# Machine learning enables accurate prediction of breast cancer five-year survival using somatic genomic variants

Xiaosen Jiang[1,2], Laizhi Zhang[3], Guangshuo Cao[4], Jia Li[5,6], Yong Bai[2*]

1 College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049, China.

2 BGI-Shenzhen, Shenzhen, China.

3 School of Basic Medicine, QINHGDAO University.

4 State Key Laboratory of Public Big Data, College of Computer Science and Technology, Guizhou University, Guiyang, 550025, China.

5 BGI Genomics, BGI-Shenzhen, Shenzhen, 518083, China.

6 Hebei Industrial Technology Research Institute of Genomics in Maternal & Child Health, Shijiazhuang BGI Genomics Co., Ltd, Shijiazhuang 050000, Hebei Province, China.

* Correspondence:

Yong Bai

baiyong@genomics.cn

## Abstract

Breast cancer is one of the most common cancers, accounting for about 30% of female cancers and a mortality rate of 15%. The 5-year survival rate is most commonly used to assess cancer progression and guide clinical practice. We used the CatBoost model to systematically construct a five-year mortality risk prediction model based on two independent data sets (BRCA_METABRIC, BRCA_TCGA). The model input data are the somatic genomic variants (copy number variation, SNP locus, cumulative mutation number of genes) and phenotype data of cancer samples. The optimal model combined all the above characteristics, and the AUC reached 0.70 in an independent external data set. At the same time, we also conducted a biological analysis of the characteristics of the model and found some potential biomarkers

29    (TP53, DNAH11, MAP3K1, PHF20L1, etc.). The results of model risk stratification

30    can be used as a guide for the prognosis of breast cancer.

## Introduction

32        Breast cancer has overtaken lung cancer to become the most common malignancy

33    and the first leading lethal cancer among women, accounting for approximately 11.7%

34    of all cancer cases diagnosed (2.3 million) in 2020 around the world[1]. In spite of

35    significant medical improvements in early diagnosis and modern therapy [2], breast

36    cancer still poses increasing burden globally and the overall survival outcomes remain

37    unsatisfactory given the mortality of 1 in 6 female cancers (685,000 deaths)[1],

38    reflecting high biological complexity and genetic heterogeneity of the disease[3-5].

39    Consequently, identification of genetic prognostic indicators plays crucial roles in

40    understanding inter-individual differences in pathogenesis between breast cancer

41    patients, providing better insight into therapeutic decision-making and optimizing

42    personalized precise treatment. Additionally, the prognosis of breast cancer, which has

43    a 5-year survival rate greater than 85%, is better than other cancers[6]. Since the

44    threshold of five years is most commonly used to assess the process of the cancer, we

45    can anticipate the 5-year mortality of breast cancer to guide clinical intervention.

46        One of the earliest survival prognostic models, Nottingham Prognostic

47    Index(NPI)[7], was constructed based on clinical factors using Cox regression[8].

48    Since then, more variables had been considered to improve the accuracy of NPI[9]. In

49    2001, Adjuvant, an internet-based tools was developed and widely applied in

50    prognosis analysis in breast cancer[10]. In 2010, two independent models,

51  OPTIONS[11] based on parametric regression and PREDICT based on Cox

52  proportional hazards, were established. But both Adjuvant and PREDICT showed low

53  accuracy of estimating mortality in different sub-groups, especially for young breast

54  cancer patients[12]. Additionally, these statistically constructed models based on

55  clinical data are limited by the prolonged process of data collecting and poor

56  timeliness. These drawbacks entail the need for prediction using data collected right

57  after diagnosis. With the advent of microarray-based gene expression profiling, some

58  gene-related studies demonstrate the impact of genetic factors on prognostic and

59  survival prediction of breast cancer and lots of predictive signatures have been

60  found[13]: such as MammaPrint[14], Oncotype DX[15], Endopredict[16]. Although

61  these models are applied to sub-groups of breast cancer patients, their signatures

62  cannot sensibly interpret their relationships with breast cancer outcomes, which are

63  called black-box models[17]. Thus models with high accuracy and high

64  interpretability need to be further developed.

65      Machine learning(ML) is a feasible method, since ML can extract features from

66  large genetic datasets and perform risk scoring and classification[18]. Genetic

67  signature copy number alteration(CNA) has a strong correlation with the prognosis

68  and mortality of cancer[19]. However, because the genomic CNA dataset is large and

69  relatively sparse, traditional models based on single or several CNA signatures are not

70  explicable. Furthermore, Somatic mutations can also be used to construct predictive

71  models for risk scoring and survival prediction. Nguyen et al. selected multi-features

72  with Random Forest(RF), largely improving the accuracy of the predictive model[20].

73    Support Vector Machine(SVM), Artificial Neural Network(ANN) and

74    semi-supervised learning methods were employed to construct predictive models for

75    assessing the survivability of breast cancer patients[21]. Compared to the integrated

76    model, the models based on somatic mutations alone have lower accuracies, and

77    integrated models are limited by small samples and may have overfitting problems[22,

78    23]. We found that most of the reported 5-year survival prediction models for breast

79    cancer have considered data preprocessing, feature selection, class imbalance

80    processing, and model validation.[24]. We only find two studies that were further

81    verified externally[25, 26]. Both of the studies used the Molecular Taxonomy of

82    Breast Cancer International Consortium(BRCA_METABRIC) dataset for training and

83    internal validation, and The Cancer Genome Atlas(BRCA_TCGA) dataset was used

84    for further validation. However, after scrutinizing these two studies, we found that

85    their independent dataset, which should be used for external test, was fed into their

86    models for training and internal testing again. External validation is necessary because

87    it can reflect the generalization ability of the predictive model. So far as we know,

88    there is no breast-survival-predictor that undertakes external test using independent

89    cancer datasets.

90       In our study, we developed a CatBoost-based machine learning model that

91    integrates multi-dimensional data including single-nucleotide variants(SNV),

92    cumulative number of gene mutations(CNGM) ,copy number alteration(CNA) and

93    phenotype data. BRCA_METABRIC dataset was employed for training and internal

94    validation and BRCA_TCGA dataset for external testing. The final result of the model

95    has good generalization ability, and the AUC in the external test set is 0.70. In

96    addition, the feature interpretation of the model found that the model has a high

97    learning ability, and some features that have been reported to be highly related to

98    cancer were found in the model. The code required by the model can be viewed in

99    github: https://github.com/jxs1996/Breast_cancer-5-year-survival-prediction

## Materials and methods

**Data preparation**

102    We obtained two breast cancer data sets from the public database

103    cBioPortal(BRCA_TCGA(n=1108                                    samples)[27],

104    https://www.cbioportal.org/study/summary?id=brca_tcga;

105    BRCA_METABRIC(n=2509                    samples)                    [28-30],

106    https://www.cbioportal.org/study/summary?id=brca_metabric). Overall, the median

107    age was 62□years; the 5-year survival rate was 75%. Both data sets contain clinical

108    data, somatic mutations, CNA and gene expression data. In the preliminary data

109    processing(Figure1.A), The SNV, CNGM, and CNA features in the two data sets were

110    separately counted and constructed into the input data required by the model.

