Precision Cancer Classification and Biomarker Identification from mRNA Gene Expression via Dimensionality Reduction and Explainable AI

Farzana Tabassum^{a,e}, Sabrina Islam^{a,e}, Siana Rizwan^{a,c}, Masrur Sobhan^{a,b}, Tasnim Ahmed^{a,c}, Sabbir Ahmed^a, Tareque Mohmud Chowdhury^{a,d}

^aDepartment of Computer Science and Engineering, Islamic University of
Technology, Board Bazar, Gazipur, 1704, Bangladesh

^bKnight Foundation School of Computing and Information Sciences (KFSCIS), Florida
International University, Miami, FL 33199, Florida, USA

^cSchool of Computing, Queen's University, Kingston, K7L 3N6, Ontario, Canada

^dtareque@iut-dhaka.edu

^e(these authors contributed equally)

Abstract

Gene expression analysis is a critical method for cancer classification, enabling precise diagnoses through the identification of unique molecular signatures associated with various tumors. Identifying cancer-specific genes from gene expression values enables a more tailored and personalized treatment approach. However, the high dimensionality of mRNA gene expression data poses challenges for analysis and data extraction. This research presents a comprehensive pipeline designed to accurately identify 33 distinct cancer types and their corresponding gene sets. It incorporates a combination of normalization and feature selection techniques to reduce dataset dimensionality effectively while ensuring high performance. Notably, our pipeline successfully identifies a substantial number of cancer-specific genes using a reduced feature set of just 500, in contrast to using the full dataset comprising 19,238 features. By employing an ensemble approach that combines three top-performing

Email addresses: farzana@iut-dhaka.edu (Farzana Tabassum), sabrinaislam22@iut-dhaka.edu (Sabrina Islam), sianarizwan@iut-dhaka.edu, siana.rizwan@queensu.ca (Siana Rizwan), msobh002@fiu.edu (Masrur Sobhan), tasnimahmed@iut-dhaka.edu,tasnim.ahmed@queensu.ca (Tasnim Ahmed), sabbirahmed@iut-dhaka.edu (Sabbir Ahmed), tareque@iut-dhaka.edu (Tareque Mohmud Chowdhury)

classifiers, a classification accuracy of 96.61% was achieved. Furthermore, we leverage Explainable AI to elucidate the biological significance of the identified cancer-specific genes, employing Differential Gene Expression (DGE) analysis.

Keywords: RNA sequence data, Shap analysis, Differential Gene Expression, Feature Reduction, Explainable AI, Model Ensemble

1. Introduction

Cancer, a complex disease marked by the uncontrolled proliferation of cells [1, 2], remains a leading cause of mortality worldwide. In 2022, it accounted for nearly 10 million deaths [3]. Gene expression analysis has emerged as a pivotal tool in this pursuit. Researchers have directed their efforts toward identifying potential biomarkers to tackle the challenges associated with cancer diagnosis and drug discovery. In response to the challenges associated with cancer diagnosis and drug discovery, researchers have focused on identifying potential biomarkers. encompassing the transcription of DNA into messenger RNA (mRNA) and the subsequent translation into proteins [4, 5]. In recent years, Machine Learning (ML) methods have become increasingly popular for cancer cell classification due to their effectiveness in identifying important features and the availability of data from high-throughput machines [6].

In the domain of Cancer classification with RNA data, Podolsky et al. conducted an evaluation of four publicly available lung cancer datasets from various institutions, utilizing seven ML-based techniques [7]. Their findings indicated that performance varied significantly across the datasets, with the k-nearest neighbor (KNN) [8] and support vector machine (SVM) [9] methods achieving higher area under the curve (AUC) values [10]. Additionally, experiments employing deep learning (DL) methodologies have been performed on mRNA datasets. Lyu et al. implemented a DL-based approach to classify tumor types using gene expression data derived from 33 tumors in the Pan-Cancer Atlas [11]. Their classification model incorporated three convolutional layers with max-pooling and batch normalization [12], followed by three fully connected layers.

Laplante et al. developed a deep neural network model to infer cancer locations using 27 miRNA data cohorts from The Cancer Genome Atlas (TCGA), categorized into 20 anatomical sites [13]. The authors employed log transformation and MinMaxScaler for data preprocessing, followed by a classification model built upon a six-layer artificial neural network (ANN)

[14]. Hsu et al. utilized TCGA RNA-sequencing data to classify 33 distinct types of cancer [15]. Their methodology included feature selection techniques such as tree classifiers [16] and variance thresholding [17], along with two data normalization approaches: min-max scaling and standard scaling. They showed a performance comparison using a decision tree, KNN, linear SVM (SVM) [18], polynomial SVM, and ANN models where linear SVM outperformed other classifiers. Mostavi et al. introduced three Convolutional Neural Network (CNN) models [19] for the classification of tumor versus non-tumor samples [20]. Their study indicated that the 1D-CNN model exhibited superior computational efficiency. Mahin et al. proposed PanClassif, a method designed to enhance the performance of machine learning classifiers in cancer identification by utilizing a limited number of efficient genes extracted from RNA-seq data [21]. Their approach incorporated six classifiers: SVM (linear kernel), SVM (RBF kernel), Random Forest [22], Neural Network, k-Nearest Neighbor (KNN), and AdaBoost [23]. Notably, none of these studies focused on the analysis and identification of significant gene sets or biomarkers specific to each cancer type within their datasets.

Focusing on the issue of identifying biomarkers for cancers, researchers conducted experiment with TCGA RNAseq data containing 33 types of cancer created a heatmap to represent gene scores, and used guided backpropagation and Grad-CAM [24] to analyze genes [25]. A gene functional classification method was utilized to evaluate 400 biomarkers and annotate genes by function similarity. The functional analysis results were compared to the cohort cancer types' relevant pathways using p-values for correlation with major gene pathways and biomarkers. In another study, Lopez-Rincon, A. et al. proposed an ensemble approach to extract 100 miRNA for multi-class cancer classification [26]. According to their findings, Logistic Regression [27] performed best across all of their experiments. They then conducted a bibliographical meta-analysis to confirm their findings.

While several studies have aimed to identify cancer biomarkers, most have relied on ML and DL-based models as 'black boxes' [28]. Recently, various approaches have been developed to explain these black-box models, including Shapley Sampling [29], Relevance Propagation [30], LIME [31], ANCHOR[32], DeepLIFT [33], etc. These explainability analysis techniques have widely been used in a wide spectrum of ML-based and DL-based tasks to enhance model transparency and trust, enabling researchers to leverage complex models while ensuring accountability and understanding [34]. In 2017, a game-theoretic approach known as SHAP (SHapley Additive ExPlanation)

[35] was introduced to provide interpretability to black box models. Levy et al. utilized SHAP to identify significant methylation states across different cell types and cancer subtypes [36]. Yap et al. utilized SHAP to provide explanations for a DL-based model that employed RNA-sequence data for cancer tissue classification [37]. Divate et al. used SHAP to interpret the classification model and identify specific gene signatures that contributed to classifying different cancer types [38]. However, they only provided a certain number of features for the combination of all kinds of cancer.

Additionally, statistical validation in studying cancer-specific gene sets is essential to confirm their biological significance and eliminate the possibility of random chance. Differential Gene Expression (DGE) analysis is a frequently employed approach that identifies genes exhibiting significant expression level variations across various cancer types [39, 40]. Ni et al. implemented DGEs to identify potential genes associated with the pathogenesis and prognosis of lung cancer, suggesting that DGEs can identify important genes across the samples [41].

To address the limitations of existing studies, which includes the lack of focus on identifying significant gene sets for specific cancer types, dimentionality curse of RNAseq data and the black-box nature different classifiers, in this paper, we proposed an ensemble architecture to efficiently and accurately identify 33 different cancer types and their corresponding gene sets using mRNA gene expression data. We utilized a normalization and feature selection technique to enhance convergence efficiency and address the high dimensionality in gene expression data. Additionally, we implemented an interpretable machine learning tool to determine the significant global feature contribution responsible for the models' high performance and accuracy. The proposed approach demonstrates competence in achieving high accuracy with a limited number of 500 features only.

