


Lift the Veil of Breast Cancers Using 4 or Fewer Critical Genes

Zhengjun Zhang 

Department of Statistics, University of Wisconsin, Madison, WI, USA.

Cancer Informatics
Volume 21: 1–11
© The Author(s) 2022
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/11769351221076360



ABSTRACT: Known genes in the breast cancer study literature could not be confirmed whether they are vital to breast cancer formations due to lack of convincing accuracy, although they may be biologically directly related to breast cancer based on present biological knowledge. It is hoped vital genes can be identified with the highest possible accuracy, for example, 100% accuracy and convincing causal patterns beyond what has been known in breast cancer. One hope is that finding gene-gene interaction signatures and functional effects may solve the puzzle. This research uses a recently developed competing linear factor analysis method in differentially expressed gene detection to advance the study of breast cancer formation. Surprisingly, 3 genes are detected to be differentially expressed in TNBC and non-TNBC (Her2, Luminal A, Luminal B) samples with 100% sensitivity and 100% specificity in 1 study of triple-negative breast cancers (TNBC, with 54675 genes and 265 samples). These 3 genes show a clear signature pattern of how TNBC patients can be grouped. For another TNBC study (with 54673 genes and 66 samples), 4 genes bring the same accuracy of 100% sensitivity and 100% specificity. Four genes are found to have the same accuracy of 100% sensitivity and 100% specificity in 1 breast cancer study (with 54675 genes and 121 samples), and the same 4 genes bring an accuracy of 100% sensitivity and 96.5% specificity in the fourth breast cancer study (with 60483 genes and 1217 samples). These results show the 4-gene-based classifiers are robust and accurate. The detected genes naturally classify patients into subtypes, for example, 7 subtypes. These findings demonstrate the clearest gene-gene interaction patterns and functional effects with the smallest numbers of genes and the highest accuracy compared with findings reported in the literature. The 4 genes are considered to be essential for breast cancer studies and practice. They can provide focused, targeted researches and precision medicine for each subtype of breast cancer. New breast cancer disease types may be detected using the classified subtypes, and hence new effective therapies can be developed.

KEYWORDS: Direct and indirect effects, breast cancer detection, gene-gene interaction, functional effects, joint risk competing

RECEIVED: August 14, 2021. **ACCEPTED:** December 30, 2021.

TYPE: Original Research

FUNDING: The author disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: The project was partially supported by NSF grant DMS-2012298.

DECLARATION OF CONFLICTING INTERESTS: The author declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

CORRESPONDING AUTHOR: Zhengjun Zhang, Department of Statistics, University of Wisconsin, Madison, WI 53706, USA. Email: zjz@stat.wisc.edu

Introduction

Breast cancer has been an unconquered plague for centuries. It has had the highest death rate among all cancers women have had for many years. It has caused enormous economic losses and costs. To save lives and protect women from breast cancers, enormous research efforts and money have been investigated. Although there have been some considerable signs of progress in breast cancer diagnoses and therapies, many women still suffer from being diagnosed with breast cancer and lost their lives every year. No apparent clues or research results show the most critical genetic causality in breast cancer formation. The most hopeful direction, finding critical genes, or primary differentially expressed genes related to breast cancer formation, has been drawing much attention in breast cancer studies. The most recent editorial summary by Narod¹ states “Results of two large case-control studies that analyzed the associations between a number of putative cancer susceptibility genes and breast cancer risk are now reported in the Journal. The study by Dorling et al² included 34 genes and 113 000 women from 25 countries, and the study by Hu et al³ included 28 genes and 64 000 women from the United States. Variants in 8 genes—BRCA1, BRCA2, PALB2, BARD1, RAD51C, RAD51D, ATM, and CHEK2—had a significant association with breast cancer risk in both studies.” However, a significant association does not mean the corresponding gene is truly informative. For example, it has been reported by Berger⁴ that the risk of

developing breast cancer was 40% to 60% greater among women with the PALB2 mutation. On the other hand, the study by de Magalhães⁵ shows every gene can (and possibly will) be associated with cancer, see also an interview report by Robitzski⁶ in *The Scientist*. It becomes clear that having a lot of genes associated with a disease doesn't mean they're important. This paper intends to identify a truly important small subset of breast cancer risk genes.

Differential expression analysis between tumor and non-tumor cells helps breast cancer prognosis prediction at a relatively early stage, identifying some clear patterns from patients to patients, recommending different precision therapies according to breast cancer subtypes. Efforts have been made in identifying genes associated with breast cancer symptoms. We now give a brief review of some of the most recent studies. In a systems biology comprehensive analysis on breast cancer to identify key gene modules and genes associated with TNM-based clinical stages,⁷ the authors have identified various numbers of genes that can be key genes related to breast cancers at different cancer stages. Malvia et al⁸ studied gene expression profiles of breast cancers in Indian women, obtained 2413 differentially expressed genes, and demonstrated the existence of molecular subtypes in Indian women. Lv et al⁹ aimed to explore some novel genes and pathways related to TNBC prognosis through bioinformatics methods as well as potential initiation and progression mechanisms. Seven hundred



