpoint was the primary endpoint for the competing risks analysis.

Data preprocessing

We excluded 79 patients with missing data in any of the selected clinicopathological variables, endpoints or follow-up data. Four patients were removed because they didn't have gene expression or mutational profile data. This resulted in a total of 387 patients which were included in this study.

For RNA expression data, the FPKM (fragments per kilobase of exon model per million reads mapped) value was used to represent gene expression level. We kept genes by the following criteria: 1) mean of FPKM across all patients is greater than 1; 2) have non-zero expression in more than 60% patients (232 patients); 3) standard deviation of log-transformed FPKM across samples is greater than 0.5. These resulted in a total of 5390 genes for further analysis. For mutational profiles, we selected 772 genes out of 17509 genes which have at least one somatic mutation. Each of the selected genes had at least one somatic mutation in more than 30 samples. We summarized the number of these three types of data before and after feature preprocessing in [Table S1](#).

The clinicopathological variables of the enrolled patients were summarized in [Table S2](#), in which we also computed the p values of asso-

ciation analysis with recurrence or death status. The methods we used for association analysis were as follows: 1) one-way ANOVA for continuous variables if following a normal distribution; 2) Wilcoxon test for continuous variables if not normally distributed; 3) chi-square test for discrete variables which contained a number of samples greater than 8 in all categories; 4) Fisher exact test for discrete variables if there's any category whose number of samples is smaller than 9.

Model development framework

For the gene expression-based and gene mutation-based models, we first applied univariate Cox proportional hazard model to screen for potentially important features. The Wald's test p value of model fitting coefficient was used to rank the importance of features, and only a certain number of features with the smallest p -values were selected. To select the optimal number of features, we applied a "nested-cross validation" approach. We first randomly split the 387 input patients into 3 groups, two of which were used for inner 3-fold cross validation. In this inner cross validation, we aimed to select the best penalty parameter inherent in the regularized survival model, while in the outer cross validation, we used the one-portion left-out data to validate the model performance given a specific number of input features. This procedure was repeated 5 times to improve the randomness of data split. We used mean and standard deviation of C-index to evaluate the accuracy of model validation.

We compared the performance of different combinatorial modeling based on a random sample of 129 patient's (1/3 of 387 patients) data. The remaining 258 patients were used for 3-fold cross-validation to determine the best model within each data type or data integration type.

Survival modeling

For competing risks modeling, the cause-specific hazard was fitted to train a risk score. Then a Fine and Gray model was fitted by setting this score as the single covariate. The cause-specific hazard model was combined with Lasso (Least absolute shrinkage and selection operator) method to deal with high-dimensional data as previously described [15].

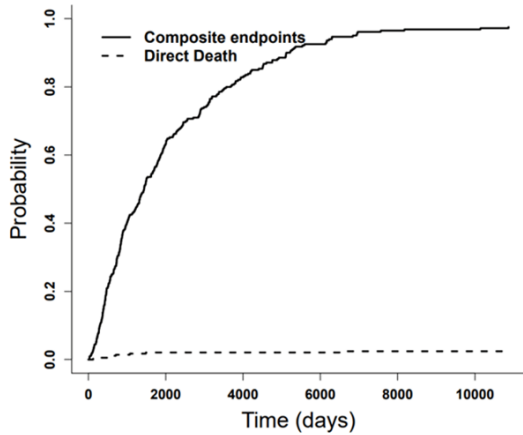


Figure 2. The cumulative incidence function of composite endpoints and death without recurrence.

Lasso is an approach to regularize model and select features concurrently. It constrains the regression coefficients by adding L_1 norm term to cost function. For Cox PH model, the lasso form of optimal function is given as [16, 17]:

$$\operatorname{argmax}_{\beta} \left\{ \sum_{r \in D} \ln \frac{\exp(\beta^T x^r)}{\sum_{i \in R_r} \exp(\beta^T x^i)} - \lambda \sum_j |\beta_j| \right\}$$

where the β is the regression coefficients, x is the covariates, D is the set of indices of observed events, R_r denotes the individuals set at risk at time t_r and λ is the coefficient controlling the penalty strength. In this study, the optimal λ was selected by a 3-fold cross validation.

The Fine and Gray model [18] regresses directly on cumulative incidence function (CIF). The CIF, denoted as $I_k(t)$ of cause k , can also be interpreted as the cumulative incidence probability $\Pr(T \leq t, D = k)$. It was defined as:

$$I_k(t) = \int_0^t \lambda_k(s) S(s) ds$$

where $\lambda_k(t)$ is the hazard of cause k at time t , $S(t) = \exp[-\sum_{k=1}^K \int_0^t \lambda_k(s) ds]$ is the survival function. To incorporate covariates information, Fine and Gray imposes a proportional hazards assumption on the subdistribution hazards:

$$\bar{\lambda}_k(t | Z) = \bar{\lambda}_{k,0}(t) \exp(\beta_k^T Z)$$

In the cause-specific hazard analysis, the risk set decreases at each time point at which there is a failure of another cause. For $\bar{\lambda}_k(t)$ estima-

tion, individuals who fail from another cause remained in the risk set [18].

To reduce bias in data splitting, we repeated the model training and testing 10 times. The time-dependent AUC curves [19] were computed to evaluate the prediction performance. All analyses were performed by R software and packages including “survival”, “riskRegression”, “glmnet”, “caret” and “survminer”.