2. Methods

Our study comprised two main portions: 1) Classification and 2) Identification of cancer-specific gene sets for each cancer type. In the classification stage, we addressed the challenge of high dimensionality in the dataset and aimed to reduce noise before applying classifiers. To achieve this, we explored various combinations of feature selection and normalization techniques. Additionally, we evaluated the performance of different standalone classifiers and implemented two ensemble techniques to enhance classification accuracy. To

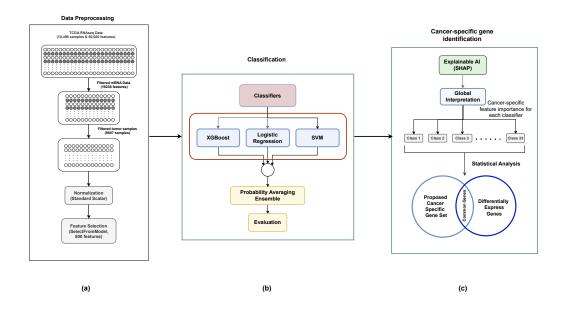


Figure 1: Overview of the proposed pipeline. (a) Data preprocessing steps include selecting mRNA data from the TCGA dataset and reducing it to the top 500 mRNA genes using data normalization and feature selection. (b) Performed probability averaging ensemble technique by combining top 3 performed classifiers. (c) Analyzed model's performance using explainable AI (XAI) to identify cancer-specific gene sets and validated it through statistical validation.

ensure a robust performance analysis for cancer classification, we employed Stratified 5-Fold cross-validation on the dataset. This approach was used to minimize bias and enhance the reliability of our results. In the second stage of our study, SHAP [42] was utilized to investigate the model explainability and identify significant genes specific to each cancer type. To assess the biological significance of our identified gene sets, we performed Differential Gene Expression (DGE) [43] analysis.

After a thorough analysis, we proposed a final pipeline that yields the best results and a substantial number of biologically significant genes. The pipeline is visually represented in Fig. 1.

2.1. Dataset

We conducted the experiments on the widely used pan-cancer dataset which was downloaded from the TCGA (The Cancer Genome Atlas Program) Data Portal [44]. TCGA gene expression data generated using RNA-Seq is quantified using the FPKM (Fragments Per Kilobase of transcript per Million mapped reads) [45] metric and tools like TOIL (Transcriptome-Oriented Incremental Learning) and RSEM (RNA-Seq by Expectation Maximization) [46] are used to analyze and interpret this data. TOIL, RSEM, and FPKM data estimate the expression level of genes based on the sequencing reads of mRNA transcripts. In this study, we worked with RNA-seq FPKM data. There are a total of 10,496 patient samples along with their 60,499 genes. The Pareto principle [47], also known as the 80/20 rule or the principle of factor sparsity, was applied to divide the sample data into two sets: 80% for training and 20% for testing [48]. A validation set, which constitutes 20% of the training set, is employed for the entire dataset. Tab. 1 presents the total number of samples in each cancer type.

Table 1: Number of samples in each cancer type from the TCGA Dataset [44]

#	Cancer Type	Count
1	Adrenocortical carcinoma (ACC)	77
2	Bladder Urothelial Carcinoma (BLCA)	426
3	Breast invasive carcinoma (BRCA)	1211
4	Cervical squamous cell carcinoma endocervical adenocarcinoma (CESC)	309
5	Cholangiocarcinoma (CHOL)	45
6	Colon adenocarcinoma (COAD)	329
7	Diffuse large B cell lymphoma (DLBCL)	47
8	Esophageal carcinoma (ESCA)	195
9	Glioblastoma multiforme (GBM)	165
10	Head and Neck squamous cell (HNSC)	564
11	Kidney Chromophobe (KICH)	91
12	Kidney renal clear cell carcinoma (KIRC)	603
13	Kidney renal papillary cell carcinoma (KIRP)	321
14	Acute myeloid leukemia (LAML)	173
15	Brain Lower Grade Glioma (LGG)	522
16	Liver hepatocellular carcinoma (LIHC)	421
17	Lung adenocarcinoma (LUAD)	574
18	Lung squamous cell carcinoma (LUSC)	548
19	Mesothelioma (MESO)	87
20	Ovarian serous cystadenocarcinoma (OV)	427
21	Pancreatic adenocarcinoma (PAAD)	183
22	Pheochromocytoma and Paraganglioma (PCPG)	185
23	Prostate adenocarcinoma (PRAD)	548
24	Rectum adenocarcinoma (READ)	102

Table 1: Number of samples in each cancer type from the TCGA Dataset [44]

#	Cancer Type	Count
25	Sarcoma (SARC)	264
26	Skin Cutaneous Melanoma (SKCM)	470
27	Stomach adenocarcinoma (STAD)	450
28	Testicular Germ Cell Tumors (TGCT)	137
29	Thyroid carcinoma (THCA)	571
30	Thymoma (THYM)	121
31	Uterine Corpus Endometrial Carcinoma (UCEC)	194
32	Uterine Carcinosarcoma (UCS)	57
33	Uveal Melanoma (UVM)	79
	Total	10,496

2.2. Data Preprocessing

This section explores the data preprocessing steps used for the classifications of 33 types cancers. This study utilized the cancerous mRNA samples to capture the genetic profile of malignancies, normalization techniques to address differences in data magnitude, and feature selection methods aimed at reducing dimensionality and identifying the most relevant genes for cancer classification.

2.2.1. Selecting Gene Identifiers and Excluding Normal Tissue Samples

Initially, the dataset contained 10,496 patient tissue samples and each sample contained 60,499 gene expressions. In this study, the examination was restricted solely to mRNA characteristics. Following the exclusion of non-mRNA genes, the dataset comprised a total of 19,238 mRNA genes. Additionally, as one of the main objectives of this study was to identify cancer-specific gene sets, genes that exhibited predominant expression in tumor tissues while being absent or minimally expressed in normal tissue samples were considered. It was imperative to mitigate any potential impact that the existence of non-cancerous tissues may exert on the genetic profile of particular malignancies. The exclusion of normal tissue samples from the analysis reduced the likelihood of encountering false positive or negative results. After excluding all normal tissue samples, the total number of samples was 9807.

2.2.2. Normalization

When data exhibit significant differences in magnitudes, it can adversely affect the training of models. This can result in slower convergence or even complete failure of the model to converge. To address this issue, data normalization techniques can be employed [49, 50, 51]. Previous studies have utilized various approaches such as standard scaler, z-score, min-max scaler, and feature vector data normalization approaches [52, 53, 54]. In our study, we specifically implemented two data normalization techniques: Min-Max Scaler [55, 56] and Standard Scaler [57, 58, 53] due to their superior performance [59, 60].

The min-max scaler involves scaling data to a predetermined range, often spanning from 0 to 1. It is susceptible to outliers and may lead to a reduction of data integrity in the tails of the data distribution. Conversely, the standard scaler is a normalization technique that converts the data to possess a mean of zero and a variance of one. This normalization technique exhibits greater resilience in handling data with unknown ranges and is comparatively less susceptible to the influence of outliers. These techniques were applied as preprocessing steps to the dataset before performing classification tasks. Our objective was to assess the impact of these techniques on the performance of the classification models and choose the most suitable one.

2.2.3. Feature Selection

The large dimensionality of the data makes it essential to explicitly specify features while classifying cancer [61]. Hence, in previous studies they have implemented different feature selection approaches to reduce data dimensionality by identifying pertinent genes that are linked to cancer and enhancing the efficacy of classification models as well [26, 62, 63, 64, 65, 61]. In our study, we employed four distinct embedded feature selection methods [66], namely SelectKBest [67], ElasticNet [68], Lasso [63, 69], and SelectFromModel [62, 70], to ascertain the most pertinent features for the categorization of 33 diverse cancer types. It is imperative to clarify that the term "feature" in this context pertains to the mRNA genes present in the dataset.

The SelectKBest method selects the top k features using univariate statistical measures like chi-square or mutual information. The study used the $ANOVA\ F-value[71]$ scoring function for selection (Eq. 1).

$$F = \frac{\sum_{i} n_{i} (\overline{y_{i}} - \overline{y})^{2} / (k - 1)}{\sum_{ij} (y_{ij} - \overline{y_{i}})^{2} / (N - k)}$$

$$(1)$$

where.

k = total number of groups,

N = overall sample size

The ElasticNet (Eq. 2) algorithm is popular in bioinformatics due to its ability to create parsimonious models with minimal non-zero weights [72]. It balances precision and weight magnitude through L1 and L2 regularization [73, 74]. In our experiment, we set the L1 ratio to 0.4 and selected top genes based on the number of features.

$$F = \frac{1}{m} \left[\sum_{l=1}^{m} (y^{(i)} - h(x^{(i)}))^2 + \lambda_1 \sum_{j=1}^{n} w_j + \lambda_2 \sum_{j=1}^{n} w_j^2 \right]$$
 (2)

where,

 w_i : the weight for the j^{th} feature,

n: the number of features in the dataset,

 λ_1 : the regularization strength for the L_1 norm,

 λ_2 : the regularization strength for the L_2 norm.

The LASSO (Eq. 3) technique is a regression analysis approach that incorporates variable selection and regularization to enhance the predictive accuracy and interpretability of the resulting statistical model [75]. We utilized an *alpha* value of 0.5 to identify the specified number of genes.

Objective Function = Mean Squared Error +
$$\alpha \sum_{i} |w_{i}|$$
 (3)

where,

 α : the hyperparameter controlling the strength of the L1 regularization term,

 w_i : weights (coefficients) assigned to each feature.