fifty-five differentially expressed overlapping mRNAs were detected between TNBC/non-TNBC samples and normal tissue. The authors found 8 hub genes associated with the cell cycle pathway highly expressed in TNBC. Additionally, a novel 6-gene (TMEM252, PRB2, SMCO1, IVL, SMR3B, and COL9A3) signature from the 755 differentially expressed mRNAs were constructed and significantly associated with prognosis as an independent prognostic signature. Zhong et al¹⁰ conducted a robust rank aggregation (RRA) analysis based on genome-wide gene expression datasets involving TNBC patients from the Gene Expression Omnibus (GEO) database to identify key genes associated with TNBC. A total of 194 highly ranked differentially expressed genes (DEGs) were identified in TNBC versus non-TNBC. Gene ontology (GO) and Kyoto Encyclopedia of Genes and Genomes pathway (KEGG) enrichment analysis was utilized to explore the identified genes' biological functions. The authors also found that some genes are positively correlated to the life expectancy ($P < .05$) of TNBC patients. Lin et al¹¹ identified potential key genes for HER-2 positive breast cancer based on bioinformatics analysis. A total of 54 up-regulated DEGs and 269 down-regulated DEGs were identified. Among them, 10 hub genes including CCNB1, RAC1, TOP2A, KIF20A, RRM2, ASPM, NUSAP1, BIRC5, BUB1B, and CEP55 demonstrated by connectivity degree in the PPI network were screened out. Chen et al¹² systematically searched the electronic databases of MEDLINE (PubMed), Embase, and Cochrane Library to identify relevant publications from April, 1959 to November, 2017. identified 16 qualified studies from 527 publications with 46,870 breast cancer patients including 868 BRCA1 mutations carriers, 739 BRCA2 mutations carriers, and 45,263 non-carriers. The results showed that breast cancer patients with BRCA1Mut carriers were more likely to have TNBC than those of BRCA2Mut carriers (OR: 3.292; 95% CI: 2.773-3.909) or non-carriers (OR: 8.889; 95% CI: 6.925-11.410). Deng et al¹³ identified potential crucial genes and key pathways in breast cancer using bioinformatic analysis. Two hundred three up-regulated and 118 down-regulated DEGs were identified. Six hub genes were selected and validated in clinical sample for further analysis due to the high degree of connectivity, including CDK1, CCNA2, TOP2A, CCNB1, KIF11, and MELK. They were all correlated to worse overall survival (OS) in breast cancer. Zhu et al¹⁴ identified some key genes and pathways associated with irradiation in breast cancer tissue and breast cancer cell lines. A total of 82 DEGs (74 up-regulated and 8 down-regulated genes) were identified. Two characteristic subnetworks and 3 hub genes (FOS, CCL2, and CXCL12) were strongly distinguished in PPI network. Dong et al¹⁵ aimed to identify the key pathways and genes and find the potential initiation and progression mechanism of TNBC. Fifty-six up-regulated and 151 down-regulated genes were listed, and the gene ontology (GO) and Kyoto Encyclopedia of Genes and Genomes pathway (KEGG) enrichment analysis was

performed. The authors found that SOX8, AR, C9orf152, NRK and RAB30, and other key genes and pathways might be promising targets for the TNBC treatment. Lu et al¹⁶ identified 5 hub genes (PHLPP1, UBC, ACACB, TGFB1, and ACTB) associated with HER2+BC with brain metastasis. The GSEA analysis revealed that the ribosomal pathway seems to play a very important role in the pathogenesis of HER2+BC with brain metastasis. Among these studies with various study designs, many genes are linked to breast cancer, which provides additional evidence stated by de Magalhães⁵. As a result, many efforts are needed in finding vital genes with the highest possible accuracy, for example, 100% accuracy and convincing causal patterns.

The reported genes in the published work point out some promising directions in breast cancer research and treatments. But it is not clear whether or not they are fundamental causes or direct causes of breast cancer. The problem is mainly due to the following 3 main limitations. (1) The number of human genes is ultra-large compared to the number of patients in affordable study designs. Identifying a few key (single digit) genes that are uniformly optimal across different trials, different study purposes, different measurement methods, and different cohorts is rather challenging. From the aforementioned research outcomes, we can see there are many different genes are identified. As a result, it's impossible to see which one is the most important one, which can be a driver of breast cancer disease. (2) The inefficient detecting power of existing analysis methods due to restricted model assumptions cannot deal with heterogeneous populations (different breast cancer subtypes). As a result, the sensitivity and specificity of many published gene classifiers are not satisfactory. (3) It isn't easy to extract informative messages from existing models and analysis methods. Also, many gene-related classifiers are not interpretable as gene-gene inter-relationships, and functional effects are hardly expressed. As a result, scientific research progress in breast cancer studies is still limited. Much literature attention has been focused on individual genes and their expression levels, that is, not gene-gene interactions, genes-subtypes (of breast cancers) interactions, and functional effects. As a result, the fundamental genetic causes of breast cancer formations can be masked by those suboptimal focuses, and the researches can still be in a primitive state. Many unknown factors exist. They can be essential to conquer the breast cancer plague, and therefore there is an urgent need for identifying critical DEGs with the highest possible sensitivity and specificity for breast cancer detection.

This work aims to lift the veil of breast cancers by discovering the joint functional effects of 4 or fewer critical DEGs that show the highest detecting power of breast cancer in 4 gene expression RNA-seq datasets. According to our analysis, these 4 genes and their functional effects describe breast cancers' overall features at the genomic level, with the highest possible sensitivity (up to 100%) and specificity (up to 100%) for breast

cancer detection. In addition, they are invariance preserving with the same group of patients but measured in different scales, and they are robust from 1 trial to another trial.

Statistical Methodology

In the medical literature and practice, logistic regression has been widely used in studying the disease types and risk probabilities. Recently, Teng and Zhang¹⁷ pointed out a limitation of the classical logistic regression model: it can only model absolute treatment effects in medical data modeling, that is, it does not model relative treatment effects. As a result, many well-designed trial studies were tested to be insignificant due to the lack of detecting power of the classical logistic regression. In their paper, Teng and Zhang¹⁷ introduced relative treatment effects in their enhanced logistic regression model (AbRelaTEs) and demonstrated its better modeling capability using 4 clinical trials studies.

When data are drawn from a homogeneous population with 1 disease type, the classical logistic regression and the AbRelaTEs model are applicable. However, when data are drawn from a heterogeneous population, we need a different modeling framework to deal with competing risks, for example, TNBC, Her2, Luminal A, Luminal B in breast cancer. The most recently developed max-linear competing factor models,¹⁸ max-linear regression models,¹⁹ and max-linear logistic models^{20,21} have proven to be powerful models and analysis approaches to study heteroscedastic populations and competing risks and resources. The theoretical foundations of these models have been established in Cui and Zhang,¹⁸ Cui et al,¹⁹ Malinowski et al,²² Xu,²⁰ and Zhang.^{21,23} The difference between the max-linear competing models and the classical statistical models is that the original linear combination

of predictors is replaced by the maximum of a set of linear combinations of predictors, called competing factors or competing-risk factors. The max-linear competing factor models are different from existing popular classification models such as random forest, support vector machine, group lasso-based machine learning methods, and deep learning methods. However, the max-linear competing factor models are interpretable and outperform existing methods.¹⁹ This study implements the max-linear logistic regression model to build a competing factor breast cancer classifier. For completeness, the model is stated as follows.