Results

Workflow and patient characteristics

The flowchart of model comparison analysis was shown in **Figure 1**. The model input data included clinicopathological variables, expressional and mutational profiles. We developed prognostic models using these three types of data or their combinations. The prognostic model we developed here was a competing risks model, by which we predicted the risk of recurrence occurrence after initial treatment while competing with the risk of death without recurrence.

A total of 387 patients with complete information were enrolled in our study, in which 308 patients (79.6%) developed recurrence. The composite endpoints consisted of 87 cases of recurrence, 118 deaths, and 103 last follow-up events. The median survival time of composite endpoints was 1244 days. The median survival time of those who died without experiencing recurrence (8 patients, 2.1% of cohort) was 662 days. The CIF curve of recurrence onset was shown in **Figure 2**. Most patients experienced composite events before 6000 days, and the hazard of composite events was relatively high until about 2000 days. We prepared 5 sets of features from single type of data (STF) and 4 sets of integrative features (ITF) as the model input. The composition of each feature set used to develop the competing risks model was shown in **Table S3**. The mean of median failure time of 10 repetitions for the development and testing dataset were 907 days and 884 days, respectively.

Modeling with a single type of data

Overall, under the context of competing risks, the clinicopathological variables showed the best prediction performance. Its mean time-

Prognosis of cutaneous melanoma

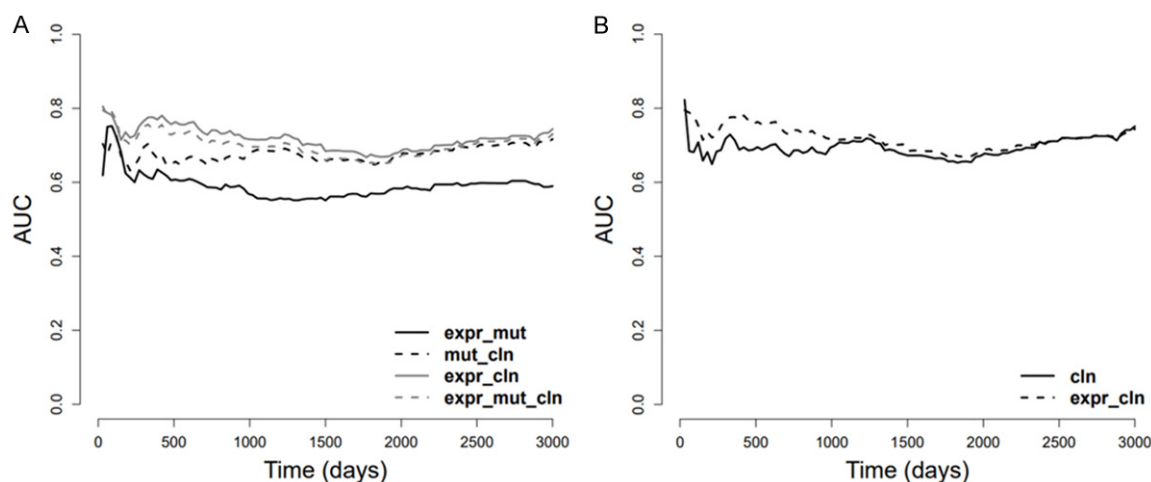


Figure 3. Time-dependent AUC curve of integrative models (A) and comparison of clinicopathological variables-based model and the best integrative model (B). Abbreviations: cln, models based on clinicopathological variables; expr_mut, models based on RNA expressional and gene mutational profiles; mut_cln, models based on clinicopathological variables and gene mutational profiles; expr_cln, models based on clinicopathological variables and RNA expressional profiles; expr_mut_cln, models based on the three types of data.

dependent AUC within 3000 days after initial treatment was 0.694 (sd = 0.023). The gene expression-based model showed a better prognostic ability as compared to that of gene mutation with (mean time-dependent AUC = 0.591 vs. 0.536) or without (mean time-dependent AUC = 0.604 vs. 0.545) performing univariate feature selection.

We next investigated whether adding univariate feature selection step benefits prognostic prediction. For RNA expression-based model, the mean time-dependent AUC increased by 0.013 with the extra feature selection step. The performance was significantly improved in the time interval before about 1000 days and after about 2300 days (mean time-dependent AUC = 0.619 vs. 0.596, p value = $1.54e-5$), although such improvement was limited (Figure S2C). For the gene mutation-based model, the overall accuracy within 3000 days increased by 0.009 after the feature selection, but the improvement became more obvious (mean time-dependent AUC = 0.557 vs. 0.531, p value = $2.2e-16$) after about 1200 days (Figure S2C), suggesting the pre-selection strategy helped increase the accuracy within specific time periods.

Modeling with integrated data

We evaluated the benefit of combinatorial modeling by combining clinicopathological variables, gene expressional or mutational profiles

as model input. As shown in Figure 3A, models combining expressional features and clinicopathological covariates performed best across all time periods. The mean of its time-dependent AUC curve within 3000 days was 0.719 (sd = 0.030), which is significantly higher than that of the second best model based on a combination of three types of data (mean time-dependent AUC = 0.702, sd = 0.031, p value < 0.001).

We next compared the best integrative model to the clinicopathological variable-based model as described above. Their time-dependent AUC curves were shown in Figure 3B. There was a slight improvement of the models (increased 2.6%) with the addition of gene expressional profiles to clinicopathological variables. Such improvement was much more obvious before 1000 days (mean time-dependent AUC = 0.750 vs. 0.692, p value = $3.49e-13$).