The SelectFromModel feature selection approach relies on a specified estimator to identify the most significant features. This approach has the advantage of autonomously identifying significant features without requiring extensive prior expertise. The function receives an estimator and a pre-defined limit on the number of features to be selected.

The techniques were applied to the dataset that underwent standard scaler normalization to determine the top 100, 250, 500, 750, and 1000 genes that are the most significant to identify each cancer type. Subsequently, the chosen genes were employed in a classification framework to categorize the 33 distinct cancer types present in the dataset.

2.3. Standalone Classifiers

In the case of cancer classification, the choice of appropriate classifiers can reduce the dimensionality curse and achieve higher accuracy [21, 25, 76, 77, 78, 79, 80]. In order to identify the optimal classifiers for our dataset in cancer classification, we explored seven distinct methodologies: Random Forest [22], XGBoost [81], Logistic Regression [27], SVM [9], MLP [82], 1-D CNN [83], and TabNet [84]. This approach allowed us to investigate different sectors, including traditional machine learning approaches, decision tree-based approaches, deep learning, and transfer learning. By doing so, we aimed to determine which type of classifier performs best on our dataset. To maintain consistency in the comparison, the parameters used for model fitting were held constant throughout all experiments.

Logistic Regression: Logistic regression estimates class probabilities by fitting binary regression models and learned feature weights, predicting the class with the highest probability. In our logistic regression model, we implemented 100 maximum iterations to optimize the loss function and achieve convergence and L2 or ridge regularization [85] to handle outliers and model complexity. Additionally, to avoid overfitting, regularization [86] strength is set to 100.

Support Vector Machines (SVM): SVM is applied to accurately categorize gene expression data by determining an optimal hyperplane with a maximum margin and identifying support vectors closest to the decision boundary [87, 88]. In our study, a tolerance of 1e-5 was set to ensure precise convergence during the optimization process. The "one-vs-one" [89] decision function was implemented to create binary classifiers for each class pair and use a voting technique for final predictions.

XGBoost: XGBoost exhibits exceptional performance in structured tabular data classification by leveraging a gradient-boosting [90] algorithm. It optimizes performance through parallel processing, tree pruning, handling missing values, and regularization, effectively mitigating bias and overfitting. For our experimental setup, we limited the maximum depth of each decision tree to four. A learning rate of 0.1 was used to control the step size during the

optimization process. Additionally, we configured the number of estimators to be 1000. When the difference between accuracy and loss remained consistent for 10 consecutive epochs, we stopped the training procedure. These values were employed to minimize overfitting, reduce training time, and improve accuracy in the model.

Random Forest: The Random Forest algorithm constructs multiple decision trees by randomly selecting subsets of features and samples. The number of estimators was set to 20 in our classification step to achieve optimal accuracy while avoiding overfitting the model.

Multilayer Perceptron (MLP): MLP is a feedforward [91] neural network with input, hidden, and output layers. In gene expression analysis, the input is gene expression profiles, and the output layer predicts class probabilities [92, 93, 94]. In our study, we applied three hidden layers with 100 neurons each. An alpha value of 0.001 was used to balance the capturing of data patterns and prevent overfitting. A learning rate of 0.001 ensured gradual convergence. The Rectified Linear Unit (ReLU) [95] activation function captured non-linear associations. Weight optimization was performed using the 'Adam' optimizer [96]. These hyper-parameter choices aimed to balance model complexity and generalization to new data.

1D-CNN (1-Dimensional Convolutional Neural Network): For predicting cancer types, various convolutional neural network (CNN) models have been proposed [20, 97]. Through experimentation, we found that using three fully connected layers with 512, 256, and 128 nodes respectively helped us achieve global minima in terms of loss. Batch Normalization was added between layers to standardize input and stabilize learning, reducing training epochs for faster convergence. The ReLU activation function was used to improve optimization and generalization. A dropout layer between the two layers worked as a regularizer to avoid overfitting. To downsample data average pooling was applied. The final output layer of the model was a densely connected layer with a softmax activation function [98] and had 33 nodes representing the class labels. Sparse categorical cross-entropy loss, Adam optimizer with a learning rate of 0.001, and accuracy as the evaluation metric were employed.

Attentive Interpretable Tabular Learning neural network (Tab-Net): The effectiveness of the transformer architecture for cancer classification [99, 100, 101] is due to its ability to capture long-range dependencies and contextual relationships from gene expression data. For this experiment, we used the TabNet classifier which incorporates an attention mechanism

for working with tabular data [102, 103, 104]. Its dynamic feature selection process enhances interpretability which improves generalization and reduces overfitting. In our study, the model used the Adam optimizer along with a decay rate (gamma) of 0.9. StepLR scheduler with a multiplication factor of 10 was employed to adjust the learning rate during training. Early stopping was implemented with a training threshold of 150 epochs. Batch size and virtual batch size were set to 512, and equal weighting was given to all training instances.

2.4. Ensemble Approach

An ensemble method is predicated on the notion that collective decisions are frequently superior to individual decisions [105, 106]. Additionally, by reducing the impact of outliers or noisy data points, the use of ensemble techniques has the potential to improve stability and dependability. In this study, we implemented two types of ensemble approaches— Max-Voting [107] and Probability-Averaging [108]. In the case of max voting, each model generates a prediction and the class label with the highest number of votes is chosen as the ultimate prediction. It is utilized to obtain a more precise estimation of class probabilities, rather than relying on the majority vote. This methodology incorporates model confidence and assigns greater weight to dependable predictions, leading to improved accuracy and calibrated probability estimation.

We selected an odd quantity of classification models to employ ensemble techniques. An odd number of models in an ensemble guarantees a majority class when voting based on their predicted class labels. This can facilitate decision-making and mitigate the occurrence of ties, which may arise when an even number of models are evaluated. Also, ensembling an odd number of models can effectively ignore outlier predictions from a single model and rely on the majority of predictions.

2.5. Cancer-specific gene set identification

Understanding the significant features of each cancer type is crucial for unraveling molecular mechanisms, identifying novel biomarkers, and discovering potential therapeutic targets. Explainable AI methods offer a powerful approach for extracting feature importance and understanding the contribution of individual genes in the prediction of cancer samples [35]. In this study, we utilized SHAP [42] as the explainable AI method. SHAP is based on game theory principles and offers a mathematical approach to elucidate

machine learning model predictions by calculating the individual impact of each feature. To calculate feature importance using SHAP, we passed the test data to the explainer along with the trained model and the explainer produced SHAP values, indicating the impact of each feature on the output prediction for each instance in the test data.

2.5.1. Identifying the correctly predicted samples

In the analysis of feature importance using SHAP scores from a model, we focused exclusively on the values from the correctly predicted samples. This is necessary because the SHAP scores of features from correctly predicted samples can provide genuine insights into the importance of each feature for a specific cancer. Thus, the genes that consistently and positively influence the model's ability to make accurate predictions for each cancer type can be prioritized [109]. Including SHAP scores from incorrectly predicted samples can introduce misleading information, as those scores may not be relevant to that cancer and can lead to misinterpretation of the gene contributions. By considering only the SHAP scores from correctly predicted samples, we ensured more accurate and reliable identification of the significant features that are more likely to be biologically relevant and specific to each cancer type [37].

2.5.2. Global interpretation of features using SHAP:

SHAP provides feature attribution by calculating scores of each feature of all samples across various types of cancers. To focus on a specific cancer type, we identified the samples belonging to that particular cancer and considered only the SHAP scores of those samples for analysis. In order to evaluate the global importance of each feature for a specific cancer, we computed the median value for each feature individually across the identified samples specific to that particular cancer type. This process was repeated for all 33 types of cancers, enabling us to determine the median score for each feature across the samples associated with each cancer type. This approach allowed us to precisely assess the significance of each feature in the context of different cancer types and identify cancer-specific significant genes.

We utilized specific SHAP explainers tailored to each algorithm to calculate the SHAP scores of features for models trained by different classifiers. The LinearExplainer from SHAP was employed for Logistic Regression and Support Vector Machine (SVM) models. For XGBoost and Random Forest models, we utilized the TreeExplainer. Lastly, the DeepExplainer was used for Multilayer Perceptron (MLP) and 1D Convolutional Neural Network (1D-CNN) models.

2.6. Statistical Validation using Differential Gene Expression

In addition to employing different approaches for cancer classification, researchers utilized statistical tools such as DESeq2 [43], edgeR [110], LIMMA [111], etc. to identify Differential Gene Expression (DGE). Differential gene expression (DGE) analysis facilitates the detection of genes with significant expression level variations across different cancer types [40, 112]. Sobhan et al. and Hossain et al. implemented DESeq2 for statistical genomic analysis and quantification of differential gene expression [113, 114, 39]. DESeq2 is a practical and widely used tool for analyzing differential gene expression in RNA-seq data. It employs negative binomial generalized linear models and applies empirical Bayes approaches to estimate priors for log fold change and dispersion, providing posterior estimates for these values [115].