111    Predictive features were expressed as follows: (1)SNV: if there is a mutation, it is

112    marked as "1", if it is not marked as "0"; CNGM: each additional SNV, the value

113    increases by 1, and each additional insertion or deletion, the value increases by 10

114    (This is because we assume that insertions and deletions have a greater accumulation

115    of mutations and a greater impact on genes than SNV); CNA: "-2" for homozygous

116    deletion, "-1" for hemizygous deletion, "0" for neutral / no change, "1" for gain, "2"

117  for high level amplification. The missing values of all features are filled with "0".

118  Gene expression data is not included in the feature representation. There are two main

119  reasons. 1) The data obtained from the two datasets are standardized with Z-score,

120  which means they are not in the same spatial dimension. Training results in one data

121  cannot be applied in another dataset; 2) If genomic data alone yields better prediction

122  results.This can significantly reduce the cost of the application. Nevertheless, gene

123  expression data will be used in subsequent analysis to observe the performance of the

124  model features at the transcriptome level. In order to ensure the consistency of

125  features, the intersection of SNV, CNGM and CNA is obtained in two independent

126  datasets. Next, we labeled samples as survival($OS\_Months > 60$) or

127  death($OS\_MONTHS < 60$ and $OS\_Status$ = deceased). Some sample data were

128  discarded($OS\_MONTHS < 60$ and $OS\_STATUS$ = LIVING ) .

129      We retained a total of CNA - 18533, SNV - 215, CNGM - 170 in both datasets.

130  BRCA_METABRIC dataset retained 1904 individuals(survival - 1432, death - 412,

131  discarded - 60), BRCA_TCGA dataset retains 513 individuals(survival-148, death-46,

132  discard-319). In addition, we screened the clinical data shared by the two

133  datasets(because we hope to establish an early risk prediction model, the data of

134  intervention treatment will not be considered), and finally obtained age, gender,

135  number of positive limph nodes and menopausal state. Most patients are female (there

136  are only three males in the BRCA_TCGA dataset), so the sample is no longer grouped

137  by gender. The statistical results of other phenotypes are shown in Supplementary

138  Figure S1. The average age of all breast cancer patients is 60.65 years, and the

139      average lymph node is 2.02. 443 are not yet menopausal, and 1548 are in menopause.

140      If there is no measurement data in the phenotype, it will be represented by -9 and will

141      not participate in the mean calculation.

142      **Machine learning analysis process development**

143      As shown in Figure 1.B, Catboost, a high-performance open source library for

144      gradient boosting on decision trees, was developed to predict the five-year mortality

145      risk of patients. The analysis process is systematically constructed using the machine

146      learning framework scikit-learn(https://scikit-learn.org/stable/, version=0.24.2). The

147      BRCA_METABRIC data(1844 samples: survival-1432, death-412) set was split into

148      training set(80%) and testing set(20%) by random stratified sampling. The

149      independent external data set BRCA_TCGA (194 samples: survival-148, death-46)

150      will be used for model evaluation. For the three features(SNV, CNA, CNGM),

151      separate models are established to evaluate their effects on prediction. After that, we

152      extract the model features constructed by SNV, CNA, and CNGM and merge them to

153      construct a new multi-dimensional feature set for training and evaluation. Finally,

154      phenotypic characteristics (age, number of positive lymph nodes, menopausal status)

155      are also integrated to further improve the accuracy of the model.

156      **Single-dimensional feature model construction**

157      Separate models were constructed for CNV, SNV, and CNGM characteristics to

158      explore their impact on the five-year mortality risk prediction. As shown in Figure 1.B,

159      For CNA (18533 features), the training set is first standardized (StandardScaler

160      method), the average value and standard deviation are retained and then applied to the

161    corresponding features in the test phase, and feature selection is performed on the

162    processed data (described in Feature selection part). Next, use the CatBoost model to

163    select the hyperparameters of the model (described in Hyperparameter selection part).

164    After fixing the hyperparameters, perform model training. Five-fold cross-validation

165    is used to evaluate the stability of the model, and finally tested on the test set and

166    independent external data set. The processing of SNV and CNGM is similar to CNA,

167    but due to the small number of SNV features (215) and onehot-encoded, data

168    standardization and feature selection are not performed; CNGM uses log first and then

169    logMinMaxScaler method during standardization. Since there are only 170 features,

170    feature selection is also omitted.

171    **Multi-dimensional feature model construction**

172    Combine the feature selection results of the single-dimensional feature model and

173    perform hyperparameter selection (described in Hyperparameter selection part). After

174    the hyperparameter results are fixed, perform training and evaluation. In addition, we

175    combined the phenotypic data (age, gender, number of positive lymph nodes, and

176    menopausal status) with the feature selection results of the single-dimensional feature

177    model to observe whether the phenotypic data can improve the performance of the

178    model.

179    **Feature selection**

180    For CNA(18533 features), irrelevant features may decrease the performance of

181    the model. We propose a hybrid feature selection method to subtract features. In this

182    method, mutual information (MI) technology[31], recursive feature elimination (RFE)

183   algorithm[32] and Boruta algorithm[33] are used to obtain the relevant subset of the

184   raw features. MI calculates feature weights based on the relationship between features

185   of mutual information; RFE selects features by recursively considering smaller and

186   smaller feature sets, and this method can obtain all feature rankings. The Boruta

187   algorithm is a packaging method that selects a subset of features based on a random

188   forest machine learning algorithm, which can be used to measure the importance of

189   features. Respectively use the above methods to obtain the feature ranking and retain

190   the top 3% features (extract the most effective features and maintain a balance with

191   the SNV and CNGM feature numbers). The features selected by any two methods will

192   be retained eventually.

**Hyperparameter selection**

194      Due to the imbalance between death and survival samples (~1:3), when 80% of

195   the training set is used for hyperparameter training, a small number of samples are

196   randomly sampled to make the ratio of positive and negative samples reach 1:1. We

197   implemented a basic grid search algorithm with 5-fold cross-validation to optimize

198   the Catboost model parameters while maximizing the weighted F1 score.

**Model comparison**

200      After the model training is completed, we will use the five-fold cross-validation

201   data set, test set and independent external data set to evaluate the model. The specific

202   evaluation indicators are as follows:

203      *TP:* True Positive. In the samples judged to be positive, the number of correct

204   judgments.

205  *FP:* False Positive. In the samples judged as positive, the number of judgment

206 errors.

207  *TN:* True Negative. Among the samples judged as negative, the correct number is

208 judged.

209  *FN:* False Negative. In the samples judged as negative, the number of judgment

210 errors.

211  $Accuracy = \dfrac{TP + TN}{TP + TN + FN + FP}$

212 (1)

213  $Precision = \dfrac{TP}{TP+FP}$            (2)

214  $Recall = \dfrac{TP}{TP+FN}$

215 (3)

216  *AUC:* The area under the receiver operating characteristic(ROC) curve, is

217 currently considered to be the standard method to assess the accuracy of predictive

218 distribution models[34]. with AUC = 1 represents perfect performance and 0.5 means

219 random guess.

220  *F1-score:* The harmonic mean of the precision(2) and recall(3). The highest

221 possible value of an F-score is 1.0, indicating perfect precision and recall, and the

222 lowest possible value is 0, if either the precision or the recall is zero.