Validation with RFS and OS analysis

To further evaluate the results developed based on composite endpoints, we carried out OS and RFS analysis without considering competing risks. The OS analysis was based on the 387 patients described above and the RFS analysis was developed on 166 patients after removing 221 patients who were missing the data of time of recurrence. We used C-index to evaluate the prediction performance. The boxplot summary of C-index was shown in

Prognosis of cutaneous melanoma

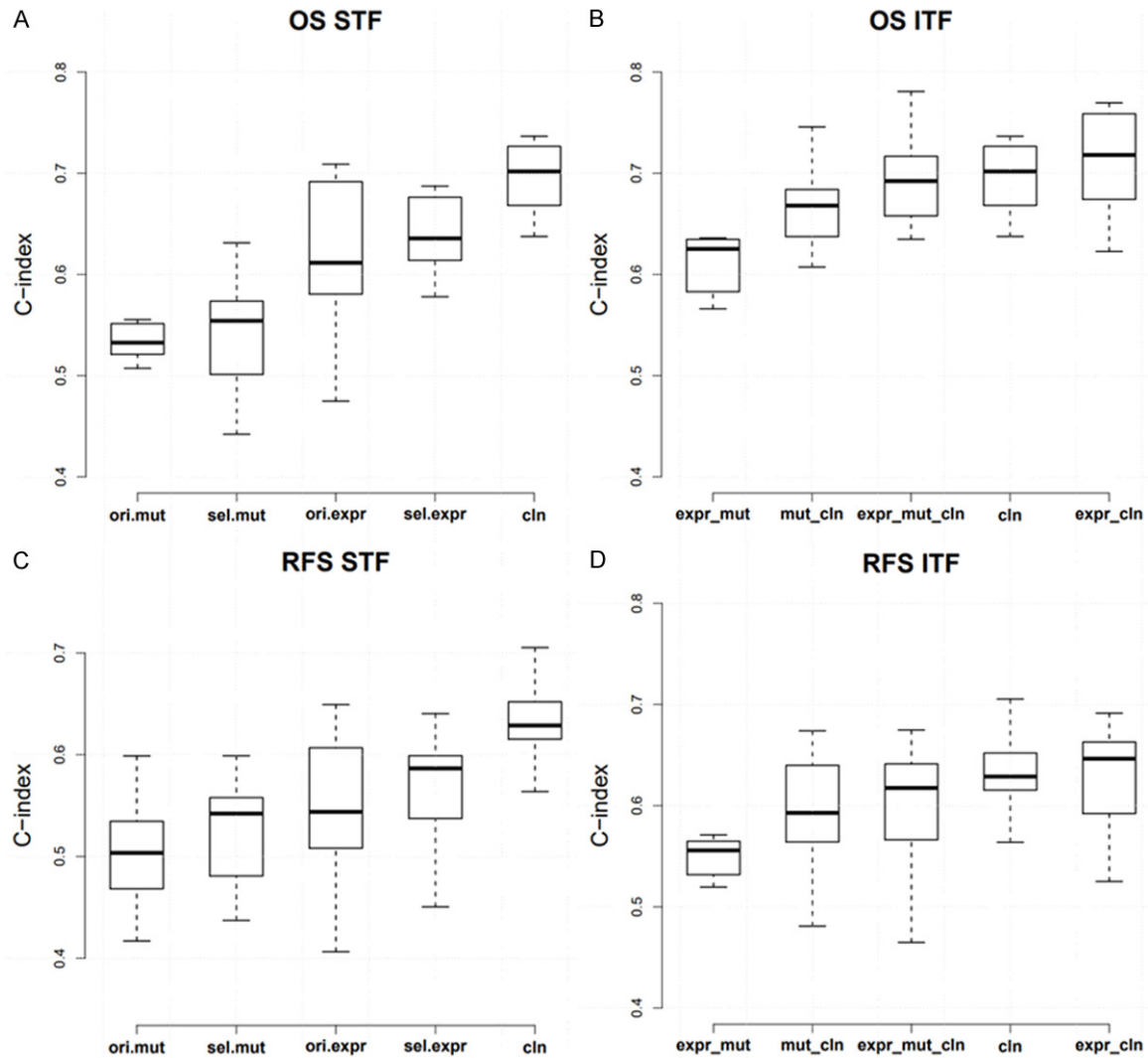


Figure 4. Boxplots of C-index for OS and RFS analysis. A. Shows the model performance based on single type data for OS; B. Represents integrate clinicopathological models for OS; C. For single type data for RFS; D. For integrative models for RFS.

Figure 4 and the mean of each box was summarized in [Table S3](#). For models based on single types of data, the clinicopathological variables performed best in both OS (mean/sd C-index: 0.697/0.033) and RFS (mean/sd C-index: 0.629/0.041) analysis. For models based on integrated data, those combining RNA expression data and clinicopathological variables performed best in both OS (mean/sd C-index: 0.709/0.053) and RFS (mean/sd C-index: 0.631/0.052) analysis.

Integrative prognostic score and interaction analysis

To illustrate the potential usage of the integrative model in clinical setting, we computed a

prognostic score by summing over the product of features and their coefficients ([Table S4](#)) from the competing risk model developed above. By setting the median of the score as the threshold, the patients were categorized into higher and lower risk subgroups. The OS or RFS of these two subgroups were all significantly different under Kaplan-Meier analysis ([Figures 5A, S3A](#)). To see whether the prognostic score could provide higher-resolution risk stratifications, we performed Kaplan-Meier analysis within locoregional or metastatic subgroup ([Figure 5B, 5C](#)). The score can further differentiate higher risk from lower risk patients within each subgroup for either OS ([Figure 5](#)) or RFS ([Figure S3](#)).