The differential expression analysis with DESeq2 involves several steps. It utilizes normalization factors (size factors) to model raw counts and estimates gene-wise dispersion, enhancing the accuracy of dispersion estimates for count simulation. For our analysis, we utilized raw counts of gene expression values from both tumor samples and healthy tissue samples [116, 117]. Low-count genes were removed, retaining rows with at least 10 reads. The factor level was set to "healthy tissue". Our analysis was performed on individual samples of each cancer type to identify cancer-specific genes with a significant impact on each type of cancer. These genes capture the overall behavior of the population concerning each cancer type and can be considered as global features. We set the threshold for differential expression as |log2Fold-change|>3 and an adjusted p-value < 0.001, ensuring a stringent selection of genes exhibiting substantial expression changes that are statistically significant. Tab. 2 depicts a comprehensive summary of the sample sizes for both tumor and healthy tissues, along with the aggregate count of genes that exhibit differential expression across 17 distinct cancer types.

3. Results

3.1. Experimental Setup

Both Python and R programming languages were employed at different stages of the study. Python facilitated the majority of the experimental processes, while RStudio was utilized specifically to apply DESeq2, a widely

Table 2: Total number of differentially expressed genes using DESeq2 for each cancer class

#	Cancer type	Tumor sample count	Healthy tissue sample count	Total differentially expressed gene
1	BLCA	364	19	802
2	BRCA	984	112	707
3	CHOL	33	11	1904
4	COAD	287	43	878
5	ESCA	185	13	818
6	HNSC	462	44	889
7	KICH	62	27	1449
8	KIRC	477	74	821
9	KIRP	238	31	768
10	LIHC	297	50	737
11	LUAD	505	61	1026
12	LUSC	491	53	1769
13	PRAD	428	50	266
14	READ	89	12	883
15	STAD	382	35	637
16	THCA	443	55	445
17	UCEC	143	25	1205

used R package, for the identification of differentially expressed genes (DEGs) in individual cancer types. All experiments were conducted on a system equipped with an Intel® $Core^{TM}$ i9-12900K CPU and dual NVIDIA GeForce GPUs, each with 24GB of memory.

For evaluating our classification models and proposed pipeline, we employed four standard metrics: Accuracy, Precision, Recall, and F1-score. These metrics provided a comprehensive evaluation of model effectiveness.

Accuracy quantifies the proportion of correct predictions relative to the total number of predictions made by the model. However, in the context of imbalanced datasets, this metric may yield misleading results. Precision measures the model's ability to correctly identify positive instances, while Recall evaluates the model's capacity to detect true positive instances among all actual positive cases. It is important to recognize that these evaluation metrics may not fully capture model performance in imbalanced datasets, where one class may significantly dominate the others. Such metrics can exhibit bias toward the majority class, potentially resulting in deceptive conclusions.

Accuracy measures the ratio of correct predictions made by a model to the total number of predictions. However, imbalanced datasets can lead to misleading results. Precision refers to the model's ability to accurately identify positive instances. Recall assesses the model's capability to accurately detect positive instances among all the true positive instances. It is important to note that these evaluation metrics may not provide a complete assessment of the model's performance in the case of an imbalanced dataset, where one class is much more dominant than the other. These metrics are biased toward the majority class, which can result in misleading outcomes. The F1 score is a more informative metric for assessing model performance in such situations. It is important in cancer classification as it considers false positives and false negatives, which are crucial in medical diagnosis. The F1 score is more helpful in the case of imbalanced datasets where the minority class is of interest. It takes into account both recall and precision to offer a balanced assessment of a model's performance on imbalanced datasets by weighing the trade-off between two metrics.

3.2. Performance on the Baseline Dataset

The dataset consisting of all the 19,238 mRNA features and 9807 tumor samples was considered as the baseline dataset. Initially, we evaluated the performance of the seven classifiers on this baseline dataset. The results are summarized in Tab. 3.

Among all the classifiers, Logistic Regression achieved the highest accuracy of 96.43% and an F1 score of 0.9391. Additionally, SVM and XGBoost also demonstrated strong performances, surpassing 95% accuracy with satisfactory F1 scores. It is worth noting that although the DL models (MLP and 1D CNN) achieved high accuracies (95.34% and 94.50% respectively), their F1 scores were not as impressive. It indicates that while the model is good at making overall predictions, it struggles to account for class disparities and misclassifies certain minority classes.

To further enhance the overall performance, we adopted Probability averaging and Max-Voting ensemble approaches. For the ensemble method, we combined the predictions of the top three performing models, namely Logistic Regression, SVM, and XGBoost. These three models were chosen as they all achieved accuracies above 95% with good F1 scores. After implementing the ensemble approaches, Probability averaging outperformed Max-Voting, resulting in a slight improvement in accuracy to 96.45%. It also surpassed

Table 3: Performance analysis of the standalone classifiers and the two ensemble techniques on the dataset containing 19238 features and the performance of our proposed pipeline (implementing Standard Scaler normalization and SelectFromModel feature selection techniques and Probability Averaging Ensemble combining Logistic Regression, SVM and XGBoost) on 500 selected features.

Number of Features	Classifier Models	Acc (%)	F1- Score
	Random Forest	92.17	0.9033
	1D-CNN	94.5	0.844
	Tabnet	94.56	0.9129
19238	MLP	95.34	0.8673
	XGBoost	95.52	0.9203
	SVM	96.3	0.9387
	Logistic Regression	96.43	0.9391
	Max Voting Ensemble (Logistic Regression,	96.33	0.9386
	SVM, XGBoost)		
	Probability Averaging Ensemble (Logistic Regression, SVM, XGBoost)	96.45	0.9399
	XGBoost	95.10	0.9315
500	SVM	95.87	0.9344
500	Logistic Regression	96.31	0.9382
	Ours (using probability averaging ensemble)	96.61	0.9415

the individual accuracy of the Logistic Regression model. Hence, the highest accuracy achieved on the baseline dataset was 96.45%.

We conducted a comprehensive evaluation using feature selection methods in order to tackle the high dimensionality issue posed by the dataset's 19,238 mRNA features. Our pipeline involved applying standard scaler normalization technique followed by the SelectFromModel feature selection technique to identify the most informative 500 features out of the original 19,238.

Remarkably, despite the significant reduction in feature space, Logistic Regression displayed great performance, achieving an impressive accuracy of 96.31%. This accuracy level not only outperformed the other classifiers but also came very close to the accuracy obtained on the baseline dataset with all 19,238 features. SVM and XGBoost classifiers also managed to achieve accuracies higher than 95% with the reduced set of features. It is important to highlight that the accuracy values were averaged across 5 folds, ensuring the consistency of the results. Additionally, we utilized ensemble techniques by combining the predictions from the top-performing classifiers to enhance the

overall performance. The results showed that Probability Averaging led the way, achieving the highest accuracy of 96.61%. F1-score of 0.9415 indicates a robust and balanced performance of the model in correctly identifying all the different cancer classes. Notably, this accuracy even surpassed the highest accuracy attained on the baseline dataset.

A particularly impressive aspect of this achievement is the drastic reduction in computational expenses. By utilizing selective features obtained through the feature selection technique, we were able to streamline the computational burden without compromising on accuracy.

3.3. Ablation study

The results of all the experiments conducted were meticulously analyzed to determine the optimal combination of normalization techniques, feature selection methods, and classifier models that yield the best performance. In this analysis, the accuracies and the F1 scores of the experiments were carefully examined.

Table 4: Performance analysis of the classifiers (in terms of Accuracy) after applying Standard Scaler and MinMax Scaler normalization techniques using reduced feature set

Classifier	Normalization Technique		
C 100021101	Standard Scaler	MinMax Scaler	
Logistic Regression	96.31	96.17	
SVM	96.03	96.02	
XGBoost	$\boldsymbol{95.52}$	95.47	
MLP	94.79	84.46	
1D-CNN	94.6	94.96	
Random Forest	91.65	92.23	

3.3.1. Choice of Normalization technique

After obtaining the performance on the baseline dataset, we applied different techniques to address the high dimensionality issue while maintaining the accuracy achieved with all the features. An ablation study was conducted to assess the impact of various normalization techniques, feature selection methods, and classifier models on overall performance. We evaluated various combinations to identify the most effective combination of these techniques for achieving the best performance with minimal features. At first, we applied the two normalization techniques. The results are shown in Tab. 4. We

observed that Standard Scalar enhances the performance better than MinMax Scaler for 4 out of the 6 classifiers. So we have chosen Standard Scalar as our normalization technique and proceeded for further implementations.