Suppose there are $i=1, \dots, n$ patients with breast cancer status label $Y_i=1$ for cancer and $Y_i=0$ for cancer-free, and Y_i is related to G groups of genes by

$$\Phi_{ij} = (X_{i,j_1}, X_{i,j_2}, \dots, X_{i,j_{g_j}}), j=1, \dots, G, g_j \geq 0 \quad (1)$$

where i is the i th individual in the sample, g_j is the number of genes in j th group. The competing (risk) factor classifier for the i th outcome variable is defined as

$$\log\left(\frac{p_i}{1-p_i}\right) = \max(\beta_{0j} + \Phi_{ij}\beta_j, \beta_{02} + \Phi_{i2}\beta_2, \dots, \beta_{0G} + \Phi_{iG}\beta_G) \quad (2)$$

where β_{0j} 's are intercepts, Φ_{ij} is a $1 \times g_j$ observed vector, β_j is a $g_j \times 1$ coefficient vector which characterizes the contribution of each predictor to the outcome variable Y in the j th group to the risk, and $\beta_{0j} + \Phi_{ij}\beta_j$ is called the j th competing risk factor, that is, j th signature. The unknown parameters are estimated from

$$\begin{aligned} (\hat{\beta}, \hat{S}, \hat{G}) = \operatorname{argmin}_{\beta, S_j \in S, j=1, 2, \dots, G} \{ & (1 + \lambda_1 + |S_u|)^{\sum_{i=1}^n (I(p_i \leq 0.5)I(Y_i=1) + I(p_i > 0.5)I(Y_i=0))} \\ & + \lambda_2 (|S_u| - \frac{|S_u| + G - 1}{(|S_u| + 1) \times G - 1}) \} \end{aligned} \quad (3)$$

where 0.5 is a probability threshold value that is commonly used in machine learning classifiers, $I(\cdot)$ is an indicate function, p_i is defined in the equation (2). $S = \{1, 2, \dots, 54675\}$ is the index set of all genes, $S_1 = \{1_1, 1_2, \dots, 1_{g_1}\}$, $S_2 = \{2_1, \dots, 2_{g_2}\}, \dots, S_G = \{G_1, \dots, G_{g_G}\}$ are index sets corresponding to (1), and \hat{G} and $\hat{S} = \{j_{j_1}, \dots, j_{j_{g_j}}, j=1, \dots, \hat{G}\}$ are the final gene set selected in the final classifiers.

The goal is to identify the clearest patterns of gene-gene interactions and functional effects related to breast cancer samples and non-tumor samples. We start with 3 competing factors in max-linear logistic regression models, with each factor having only 3 genes randomly drawing from 54675, 54673, or 60483 genes. Then, a Monte Carlo method with extensive computation is used to find the final model with the best performance of sensitivity and specificity and the smallest number

of genes. Finally, the complete computing description is listed in Zhang²¹ in which 5 Covid-19 critical genes and 7 subtypes were identified, and the validity of (3) is shown theoretically in Zhang.²⁴

Data Descriptions

There are 4 datasets used in this study. The first dataset is triple-negative breast cancer (TNBC, North American cohort) study conducted by Burstein et al²⁵ and den Hollander et al²⁶ with 54675 genes, 198 TNBC tumor samples, and 67 not TNBC (Her2, Luminal A, Luminal B) samples. The data link and descriptions are <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE76275>. The platforms are GPL570 [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array. The expression values are $\log_2(\text{RMA signal})$. The

second dataset is a European cohort with 55 TNBC samples and 11 normal breast tissue samples, studied by Maire et al.^{27,28} and Maubant et al.²⁹ The data link and description are <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE65194>. The platforms are GPL570 [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array. The number of genes is 54673. The expression values are $\log_2(\text{GCRMA signal from Affy cdf})$. The third dataset is gene expression profiling of 104 breast cancer and 17 normal breast biopsies by Clarke et al.³⁰ It is from a European cohort. The data link and descriptions are <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE42568>. The platforms are GPL570 [HG-U133_Plus_2] Affymetrix Human Genome U133 Plus 2.0 Array. The expression values are $\log_2(\text{GC-RMA signal intensity})$. The fourth dataset is GDC TCGA Breast Cancer cohort by Genomic Data Commons. The dataset contains 60484 identifiers (genes) and 1217 (1104 tumors and 113 tumor-free) samples. Data from the same sample but from different vials/portions/analytes/aliquotes is averaged; data from different samples are combined into genomicMatrix; all data is then $\log_2(\text{fpkm} + 1)$ transformed. The platform is Illumina. The type of data is gene expression RNAseq. The data link and descriptions are https://xenabrowser.net/datapages/?dataset=TCGA-BRCA.htseq_fpkm.tsv&host=https%3A%2F%2Fgdc.xenahubs.net&removeHub=https%3A%2F%2Fxcena.treehouse.gi.ucsc.edu%3A443.

Results and Interpretations

In medical studies, sensitivity and specificity are 2 main indexes to evaluate treatment effectiveness and disease classification capability in diagnoses. If the intention is to rule out disease, a test with high sensitivity is demanded. If it is desired to confirm a diagnosis or find evidence of disease, a test with high specificity is required. We adopt these 2 metrics in this study. The aim is to identify the smallest number of genes that lead to the highest sensitivities and specificities and establish mathematical equivalence and biological equivalence between the chosen genes and the disease types. Meanwhile, we also present graphical illustration tools for practical doctors to use in their daily medical practice.

Using a probability higher than 50% as the threshold, we identify 3 critical DEGs: RBM22 (RNA binding motif protein 22), RNF213 (ring finger protein 213), and CACNG4 (Calcium Voltage-Gated Channel Auxiliary Subunit Gamma 4), which lead to 100% sensitivity and 100% specificity of classifying all 265 samples in their respective groups in the first TNBC dataset; 4 critical DEGs: MYCT1 (MYC Target 1), NUA2 (NUAK Family Kinase 2), NAT8L (N-Acetyltransferase 8 Like), and CACNG4, which lead to 100% sensitivity and 100% specificity of classifying all 66 samples in their respective groups in the second TNBC dataset; 4 critical DEGs: MYCT1, UNC5B (Unc-5 Netrin Receptor B), NUA2, and NAT8L, which also lead to 100% sensitivity and 100% specificity of classifying all 121 samples

in their respective groups in the third breast cancer dataset; and the same 4 critical DEGs as in the third dataset, which leads to 100% sensitivity and 96.5% specificity of classifying all 1217 samples in their respective groups in the fourth breast cancer dataset.

