

Machine Learning-Based Prediction of Survival Prognosis in Cervical Cancer

Dongyan Ding

Chongqing University

Tingyuan Lang

Chongqing University

Dongling Zou

Chongqing University

Jiawei Tan

Changchun University of technology

Jia Chen

Changchun University of technology

Dong Wang

Chongqing University

Rong Li

Chongqing University

Yunzhe Li

Chongqing University

Jingshu Liu

Chongqing University

Cui Ma

Jilin University

Qi Zhou (✉ cqzl_zq@163.com)

Chongqing University

Lei Zhou

Singapore Eye Research Institute

Research Article

Keywords: Cervical cancer, microRNAs, machine learning, survival prediction.

Posted Date: December 31st, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-134659/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Abstract

Background: Accurately forecasting the prognosis could improve therapeutic management of cancer patients, however, the currently used clinical features are difficult to provide enough information. The purpose of this study is to develop a survival prediction model for cervical cancer patients with big data and machine learning algorithms.

Results: The cancer genome atlas cervical cancer data, including the expression of 1046 microRNAs and the clinical information of 309 cervical and endocervical cancer and 3 control samples, were downloaded. Missing values and outliers imputation, samples normalization, log transformation and features scaling were performed for preprocessing and 3 control, 2 metastatic samples and 707 microRNAs with missing values $\geq 20\%$ were excluded. By Cox Proportional-Hazards analysis, 55 prognosis-related microRNAs (20 positively and 35 negatively correlated with survival) were identified. K-means clustering analysis showed that the cervical cancer samples can be separated into two and three subgroups with top 20 identified survival-related microRNAs for best stratification. By Support Vector Machine algorithm, two prediction models were developed which can segment the patients into two and three groups with different survival rate, respectively. The models exhibit high performance : for two classes, Area under the curve = 0.976 (training set), 0.972 (test set), 0.974 (whole data set); for three classes, AUC = 0.983, 0.996 and 0.991 (group1, 2 and 3 in training set), 0.955, 0.989 and 0.991 (group 1, 2 and 3 in test set), 0.974, 0.993 and 0.991 (group 1, 2 and 3 in whole data set) .

Conclusion: The survival prediction models for cervical cancer were developed. The patients with very low survival rate ($\leq 40\%$) can be separated by the three classes prediction model first. The rest patients can be identified by the two classes prediction model as high survival rate ($\approx 75\%$) and low survival rate ($\approx 50\%$).

Background

Cervical cancer is one of the leading causes of cancer-related women deaths worldwide which accounts for more than 520,000 new cases and 260,000 deaths each year.¹ Although vaccines against the most important carcinogenic human papilloma virus (HPV) types, HPV16 and HPV 18 for example, are available commercially, the number of women receiving the vaccine is still low, especially in developing countries.² Furthermore, despite effective treatment of early cervical cancer with surgery and radiation therapy, it will be uncontrollable when the cancer cells become metastatic.³

Survival prediction is important for both health care professionals and patients; accurately forecasting the prognosis could improve therapeutic management, including treatment decision and recommendations.⁴ Furthermore, accurate prediction of the survival would be useful for prevention of disease progression after first diagnosis and treatment.⁵ On the other hand, patients and their families can also set appropriate goals base on an accurate prediction.⁶ Inaccurate predictions may lead to worse treatment decision, such as over-treatment or late palliative care.^{4,5}

However, accurate prediction of the survival of cervical cancer patients is difficult as the currently used clinical characteristics, such as clinical stage, could not provide enough information involved in the mechanistic details of the tumor cells.^{6,7} The molecular information of the cervical cancer cells, such as expression level of gene, microRNA (miRNA) and long non-coding RNA (lncRNA) and the status of gene mutation and methylation, has not been well analyzed for survival prediction of cervical cancer yet. Meanwhile, the rapidly developing deep learning and machine learning technologies have been widely used for analysis of the medical data and discriminative model development,⁸⁻¹⁰ which raises the probabilities of development of survival prediction model for cervical cancer with big molecular data, The cancer genome atlas (TCGA) data for example.¹¹

Thus, the objective of this study is to develop a survival prediction model for cervical cancer patients with big data and machine learning algorithms.

Results

Workflow and data preprocessing

The workflow of the study was described in Fig. 1a. The TCGA data of 309 cervical and endocervical cancer samples and 3 control samples were downloaded by Firehose online tool, which includes the expression data of 1046 miRNAs and clinical information. According to experimental design, 3 control and 2 metastatic samples and 707 miRNAs with missing values $\geq 20\%$ were excluded. The rest missing values and identified outliers were replaced by K-Nearest Neighbor (KNN) method. The data was then normalized and scaled as described in method section. The preprocessed miRNA expression data in 307 primary cervical cancer samples was given in Supplementary Table 1 and the expression profiles were presented as heatmap (Supplementary Fig. 1). The preprocessed data were subjected to further analysis.

Association of clinical features and the survival rate of cervical cancer patients

To assess the predictive value of clinical features for survival in cervical cancer, the clinical characteristics were summarized in Supplementary Table 2. The results from Chi-Square analysis showed that only clinical stage is significantly associated with the 3-years survival ($P = 0.040$) (Figure. 1b). The subsequent Kaplan-meier analysis showed that patients in clinical stage I-IIa and IIb-IVb have statistically different 3-year survival rate with Log-rank P value of 0.020 (Fig. 1c). While, as mentioned above, a better prediction model is needed, we thus intend to achieve this goal by machine learning and molecular characteristics.