3.3.2. Choice of Feature Selection technique

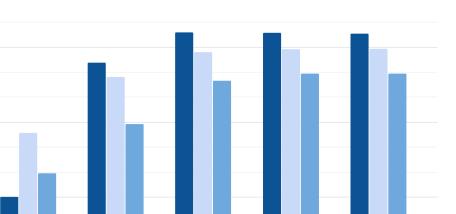
After applying the normalization technique, we have implemented four different feature selection techniques with 1000 features. We chose 1000 feature count as a starter as working with around 5% features reduces the computational cost greatly. The results after applying different feature selection methods of the three top performing models are shown in Tab. 5. From the results, it is observed that SelectFromModel feature selection methods performs the best among the four feature selection methods. For this reason, we chose SelectFromModel as our feature selection method.

Table 5: Performance analysis of the three best-performing classifiers with 1000 features chosen by different feature selection methods using 'Standard Scaler' feature normalization technique.

Feature Selection Technique	Classifier	$\mathrm{Acc}(\%)$
SelectFromModel	Logistic Regression SVM XGBoost	96.49 95.87 95.10
Select-K-Best	Logistic Regression SVM XGBoost	95.51 95.14 94.16
Lasso	Logistic Regression SVM XGBoost	95.47 95.38 94.87
ElasticNet	Logistic Regression SVM XGBoost	94.28 94.29 93.58

3.3.3. Selecting the appropriate number of features

Experiments were performed to find out how different feature selection affects the overall performance. Using the feature selection approaches, it was observed that using 100 and 250 features did not yield satisfactory outcomes, while performance significantly improved with 500, 750, and 1000 features.



■ Logistic Regression ■ SVM ■ XGBoost

Figure 2: The performance of the three best-performing classifiers with 100, 250, 500, 750 and 1000 features chosen by 'SelectFromModel' feature selection method using 'Standard Scaler' feature normalization

Interestingly, the top three models (Logistic regression, SVM, XGBoost) exhibited nearly identical accuracies across the 500, 750, and 1000 feature sets (Fig. 2). Consequently, the decision was made to proceed with 500 features, striking a balance between computational efficiency and satisfactory performance.

3.4. Class-wise performance analysis

96.00%

94.009

90.00%

It is worth mentioning that in total 29 out of 33 cancer types achieved over 90% accuracy in identification (Tab. 6) with only 500 gene set. Specially, BLCA, BRCA, DLBC, GBM, HNSC, KIHC, LAML, MESO, OV, PCPG, PRAD, SKCM, TGCT, THCA, THYM, and UVM, achieved accuracy over 98% and exhibited distinct characteristics that allowed for their identification from other cancer types. The proposed architecture proved successful in handling cancer types with limited patient samples. DLBC attained a 97.78% accuracy rate despite having a sample size of only 47. KICH, MESO, and UVM achieved accuracy rates of 93.96%, 98.7%, and 98.67% respectively, with sample sizes of 66, 87, and 79. This demonstrates that our proposed

architecture is capable of handling class imbalance and performs well even with limited sample sizes. However, cancers originating from the same tissue origins posed greater difficulty in distinguishing them from each other compared to those originating from different lineages. For instance, READ and COAD were more challenging to differentiate than BRCA and COAD. Notably, more than half of the READ samples were misclassified as COAD samples, but combining these samples improved the accuracy from 96.39% to 97.36% in logistic regression.

Table 6: Accuracy of each cancer type of the proposed pipeline.

#	Cancer Type	$\mathrm{Acc}(\%)$
1	Adrenocortical carcinoma (ACC)	96.00
2	Bladder Urothelial Carcinoma (BLCA)	98.53
3	Breast invasive carcinoma (BRCA)	99.55
4	Cervical squamous cell carcinoma endocervical adenocarcinoma (CESC)	94.46
5	Cholangiocarcinoma (CHOL)	81.07
6	Colon adenocarcinoma (COAD)	89.23
7	Diffuse large B cell lymphoma (DLBCL)	97.78
8	Esophageal carcinoma (ESCA)	86.28
9	Glioblastoma multiforme (GBM)	98.79
10	Head and Neck squamous cell (HNSC)	98.46
11	Kidney Chromophobe (KICH)	93.96
12	Kidney renal clear cell carcinoma (KIRC)	96.42
13	Kidney renal papillary cell carcinoma (KIRP)	93.77
14	Acute myeloid leukemia (LAML)	100.00
15	Brain Lower Grade Glioma (LGG)	98.66
16	Liver hepatocellular carcinoma (LIHC)	97.04
17	Lung adenocarcinoma (LUAD)	95.34
18	Lung squamous cell carcinoma (LUSC)	94.97
19	Mesothelioma (MESO)	98.71
20	Ovarian serous cystadenocarcinoma (OV)	99.76
21	Pancreatic adenocarcinoma (PAAD)	97.21
22	Pheochromocytoma and Paraganglioma (PCPG)	98.90
23	Prostate adenocarcinoma (PRAD)	100.00
24	Rectum adenocarcinoma (READ)	42.34
25	Sarcoma (SARC)	93.93
26	Skin Cutaneous Melanoma (SKCM)	98.51
27	Stomach adenocarcinoma (STAD)	93.49
28	Testicular Germ Cell Tumors (TGCT)	99.26
29	Thyroid carcinoma (THCA)	100.00
30	Thymoma (THYM)	98.33

Table 6: Accuracy of each cancer type of the proposed pipeline.

#	Cancer Type	Acc(%)
31	Uterine Corpus Endometrial Carcinoma (UCEC)	92.28
32	Uterine Carcinosarcoma (UCS)	78.79
33	Uveal Melanoma (UVM)	98.67

In conclusion, the experiments highlighted the significance of appropriate preprocessing approaches in improving model performance while considering computational constraints. The adoption of StandardScaler normalization, SelectFromModel with 500 features, and ensembling of the top three classifiers struck a balance between accuracy and resource efficiency.

3.5. Performance Comparison with State-of-the-Art Approaches

A comparative analysis was conducted with various state-of-the-art [15, 11, 25] techniques that utilized the TCGA pan-cancer dataset to perform classification tasks on 33 distinct cancer categories. Tab. 7 displays the performance evaluation for each of the architectures, including the architecture proposed by us. To the best of our knowledge, the proposed architecture outperforms all existing methodologies.

Table 7: Performance comparison with state-of-the-art works on 33 types cancer classification using RNA sequence data.

Approach	$\begin{array}{c} \textbf{Number of} \\ \textbf{Genes} \end{array}$	$\begin{array}{c} \textbf{Avg.} \\ \textbf{Acc}(\%) \end{array}$
Linear SVM [15]	9900	94.98
Deep learning based algorithm [11]	10381	95.59
Deep learning based algorithm [25]	20531	95.65
Ours	500	96.61

One crucial advantage of the proposed pipeline is its ability to achieve superior performance while utilizing a significantly reduced number of features. Our pipeline utilized only 500 features, which is considerably less than other methods. Reducing the number of features greatly reduces the computational resources required for classification tasks. The proposed architecture demonstrates high-performance accuracy and efficient usage of computational resources, making it a promising solution for multi-class cancer classification.

3.6. Cancer Specific Gene Set

For the identification of gene biomarkers responsible for specific cancers, we considered the SHAP values of the three best-performing classifiers (XGBoost, Logistic Regression, and SVM).

As discussed previously, we employed a 5-fold cross-validation approach and utilized SelectFromModel to select 500 features in each fold for each classifier. So it is important to note that the 500 selected features in each fold were not exactly identical across all the 5 folds, as they were determined based on the specific data partition in each iteration. Consequently, we obtained SHAP values of features from the 5 folds for each classifier. To obtain a comprehensive global interpretation from each classifier, we combined the SHAP values from the 5 folds. It resulted in a collection of 1008, 1595, and 723 distinct features and their SHAP values for Logistic Regression, XGBoost, and SVM, respectively. After the processing of these values as described in Sec. 2.5.2, we got 33 distinct sets of feature attribution values specific to the 33 types of cancers from each of the classifiers. We sorted the genes in descending order based on their significance within each cancer type to get the most important genes.

In order to validate the significance of the genes identified through global interpretation, we conducted a comparison with the differentially expressed genes (DEGs) available for 17 cancer types obtained from Deseq2 analysis. For each cancer, we took the top genes same as the number of the DEGs from our obtained gene sets. From the analysis, we observed that Logistic Regression, XGBoost, and SVM were able to identify a significant number of important genes for most of the cancers. The comparison results are shown in Tab. 8. For instance, in the case of BLCA cancer, the Logistic Regression model identified 97 common genes, XGBoost identified 49 common genes, and SVM identified 122 common genes, all of which were shared with the corresponding DGE. With these findings, we can consider that the genes identified as common with the DEGs are of utmost significance for each cancer type.

We further validated the specificity of the class-specific genes obtained in our study. Typically, if a gene set is truly specific to a particular cancer, it should not exhibit any overlap with genes from a different cancer set. In other words, the number of overlapping genes between different cancer types should be minimal, approaching zero. Our findings indicate that there is indeed a small number of overlapping genes among different cancer types.