223 **Feature analysis**

224  We will select the model with the best comprehensive score and use the

225 SHAP(Shapley Additive exPlanations) tool to analyze the characteristics of the final

226 model[35]. SHAP is a unified method to explain machine learning predictions based

227 on the optimal Shapley value of game theory. SHAP computed the contribution of

228 each feature to the prediction, which was quantified using Shapley values from

229 coalitional game theory. The Shapley value was represented as an additive feature

230 attribution method, providing the average of the marginal contributions across all

231 permutations of features and distribution of model prediction among features. As an

232 alternative to permutation feature importance, SHAP feature importance was based on

233 magnitude of feature attributions. The absolute Shapley values per feature across the

234 data was further averaged as the global importance was needed. We ranked the

235 features importance in descending order and picked the top 30 most important

236 features. The SHAP value can be plotted for each sample corresponding to the first 30

237 features. We used the Python library to implement the SHAP algorithm

238 (https://github.com/slundberg/shap).

239 For features, the enrichment analysis in CLINVAR, KEGG, GO, and Reactome

240 will also be performed using ClueGO[36]. In addition, the genes corresponding to the

241 optimal model features were extracted, and the Kruskal-Wallis test was used in the

242 BRCA_METABRIC gene expression data set (Bonferroni correction of the results,

243 adjusted P value <=0.05) to calculate the difference between survival and death

244 groups. For genes with significant differences, use the limma tool[37] to calculate the

245 expression fold difference.

**Risk stratification analysis**

247 The original output result of the model is a probability value (between 0 and 1).

248 Based on the optimal model result, We will divide all samples into high, medium and

249   low risk groups (BRCA_METABRIC, BRCA_TCGA), and draw Kaplan-Meier (K-M)

250   curve.

## Result

**Comparison of performance of different machine learning models**

253   The optimization process of the five models (CNA, SNV, CNGM, SNV+CNGM+CNA(combined variants) , combined variants+phenotype)was shown in Supplementary Figure S2. The number of features and AUC values corresponding to the optimal model were: SNV(AUC:0.56; features: 93), CNGM(AUC:0.63; features: 4), CNA(AUC:0.64; features: 75), combined variants (AUC:0.72; features: 353), combined variants + phenotype (AUC:0.81; features: 172).

259   We have drawn the Precision-Recall and ROC curves for the above optimal models using the 5-fold cross-validation method in the BRCA_METABRIC, internal test set and external test set (Figure 2). Taking the test result of the external data set BRCA_TCGA as the final evaluation index, the indexes of each model are as follows: SNV(AUC:0.53, APS:0.25); CNGM(AUC:0.54, APS:0.26); CNA(AUC:0.62, APS:0.42); combined variants (AUC:0.61, APS:0.35); combined variants+phenotype(AUC:0.70, APS:0.43); More model evaluation indicators can be viewed in Table 1. The best comprehensive score was the combined variants+phenotype model, which performed best in both the internal test set (AUC:0.81, APS:0.55) and the independent external data set (AUC=0.70, APS:0.43)(Table 1).

**Optimal model feature ranking**

271  The combined variants + phenotype model comprised a total of 172 features,

272  including 121 CNA, 45 CNGM, 4 SNV, and 2 phenotypes (age, number of positive

273  limph nodes). We used shap to analyze the importance of the predictive characteristics

274  of the model. As shown in Figure 3, among the 172 features of the model, we

275  extracted the top 30 most important features. The phenotypic characteristics age and

276  number of positive lymph nodes ranked first and second, and showed positive

277  correlation with death within five years. The remaining 28 features included 18 CNA

278  (ZNF720, TBC1D13, SCAF4, CDRT15, TMED6, OR4M2, C17orf102, TAS2R10,

279  PHF20L1, RNF187, STIM2, CCDC136, TTI2, MTBP, FAM24B, TMEM26, OR4F15,

280  PDCL2), 9 CNGM (TP53, DNAH11, DNAH2, PIK3CA, MAP3K1, GATA3, CDH1,

281  PDE4DIP, 80273), 1 SNV(chr3:178936091:G:A). Some characteristics also showed

282  positive correlation with mortality within five years, such as CNGM-TP53,

283  CNGM-DNAH2, CNGM-PIK3CA, CNA-SCAF4, CNGM-CDH1, etc. There were

284  also some opposite manifestations, such as CNA-ZNF720, CNGM-DNAH11,

285  CNA-TMED6, CNGM-MAP3K1, SNP-3-178936091-G-A, etc.

286  **Enrichment analysis of optimal model features**

287  We used ClueGO to perform enrichment analysis on the genes corresponding to

288  172 features. The selected data sets included CLINVAR, KEGG, GO, and Reactome

289  pathways. The enrichment results were corrected by bonferoni multiple test. After

290  correction, the pathways with adjusted P value less than 0.05 were selected( Figure

291  4.A). A total of 33 records were obtained. In CLINVAR and KEGG, the features were

292  enriched in pan-cancer or breast cancer-related pathways (C0006142, C1458155,

293 KEGG:05212, KEGG:05222, KEGG:05224). In GO biological process pathways,

294 these genes were over-represented in some pathways related to cell cycle and cell

295 proliferation (GO:0048103, GO:1904030, GO:0000079, GO:0061982, etc.). Two

296 REACTOME pathways reached statistical significance, both of which are related to

297 the NOTCH signaling pathway (R-HSA: 350054, R-HSA: 1980143)

298 **Difference analysis of features at the transcriptoional level**

299 　　In the optimal model, we extracted the genes corresponding to 170 features

300 (excluding two phenotypic features). In BRCA_METABRIC, a total of 17 genes were

301 differentially expressed between the living and dead breast cancer patient groups

302 (adjusted P value < 0.05 for all cases, Kruskal-Wallis test), such as TP53, DNAH11,

303 MAP3K1, PHF20L1, etc. (Figure 4.B). TP53 (No. 3), DNAH11 (No. 5), MAP3K1

304 (No. 12), PHF20L1 (No. 20)　ranked in the top 30 of the model feature weights. The

305 limma results showed that none of these genes had a significant fold change in

306 expression, between the living and dead breast cancer patient groups.

307 **Results of risk stratification**

308 　　According to the model prediction results of all samples, we assigned samples

309 with probability values less than 0.1(TPR>0.93) to the low-risk group (1473 samples),

310 samples with probability values greater than 0.9(TNR>0.99) to high-risk group (363

311 samples), and others to medium-risk group (202 samples). The stratification results

312 are shown in Figure 5.A. The Kaplan-Meier survival curves corresponding to the

313 three sets of results are shown in Figure 5.B. The results showed that the three groups

314 of patients had significantly different survival outcomes. This has clinical implications.

315  For high-risk patients, more clinical intervention and active treatment may be

316  required.

## Discussion

**Model evaluation**

319  The CatBoost algorithm model is used. Different models performed similarly in

320  the training set and the test set without serious overfitting and strong generalization

321  ability.  The best model result is the combined variants+phenotype (AUC: 0.70). For

322  a single feature such as SNV, CNGM has a lower AUC In an independent external

323  data set  (SNV: AUC=0.53, CNGM: AUC=0.54). CNA, as a single-dimensional

324  feature, is similar to the combined variants model's result in external data

325  set(AUC=0.62), and compared to CNGM and SNV, CNA has better generalization

326  capabilities. But this is also related to the small number of CNA and SNV features,

327  and more comprehensive data needs to be collected for further verification. In

328  addition, the addition of phenotype, especially age and the number of lymph nodes,

329  has a very large impact on death, as can be seen from the feature weights of the

330  optimal model (Figure 3).