Our final classifiers are combined classifiers of 3 competing factor ($CF_i, i = 1, 2, 3$) classifiers expressed as:

For the first TNBC (North American cohort) dataset:

$$\text{Data-1-CF1: } 19.0107 + 3.1105 * \text{RNF213} - 3.6692 * \text{CACNG4}$$

$$\text{Data-1-CF2: } -0.4312 + 8.0992 * \text{RNF213} - 9.5921 * \text{RBM22}$$

$$\text{Data-1-CFmax: } \max(\text{Data-1-CF1}, \text{Data-1-CF2})$$

For the second TNBC (European cohort) dataset:

$$\text{Data-2-CF1: } 39.8651 - 1.6945 * \text{NAT8L} - 3.5933 * \text{CACNG4}$$

$$\text{Data-2-CF2: } 9.8676 - 5.1333 * \text{MYCT1} + 0.4595 * \text{NUAK2}$$

$$\text{Data-2-CFmax: } \max(\text{Data-2-CF1}, \text{Data-2-CF2})$$

For the third (European cohort) dataset:

$$\text{Data-3-CF1: } 25.1089 - 10.1863 * \text{MYCT1} + 3.1654 * \text{NUAK2} - 2.0708 * \text{NAT8L}$$

$$\text{Data-3-CF2: } 2.4425 + 2.0119 * \text{UNC5B} - 4.1677 * \text{NUAK2} - 0.8255 * \text{NAT8L}$$

$$\text{Data-3-CFmax: } \max(\text{Data-3-CF1}, \text{Data-3-CF2})$$

For the fourth (Genomic Data Commons) dataset:

$$\text{Data-4-CF1: } 5.7644 - 2.5133 * \text{MYCT1} + 2.3383 * \text{NUAK2} - 1.2537 * \text{NAT8L}$$

$$\text{Data-4-CF2: } -9.5458 + 3.1219 * \text{UNC5B} + 0.7849 * \text{NUAK2}$$

$$\text{Data-4-CF3: } 7.0281 - 2.9389 * \text{MYCT1} + 4.3574 * \text{NUAK2} - 2.8591 * \text{NAT8L}$$

$$\text{Data-4-CFmax: } \max(\text{Data-4-CF1}, \text{Data-4-CF2}, \text{Data-4-CF3})$$

We note that the presentation of the final models above for different cohorts intends to visually show common genes in column-wise clusters.

The risk probabilities are calculated using the logistic function of $\exp(\text{Data-}i\text{-CFmax}) / (1 + \exp(\text{Data-}i\text{-CFmax}))$ for the combined classifiers in each dataset, or $\exp(\text{Data-}i\text{-CF}_j) / (1 + \exp(\text{Data-}i\text{-CF}_j))$ for each individual classifier $i = 1, 2, 3, 4, j = 1, 2, 3$.

In the first 3 cohorts, multiple ID-ref subtype genes correspond to a gene symbol. The following ID-ref subtype genes are used in the classifiers: 236872_at (RBM22), 241480_at (RNF213), 62987_r_at (CACNG4), 220471_s_at (MYCT1), 220987_s_at (NUAK2), 228880_at (NAT8L), 226899_at (UNC5B).

Table 1 lists gene expression values, individual classifiers' computed values, the combined classifier's computed values, and the risk probabilities. Figure 1 plots all patients' risk probabilities with circles for breast cancer samples and asters for non-breast cancer samples. Figure 2 is a Venn diagram that plots individual classifiers' performance.

This study is the first time TNBC and other breast cancer types can be further classified into subtypes based on critical genes' functions. This new classification opens a new research

Table 1. Gene information, expression values, competing factors, risk probabilities.

THE FIRST (TNBC) DATASET									
ID	TNBC/NO	236872_AT	241480_AT	62987_R_AT	CF1	CF2	CFMAX	PMAX	
GSM1974566	1	6.09	7.42	10.55	3.37	1.19	3.37	0.97	
GSM1974567	1	5.80	7.35	10.23	4.32	3.47	4.32	0.99	
.....									
GSM1974763	1	5.82	7.25	10.20	4.15	2.52	4.15	0.98	
GSM1978883	0	6.15	6.81	11.50	-2.01	-4.26	-2.01	0.12	
.....									
GSM1978948	0	6.17	6.69	11.92	-3.93	-5.47	-3.93	0.02	
GSM1978949	0	6.21	6.73	11.73	-3.11	-5.45	-3.11	0.04	
THE SECOND (TNBC) DATASET									
ID	TNBC/NO	220471_S_AT	220987_S_AT	228880_AT	62987_R_AT	CF1	CF2	CFMAX	PMAX
GSM1588970	1	2.37	5.26	5.47	8.08	1.58	0.13	1.58	0.83
GSM1588971	1	2.36	5.40	4.86	8.25	1.97	0.25	1.97	0.88
.....									
GSM1589150	0	2.42	4.98	8.83	8.64	-6.16	-0.26	-0.26	0.44
GSM1589151	0	2.46	5.02	8.95	8.57	-6.09	-0.47	-0.47	0.38
THE THIRD BREAST CANCER DATASET									
ID	BC/NOBC	220471_S_AT	226699_AT	220987_S_AT	228880_AT	CF1	CF2	CFMAX	PMAX
GSM1045191	0	3.21	6.95	2.94	5.65	-10.01	-0.49	-0.49	0.38
GSM1045192	0	2.90	7.88	3.57	4.27	-2.00	-0.11	-0.11	0.47
GSM1045193	0	3.16	6.57	3.19	4.67	-6.67	-1.49	-1.49	0.18
.....									
GSM1045207	0	2.31	6.68	3.41	7.61	-3.40	-4.63	-3.40	0.03
GSM1045208	1	2.31	8.56	4.02	4.43	5.08	-0.73	5.08	0.99
.....									
GSM1045310	1	2.31	7.84	5.54	4.67	9.40	-8.70	9.40	1.00
GSM1045311	1	2.31	9.91	4.09	4.85	4.43	1.35	4.43	0.99

(Continued)

Table 1. (Continued)

THE FOURTH BREAST CANCER DATASET												
ID	BC/NOBC	MYCT1	UNC5B	NUAK2	NAT8L	CF1	CF2	CF3	CFMAX	PMAX		
TCGA-E9-A1NI-01A	1	1.82	3.40	2.24	1.08	5.09	2.82	8.38	8.38	1.00		
TCGA-A1-A0SP-01A	1	1.48	2.74	1.27	3.16	1.05	-0.01	-0.82	1.05	0.74		
TCGA-BH-A1EU-11A	0	3.34	2.01	1.12	2.21	-2.79	-2.39	-4.23	-2.39	0.08		
.....												
TCGA-BH-A0BW-11A	0	2.86	2.85	2.12	1.98	1.05	1.03	2.20	2.20	0.90		
.....												
TCGA-BH-A0DK-11A	0	2.08	1.77	0.56	0.94	0.68	-3.57	0.69	0.69	0.67		
.....												
TCGA-E2-A1LH-11A	0	2.27	1.92	1.46	2.16	0.76	-2.39	0.53	0.76	0.68		
.....												
TCGA-AC-A2FF-11A	0	2.91	3.18	1.18	1.33	-0.46	1.31	-0.19	1.31	0.79		
.....												
TCGA-A7-A5ZW-01A	1	2.42	3.96	1.52	0.42	2.69	4.01	5.31	5.31	1.00		
TCGA-BH-A203-01A	1	1.74	2.87	1.44	0.79	3.77	0.54	5.94	5.94	1.00		