Table 8: Number of common gene biomarkers between the gene sets obtained from our pipeline and the Differential Gene Expression from DESeq2.

#	Cancer	Number of common features between our selected features using SHAP analysis and DEGs		
		Logistic Regression	XGBoost	Support Vector Machine (SVM)
1	BLCA	97	49	122
2	BRCA	60	31	87
3	CHOL	112	233	203
4	COAD	112	79	159
5	ESCA	90	57	132
6	HNSC	120	58	133
7	KICH	139	159	183
8	KIRC	82	57	108
9	KIRP	74	48	99
10	LIHC	80	43	86
11	LUAD	155	97	200
12	LUSC	239	240	246
13	PRAD	19	7	25
14	READ	102	76	132
15	STAD	71	40	112
16	THCA	19	23	52
_17	UCEC	145	135	178

This outcome serves as validation for the cancer-specific nature of our gene sets.

4. Conclusion

Gene expression analysis from mRNA data is a valuable tool for cancer classification, leveraging the altered expression levels of specific genes to differentiate between different types of cancer [88]. In our study, we introduced a method that efficiently and accurately identifies cancer types and the corresponding gene sets based on mRNA gene expression data. By employing appropriate feature selection techniques, we achieved a significant reduction in computational cost while maintaining a high accuracy rate of 96.61%. This was accomplished through an ensemble approach incorporating three well-performing classifiers: Logistic regression, Support Vector Machine (SVM), and XGBoost. Furthermore, we computed SHAP scores for each feature to identify and prioritize the most important biomarker genes specific to

each cancer type. Comparing these genes with those obtained from DGE analysis underscored the effectiveness of our approach in accurately capturing cancer-specific genes. Overall, our method offers a precise and rapid means of cancer classification while highlighting biologically relevant genes associated with each cancer type. This study forms the groundwork to identify biomarker genes for individual cancer patients. Analyzing patient-specific gene sets can play an important role in facilitating personalized and targeted treatment approaches.

References

- [1] J. Wang, X. Zhang, W. Chen, X. Hu, J. Li, C. Liu, Regulatory roles of long noncoding rnas implicated in cancer hallmarks, International journal of cancer 146 (4) (2020) 906–916.
- [2] S. Bekisz, L. Geris, Cancer modeling: From mechanistic to data-driven approaches, and from fundamental insights to clinical applications, Journal of Computational Science 46 (2020) 101198.
- [3] Global cancer burden growing, amidst mounting need for services, https://www.who.int/news/item/
 01-02-2024-global-cancer-burden-growing-amidst-mounting-need-for-services,
 [Accessed 07-10-2024].
- [4] F. Alharbi, A. Vakanski, Machine learning methods for cancer classification using gene expression data: A review, Bioengineering 10 (2) (2023) 173.
- [5] Y. Yuan, F. Gao, Y. Chang, Q. Zhao, X. He, Advances of mrna vaccine in tumor: a maze of opportunities and challenges, Biomarker research 11 (1) (2023).
- [6] K. Kourou, T. P. Exarchos, K. P. Exarchos, M. V. Karamouzis, D. I. Fotiadis, Machine learning applications in cancer prognosis and prediction, Computational and structural biotechnology journal 13 (2015) 8–17.
- [7] M. D. Podolsky, A. A. Barchuk, V. I. Kuznetcov, N. F. Gusarova, V. S. Gaidukov, S. A. Tarakanov, Evaluation of machine learning algorithm utilization for lung cancer classification based on gene expression levels, Asian Pacific journal of cancer prevention 17 (2) (2016) 835–838.

- [8] G. Guo, H. Wang, D. Bell, Y. Bi, K. Greer, Knn model-based approach in classification, in: On The Move to Meaningful Internet Systems 2003: CoopIS, DOA, and ODBASE, Springer Berlin Heidelberg, 2003, pp. 986–996.
- [9] L. Wang, Support vector machines: theory and applications, Vol. 177, Springer Science & Business Media, 2005.
- [10] J. Huang, C. X. Ling, Using auc and accuracy in evaluating learning algorithms, IEEE Transactions on knowledge and Data Engineering 17 (3) (2005) 299–310.
- [11] B. Lyu, A. Haque, Deep learning based tumor type classification using gene expression data, in: Proceedings of the 2018 ACM international conference on bioinformatics, computational biology, and health informatics, 2018, pp. 89–96.
- [12] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: International conference on machine learning, pmlr, 2015, pp. 448–456.
- [13] J.-F. Laplante, M. A. Akhloufi, Predicting cancer types from mirna stem-loops using deep learning, in: 2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC), IEEE, 2020, pp. 5312–5315.
- [14] A. K. Jain, J. Mao, K. M. Mohiuddin, Artificial neural networks: A tutorial, Computer 29 (3) (1996) 31–44.
- [15] Y.-H. Hsu, D. Si, Cancer type prediction and classification based on rna-sequencing data, in: 2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), IEEE, 2018, pp. 5374–5377.
- [16] G. Stein, B. Chen, A. S. Wu, K. A. Hua, Decision tree classifier for network intrusion detection with ga-based feature selection, in: Proceedings of the 43rd annual Southeast regional conference-Volume 2, 2005, pp. 136–141.
- [17] M. A. F. A. Fida, T. Ahmad, M. Ntahobari, Variance threshold as early screening to boruta feature selection for intrusion detection system, in:

- 2021 13th International Conference on Information & Communication Technology and System (ICTS), IEEE, 2021, pp. 46–50.
- [18] Y.-W. Chang, C.-J. Lin, Feature ranking using linear sym, in: Causation and prediction challenge, PMLR, 2008, pp. 53–64.
- [19] S. Albawi, T. A. Mohammed, S. Al-Zawi, Understanding of a convolutional neural network, in: 2017 international conference on engineering and technology (ICET), Ieee, 2017, pp. 1–6.
- [20] M. Mostavi, Y.-C. Chiu, Y. Huang, Y. Chen, Convolutional neural network models for cancer type prediction based on gene expression, BMC medical genomics 13 (2020) 1–13.
- [21] K. F. Mahin, M. Robiuddin, M. Islam, S. Ashraf, F. Yeasmin, S. Shatabda, Panclassif: Improving pan cancer classification of single cell rna-seq gene expression data using machine learning, Genomics 114 (2) (2022) 110264.
- [22] L. Breiman, Random forests, Machine learning 45 (1) (2001) 5–32.
- [23] Y. Freund, R. Schapire, N. Abe, A short introduction to boosting, Journal-Japanese Society For Artificial Intelligence 14 (771-780) (1999) 1612.
- [24] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 618–626.
- [25] J. M. de Guia, M. Devaraj, C. K. Leung, Deepgx: deep learning using gene expression for cancer classification, in: Proceedings of the 2019 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining, 2019, pp. 913–920.
- [26] A. Lopez-Rincon, M. Martinez-Archundia, G. U. Martinez-Ruiz, A. Schoenhuth, A. Tonda, Automatic discovery of 100-mirna signature for cancer classification using ensemble feature selection, BMC bioinformatics 20 (1) (2019) 1–17.

- [27] H.-H. Huang, X.-Y. Liu, Y. Liang, Feature selection and cancer classification via sparse logistic regression with the hybrid l1/2+ 2 regularization, PloS one 11 (5) (2016) e0149675.
- [28] P. Linardatos, V. Papastefanopoulos, S. Kotsiantis, Explainable ai: A review of machine learning interpretability methods, Entropy 23 (1) (2020) 18.
- [29] E. Strumbelj, I. Kononenko, Explaining prediction models and individual predictions with feature contributions, Knowledge and information systems 41 (2014) 647–665.
- [30] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, PloS one 10 (7) (2015) e0130140.
- [31] M. T. Ribeiro, S. Singh, C. Guestrin, Why should i trust you?": Explaining the predictions of any classifier. corr abs/1602.04938 (2016), arXiv preprint arXiv:1602.04938 (2016).
- [32] M. T. Ribeiro, S. Singh, C. Guestrin, Anchors: High-precision model-agnostic explanations, in: Proceedings of the AAAI conference on artificial intelligence, Vol. 32, 2018.
- [33] A. Shrikumar, P. Greenside, A. Kundaje, Learning important features through propagating activation differences. 34th int. conf, Mach. Learn. ICML 7 (2017) (2017) 4844–4866.
- [34] Z. Zhang, C. Sun, Z.-P. Liu, Discovering biomarkers of hepatocellular carcinoma from single-cell rna sequencing data by cooperative games on gene regulatory network, Journal of Computational Science 65 (2022) 101881.
- [35] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, Advances in neural information processing systems 30 (2017).
- [36] J. J. Levy, A. J. Titus, C. L. Petersen, Y. Chen, L. A. Salas, B. C. Christensen, Methylnet: an automated and modular deep learning approach for dna methylation analysis, BMC bioinformatics 21 (2020).