Survival-related miRNA features selection

To cluster the cervical cancer patients with different prognosis, the survival-related features were identified; miRNAs were used in this study as its important roles in regulating signaling networks involved in tumorigenesis.¹² The survival-related miRNAs were identified by Cox Proportional-Hazards (Cox-PH)

regression. In total, 55 prognosis-related miRNAs (20 positively and 35 negatively correlated with survival) (with $P < 0.05$) were identified (Supplementary Table 3). To verify these prognosis-related miRNAs, the average expression levels of prognosis-related miRNAs in samples grouped base on 3-years survival were shown in Fig. 2a and 2b. We found that all identified miRNAs were differently expressed in the two groups, and interestingly, tumor samples collected from patients who loss to follow up < 3 years have the similar miRNA expression profiles with samples of patients with high 3-years survival rate. These data demonstrated that the survival-related miRNAs of cervical cancer samples were correctly identified.

Stratification of the patients by survival-related features

K-means clustering algorithm was used to stratify the patients. We tried all 55, top 30 and top 20 survival-related miRNAs as features and the parameter K was set to 2 to 6. The clustered patients (Supplementary Table 4) were subsequently subjected to Kaplan-meier analysis. We found that clustering by top 20 survival-related miRNAs obtained the best stratification (Fig. 3) and the patients could be classified into 2 and 3 groups ($P = 0.0092$ for $K = 2$, $P = 0.0045$ for $K = 3$, Log-rank test). We therefore determined these parameters for Support Vector Machine (SVM) model development.

SVM model development and validation

We next used the labeled patients and top 20 survival miRNA features to develop classification model by SVM algorithm and 70%/30% split was used. Receiver operating characteristics (ROC) curve and Area under the curve (AUC) were employed to evaluate the performance of the classification model. Kaplan-meier analysis with Log-rank p value was used to assess the accuracy of survival subtype prediction; to avoid the sampling error, all samples were used. The whole strategy for model development and verification is provided in Fig. 4a. As shown Fig. 4b and 4c, all training, test and whole dataset generated high AUC values (for two classes: AUC = 0.976 (training set), 0.972 (test set), 0.974 (whole data set); for three classes: AUC = 0.983, 0.996 and 0.991 (group1, 2 and 3 in training set), 0.955, 0.989 and 0.991 (group 1, 2 and 3 in test set), 0.974, 0.993 and 0.991 (group 1, 2 and 3 in whole data set)), and significant Log-rank p values (for two classes: Log-rank $P = 0.025$; for three classes: Log-rank $P = 0.00078$ (Group 1 vs 2 $P \leq 0.001$, Group 1 vs 3 $P = 0.015$, Group 2 vs 3 $P = 0.074$)). These results demonstrated that the two SVM models are robust to predict survival-specific clusters. Although three classes prediction model failed to separate group 2 and 3, it can be used to identify the patients with very low survival rate (3-years survival rate $\leq 40\%$) accurately. The rest patients can be further divided into two groups (high survival rate (3-years survival rate $\approx 75\%$) and low survival rate (3-years survival rate $\approx 50\%$)) with the two classes prediction model.

Bioinformatic analysis of features used in prediction model

Finally, the pathways analysis of target genes predicted with top 20 survival-related miRNAs showed that pathways involved in presynaptic depolarization and calcium channel opening, PPAR α , cell-extracellular matrix interactions, L1 and Ankyrins, apoptosis, FOXO-mediated transcription of cell death genes, NOTCH1, MECP2, type I IFN production and TP53 degradation, etc. were primarily impacted (Supplementary Table 5 and Figure. 5).

Discussion

In this study, the survival prediction models for cervical cancer were developed. The survival-related miRNAs were used as the features. The patients can be divided into three groups: high survival rate ($\approx 75\%$), low survival rate ($\approx 50\%$) and very low survival rate ($\leq 40\%$).

Cervical cancer is still one of the leading causes of cancer-related women deaths worldwide.¹⁻³ Despite the development of HPV vaccine, this disease remains high incidence rate, especially in developing countries.^{1,2} The metastatic cervical cancer is incurable, and thus accurately forecasting the survival is important for both health care professionals and patients to improve the therapeutic management and satisfaction of patients and their families. To achieve this goal, the big data and the machine learning technologies were employed in this study.

MiRNAs are small endogenous non-coding RNAs that regulates gene expression.¹³⁻¹⁵ Evidence have demonstrated that the miRNAs plays important roles in human cancer cells, including activating proliferation, invasion, metastasis and angiogenesis and suppressing tumor suppressors.^{16,17} Accumulating studies have shown that miRNAs can be served as biomarkers and therapeutic targets.^{18,19} Thus, given the important roles of miRNAs in cancer cell regulation, we used miRNAs as the features for discriminative model development.

In this study, we employed the results of K-means clustering analysis to determine the number of survival-related features that would be used for SVM model development. We found that the number of the features would be carefully chosen according to the discriminative ability of the features and the sample size of the study. In this case, the top 20 survival-related miRNAs obtained the best stratification (Fig. 3) and the patients could be classified into 2 and 3 groups ($P = 0.0092$ for $K = 2$, $P = 0.0045$ for $K = 3$, Log-rank test).