331  In general, we comprehensively assessed the impact of different characteristics on

332  the five-year mortality risk. In the process of model evaluation, we found that a single

333  feature has poorer performance than the feature fusion model. CNA accounts for a

334  relatively large number of model features due to the large number of original features,

335  but there are still more CNGM features in the top 30 features. The contribution of

336  SNV features in risk prediction is low. The addition of phenotypic information such as

337    age and number of lymph nodes can increase the accuracy of the model.

338    **Discovery of biomarkers associated with five-year mortality risk**

339    In the optimal combined variants+phenotype model, in addition to phenotypic

340    features, some genomics features that have a greater contribution to the model have

341    also been found. And these features still have significant differences between the

342    survival and death groups at the transcriptome level, although there is no large fold

343    difference. For example, TP53, DNAH11, MAP3K1, PHF20L1. As a very complex

344    biomarker, TP53 acts as a tumor suppressor in many tumor types; induces growth

345    arrest or apoptosis depending on the physiological circumstances and cell type which

346    has been widely reported. Its mutations are widely present in various cancers[38-41].

347    IARC TP53 Database (https://p53.iarc.fr/) records all the resources of TP53

348    mutations[41]. They pointed out that there are 28 mutations that lead to a poor

349    prognosis (https://p53.iarc.fr/SomaticPrognosisStats.aspx). In our model, the TP53

350    feature comes from the CNGM feature dimension. The model results indicate that the

351    greater the cumulative number of TP53 mutations, the greater the probability of death

352    within five years (Figure. 3). The DNAH11 gene mutation rs2285947 is considered a

353    potential risk factor for ovarian cancer and breast cancer[42], and there is no clear

354    report related to prognosis. MAP3K1 is a component of a protein kinase signal

355    transduction cascade, which has dual regulatory effects on cell survival and apoptosis,

356    and its regulatory mechanism is not yet clear[43, 44]. These characteristics have been

357    reported to be related to cancer, and our study further verified their relationship with

358    the five-year mortality risk.

359     Through feature selection and multi-dimensional feature fusion, the optimal model

360     features are concentrated on pathways related to cancer, cell division, and

361     proliferation without adding additional prior information. This reflects that the design

362     of the model is relatively reliable, and the model can eliminate features that are not

363     related to the training target from a large amount of input data. The genes

364     corresponding to the features retained by the model are potential biomarkers for

365     prognostic analysis and drug development.

## Conclusion

366

367     In general, in this article, based on the CatBoost algorithm, we use independent

368     data sets of BRCA_METABRIC and BRCA_TCGA to conduct systematic model

369     training on features of different dimensions. The effects of different dimensional

370     features at the genome level on the prediction results of the model are compared. Our

371     best model combines all the features, and the AUC in the external independent

372     BRCA_TCGA is 0.70. In addition, the risk stratification results of all samples showed

373     significant differences between different populations. For high-risk groups classified

374     by the model, active clinical treatment is very necessary. This is the first five-year

375     breast cancer death analysis based on genomic data and using external independent

376     data for evaluation. And compared with other studies, the model based on somatic

377     genomic variants data and phenotypic data (age, number of lymph nodes) is more

378     prospective, and the patient's condition can be evaluated before clinical intervention,

379     providing guidance for follow-up treatment

380     Nevertheless, the research still has limitations. When selecting the features that

381   the two data sets contain in common, the SNP and CNGM features only get very little

382   intersection, which may lead to the underestimation of the role of SNP and CNGM.

383   Deep learning algorithms have not been used and compared. We will continue to

384   conduct in-depth research, collect more comprehensive data, design and develop new

385   algorithms based on existing experience, and further compare the performance

386   differences between machine learning and deep learning. In addition, we will also try

387   to collect other cancer data, conduct migration learning, and develop a five-year

388   mortality risk model for pan-cancer.

## Acknowledgments

## Reference

1.   Sung, H., et al., *Global cancer statistics 2020: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries.* CA: a cancer journal for clinicians, 2021. **71**(3): p. 209-249.

2.   Loibl, S., et al., *Breast cancer.* Lancet, 2021.

3.   Garcia-Closas, M., et al., *Heterogeneity of breast cancer associations with five susceptibility loci by clinical and pathological characteristics.* PLoS genetics, 2008. **4**(4): p. e1000054.

4.   Lüönd, F., S. Tiede, and G. Christofori, *Breast cancer as an example of tumour heterogeneity and tumour cell plasticity during malignant progression.* British Journal of Cancer, 2021. **125**(2): p. 164-175.

5.   McClellan, J. and M.-C. King, *Genetic heterogeneity in human disease.* Cell, 2010. **141**(2): p. 210-217.

6.   Allemani, C., et al., *Global surveillance of trends in cancer survival 2000-14 (CONCORD-3): analysis of individual records for 37 513 025 patients diagnosed with one of 18 cancers from 322 population-based registries in 71 countries.* Lancet, 2018. **391**(10125): p. 1023-1075.

7.   Haybittle, J.L., et al., *A prognostic index in primary breast cancer.* British Journal of Cancer, 1982. **45**(3): p. 361-366.

8.   Cox, D.R., *Regression Models and Life-Tables.* Journal of the Royal Statistical Society. Series B (Methodological), 1972. **34**(2): p. 187-220.

9.   Kattan, M.W., et al., *A tool for predicting breast carcinoma mortality in women who do not receive adjuvant therapy.* Cancer, 2004. **101**(11): p. 2509-15.

10.   Ravdin, P.M., et al., *Computer program to assist in making decisions about adjuvant therapy for women with early breast cancer.* J Clin Oncol, 2001. **19**(4): p. 980-91.