Prognosis of cutaneous melanoma

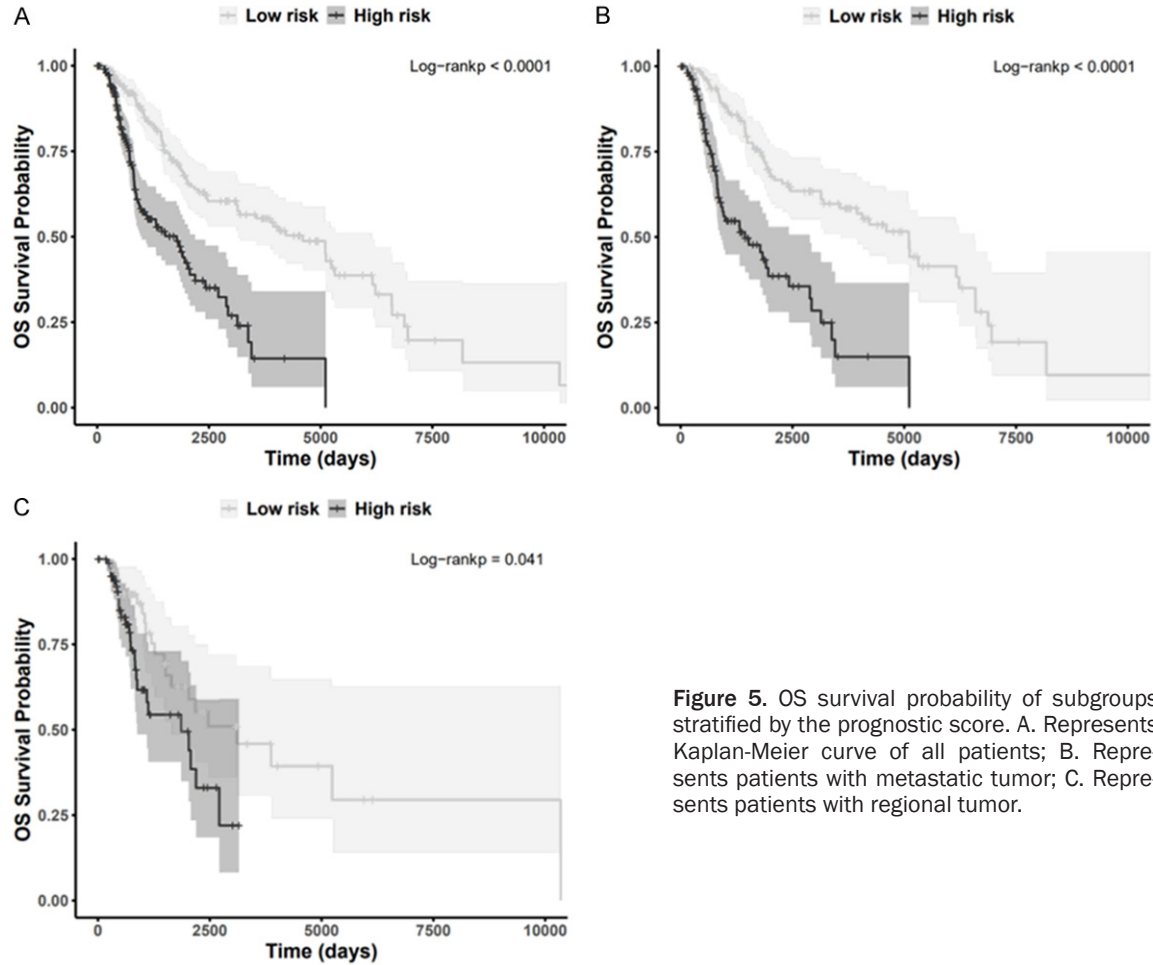


Figure 5. OS survival probability of subgroups stratified by the prognostic score. A. Represents Kaplan-Meier curve of all patients; B. Represents patients with metastatic tumor; C. Represents patients with regional tumor.

Table 1. Summary of prediction accuracy (mean C-index/standard deviation) on subgroups stratified by tumor metastatic status

Feature set	OS		RFS	
	Locoregional	Metastatic	Locoregional	Metastatic
ori.expr	0.521/0.070	0.586/0.053	0.522/0.122	0.463/0.091
sel.expr	0.561/0.102	0.599/0.040	0.658/0.103	0.492/0.079
ori.mut	0.481/0.064	0.499/0.036	0.513/0.058	0.523/0.064
sel.mut	0.484/0.065	0.525/0.032	0.568/0.058	0.543/0.062
Cln	0.551/0.076	0.719/0.054	0.612/0.057	0.666/0.079
expr_cln	0.564/0.084	0.725/0.038	0.680/0.091	0.589/0.129
mut_cln	0.449/0.055	0.701/0.051	0.610/0.082	0.614/0.086
expr_mut	0.530/0.078	0.589/0.050	0.619/0.081	0.552/0.085
expr_mut_cln	0.526/0.082	0.712/0.048	0.631/0.087	0.567/0.076