Most of survival prediction models developed so far separate the patients into two groups (high and low survival rate).²⁰⁻²² However, by combination of the stratification models, the patients can be identified more accurately. In our study, we developed two prediction models that divide the patients into two and three groups, respectively (Figs. 3 and 4). Although three classes prediction model failed to separate group 2 and 3, it can be used to identify the patients with very low survival rate (3-years survival rate $\leq 40\%$) accurately first. Then, the rest patients can be further identified as patients with high survival rate (3-years survival rate $\approx 75\%$) and low survival rate (3-years survival rate $\approx 50\%$).

Conclusion

In summary, the survival prediction models for cervical cancer were developed in this study. The patients with very low survival rate ($\leq 40\%$) can be separated by the three classes prediction model first. The rest patients can be identified by the two classes prediction model as high survival rate ($\approx 75\%$) and low survival rate ($\approx 50\%$).

Material And Methods

Datasets and workflow

We used the TCGA data to develop the prognosis prediction program. The TCGA data of 309 cervical and endocervical cancer samples and 3 control samples were downloaded by Firehose online tools, which includes the expression data of 1046 miRNAs (file name: `gdac.broadinstitute.org_CESC.miRseq_Preprocess.Level_3.2016012800.0.0`) and clinical information (file name: `gdac.broadinstitute.org_CESC.Merge_Clinical.Level_1.2016012800.0.0`). The survival-related miRNAs identified by Cox-PH model were used for clustering the samples. The labeled samples and related survival-related miRNA features were used for SVM model development. The details of workflow were presented in Fig. 1a.

Survival-related clinical features identification

Frequency distribution table containing selected clinical features was made according to the standard protocol. The patients were stratified into survival ≤ 3 years and survival > 3 years. R command (`chisq.test()`) was used to implement Chi-square test and $P < 0.05$ was regarded as statistically significant.

Data preprocessing

According to experimental design (Fig. 1a), 3 control samples, 2 metastasis samples and 707 miRNAs with missing values $\geq 20\%$ were excluded. MetaboAnalyst 4.0 online software was used for missing values replacement, outliers identification, samples normalization and features scaling. Briefly, the missing values were replaced by KNN method. The outliers were identified by Interquartile Range (IQR) method. Median normalization was used for sample normalization. Autoscaling algorithm was used for features scaling. The expression profiles of preprocessed miRNAs in cervical cancer patients were presented as heatmap.

Survival-related miRNAs identification

Cox-PH model was used to identification of survival-related miRNAs. R package “survival” was used to calculate of p value for each miRNA feature. The expression profiles of survival-related miRNAs in patients were presented as heatmap.

K-means clustering

R command, `kmeans()` was used to perform stratify the patients according to survival. The expression data of top 20, top 30 or total survival-related miRNAs were input into K-means program. The K parameter was set to 2 to 6.

Kaplan-meier analysis

Kaplan-meier analysis was used to calculate and plot the survival rate of stratified patients. `Survfit()` command in Survival R package and `ggsurvplot()` command in Survminer R package were used to plot Kaplan-meier curve. Log-rank p value for each analysis was given.

Supervised classification model development

The supervised classification model was developed by SVM algorithm with labeled samples and the expression data of top 20 survival-related miRNAs. The SVM model was developed by splitting the samples 70%/30% to training and held-out testing data. The R command `svm()` in e1071 R package was used to SVM algorithm implementation.

ROC curve

The R command `roc()` in pROC R package was used to plot ROC curve.

Bioinformatic analysis

The mature miRNAs of survival-related stem-loop miRNAs were checked by miRbase online database. The predicted targets of mature miRNAs derived from top 20 survival-related stem-loop miRNAs were analyzed by miRDB online tool. Only predicted targets with target score ≥ 90 were recorded and the top 20 targets were included for bioinformatic analysis. Pathway analysis was performed by reactome online software.

Abbreviations

HPV

human papilloma virus; miRNA:microRNA; lncRNA:long non-coding RNA; TCGA:The cancer genome atlas; KNN:K-Nearest Neighbor; Cox-PH:Cox Proportional-Hazards; SVM:Support Vector Machine; ROC:Receiver operating characteristics; AUC:Area under the curve; IQR:Interquartile Range

Declarations

Acknowledgments

This study was funded by Chongqing Science & Technology Commission cstc2019jscx-msxmX0174 and cstc2019jscx-msxmX0106. We thank Professor Guangwu Tang from Chongqing research and Design Institute of Transportation for the idea of artificial intelligence. We thank Mr. Tingxiu Lang and Jin Chen from Shenzhen Forms Syntron Information Co., Ltd for the help in computer technology. We thank Professor Qing Zeng from Department of Public Health of Chongqing Medical University for statistics idea.

Authors' contributions

DYD, TYL performed all the experiments. DLZ, JT interpreted the data. JC, LZ participated in statistical method development, DW, RL, YZL, JSL, CM performed data analysis involved in clinical information. TYL

and DYD wrote the manuscript. TYL and QZ designed and supervised the study.

Funding

The study is funded by Chongqing Science & Technology Commission cstc2019jscx-msxmX0174 and cstc2019jscx-msxmX0106. The funder play no role in study design and excutation.

Availability of data and materials

All data in this study are included in this published article and the data analysis codes are deposited in GitHub (<https://github.com/dingdongyan/CESC>).

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable

Competing interests

The author reports no conflicts of interest in this work.