- [37] M. Yap, R. L. Johnston, H. Foley, S. MacDonald, O. Kondrashova, K. A. Tran, K. Nones, L. T. Koufariotis, C. Bean, J. V. Pearson, et al., Verifying explainability of a deep learning tissue classifier trained on rna-seq data, Scientific reports 11 (1) (2021) 2641.
- [38] M. Divate, A. Tyagi, D. J. Richard, P. A. Prasad, H. Gowda, S. H. Nagaraj, Deep learning-based pan-cancer classification model reveals tissue-of-origin specific gene expression signatures, Cancers 14 (5) (2022) 1185.
- [39] J.-m. Xue, Y. Liu, L.-h. Wan, Y.-x. Zhu, Comprehensive analysis of differential gene expression to identify common gene signatures in multiple cancers, Medical science monitor: international medical journal of experimental and clinical research 26 (2020) e919953–1.
- [40] A. Stupnikov, C. McInerney, K. Savage, S. McIntosh, F. Emmert-Streib, R. Kennedy, M. Salto-Tellez, K. Prise, D. McArt, Robustness of differential gene expression analysis of rna-seq, Computational and structural biotechnology journal 19 (2021) 3470–3481.
- [41] M. Ni, X. Liu, J. Wu, D. Zhang, J. Tian, T. Wang, S. Liu, Z. Meng, K. Wang, X. Duan, et al., Identification of candidate biomarkers correlated with the pathogenesis and prognosis of non-small cell lung cancer via integrated bioinformatics analysis. front genet. 2018; 9: 469 (2018).
- [42] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, ArXiv abs/1705.07874 (2017). URL https://api.semanticscholar.org/CorpusID:21889700
- [43] E. Porcu, M. C. Sadler, K. Lepik, C. Auwerx, A. R. Wood, A. Weihs, M. S. B. Sleiman, D. M. Ribeiro, S. Bandinelli, T. Tanaka, et al., Differentially expressed genes reflect disease-induced rather than disease-causing changes in the transcriptome, Nature Communications 12 (1) (2021) 5647.
- [44] J. Weinstein, Tcgar network, ea collisson et al., "the cancer genome atlas pan-cancer analysis project,", Nature Genetics 45 (10) (2013) 1113–1120.
- [45] Y. Zhao, M.-C. Li, M. M. Konaté, L. Chen, B. Das, C. Karlovich, P. M. Williams, Y. A. Evrard, J. H. Doroshow, L. M. McShane, Tpm, fpkm,

- or normalized counts? a comparative study of quantification measures for the analysis of rna-seq data from the nci patient-derived models repository, Journal of translational medicine 19 (1) (2021) 1–15.
- [46] B. Li, C. N. Dewey, Rsem: accurate transcript quantification from rna-seq data with or without a reference genome, BMC bioinformatics 12 (2011) 1–16.
- [47] N. Bunkley, Joseph juran, pioneer in quality control, dies, The New York Times 103 (2008).
- [48] V. R. Joseph, Optimal ratio for data splitting, Statistical Analysis and Data Mining: The ASA Data Science Journal 15 (4) (2022) 531–538.
- [49] H. Henderi, T. Wahyuningsih, E. Rahwanto, Comparison of min-max normalization and z-score normalization in the k-nearest neighbor (knn) algorithm to test the accuracy of types of breast cancer, International Journal of Informatics and Information Systems 4 (1) (2021) 13–20.
- [50] D. Singh, B. Singh, Investigating the impact of data normalization on classification performance, Applied Soft Computing 97 (2020) 105524.
- [51] D. Borkin, A. Némethová, G. Michal'čonok, K. Maiorov, Impact of data normalization on classification model accuracy, Research Papers Faculty of Materials Science and Technology Slovak University of Technology 27 (45) (2019) 79–84.
- [52] M. V. Polyakova, V. N. Krylov, Data normalization methods to improve the quality of classification in the breast cancer diagnostic system, tic 5 (1) (2022) 55–63.
- [53] V. G. Raju, K. P. Lakshmi, V. M. Jain, A. Kalidindi, V. Padma, Study the influence of normalization/transformation process on the accuracy of supervised classification, in: 2020 Third International Conference on Smart Systems and Inventive Technology (ICSSIT), IEEE, 2020, pp. 729–735.
- [54] Z. Swiderska-Chadaj, T. de Bel, L. Blanchet, A. Baidoshvili, D. Vossen, J. van der Laak, G. Litjens, Impact of rescanning and normalization on convolutional neural network performance in multi-center, whole-slide classification of prostate cancer, Scientific Reports 10 (1) (2020) 1–14.

- [55] S. K. Panda, S. Nag, P. K. Jana, A smoothing based task scheduling algorithm for heterogeneous multi-cloud environment, in: 2014 International Conference on Parallel, Distributed and Grid Computing, IEEE, 2014, pp. 62–67.
- [56] S. Patro, K. K. Sahu, Normalization: A preprocessing stage, arXiv preprint arXiv:1503.06462 (2015).
- [57] M. M. Ahsan, M. P. Mahmud, P. K. Saha, K. D. Gupta, Z. Siddique, Effect of data scaling methods on machine learning algorithms and model performance, Technologies 9 (3) (2021) 52.
- [58] E. Bisong, E. Bisong, Introduction to scikit-learn, Building Machine Learning and Deep Learning Models on Google Cloud Platform: A Comprehensive Guide for Beginners (2019) 215–229.
- [59] N. K. Chauhan, K. Singh, Performance assessment of machine learning classifiers using selective feature approaches for cervical cancer detection, Wireless Personal Communications 124 (3) (2022) 2335–2366.
- [60] A. Arafa, M. Radad, M. Badawy, N. El-Fishawy, Regularized logistic regression model for cancer classification, in: 2021 38th National Radio Science Conference (NRSC), Vol. 1, IEEE, 2021, pp. 251–261.
- [61] B. Yesilkaya, M. Perc, Y. Isler, Manifold learning methods for the diagnosis of ovarian cancer, Journal of Computational Science 63 (2022) 101775.
- [62] M. Huljanah, Z. Rustam, S. Utama, T. Siswantining, Feature selection using random forest classifier for predicting prostate cancer, in: IOP Conference Series: Materials Science and Engineering, Vol. 546, IOP Publishing, 2019, p. 052031.
- [63] V. Fonti, E. Belitser, Feature selection using lasso, VU Amsterdam research paper in business analytics 30 (2017) 1–25.
- [64] H. Liu, J. Li, L. Wong, A comparative study on feature selection and classification methods using gene expression profiles and proteomic patterns, Genome informatics 13 (2002) 51–60.

- [65] Z. M. Hira, D. F. Gillies, A review of feature selection and feature extraction methods applied on microarray data, Advances in bioinformatics 2015 (2015).
- [66] J. Brownlee, An introduction to feature selection, Machine learning process 6 (2014).
- [67] A. Chauhan, et al., Detection of lung cancer using machine learning techniques based on routine blood indices, in: 2020 IEEE international conference for innovation in technology (INOCON), IEEE, 2020, pp. 1–6.
- [68] H. Zou, T. Hastie, Regression shrinkage and selection via the elastic net, with applications to microarrays, JR Stat Soc Ser B 67 (2003) 301–20.
- [69] R. Muthukrishnan, R. Rohini, Lasso: A feature selection technique in predictive modeling for machine learning, in: 2016 IEEE international conference on advances in computer applications (ICACA), IEEE, 2016, pp. 18–20.
- [70] A. Agaal, M. Essgaer, Influence of feature selection methods on breast cancer early prediction phase using classification and regression tree, in: 2022 International Conference on Engineering & MIS (ICEMIS), IEEE, 2022, pp. 1–6.
- [71] T. K. Kim, Understanding one-way anova using conceptual figures, Korean journal of anesthesiology 70 (1) (2017) 22–26.
- [72] J. Friedman, T. Hastie, R. Tibshirani, Regularization paths for generalized linear models via coordinate descent, Journal of statistical software 33 (1) (2010) 1.
- [73] A. Sokolov, D. E. Carlin, E. O. Paull, R. Baertsch, J. M. Stuart, Pathway-based genomics prediction using generalized elastic net, PLoS computational biology 12 (3) (2016) e1004790.
- [74] A. Basu, R. Mitra, H. Liu, S. L. Schreiber, P. A. Clemons, Rwen: response-weighted elastic net for prediction of chemosensitivity of cancer cell lines, Bioinformatics 34 (19) (2018) 3332–3339.