11.   Campbell, H.E., et al., *Estimation and external validation of a new prognostic model for*

| | | |
|---|---|---|
| 415 | | *predicting recurrence-free survival for early breast cancer patients in the UK.* Br J Cancer, |
| 416 | | 2010. **103**(6): p. 776-86. |
| 417 | 12. | Engelhardt, E.G., et al., *Accuracy of the online prognostication tools PREDICT and Adjuvant!* |
| 418 | | *for early-stage breast cancer patients younger than 50 years.* Eur J Cancer, 2017. **78**: p. |
| 419 | | 37-44. |
| 420 | 13. | Reis-Filho, J.S. and L. Pusztai, *Gene expression profiling in breast cancer: classification,* |
| 421 | | *prognostication, and prediction.* Lancet, 2011. **378**(9805): p. 1812-23. |
| 422 | 14. | van 't Veer, L.J., et al., *Gene expression profiling predicts clinical outcome of breast cancer.* |
| 423 | | Nature, 2002. **415**(6871): p. 530-6. |
| 424 | 15. | Paik, S., et al., *A multigene assay to predict recurrence of tamoxifen-treated, node-negative* |
| 425 | | *breast cancer.* N Engl J Med, 2004. **351**(27): p. 2817-26. |
| 426 | 16. | Filipits, M., et al., *A new molecular predictor of distant recurrence in ER-positive,* |
| 427 | | *HER2-negative breast cancer adds independent information to conventional clinical risk* |
| 428 | | *factors.* Clin Cancer Res, 2011. **17**(18): p. 6012-20. |
| 429 | 17. | Manjang, K., et al., *Prognostic gene expression signatures of breast cancer are lacking a* |
| 430 | | *sensible biological meaning.* Sci Rep, 2021. **11**(1): p. 156. |
| 431 | 18. | Kourou, K., et al., *Machine learning applications in cancer prognosis and prediction.* Comput |
| 432 | | Struct Biotechnol J, 2015. **13**: p. 8-17. |
| 433 | 19. | Hieronymus, H., et al., *Tumor copy number alteration burden is a pan-cancer prognostic* |
| 434 | | *factor associated with recurrence and death.* Elife, 2018. **7**. |
| 435 | 20. | Nguyen, C., Y. Wang, and H.N. Nguyen, *Random forest classifier combined with feature* |
| 436 | | *selection for breast cancer diagnosis and prognostic.* Journal of Biomedical Science and |
| 437 | | Engineering, 2013. **Vol.06No.05**: p. 10. |
| 438 | 21. | Park, K., et al., *Robust predictive model for evaluating breast cancer survivability.* |
| 439 | | Engineering Applications of Artificial Intelligence, 2013. **26**(9): p. 2194-2205. |
| 440 | 22. | Zhang, Y., et al., *Toward the precision breast cancer survival prediction utilizing combined* |
| 441 | | *whole genome-wide expression and somatic mutation analysis.* BMC Med Genomics, 2018. |
| 442 | | **11**(Suppl 5): p. 104. |
| 443 | 23. | He, Z., et al., *Integrating Somatic Mutations for Breast Cancer Survival Prediction Using* |
| 444 | | *Machine Learning Methods.* Front Genet, 2020. **11**: p. 632901. |
| 445 | 24. | Li, J., et al., *Predicting breast cancer 5-year survival using machine learning: A systematic* |
| 446 | | *review.* PLoS One, 2021. **16**(4): p. e0250370. |
| 447 | 25. | Sun, D., M. Wang, and A. Li, *A multimodal deep neural network for human breast cancer* |
| 448 | | *prognosis prediction by integrating multi-dimensional data.* IEEE/ACM Trans Comput Biol |
| 449 | | Bioinform, 2018. |
| 450 | 26. | Arya, N. and S. Saha, *Multi-modal advanced deep learning architectures for breast cancer* |
| 451 | | *survival prediction.* Knowledge-Based Systems, 2021. **221**: p. 106965. |
| 452 | 27. | Tomczak, K., P. Czerwińska, and M. Wiznerowicz, *The Cancer Genome Atlas (TCGA): an* |
| 453 | | *immeasurable source of knowledge.* Contemporary oncology, 2015. **19**(1A): p. A68. |
| 454 | 28. | Curtis, C., et al., *The genomic and transcriptomic architecture of 2,000 breast tumours reveals* |
| 455 | | *novel subgroups.* Nature, 2012. **486**(7403): p. 346-352. |
| 456 | 29. | Pereira, B., et al., *The somatic mutation profiles of 2,433 breast cancers refine their genomic* |
| 457 | | *and transcriptomic landscapes.* Nature communications, 2016. **7**(1): p. 1-16. |
| 458 | 30. | Rueda, O.M., et al., *Dynamics of breast-cancer relapse reveal late-recurring ER-positive* |

459          *genomic subgroups.* Nature, 2019. **567**(7748): p. 399-404.

460   31.   Kraskov, A., H. Stögbauer, and P. Grassberger, *Estimating mutual information.* Physical
461          review E, 2004. **69**(6): p. 066138.

462   32.   Li, F. and Y. Yang. *Analysis of recursive feature elimination methods*. in *Proceedings of the*
463          *28th annual international ACM SIGIR conference on Research and development in*
464          *information retrieval.* 2005.

465   33.   Kursa, M.B. and W.R. Rudnicki, *Feature selection with the Boruta package.* J Stat Softw,
466          2010. **36**(11): p. 1-13.

467   34.   Lobo, J.M., A. Jimenez-Valverde, and R. Real, *AUC: a misleading measure of the*
468          *performance of predictive distribution models.* Global Ecology and Biogeography, 2008. **17**(2):
469          p. 145-151.

470   35.   Lundberg, S.M. and S.-I. Lee. *A unified approach to interpreting model predictions*. in
471          *Proceedings of the 31st international conference on neural information processing systems.*
472          2017.

473   36.   Bindea, G., et al., *ClueGO: a Cytoscape plug-in to decipher functionally grouped gene*
474          *ontology and pathway annotation networks.* Bioinformatics, 2009. **25**(8): p. 1091-3.

475   37.   Smyth, G.K., *Limma: linear models for microarray data*, in *Bioinformatics and computational*
476          *biology solutions using R and Bioconductor*. 2005, Springer. p. 397-420.

477   38.   Bertheau, P., et al., *TP53 status and response to chemotherapy in breast cancer.* Pathobiology,
478          2008. **75**(2): p. 132-139.

479   39.   Børresen☐Dale, A.L., *TP53 and breast cancer.* Human mutation, 2003. **21**(3): p. 292-300.

480   40.   Petitjean, A., et al., *TP53 mutations in human cancers: functional selection and impact on*
481          *cancer prognosis and outcomes.* Oncogene, 2007. **26**(15): p. 2157-2165.

482   41.   Petitjean, A., et al., *Impact of mutant p53 functional properties on TP53 mutation patterns and*
483          *tumor phenotype: lessons from recent developments in the IARC TP53 database.* Human
484          mutation, 2007. **28**(6): p. 622-629.

485   42.   Verma, S., et al., *Genetic variants of DNAH 11 and LRFN 2 genes and their association with*
486          *ovarian and breast cancer.* International Journal of Gynecology & Obstetrics, 2020. **148**(1): p.
487          118-122.

488   43.   Pham, T.T., S.P. Angus, and G.L. Johnson, *MAP3K1: genomic alterations in cancer and*
489          *function in promoting cell survival or apoptosis.* Genes & cancer, 2013. **4**(11-12): p. 419-426.