References

1. Graham Cancer Center & Research Institute at Christiana Care Health Services; HudsonAlpha Institute for Biotechnology; ILSbio, LLC; Indiana University School of Medicine; Institute of Human Virology; Institute for Systems Biology; et al; Integrated genomic and molecular characterization of cervical cancer. *Nature*. 2017;543(7645):378–384.
2. Roden RBS, Stern PL. Opportunities and challenges for human papillomavirus vaccination in cancer. *Nat Rev Cancer*. 2018;18(4):240–254.
3. Tyagi A, Vishnoi K, Mahata S, Verma G, Srivastava Y, Masaldan S, et al. Cervical Cancer Stem Cells Selectively Overexpress HPV Oncoprotein E6 that Controls Stemness and Self-Renewal through Upregulation of HES1. *Clin Cancer Res*. 2016;22(16):4170–84.
4. White N, Reid F, Harris A, Harries P, Stone P. A Systematic Review of Predictions of Survival in Palliative Care: How Accurate Are Clinicians and Who Are the Experts? *PLoS One*. 2016;11(8):e0161407.
5. van Amsterdam WAC, Verhoeff JJC, de Jong PA, Leiner T, Eijkemans MJC. Eliminating biasing signals in lung cancer images for prognosis predictions with deep learning. *NPJ Digit Med*. 2019;2(1):122.
6. Yousefi S, Amrollahi F, Amgad M, Dong C, Lewis JE, Song C, et al. Predicting clinical outcomes from large scale cancer genomic profiles with deep survival models. *Sci Rep*. 2017;7(1):11707.

7. Ramazzotti D, Lal A, Wang B, Batzoglou S, Sidow A. Multi-omic tumor data reveal diversity of molecular mechanisms that correlate with survival. *Nat Commun.* 2018;9(1):4453.
8. Wang Q. Cancer predisposition genes: molecular mechanisms and clinical impact on personalized cancer care: examples of Lynch and HBOC syndromes. *Acta Pharmacol Sin.* 2016;37(2):143–9.
9. Bello GA, Dawes TJW, Duan J, Biffi C, de Marvao A, Howard LSGE, et al. Deep learning cardiac motion analysis for human survival prediction. *Nat Mach Intell.* 2019;1:95–104.
10. Lee B, Chun SH, Hong JH, Woo IS, Kim S, Jeong JW, et al. DeepBTS: Prediction of Recurrence-free Survival of Non-small Cell Lung Cancer Using a Time-binned Deep Neural Network. *Sci Rep.* 2020;10(1):1952.
11. Carmona-Bayonas A, Jiménez-Fonseca P, Font C, Fenoy F, Otero R, Beato C, et al. Predicting serious complications in patients with cancer and pulmonary embolism using decision tree modelling: the EPIPHANY Index. *Br J Cancer.* 2017;116(8):994–1001.
12. ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. Pan-cancer analysis of whole genomes. *Nature.* 2020;578(7793):82–93.
13. Peng Y, Croce CM. The role of MicroRNAs in human cancer. *Signal Transduct Target Ther.* 2016;1:15004.
14. Gebert LFR, MacRae IJ. Regulation of microRNA function in animals. *Nat Rev Mol Cell Biol.* 2019;20(1):21–37.
15. Rupaimoole R, Slack FJ. MicroRNA therapeutics: towards a new era for the management of cancer and other diseases. *Nat Rev Drug Discov.* 2017;16(3):203–222.
16. Konno M, Koseki J, Asai A, Yamagata A, Shimamura T, Motooka D, et al. Distinct methylation levels of mature microRNAs in gastrointestinal cancers. *Nat Commun.* 2019;10(1):3888.
17. Wessels HH, Lebedeva S, Hirsekorn A, Wurmus R, Akalin A, Mukherjee N, et al. Global identification of functional microRNA-mRNA interactions in *Drosophila*. *Nat Commun.* 2019;10(1):1626.
18. Ghini F, Rubolino C, Climent M, Simeone I, Marzi MJ, Nicassio F. Endogenous transcripts control miRNA levels and activity in mammalian cells by target-directed miRNA degradation. *Nat Commun.* 2018;9(1):3119.
19. Dhawan A, Scott JG, Harris AL, Buffa FM. Pan-cancer characterisation of microRNA across cancer hallmarks reveals microRNA-mediated downregulation of tumour suppressors. *Nat Commun.* 2018;9(1):5228.
20. Asakura K, Kadota T, Matsuzaki J, Yoshida Y, Yamamoto Y, Nakagawa K, et al. A miRNA-based diagnostic model predicts resectable lung cancer in humans with high accuracy. *Commun Biol.* 2020;3(1):134.
21. Chaudhary K, Poirion OB, Lu L, Garmire LX. Deep Learning-Based Multi-Omics Integration Robustly Predicts Survival in Liver Cancer. *Clin Cancer Res.* 2018;24(6):1248–1259.
22. Huang Z, Johnson TS, Han Z, Helm B, Cao S, Zhang C, et al. Deep learning-based cancer survival prognosis from RNA-seq data: approaches and evaluations. *BMC Med Genomics.* 2020;13(Suppl

5):41.

23. Jiang D, Liao J, Duan H, Wu Q, Owen G, Shu C, et al. A machine learning-based prognostic predictor for stage III colon cancer. *Sci Rep.* 2020;10(1):10333.