We then further evaluated the level of importance of incorporating gene expression into prognostic analysis within pathologically-defined subgroups. We performed OS or RFS model training and testing using data from each subgroup. The mean and standard deviation

of C-index of each type of model was summarized in **Table 1**. Compared to the results based on all patients, the prediction accuracy of models combining clinicopathological variables and gene expression varied greatly between locoregional and metastatic group (mean C-index of locoregional/metastatic groups = 0.564/0.725 for OS, 0.680/0.589 for RFS). Of note, such model has significantly improved prediction accuracy (4.9%)

for the locoregional group as compared to those based on all patients in RFS analysis, or to the model simply based on clinicopathological variables for the locoregional group. We also performed predictions on subgroups stratified by AJCC tumor pathologic stage ("Stage

0~II” vs. “Stage III, IV”). As shown in [Table S5](#), the prediction accuracy of models based on clinicopathological variables and gene expression was also different across subgroups (mean C-index of “Stage 0~II”/“Stage III, IV” groups = 0.683/0.645 for OS, 0.594/0.600 for RFS).

Discussion

In this study, we extensively compared prognostic models of cutaneous melanoma from the perspective of input data types (clinicopathological, expressional and mutational profiles), endpoints (composite, recurrence or death), statistical methods (competing risk modeling or typical high-dimensional survival modeling) and evaluation criteria (time-dependent AUC or C-index). We found that combining clinicopathological variables and gene expression achieves the best overall prediction performance. We also showed that this combinatorial model can provide higher-resolution risk stratifications within non-metastatic or metastatic patients.

This risk stratification was particularly effective among locoregional patients, where the gene expression data added significantly more benefit in predicting recurrence (mean C-index = 0.680) as compared to that of clinicopathological variables-based model (mean C-index = 0.612) or other types of combinatorial models (best mean C-index = 0.631).

We noted that the best competing risks model developed in this study ([Table S4](#)) did not include any one of genes in a commercially used prognostic marker, DecisionDx-Melanoma [2] (31-GEP). One possible reason could be the different platforms used for gene expression measurements. The 31-GEP was developed using microarray data while the TCGA data we used in the study employed RNA sequencing technology. The two platforms showed certain discrepancy in previous studies [20, 21]. Another possible reason could be the different statistical learning methods used for the high-dimensional variable selection. It is still unknown about the exact advantage of 31-GEP over combinatorial or integrative models in prognostic prediction. Our finding is consistent with a previous study [13] which used different cohort and methods for model development. We all concluded that the high-throughput “omics” profiling alone does not show a

clear superiority in prognosis over commonly used clinicopathological variables, but combining gene expression data and clinicopathological variables can improve the performance.

Our study has limitations. First, the data we used has many missing values in recurrence time. The failure time of composite endpoint we defined here, which was mixed with recurrence and death events, is not clinically meaningful, but the best combination we identified based on this setting showed consistency among RFS or OS analyses. The risk score we developed from the composite endpoint-based competing risk model using training dataset can significantly stratify risk groups in both OS and RFS analyses ([Figures 5, S3](#)). Second, the dataset has a very small proportion (2.1%) of competing risks cases, that is, patients who died without experiencing recurrence. This could make our application of competing risks survival modeling unnecessary, but when the sample size becomes larger, or if the study of late-stage patients is the focus, we would expect that the modeling strategy in this study can still be applicable. Third, all the analyses in this study were developed on a single cohort dataset. Further validation on a larger cohort of data still needs to be performed.

In summary, we developed a combinatorial model and demonstrated its prediction ability over multiple combinatorial methods. The model includes expression for only 10 genes, and can be used for assessing both OS and RFS.

Acknowledgements

This work is supported by the Shenzhen Science and Technology Project, Grant JCYJ201802-28175315535.

Disclosure of conflict of interest

None.

Address correspondence to: Dr. Jialu Li, Department of Biostatistics, Huajia Biomedical Intelligence, Room 1605, Shenzhen Overseas Chinese High-Tech Venture Park, Nanshan District, Shenzhen 518057, China. E-mail: Jialu.li@huajia.bio.com

References

- [1] Gerami P, Cook RW, Russell MC, Wilkinson J, Amaria RN, Gonzalez R, Lyle S, Jackson GL, Greisinger AJ, Johnson CE, Oelschlagel KM,