490   44.   Xue, Z., et al., *MAP3K1 and MAP2K4 mutations are associated with sensitivity to MEK*
491          *inhibitors in multiple cancer models.* Cell research, 2018. **28**(7): p. 719-729.
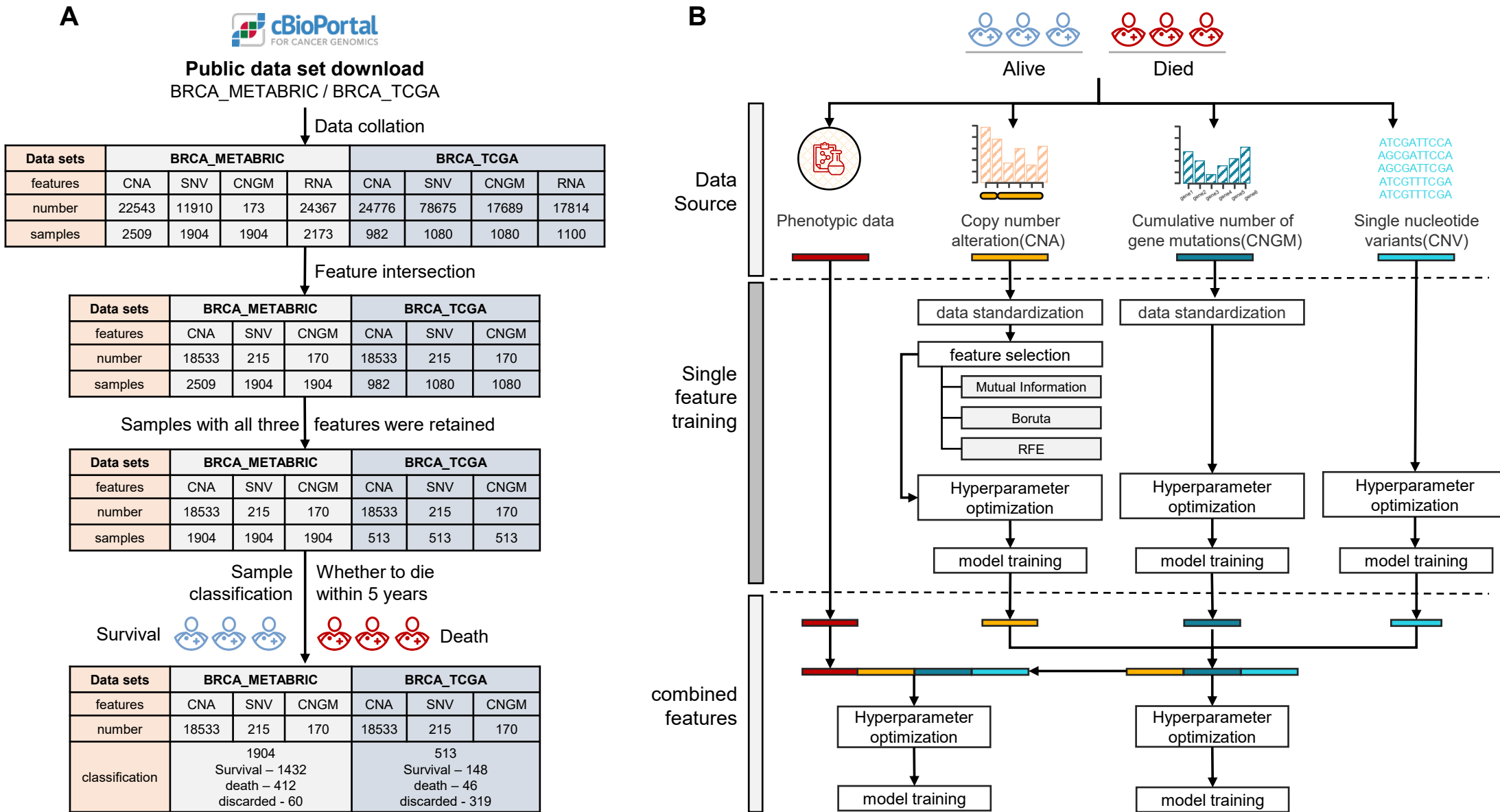
492

493

## Figure

**Figure 1. Data quality control and five-year survival prediction model building process. A)** BRCA_TCGA and BRCA_METABRIC data acquisition and quality control process; **B)** In the process of building a polygenic risk assessment model, different processing methods are adopted for different dimension characteristics.

**Figure 2. Precision-Recall and ROC curves of optimal models constructed with features of different dimensions.** The first row is the ROC curve and the second row is the Precision-Recall curve. From left to right are the results in the cross-validation set(mean), training set and test set, respectively.

**Figure 3. Optimal model feature weight analysis.** The scatter points represent the SHAP value of each feature for each sample. Features are sorted according to the sum of the magnitudes of the SHAP values of all samples. The first 30 features are shown, and the colors represent the feature values (red high, blue low). For example, as age ("AGE_AT_DIAGNOSIS") increases, the risk of death within five years of the sample will increase.

**Figure 4. A)** Optimal Model Pathway Enrichment Analysis; **B)** Transcriptome-level differential analysis of optimal model features.

**Figure 5. A)** Risk stratification for all samples based on model scoring; **B)** Plot Kaplan-Meier survival curves for three groups of stratified outcomes (high, intermediate, and low risk).

## Supplementary Figure

**Supplementary Figure S1.** Statistical results of sample distribution regarding gender, number of lymph nodes, menopause (-9 - unknown, 0 - not menopause, 1 - menopause).

**Supplementary Figure S2.** The optimization process of the five models (CNA, SNV, CNGM, SNV+CNGM+CNA(combined variants) , combined variants+phenotype).

518 **Table**

| Model | Internal test data | | | | | | External test data | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AUC | F1 | Accuracy | precision | recall | APS | AUC | F1 | Accuracy | precision | recall | APS |
| SNV | 0.56 | 0.21 | 0.72 | 0.29 | 0.17 | 0.25 | 0.53 | 0.22 | 0.70 | 0.29 | 0.17 | 0.25 |
| CNGM | 0.63 | 0.37 | 0.66 | 0.32 | 0.45 | 0.28 | 0.54 | 0.36 | 0.62 | 0.30 | 0.46 | 0.26 |
| CNA | 0.64 | 0.34 | 0.74 | 0.39 | 0.30 | 0.34 | 0.62 | 0.38 | 0.77 | 0.52 | 0.30 | 0.42 |
| SNV+CNGM+CNA | 0.72 | 0.32 | 0.76 | 0.42 | 0.26 | 0.38 | 0.61 | 0.25 | 0.73 | 0.36 | 0.20 | 0.35 |
| SNV+CNGM+CNA +Phenotype | 0.81 | 0.52 | 0.80 | 0.55 | 0.49 | 0.55 | 0.70 | 0.46 | 0.77 | 0.51 | 0.41 | 0.43 |

519 **Table 1.** The model predicts the performance indicators of breast cancer deaths within five years

520 in the internal and external test data sets.