Figures

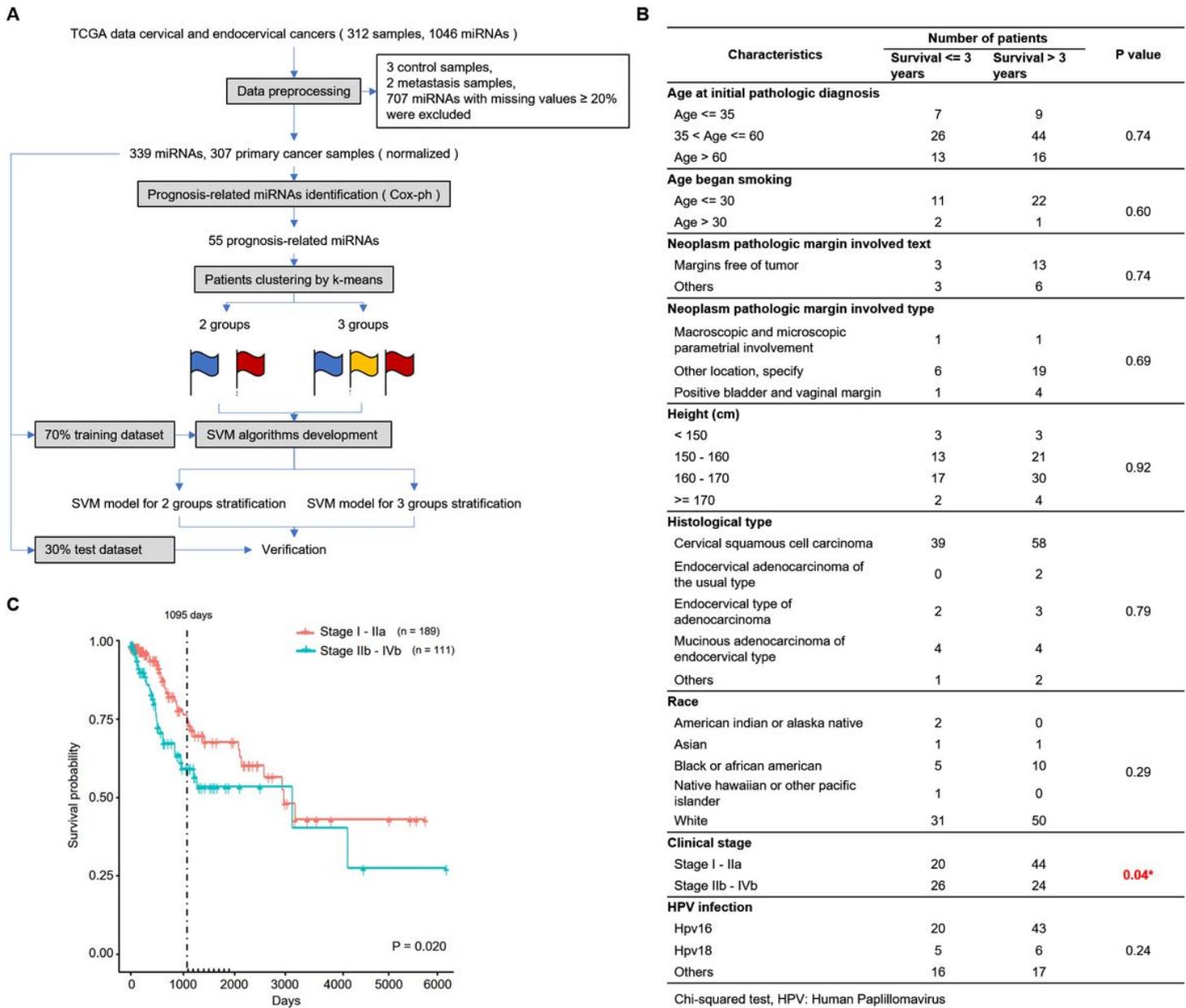


Figure 1

Workflow of the study and association of clinical characteristics and survival rate of cervical cancer patients. (A) Workflow of cervical cancer prognosis prediction model development. The TCGA (The Cancer Genome Atlas) cervical cancer miRNA expression data was preprocess and used for patients clustering by K-means algorithm; the labeled patients were then used for developing prediction model development by SVM program. (B) The association between clinical characteristics and 3-years survival

was analyzed by Chi-Square analysis ($P < 0.05$ means statistical significance). (C) The overall survival of cervical cancer patients in stage I-IIa and stage IIb-IVb was analyzed by Kaplan-meier analysis. $P < 0.05$ means significant difference.