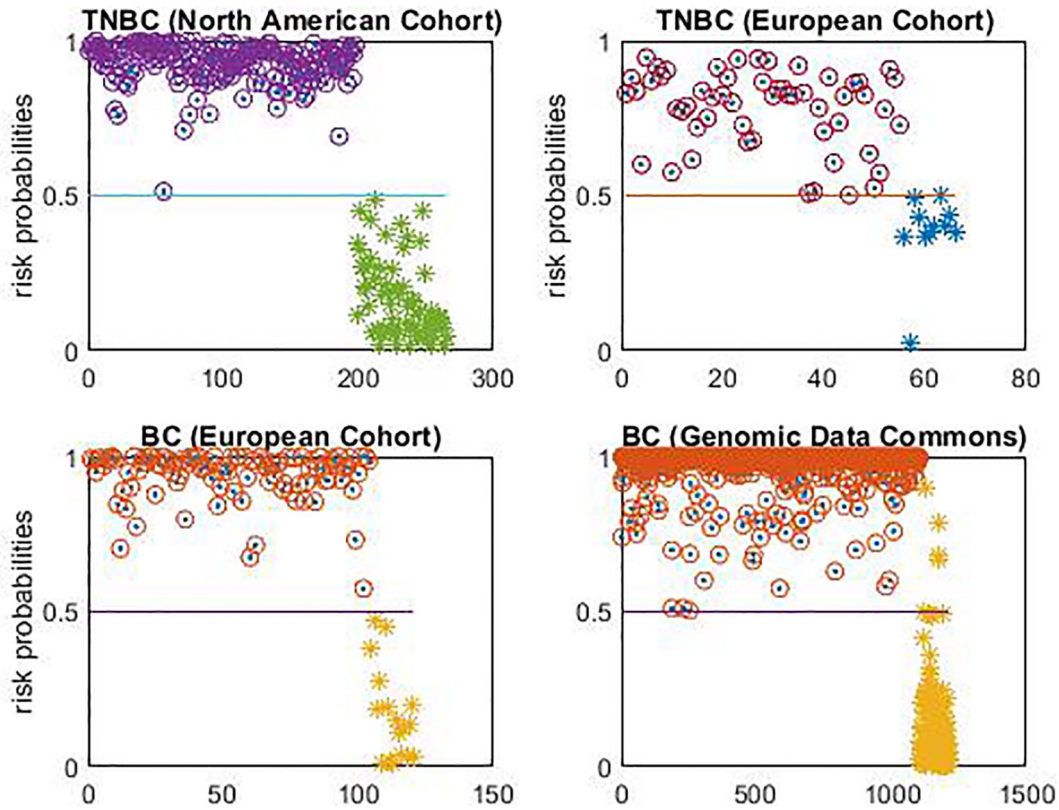


Figure 1. Risk probabilities of 4 cohorts. The circles are for patients with breast cancers. The asterisks are for tissues without breast cancers.

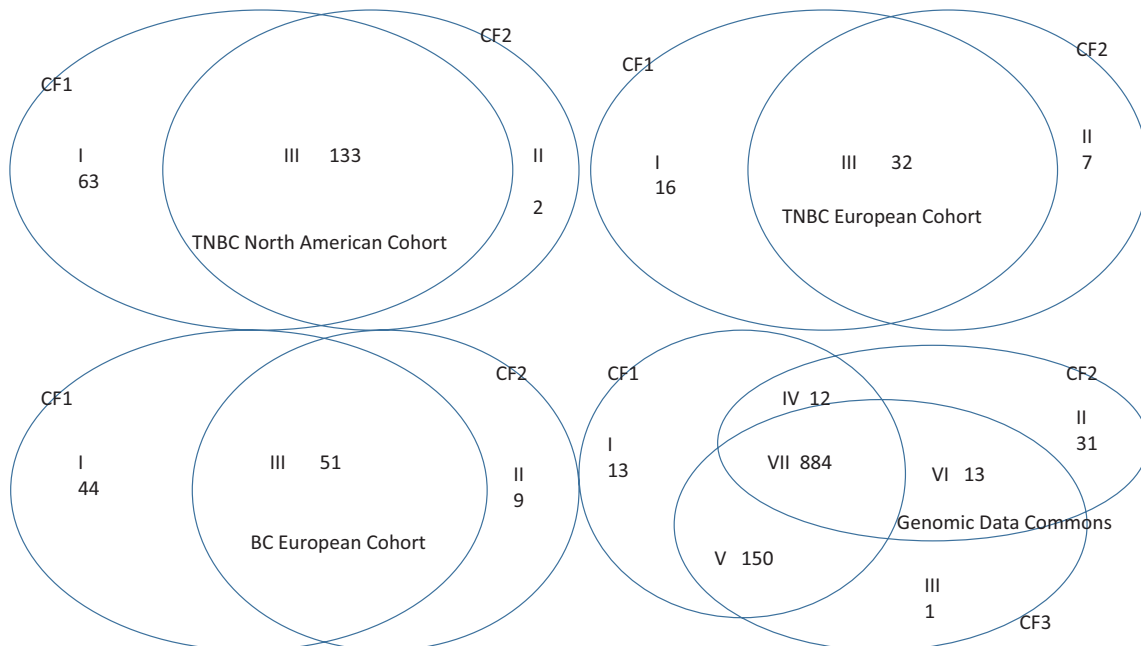


Figure 2. Venn diagrams of breast cancer subtypes. The first 3 cohorts have more than 3 subtypes. The fourth cohort has more than 7 subtypes.

direction, new drug developments, and new refined personalized therapies.

For the first TNBC (North American cohort) dataset, 3 genes (RNF213, RBM22, CACNG4) completely classify all 198 TNBC tumor samples into 3 subtypes (Figure 2) with the sensitivity of 100% and the specificity of 100%. From the individual classifiers, we can see that a decrease of RNF213 level

will reduce the risk of developing TNBC, while increases in the expression levels of RBN22 and CACNG4 will reduce the risk of developing TNBC.