521

**Figure 1. Data quality control and five-year survival prediction model building process.** **A)** BRCA_TCGA and BRCA_METABRIC data acquisition and quality control process; **B)** In the process of building a polygenic risk assessment model, different processing methods are adopted for different dimension characteristics.
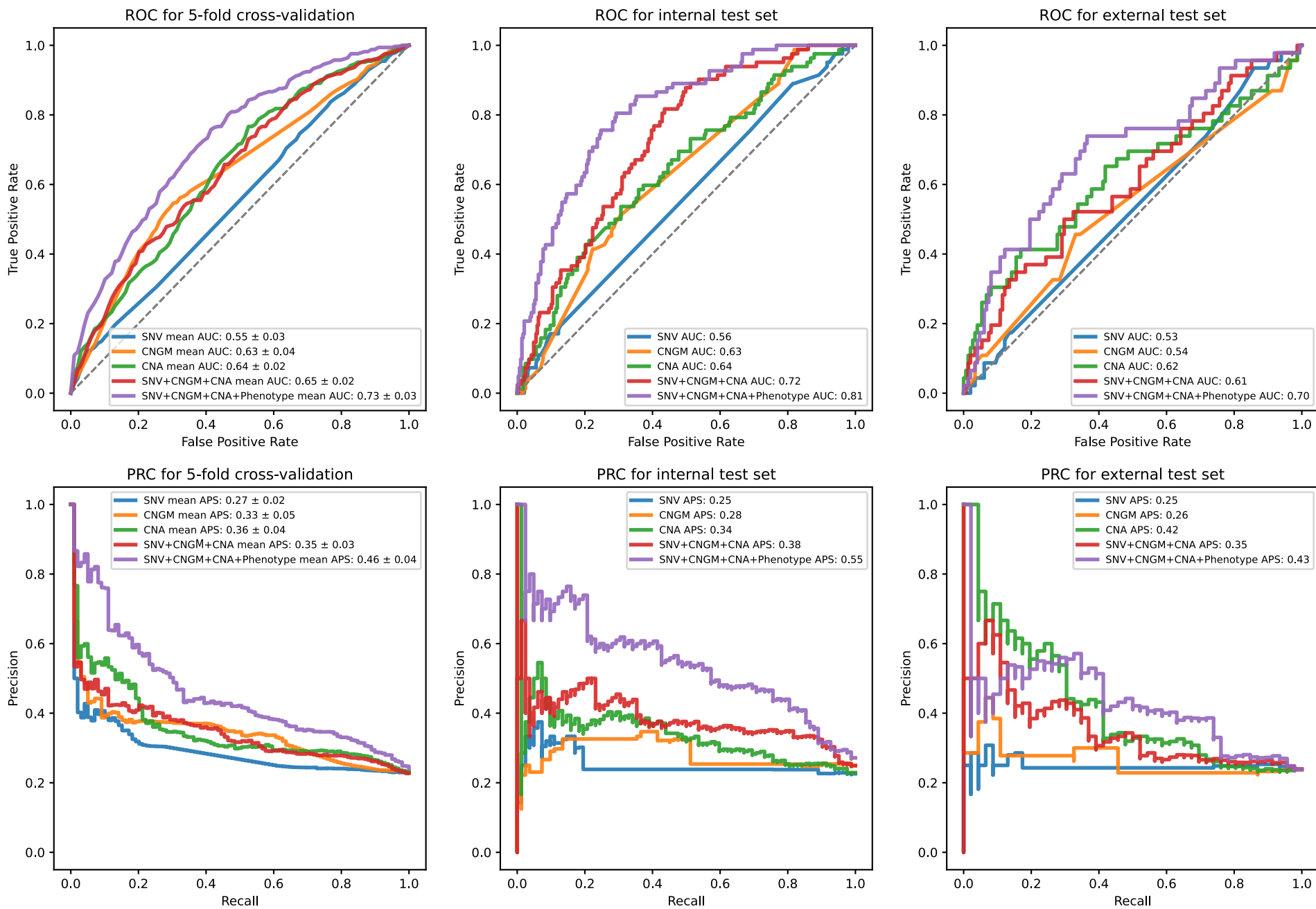
**Figure 2. Precision-Recall and ROC curves of optimal models constructed with features of different dimensions.** The first row is the ROC curve and the second row is the Precision-Recall curve. From left to right are the results in the cross-validation set(mean), training set and test set, respectively.