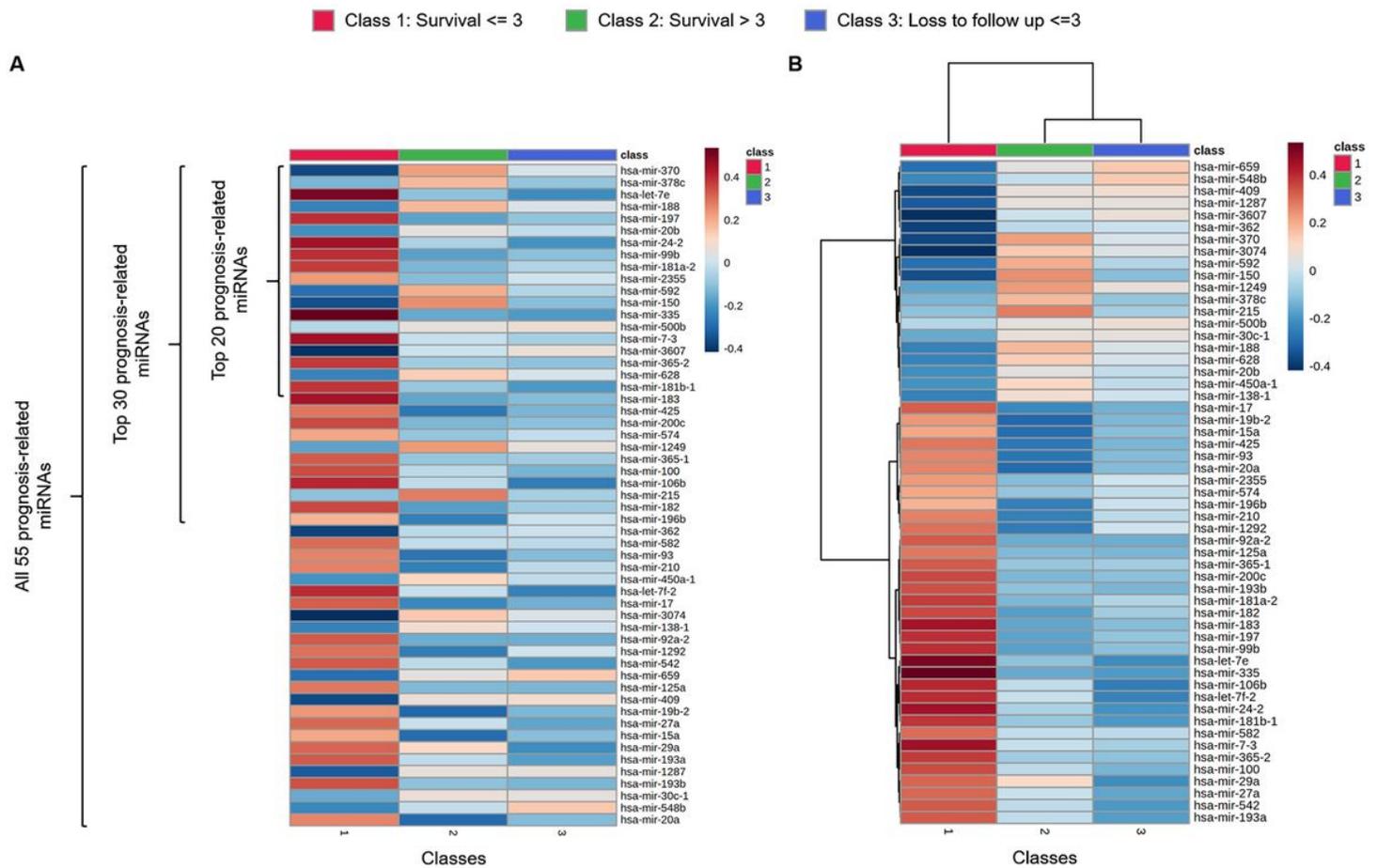


Figure 2

The expression of identified prognosis-related miRNAs in cervical cancer samples. (A) The average expression level of each prognosis-related miRNAs identified by Cox-PH regression analysis in cervical cancer samples of patients who survival ≤ 3 years (red), survival > 3 years (green), loss to follow up ≤ 3 years (blue) was presented as heatmap. (B) The prognosis-related miRNAs were clustered base on their expression in cervical cancer patients; 20 miRNAs are positively correlated with survival ability and 35 miRNAs are negatively correlated with survival ability of cervical cancer patients. Tumor samples of patients who loss to follow up < 3 years have the similar miRNA expression profiles with samples of patients with high 3-years survival rate.