For the second TNBC (European cohort) dataset, 4 genes (MYCT1, NAT8L, NUA2K, CACNG4) completely classify all 66 TNBC tumor samples into 3 subtypes (Figure 2) with the sensitivity of 100% and the specificity of 100%. From the

individual classifiers, we can see that a decrease in NUA2 level will benefit the patients, while increases in the expression levels of MYCT1, NAT8L, and CACNG4 will benefit the patients. We note that there are also Her2, Luminal A and Luminal B samples in this second dataset. After adding classifier CF3: $21.8170 - 8.8170 \cdot \text{RBM22} - 0.3047 \cdot \text{NAT8L}$, all breast cancer (TNBC, Her2, Luminal A, Luminal B) patients will again be 100% accurately classified into their respective groups.

Comparing the first and second TNBC cohorts, we see that the TNBC patients from North American and the TNBC patients from European cohorts share a common gene CACNG4 and similar coefficients (-3.6692 vs. -3.5933). Otherwise, other critical genes from these 2 cohorts are different. This observation tells that the causes, the formations, and the therapies of TNBC can be different from region to region and race to race. We want to note that based on our knowledge in the field, there does not exist any other method that can 100% accurately classify breast cancer patients and cancer-free patients into their respective groups. With 100% accuracy, regardless of how big and how small the sample is, these genes should contain basic cancer information of TNBC disease, they should be thoroughly analyzed and explored.

On the other hand, cautions should be called with any other classifiers with lower accuracy. Using genes derived/obtained from low accuracy classifiers may lead to suboptimal results and even wrong conclusions. The formulas of these 2 cohorts disclose the puzzle of TNBC as gene-gene interactions and functional effects are different. Such differences can be the most important part of studying TNBC and point out new research directions for better understanding TNBC and designing better treatments.

For the third (European cohort) dataset, 4 genes (MYCT1, NAT8L, NUA2, UNC5B) completely classify all 104 tumor samples into 3 subtypes (Figure 2) with a sensitivity of 100% and a specificity of 100%. A decrease of UNC5B level will benefit the patients in this cohort, while increases of expression levels of MYCT1 and NAT8L will benefit the patients. In addition, it can be seen that NUA2 can benefit the patients and can also harm the patients depending on the patients' breast cancer subtypes in Figure 2. These gene-gene relationships and genes-subtypes relationships tell efficient therapies to breast cancer patients depending on their subtypes' determinations.

For the fourth (Genomic Data Commons) dataset, the same 4 genes (MYCT1, NAT8L, NUA2, UNC5B) as for the third (European cohort) dataset completely classify 1104 tumor samples into 7 subtypes (Figure 2) with the sensitivity of 100% and the specificity of 96.5%. There are 4 samples among 103 normal samples being classified as tumor samples. Note that this dataset does not offer multiple ID-ref subtypes. If genes' expression values are taken the same as those ID-ref subtypes, the specificity may be improved to 100%. In this cohort, increases of MYCT1 and NAT8L levels can benefit the

patients, while decreases of UNC5B and NUA2 levels will benefit the patients.

Comparing the third and fourth breast cancer cohorts, the individual classifiers Data-3-CF1, Data-4-CF1, and Data-4-CF3 have the same component genes and coefficient signs. Data-3-CF2 has 1 more gene, NAT8L, than Data-4-CF2. However, the signs of NUA2 coefficients in these 2 individual classifiers are different. We further note that to have 2 similar individual classifiers Data-4-CF1 and Data-4-CF3 in the final classifier is completely new in machine learning literature. These observations further reveal that breast cancer formations are more complicated than simply looking at some high/low expression values of individual genes as in the literature. The most important relations in finding critical genes linked to breast cancers are gene-gene interactions, genes-individual classifiers interactions, and their functional effects.

Comparing the second TNBC cohort, the third breast cancer cohort, and the fourth cohort, we see that increasing the levels of MYCT1 and NAT8L can benefit all patients.

In Figure 2, Venn diagrams for 4 different cohorts are different, with 2 TNBC cohorts having similar patterns, while BC European cohort and Genomic Data Commons are different. It is because the numbers of component classifiers in the final classifiers for different cohorts are different. Such phenomena tell that there are commonalities among breast cancer patients and specificities from patient to patient, that is, the critical cancer informatics are expressed. Note that in Venn diagrams, the more intersections the groups, the more complex the disease, and the more difficult the treatment. Taking Genomic Data Commons as an example, patients in Group VII will be the most difficult to treat.

Figure 3 presents the gene-gene interactions, gene-subtype interactions, and functional effects of our identified competing classifiers. We can see clear signature patterns in each plot. This visualization tool provides a new way for breast cancer diagnosis.

Characteristics of studying samples

All 4 datasets are accompanied by some characteristics of patients. Here we report their inter-relationship with the competing classifiers. Table 2 displays Sex, Age, BMI, and Grade from the first dataset (TNBC, a North American cohort). Table 3 displays Age and BMI from the second dataset (TNBC samples only, A European cohort). Table 4 is for the third dataset (breast cancer samples, A European cohort). Finally, Table 5 includes disease Stage besides Age and Sex for the fourth breast sample data set (A Genomic Data Commons—TCGA).

Overall, these 4 tables show that more patients fall in groups related to more than 1 competing classifier. Obese patients can make TNBC more complex. The more the competing classifiers, the worse the grade. In Table 5, Stages (IV, X) are

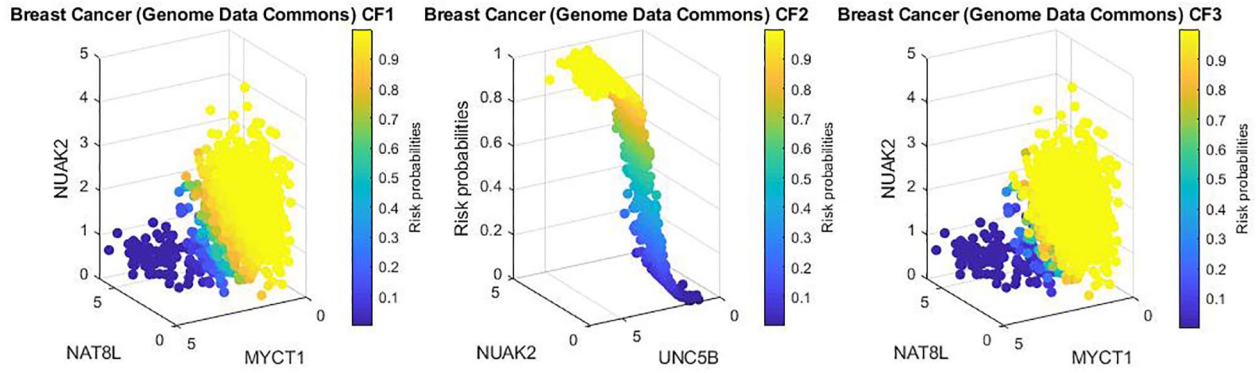


Figure 3. Four-dimensional plots for visualizing risk signature patterns from 3 competing component classifiers and the combined functional effects of gene-gene interactions and gene-subtype interactions of 4 genes.