Prognosis of cutaneous melanoma

- Stone JF, Maetzold DJ, Ferris LK, Wayne JD, Cooper C, Obregon R, Delman KA and Lawson D. Gene expression profiling for molecular staging of cutaneous melanoma in patients undergoing sentinel lymph node biopsy. *J Am Acad Dermatol* 2015; 72: 780-785, e783.
- [2] Gerami P, Cook RW, Wilkinson J, Russell MC, Dhillon N, Amaria RN, Gonzalez R, Lyle S, Johnson CE, Oelschlager KM, Jackson GL, Greisinger AJ, Maetzold D, Delman KA, Lawson DH and Stone JF. Development of a prognostic genetic signature to predict the metastatic risk associated with cutaneous melanoma. *Clin Cancer Res* 2015; 21: 175-183.
- [3] Gerstung M, Pellagatti A, Malcovati L, Giagounidis A, Porta MG, Jadersten M, Dolatshad H, Verma A, Cross NC, Vyas P, Killick S, Hellstrom-Lindberg E, Cazzola M, Papaemmanuil E, Campbell PJ and Boultonwood J. Combining gene mutation with gene expression data improves outcome prediction in myelodysplastic syndromes. *Nat Commun* 2015; 6: 5901.
- [4] Song Y, Chen D, Zhang X, Luo Y and Li S. Integrating genetic mutations and expression profiles for survival prediction of lung adenocarcinoma. *Thorac Cancer* 2019; 10: 1220-1228.
- [5] Long GV, Hauschild A, Santinami M, Atkinson V, Mandala M, Chiarion-Sileni V, Larkin J, Nyakas M, Dutriaux C, Haydon A, Robert C, Mortier L, Schachter J, Schadendorf D, Lesimple T, Plummer R, Ji R, Zhang P, Mookerjee B, Legos J, Kefford R, Dummer R and Kirkwood JM. Adjuvant dabrafenib plus trametinib in stage III BRAF-mutated melanoma. *N Engl J Med* 2017; 377: 1813-1823.
- [6] Siegel RL, Miller KD and Jemal A. Cancer statistics, 2019. *CA Cancer J Clin* 2019; 69: 7-34.
- [7] Davis LE, Shalin SC and Tackett AJ. Current state of melanoma diagnosis and treatment. *Cancer Biol Ther* 2019; 20: 1366-1379.
- [8] Dickson PV and Gershenwald JE. Staging and prognosis of cutaneous melanoma. *Surg Oncol Clin N Am* 2011; 20: 1-17.
- [9] Hyams DM, Cook RW and Buzaid AC. Identification of risk in cutaneous melanoma patients: prognostic and predictive markers. *J Surg Oncol* 2019; 119: 175-186.
- [10] Huang J, Sun SG and Hou S. Aberrant NEK2 expression might be an independent predictor for poor recurrence-free survival and overall survival of skin cutaneous melanoma. *Eur Rev Med Pharmacol Sci* 2018; 22: 3694-3702.
- [11] Cancer Genome Atlas Network. Genomic classification of cutaneous melanoma. *Cell* 2015; 161: 1681-1696.
- [12] Zhao Y, Schaafsma E, Gorlov IP, Hernando E, Thomas NE, Shen R, Turk MJ, Berwick M, Amos CI and Cheng C. A leukocyte infiltration score defined by a gene signature predicts melanoma patient prognosis. *Mol Cancer Res* 2019; 17: 109-119.
- [13] Jayawardana K, Schramm SJ, Haydu L, Thompson JF, Scolyer RA, Mann GJ, Muller S and Yang JY. Determination of prognosis in metastatic melanoma through integration of clinicopathologic, mutation, mRNA, microRNA, and protein information. *Int J Cancer* 2015; 136: 863-874.
- [14] Jiang Y, Shi X, Zhao Q, Krauthammer M, Rothberg BE and Ma S. Integrated analysis of multidimensional omics data on cutaneous melanoma prognosis. *Genomics* 2016; 107: 223-230.
- [15] Tapak L, Saidijam M, Sadeghifar M, Poorolajal J and Mahjub H. Competing risks data analysis with high-dimensional covariates: an application in bladder cancer. *Genomics Proteomics Bioinformatics* 2015; 13: 169-176.
- [16] Tibshirani R. The lasso method for variable selection in the cox model. *Stat Med* 1997; 16: 385-395.
- [17] Witten DM and Tibshirani R. Survival analysis with high-dimensional covariates. *Stat Methods Med Res* 2010; 19: 29-51.
- [18] Putter H, Fiocco M and Geskus RB. Tutorial in biostatistics: competing risks and multi-state models. *Stat Med* 2007; 26: 2389-2430.
- [19] Heagerty PJ, Lumley T and Pepe MS. Time-dependent ROC curves for censored survival data and a diagnostic marker. *Biometrics* 2000; 56: 337-344.
- [20] Wang C, Gong B, Bushel PR, Thierry-Mieg J, Thierry-Mieg D, Xu J, Fang H, Hong H, Shen J, Su Z, Meehan J, Li X, Yang L, Li H, Łabaj PP, Kreil DP, Megherbi D, Gaj S, Caiment F, van Delft J, Kleinjans J, Scherer A, Devanarayan V, Wang J, Yang Y, Qian HR, Lancashire LJ, Bessarabova M, Nikolsky Y, Furlanello C, Chierici M, Albanese D, Jurman G, Riccadonna S, Filosi M, Visintainer R, Zhang KK, Li J, Hsieh JH, Svoboda DL, Fuscoe JC, Deng Y, Shi L, Paules RS, Auerbach SS and Tong W. The concordance between RNA-seq and microarray data depends on chemical treatment and transcript abundance. *Nat Biotechnol* 2014; 32: 926-32.
- [21] Zhao S, Fung-Leung WP, Bittner A, Ngo K and Liu X. Comparison of RNA-Seq and microarray in transcriptome profiling of activated T cells. *PLoS One* 2014; 9: e78644.

Prognosis of cutaneous melanoma

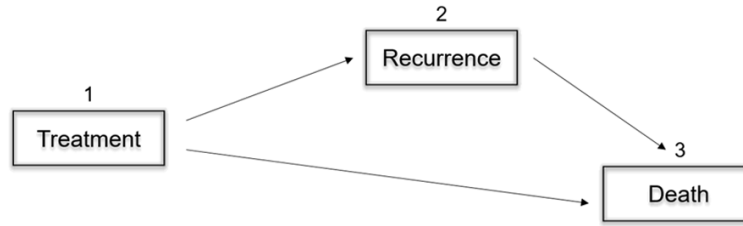


Figure S1. Sketch of disease process for patients after initial treatment. The censoring event occurred if one patient lost follow-up without experiencing recurrence or death.