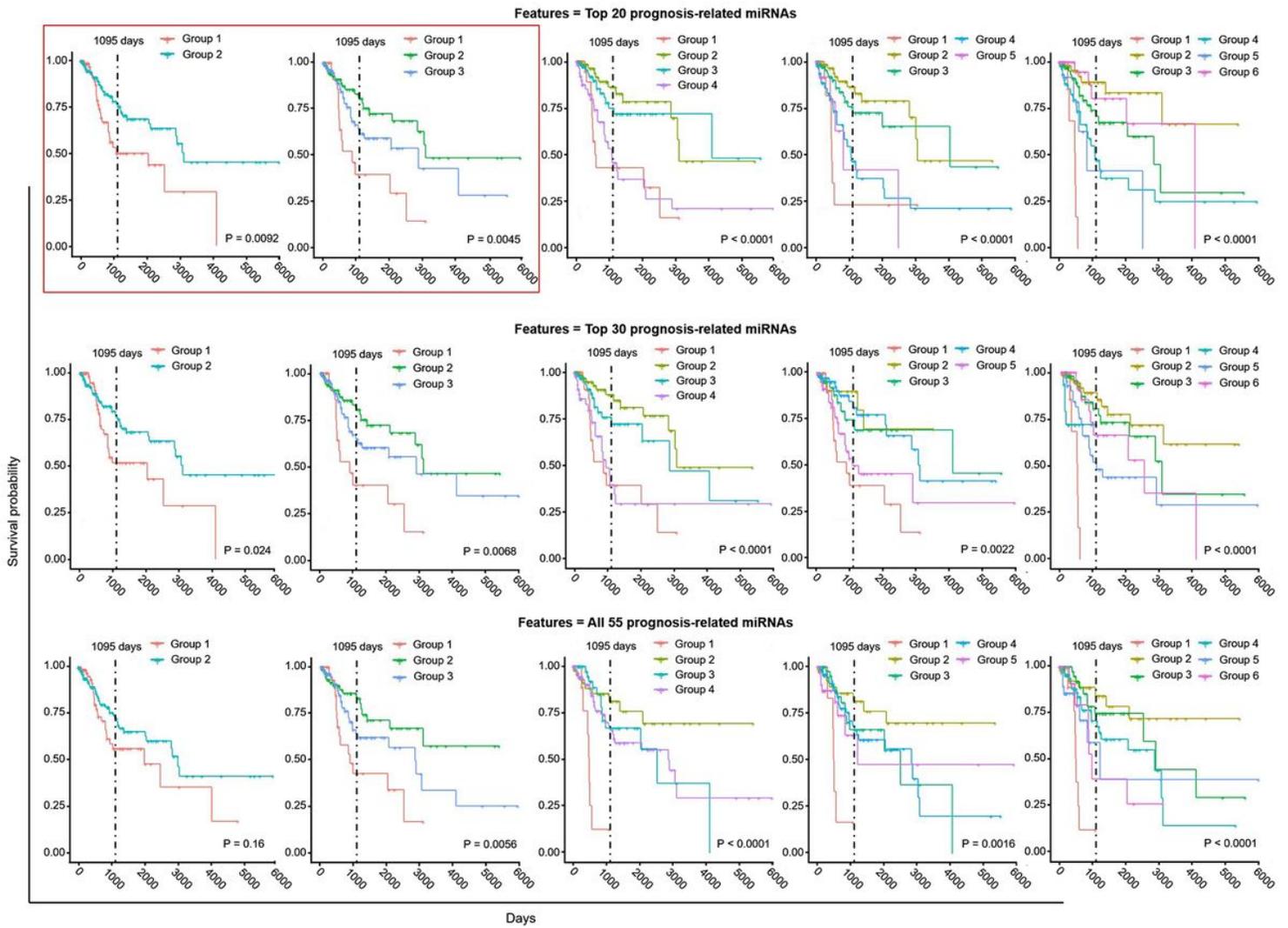


Figure 3

Kaplan-meier analysis of the survival of cervical cancer patients stratified by K-means clustering. The cervical cancer samples were clustered by K-means clustering algorithm (K = 2-6) with top 20 (top), top 30 (middle) and all 55 (bottom) identified prognosis-related miRNAs, followed by Kaplan-meier analysis of the survival of corresponding patients. Only when K=2 and 3, the patients were well stratified with statistically significant difference in survival rate ($P < 0.05$, Log-rank test).

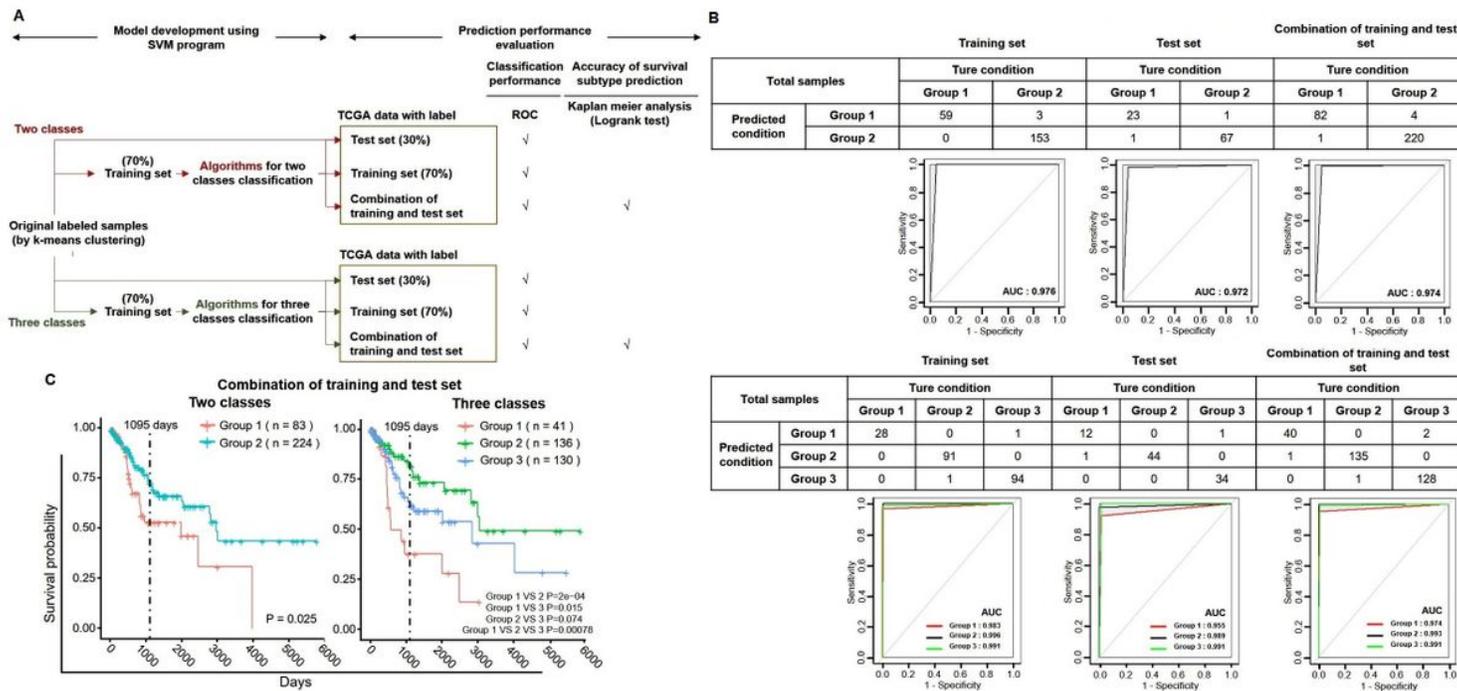


Figure 4

Validation of survival prediction models. (A) Workflow of SVM model development and verification. The SVM classification model was developed by splitting the samples 70%/30% to training and held-out testing data. The performance of model was verified by test, training and all samples dataset by calculating AUC (Area under the curve) and Log-rank P value. (B) The performance of classification for the prediction model was evaluated by ROC (Receiver operating characteristics) curve and AUC value. (C) The accuracy of survival subtype prediction was assessed by Kaplan-meier analysis and Log-rank P value.