**Figure 3. Optimal model feature weight analysis.** The scatter points represent the SHAP value of each feature for each sample. Features are sorted according to the sum of the magnitudes of the SHAP values of all samples. The first 30 features are shown, and the colors represent the feature values (red high, blue low). For example, as age ("AGE_AT_DIAGNOSIS") increases, the risk of death within five years of the sample will increase.
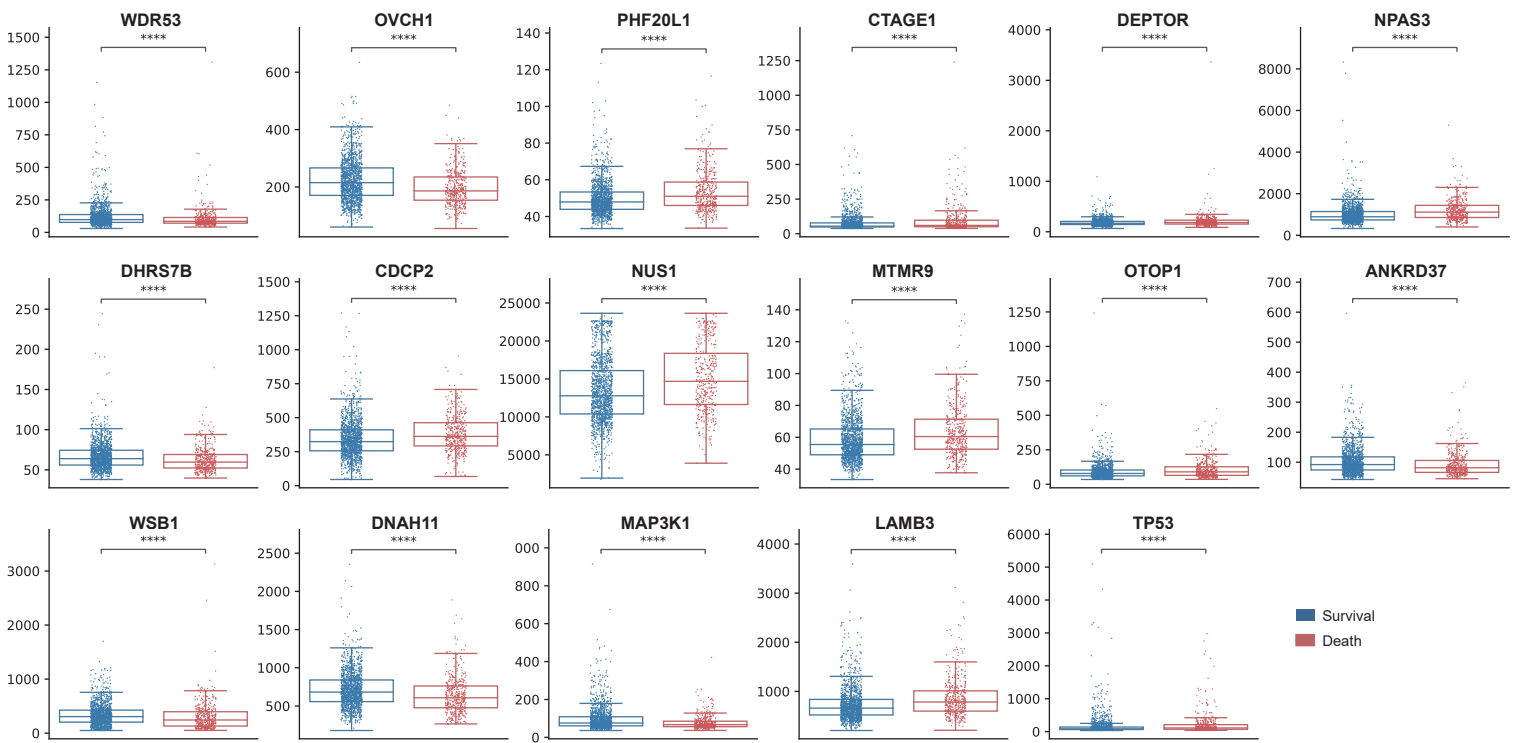
**Figure 4. A)** Optimal Model Pathway Enrichment Analysis; **B)** Transcriptome-level differential analysis of optimal model features.
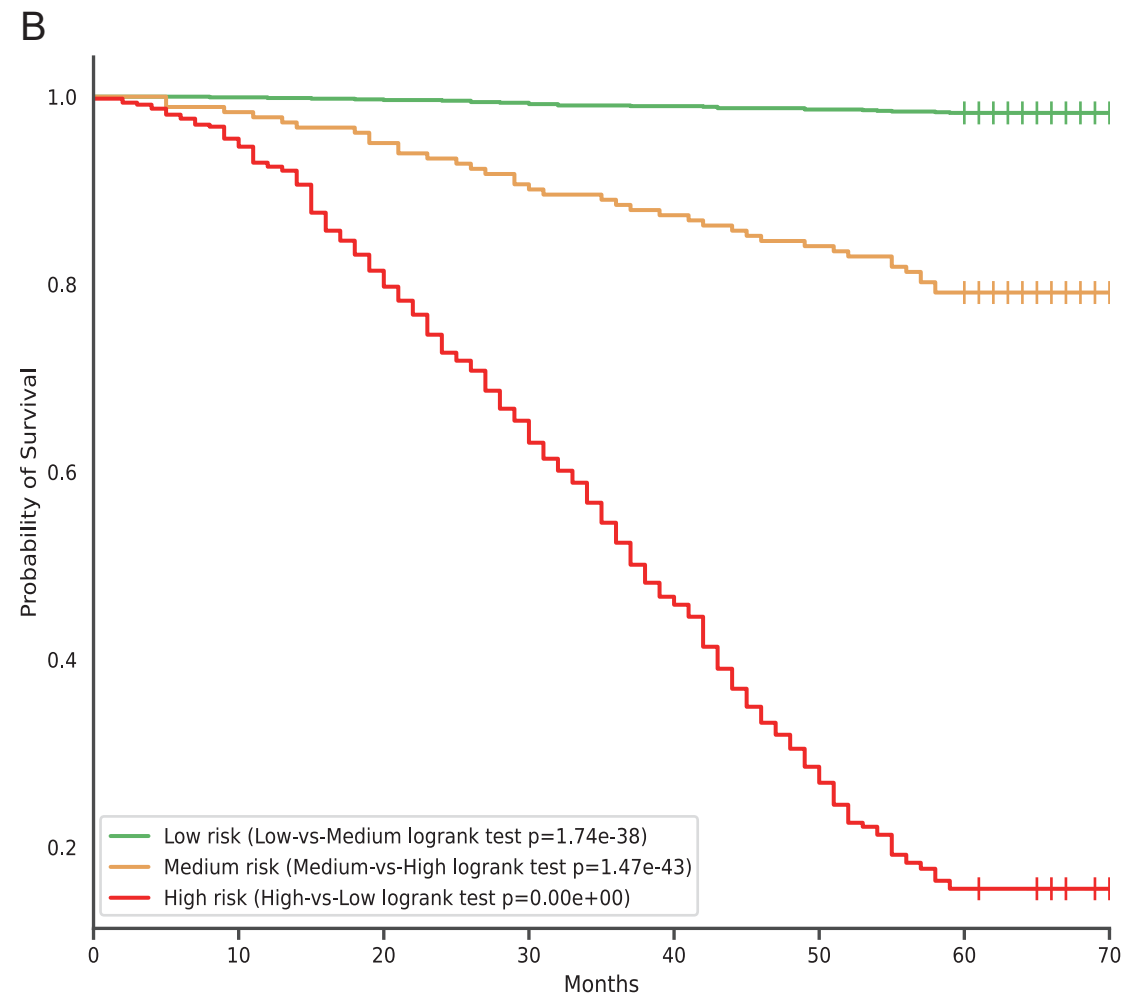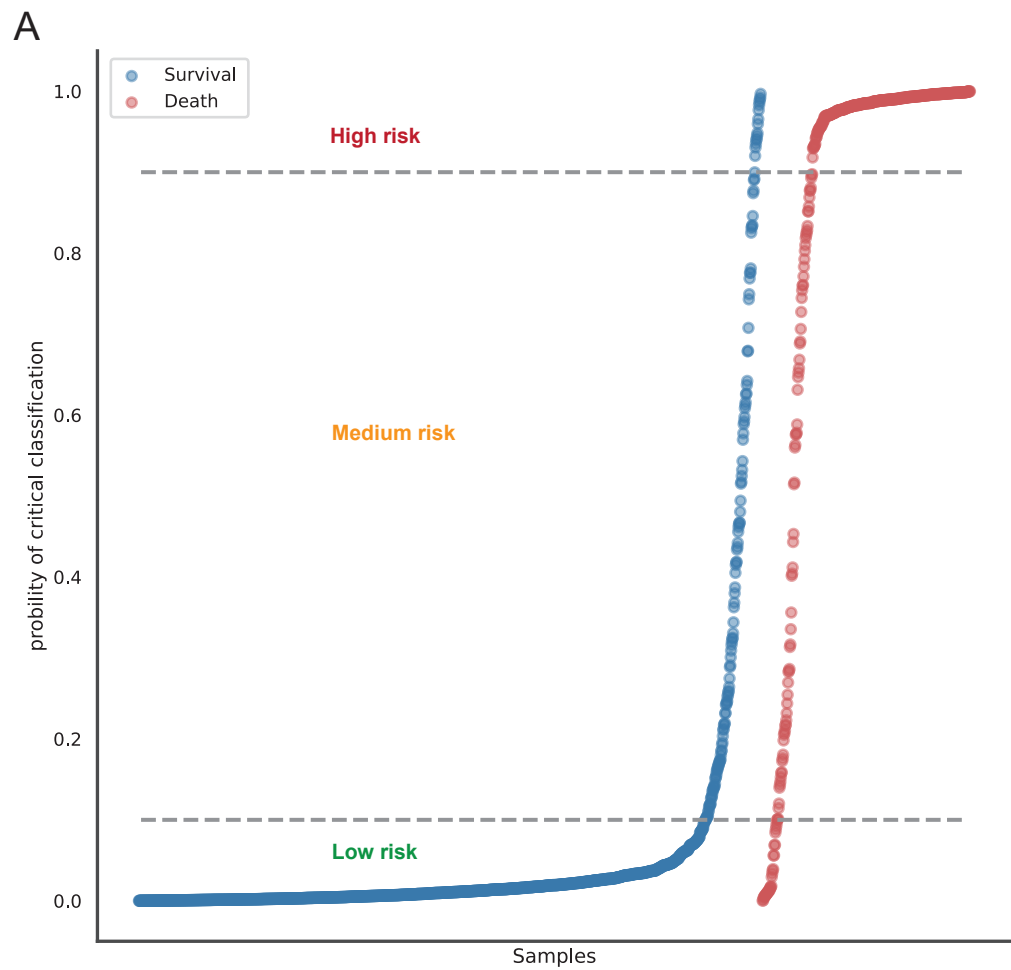
**Figure 5. A)** Risk stratification for all samples based on model scoring; **B)** Plot Kaplan-Meier survival curves for three groups of stratified outcomes (high, intermediate, and low risk).