Table S1. Number of types of data before and after feature preprocessing

	Number (patients, features)	
	Before feature preprocessing	After feature preprocessing
Clinicopathological variables	470, 70	387, 6
RNA-seq expression	468, 60483	387, 5390
Gene mutation	467, 17509	387, 772

Table S2. Summary and association analysis of clinicopathological variables of 387 patients in TCGA-SKCM

clinicopathological variables	Summary	P-Value	
		RFS status	OS status
Age at diagnosis	57.78/387 (15.45)	0.001	0.876
Male gender	241/387 (62.2%)	0.336	0.268
Height cm at diagnosis	170.31/208 (9.42)	0.043	0.029
Weight kg at diagnosis	81.61/211 (19.50)	0.023	0.345
Race White	368/380 (96.8%)	0.278	0.999
History other malignancy	16/387 (4.1%)	0.211	0.999
History neoadjuvant treatment	21/387 (5.3%)	0.272	0.064
With tumor	208/387 (53.7%)	0.044	0.074
Primary location	387	< 0.001	0.787
Regional Lymph Node	186/387 (48.0%)		
Primary Tumor	87/387 (22.4%)		
Regional Cutaneous or Subcutaneous Tissue	62/387 (16.0%)		
Distant Metastasis	52/387 (13.4%)		
Breslow thickness at diagnosis	5.62/312 (8.92)	< 0.001	0.956
Clark level at diagnosis (I, II, III)¶	86/280 (43.65%)	0.493	0.644
Primary melanoma tumor ulceration	143/271 (52.76%)	< 0.001	0.999
Primary melanoma mitotic rate	6.37/159 (6.74)	0.359	0.039
AJCC pathologic tumor stage (III, IV)†	175/387 (45.2%)	0.115	0.734
AJCC nodes pathologic stage (NO, NX)§	229/387 (59.2%)	0.083	0.999
Pharmaceutical adjuvant	79/357 (22.1%)	0.523	0.382
Radiation treatment adjuvant	73/361 (20.1%)	0.048	0.999

¶The Clark level at diagnosis included 5 categories: I, II, III, IV, V. We separated I, II, III into one group, IV and V into the other group. †The AJCC pathologic tumor stage included stage 0, I, IA, IB, I/II NOS, II, IIA, IIB, IIC, III, IIIA, IIIB, IIIC, IV. We merged stage 0~II (include sub-stages) as one group, stage III and IV as the other group. §The AJCC nodes pathologic stage included NO, X, 1, 1a, 1b, 2, 2a, 2b, 2c, 3. We combined NO and NX to a group, other categories to the other group.

Prognosis of cutaneous melanoma

Table S3. Summary of the names and prediction accuracy of RFS and OS

Feature class	Detail set	C-index (mean/std.)	
		OS	RFS
STF (Single-type features)	Clinical variables (cln)	0.697/0.033	0.629/0.041
	RNA expression before feature preprocessing (ori.expr)	0.620/0.067	0.549/0.070
	Gene mutation before feature preprocessing (ori.mut)	0.526/0.030	0.498/0.070
	RNA expression after feature selection (sel.expr)	0.639/0.033	0.578/0.071
	Gene mutation after feature selection (sel.mut)	0.543/0.051	0.528/0.051
ITF (Integrate-type features)	sel.expr and cln (expr_cln)	0.709/0.053	0.631/0.052
	sel.mut and cln (mut_cln)	0.666/0.039	0.593/0.062
	sel.expr and sel.mut (expr_mut)	0.618/0.042	0.549/0.041
	sel.expr, sel.mut and cln (expr_mut_cln)	0.694/0.042	0.595/0.069

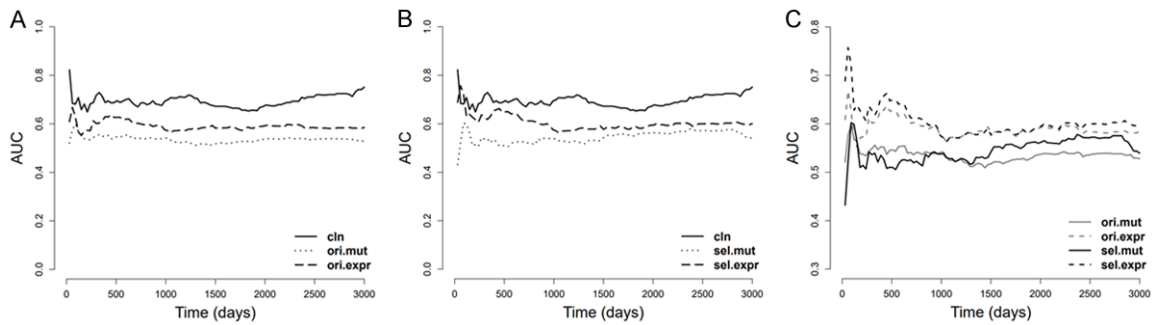


Figure S2. The prediction performance of models based on single type of data. A. Shows the models based on clinicopathological variables, RNA expression or gene mutation without performing feature selection; B. Represents those after feature selection; C. Compares the RNA expressional and gene mutational profiles before and after feature selection. Abbreviations: cln, clinicopathological variables-based model; ori.mut, models based on gene mutational profiles before feature selection; ori.expr, models based on RNA expressional profiles before feature selection; sel.mut, models based on gene mutational profiles after feature selection; sel.expr, models based on RNA expressional profiles after feature selection.