Table 2. Characteristics of the first dataset samples (TNBC, A North American Cohort).

	SEX		AGE				BMI			GRADE		
	MALE	FEMALE	≤50	(50,60]	(60,70]	>70	NORMAL	OVERWEIGHT	OBESE	POORLY	MODERATELY	WELL
CF-1	0	63	34	12	10	5	13	19	61	33	25	1
CF-2	0	2	1	0	1	0	0	1	2	1	0	0
CF-(1,2)	0	133	51	33	28	18	36	41	129	74	28	3

Table 3. Characteristics of the second dataset samples (TNBC, A European Cohort).

	AGE				BMI		
	≤50	(50,60]	(60,70]	>70	NORMAL	OVERWEIGHT	OBESE
CF-1	3	4	2	2	5	3	3
CF-2	2	2	0	1	2	2	1
CF-(1,2)	9	10	5	1	12	6	6

Table 4. Characteristics of the third dataset samples (BC, A European Cohort).

	AGE				GRADE		
	≤50	(50,60]	(60,70]	>70	POORLY	MODERATELY	WELL
CF-1	9	15	11	9	25	14	5
CF-2	3	4	2	0	4	5	0
CF-(1,2)	15	13	10	13	24	21	6

mainly related to CF-(1,2,3), which shows the classifiers are positively correlated.

Discussions

This study is the first time in the medical literature that breast cancer diseases can be classified almost 100% correctly using only a few (3 or 4) genes. The results clearly disclose the puzzle of breast cancers, including TNBC, due to the selected genes and their predicting powers through gene-gene interaction, gene-subtype interaction, and functional effects. The results also point to new treatment directions.

We note that this study does not use the primary endpoint information. It is a pure classification study. The main purpose is to identify the essential breast cancer informatics. The study has achieved 100% accuracy, which is the first in the literature. Given patients have different endpoints, the new classifier still reaches 100% accuracy, which means the classifier is robust to patients’ disease states, that is, we can conclude that the classifier is robust regardless of primary endpoints and other individual attributes.

The discovery of critical genes can motivate many new research directions and laboratory experiments. These critical

Table 5. Characteristics of the fourth dataset samples (TCGA, Genomic Data Commons).

	AGE					SEX		STAGE				
	≤50	(50,60]	(60,70]	(70,80]	>80	MALE	FEMALE	I	II	III	IV	X
CF-1	8	1	3	0	1	0	13	4	6	3	0	0
CF-2	10	9	7	2	1	1	30	6	16	9	0	0
CF-3	1	0	0	0	0	0	1	0	0	1	0	0
CF-(1,2)	5	1	4	1	1	0	12	2	7	2	0	1
CF-(1,3)	44	36	35	27	7	2	148	25	97	24	2	0
CF-(2,3)	2	5	5	1	0	0	13	3	3	5	1	1
CF-(1,2,3)	257	216	225	119	43	9	862	143	490	202	17	10

Table 6. Correlation coefficients between CFmax from the fourth data and 8 genes in the literature.

	CFMAX	BRCA1	BRCA2	PALB2	BARD1	RAD51C	RAD51D	ATM	CHEK2
CFmax	1.00	.25	.25	.24	.31	.12	.13	-.12	.21
BRCA1		1.00	.48	.52	.50	.26	.50	.14	.28
BRCA2			1.00	.47	.61	.15	.24	.28	.45
PALB2				1.00	.53	.19	.32	.30	.20
BARD1					1.00	.17	.27	.21	.41
RAD51C						1.00	.19	-.15	.17
RAD51D							1.00	.14	.14
ATM								1.00	-.15
CHEK2									1.00

genes and their derived signature patterns (individual classifiers) can be a starting point as new biomarkers for conducting gene network analysis, testing other reported genes, and finding the causal directions of gene expression in various projects. As a result, many other existing pieces of research can be enriched. It can also be hoped that new types of diseases can be discovered. Eventually, new testing procedures and therapies for breast cancer can be designed.

These critical genes enrich the biological literature of their new functions related to breast cancer from indirect relationship to direct relationship. In many scenarios, indirect effects are more significant than direct effects as direct effects can be seen and controlled, while indirect effects are hard to see and even not to say how to control.

In the introduction, 8 genes, BRCA1, BRCA2, PALB2, BARD1, RAD51C, RAD51D, ATM, and CHEK2, were discussed as they are potentially helpful. We found that in terms of detecting power in diagnoses and breast cancer risk prediction, these 8 genes are not significant (even inferior) compared with the genes presented in this paper. In Table 6 below, we use the fourth dataset (cohort) to present the linear correlation coefficients between our final classifier CFmax and each of these 8 genes and among these 8 genes.

From Table 6, we can immediately see that the correlation coefficient between each of these 8 genes and the CFmax is low, for example, between CFmax and PALB2 is 0.24. With a correlation coefficient of 0.24 and given the CFmax's super detecting power of 100% sensitivity and 96.5% specificity, it is likely PALB2 can be just around 40% to 60% of overall detecting power or even lower. Furthermore, Berger⁴ reported unlike BRCA1 and BRCA2, which are often found in the Ashkenazi Jewish population, PALB2 is not associated with the Ashkenazi group. Some studies have found a PALB2 association with Finnish and French Canadian and Greek women, but experts say more research is needed. This phenomenon is interesting; however, it highlights the more significant uncertainty of applying those 8 genes in practice. In contrast, the genes identified in this paper lead to the highest accuracy, perfect or nearly perfect. As a result, more focuses should be paid to the genes discovered in this paper.

The risk probability of a patient developing a specific type of breast cancer in her/his life is low. Among all discovered breast cancer types, growing more than 1 type of breast cancer is rare. These breast cancer types compete, and 1 type will first be diagnosed. As a result, the competing risk factor models can efficiently model multiple breast cancer types.

This study's inference/analysis approach can shed new light on all gene-related research, that is, not just the breast cancer study. Researchers can apply max-linear type models in their studies. Ultimately, our new findings may make researchers' cancer research efforts more effective and meaningful, reduce substantial research costs, and save lives and protect people.