Pahtway analysis

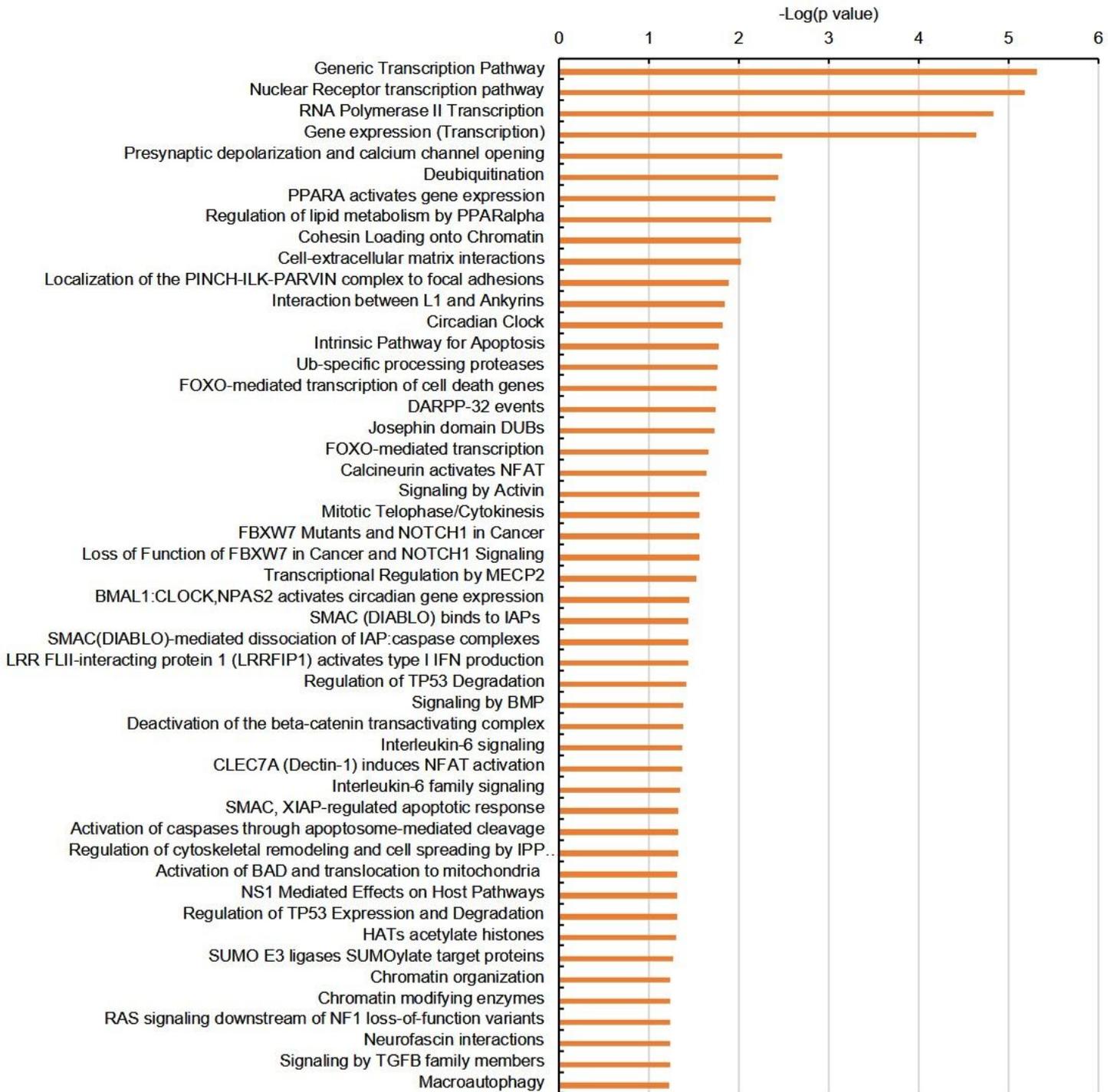


Figure 5

Pathway analysis of the predicted target genes of top 20 survival-related miRNAs. The mature miRNAs of survival-related stem-loop miRNAs were checked by miRbase online database (<http://www.mirbase.org/>). The predicted targets of mature miRNAs derived from top 20 survival-related stem-loop miRNAs were analyzed by miRDB online tool (<http://mirdb.org/index.html>). Only predicted targets with target score \geq

90 were recorded and the top 20 targets were included for bioinformatic analysis. Pathway analysis was performed by reactome online software (<https://reactome.org/>).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [dingaicc20201216SI.docx](#)
- [SupplementaryTable1.PreprocessedmiRNAdataforprimarycervicalcancersamples.csv](#)
- [SupplementaryTable3.SurvivalrelatedmiRNAsidentifiedbyCoxPHanalysis.csv](#)
- [SupplementaryTable4.Kmeansclusteringofcervicalcancerpatients.xlsx](#)
- [SupplementaryTable5.Pathwayanalysisofthepredictedtargetsof20survivalrelatedmiRNAs.csv](#)