Table S4. The features selected by the competing risks model and their coefficients

Feature name	coefficient
SLC5A3	-0.1136
MRPS6	-0.0176
HOXC10	-0.11377
NXT2	-0.05
UBE2L6	-0.11651
TRAF1	-0.12959
CSGALNACT1	-0.02233
FMNL2	-0.04611
IFITM3	-0.02595
KIT	0.028093
age at diagnosis	0.057223
primary location (distant metastasis)	-0.21144
AJCC pathologic tumor stage	0.622113

Prognosis of cutaneous melanoma

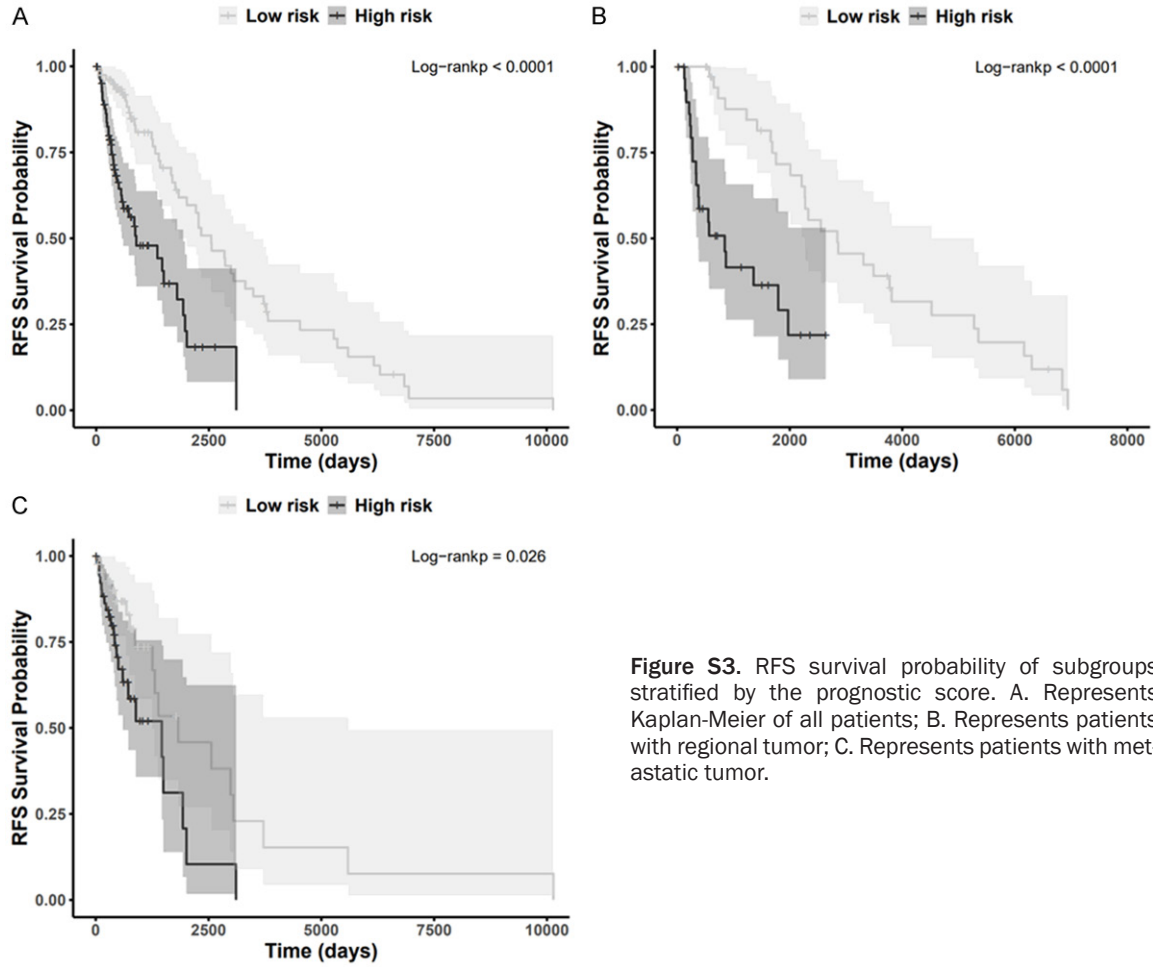


Figure S3. RFS survival probability of subgroups stratified by the prognostic score. A. Represents Kaplan-Meier of all patients; B. Represents patients with regional tumor; C. Represents patients with metastatic tumor.

Table S5. Summary of prediction accuracy (mean C-index/standard deviation) on subgroups stratified by AJCC tumor pathologic stage

Feature set	OS		RFS	
	Stage 0~II	Stage III, IV	Stage 0~II	Stage III, IV
ori.expr	0.591/0.047	0.607/0.032	0.613/0.108	0.516/0.071
sel.expr	0.601/0.037	0.616/0.039	0.625/0.083	0.601/0.074
ori.mut	0.549/0.045	0.497/0.037	0.532/0.024	0.521/0.021
sel.mut	0.547/0.037	0.478/0.036	0.498/0.065	0.547/0.053
cln	0.713/0.056	0.617/0.053	0.587/0.086	0.447/0.061
expr_cln	0.683/0.060	0.645/0.049	0.594/0.070	0.600/0.072
mut_cln	0.681/0.064	0.535/0.050	0.521/0.078	0.543/0.055
expr_mut	0.583/0.058	0.589/0.046	0.548/0.105	0.565/0.068
expr_mut_cln	0.663/0.060	0.609/0.056	0.543/0.115	0.537/0.052