Finally, we address an important medical practice issue. In this paper, all classifier formulas are explicitly expressed. Thus, the results in all tables are reproducible. Furthermore, Figures 1 and 3 show the risks of all patients. Using this paper's results, medical doctors have a powerful tool (testing kit) in their daily work, that is, in the diagnostic stage, diagnosing and analyzing patients' breast cancer risks based on the 4 or fewer critical genes' expression values and the computed risks; in the treatment stage, those signature patterns can be used to study the effectiveness of drugs and treatments, that is, conduct clinical trials, for example, survival analysis, based on classified groups; in the drug development stage, pharmaceutical companies can use the findings of critical genes to study new drugs; finally, it can be hoped that mRNA-based therapies can be introduced using the critical genes' information in the therapy stage.

Acknowledgements

The author thanks Editor-in-Chief Jimmy Efirid and two anonymous reviewers for their insightful comments and questions.

ORCID iD

Zhengjun Zhang  <https://orcid.org/0000-0003-2615-1539>

Data Availability

The datasets are publicly available. The data links are stated in Section Data Description.

Supplementary Materials

Real data and computer outputs are in a supplementary file available online and submitted together with this paper. A MATLAB® demo code for solving a final dataset example in equation (4) ($\lambda_2 = 0$) is also available.

REFERENCES

- Narod SA. Which genes for hereditary breast cancer? *New Engl J Med.* 2021;384:471-473.
- Dorling L, Carvalho S, Allen J, et al. Breast cancer risk genes — association analysis in more than 113,000 women. *New Engl J Med.* 2021;384:428-439.
- Hu C, Hart SN, Gnanaolivu R, et al. A population-based study of genes previously implicated in breast cancer. *New Engl J Med.* 2021;384:440-451.
- Berger S. This breast cancer gene is less well known, but nearly as dangerous. *New York Times.* 2021. <https://www.nytimes.com/2021/08/17/health/breast-cancer-palb2-brca.html>
- de Magalhães JP. Every gene can (and possibly will) be associated with cancer. *Trends Genet.* 2021. doi:10.1016/j.tig.2021.09.005
- Robitzski D. Q&A: Nearly every single human gene can be linked to cancer. *The Scientist.* 2021. <https://www.the-scientist.com/news-opinion/q-a-nearly-every-single-human-gene-can-be-linked-to-cancer-69365>
- Amjad E, Asnaashari S, Sokouti B, Dastmalchi S. Systems biology comprehensive analysis on breast cancer for identification of key gene modules and genes associated with TNM-based clinical stages. *Sci Rep.* 2020;10:10816.
- Malvia S, Bagadi SAR, Pradhan D, et al. Study of gene expression profiles of breast cancers in Indian women. *Sci Rep.* 2019;9:10018.
- Lv X, He M, Zhao Y, et al. Identification of potential key genes and pathways predicting pathogenesis and prognosis for triple-negative breast cancer. *Cancer Cell Int.* 2019;19:172.
- Zhong G, Lou W, Shen Q, Yu K, Zheng Y. Identification of key genes as potential biomarkers for triple-negative breast cancer using integrating genomics analysis. *Mol Med Rep.* 2020;21:557-566.
- Lin Y, Fu F, Lv J, et al. Identification of potential key genes for HER-2 positive breast cancer based on bioinformatics analysis. *Medicine.* 2020;99:e18445.
- Chen H, Wu J, Zhang Z, et al. Association between *BRCA* status and triple-negative breast cancer: a meta-analysis. *Front Pharmacol.* 2018;9:909. Published 2018 Aug 21.
- Deng J-L, Xu Y-H, Wang G. Identification of potential crucial genes and key pathways in breast cancer using bioinformatic analysis. *Front Genet.* 2019;10:695.
- Zhu C, Ge C, He J, Zhang X, Feng G, Fan S. Identification of key genes and pathways associated with irradiation in breast cancer tissue and breast cancer cell lines. *Dose Resp.* 2020;18:1559325820931252.
- Dong P, Yu B, Pan L, Tian X, Liu F. Identification of key genes and pathways in triple-negative breast cancer by integrated bioinformatics analysis. *Biomed Res Int.* 2018;2018:2760918.
- Lu X, Gao C, Liu C, et al. Identification of the key pathways and genes involved in HER2-positive breast cancer with brain metastasis. *Pathol Res Pract.* 2019;215:152475.
- Teng H, Zhang Z. Directly and simultaneously expressing absolute and relative treatment effects in medical data models and applications. *Entropy.* 2021;23:1517.
- Cui Q, Zhang Z. Max-linear competing factor models. *J Bus Econ Stat.* 2018;36:62-74.
- Cui Q, Xu Y, Zhang Z, Chan V. Max-linear regression models with regularization. *J Econom.* 2021;222:579-600.
- Xu Y. *Regression Models With Max-Linear Structure.* PhD dissertation. University of Wisconsin at Madison; 2019.
- Zhang Z. Five critical genes related to seven covid-19 subtypes: a data science discovery. *Data Sci J.* 2021;19:142-150.
- Malinowski A, Schlather M, Zhang Z. Intrinsically weighted means and non-ergodic marked point processes. *Ann Inst Stat Math.* 2016;68:1-24.
- Zhang Z. On studying extreme values and systematic risks with nonlinear time series models and tail dependence measures (with discussion). *Stat Theory Relat Fields.* Published online December 23, 2020. doi:10.1080/24754269.2020.1856590.
- Zhang Z. Functional effects of four or fewer critical genes linked to lung cancers and new subtypes detected by a new machine learning classifier. *J Clin Trials.* 2021;S14:001.
- Burstein MD, Tsimelzon A, Poage GM, et al. Comprehensive genomic analysis identifies novel subtypes and targets of triple-negative breast cancer. *Clin Cancer Res.* 2015;21:1688-1698.
- den Hollander P, Rawls K, Tsimelzon A, et al. Phosphatase PTP4A3 promotes triple-negative breast cancer growth and predicts poor patient survival. *Cancer Res.* 2016;76:1942-1953.
- Maire V, Baldeyron C, Richardson M, et al. TTK/hMPS1 is an attractive therapeutic target for triple-negative breast cancer. *PLoS One.* 2013;8:e63712.
- Maire V, Némati F, Richardson M, et al. Polo-like kinase 1: a potential therapeutic option in combination with conventional chemotherapy for the management of patients with triple-negative breast cancer. *Cancer Res.* 2013;73:813-823.
- Maubant S, Tesson B, Maire V, et al. Transcriptome analysis of Wnt3a-treated triple-negative breast cancer cells. *PLoS One.* 2015;10:e0122333.
- Clarke C, Madden SF, Doolan P, et al. Correlating transcriptional networks to breast cancer survival: a large-scale coexpression analysis. *Carcinogenesis.* 2013;34:2300-2308.