- [75] R. Tibshirani, Regression shrinkage and selection via the lasso, Journal of the Royal Statistical Society: Series B (Methodological) 58 (1) (1996) 267–288.
- [76] M. Khalsan, L. R. Machado, E. S. Al-Shamery, S. Ajit, K. Anthony, M. Mu, M. O. Agyeman, A survey of machine learning approaches applied to gene expression analysis for cancer prediction, IEEE Access 10 (2022) 27522–27534.
- [77] J. Liñares-Blanco, A. Pazos, C. Fernandez-Lozano, Machine learning analysis of tcga cancer data, Peer J. Computer Science 7 (2021) e584.
- [78] A. B. Tufail, Y.-K. Ma, M. K. Kaabar, F. Martínez, A. Junejo, I. Ullah, R. Khan, et al., Deep learning in cancer diagnosis and prognosis prediction: a minireview on challenges, recent trends, and future directions, Computational and Mathematical Methods in Medicine 2021 (2021).
- [79] V. S. Desdhanty, Z. Rustam, Liver cancer classification using random forest and extreme gradient boosting (xgboost) with genetic algorithm as feature selection, in: 2021 International Conference on Decision Aid Sciences and Application (DASA), IEEE, 2021, pp. 716–719.
- [80] L. A. V. Silva, K. Rohr, Pan-cancer prognosis prediction using multimodal deep learning, in: 2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI), IEEE, 2020, pp. 568–571.
- [81] T. Chen, C. Guestrin, Xgboost: A scalable tree boosting system, in: Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining, 2016, pp. 785–794.
- [82] D. Svozil, V. Kvasnicka, J. Pospichal, Introduction to multi-layer feed-forward neural networks, Chemometrics and intelligent laboratory systems 39 (1) (1997) 43–62.
- [83] R. S. Srinivasamurthy, Understanding 1d convolutional neural networks using multiclass time-varying signalss, Ph.D. thesis, Clemson University (2018).
- [84] S. Ö. Arik, T. Pfister, Tabnet: Attentive interpretable tabular learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 35, 2021, pp. 6679–6687.

- [85] T. Hastie, Ridge regularization: An essential concept in data science, Technometrics 62 (4) (2020) 426–433.
- [86] F. R. Bach, Bolasso: model consistent lasso estimation through the bootstrap, in: Proceedings of the 25th international conference on Machine learning, 2008, pp. 33–40.
- [87] I. Guyon, J. Weston, S. Barnhill, V. Vapnik, Gene selection for cancer classification using support vector machines, Machine learning 46 (2002) 389–422.
- [88] E. Alba, J. Garcia-Nieto, L. Jourdan, E.-G. Talbi, Gene selection in cancer classification using pso/svm and ga/svm hybrid algorithms, in: 2007 IEEE congress on evolutionary computation, IEEE, 2007, pp. 284–290.
- [89] M. Galar, A. Fernández, E. Barrenechea, H. Bustince, F. Herrera, An overview of ensemble methods for binary classifiers in multi-class problems: Experimental study on one-vs-one and one-vs-all schemes, Pattern Recognition 44 (8) (2011) 1761–1776.
- [90] J. H. Friedman, Stochastic gradient boosting, Computational statistics & data analysis 38 (4) (2002) 367–378.
- [91] T. Serre, A. Oliva, T. Poggio, A feedforward architecture accounts for rapid categorization, Proceedings of the national academy of sciences 104 (15) (2007) 6424–6429.
- [92] U. Ravindran, C. Gunavathi, A survey on gene expression data analysis using deep learning methods for cancer diagnosis, Progress in Biophysics and Molecular Biology 177 (2023) 1–13.
- [93] P. Guillen, J. Ebalunode, Cancer classification based on microarray gene expression data using deep learning, in: 2016 International Conference on Computational Science and Computational Intelligence (CSCI), IEEE, 2016, pp. 1403–1405.
- [94] F. Gao, W. Wang, M. Tan, L. Zhu, Y. Zhang, E. Fessler, L. Vermeulen, X. Wang, Deepcc: a novel deep learning-based framework for cancer molecular subtype classification, Oncogenesis 8 (9) (2019) 44.

- [95] B. Xu, N. Wang, T. Chen, M. Li, Empirical evaluation of rectified activations in convolutional network, arXiv preprint arXiv:1505.00853 (2015).
- [96] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014).
- [97] M. Mohammed, H. Mwambi, I. B. Mboya, M. K. Elbashir, B. Omolo, A stacking ensemble deep learning approach to cancer type classification based on tcga data, Scientific reports 11 (1) (2021) 1–22.
- [98] I. Goodfellow, Y. Bengio, A. C. Courville, Deep learning. das umfassende handbuch, 2018. URL https://api.semanticscholar.org/CorpusID:210451851
- [99] M. Gokhale, S. K. Mohanty, A. Ojha, Genevit: Gene vision transformer with improved deepinsight for cancer classification, Computers in Biology and Medicine 155 (2023) 106643.
- [100] A. Khan, B. Lee, Gene transformer: Transformers for the gene expression-based classification of lung cancer subtypes, arXiv preprint arXiv:2108.11833 (2021).
- [101] T.-H. Zhang, M. M. Hasib, Y.-C. Chiu, Z.-F. Han, Y.-F. Jin, M. Flores, Y. Chen, Y. Huang, Transformer for gene expression modeling (t-gem): An interpretable deep learning model for gene expression-based phenotype predictions, Cancers 14 (19) (2022) 4763.
- [102] F. B. Rahman, F. Anjum, M. H. F. Khan, Detection of lung adenocarcinoma cancer based on rna-seq gene expression data using limma and tabnet, Ph.D. thesis, Department of Computer Science and Engineering (CSE), Islamic University of ... (2022).
- [103] R. T. McLaughlin, M. Asthana, M. Di Meo, M. Ceccarelli, H. J. Jacob, D. L. Masica, Fast, accurate, and racially unbiased pan-cancer tumoronly variant calling with tabular machine learning, NPJ Precision Oncology 7 (1) (2023) 4.
- [104] A. Nasimian, M. Ahmed, I. Hedenfalk, J. U. Kazi, A deep tabular data learning model predicting cisplatin sensitivity identifies bcl2l1

- dependency in cancer, Computational and Structural Biotechnology Journal (2023).
- [105] Y. Xiao, J. Wu, Z. Lin, X. Zhao, A deep learning-based multi-model ensemble method for cancer prediction, Computer methods and programs in biomedicine 153 (2018) 1–9.
- [106] A. C. Tan, D. Gilbert, Ensemble machine learning on gene expression data for cancer classification (2003).
- [107] A. S. Assiri, S. Nazir, S. A. Velastin, Breast tumor classification using an ensemble machine learning method, Journal of Imaging 6 (6) (2020) 39.
- [108] M. Hosni, I. Abnane, A. Idri, J. M. C. de Gea, J. L. F. Alemán, Reviewing ensemble classification methods in breast cancer, Computer methods and programs in biomedicine 177 (2019) 89–112.
- [109] K. Futagami, Y. Fukazawa, N. Kapoor, T. Kito, Pairwise acquisition prediction with shap value interpretation, The Journal of Finance and Data Science 7 (2021) 22–44.
- [110] Q. Sun, X. Li, M. Xu, L. Zhang, H. Zuo, Y. Xin, L. Zhang, P. Gong, Differential expression and bioinformatics analysis of circrna in nonsmall cell lung cancer, Frontiers in Genetics 11 (2020) 586814.
- [111] N. Shriwash, P. Singh, S. Arora, S. M. Ali, S. Ali, R. Dohare, Identification of differentially expressed genes in small and non-small cell lung cancer based on meta-analysis of mrna, Heliyon 5 (6) (2019).
- [112] A. McDermaid, B. Monier, J. Zhao, B. Liu, Q. Ma, Interpretation of differential gene expression results of rna-seq data: review and integration, Briefings in bioinformatics 20 (6) (2019) 2044–2054.
- [113] M. Sobhan, A. M. Mondal, Explainable machine learning to identify patient-specific biomarkers for lung cancer, in: 2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), IEEE, 2022, pp. 3152–3159.
- [114] S. M. M. Hossain, L. Khatun, S. Ray, A. Mukhopadhyay, Pan-cancer classification by regularized multi-task learning, Scientific reports 11 (1) (2021) 24252.

- [115] M. I. Love, W. Huber, S. Anders, Moderated estimation of fold change and dispersion for rna-seq data with deseq2, Genome biology 15 (12) (2014) 1–21.
- [116] Q. Wang, J. Armenia, C. Zhang, A. V. Penson, E. Reznik, L. Zhang, T. Minet, A. Ochoa, B. E. Gross, C. A. Iacobuzio-Donahue, et al., BioRxiv (2017).
- [117] Q. Wang, J. Armenia, C. Zhang, A. V. Penson, E. Reznik, L. Zhang, T. Minet, A. Ochoa, B. E. Gross, C. A. Iacobuzio-Donahue, et al., Unifying cancer and normal rna sequencing data from different sources, Scientific data 5 (1) (2018) 1–8.