

Recurrent Neural Network Based Feature Selection for High Dimensional and Low Sample Size Micro-array Data

Shanta Chowdhury

Center for Computational Systems Biology
Prairie View A&M University
Prairie View, TX 77446, USA
schowdhury1@student.pvamu.edu

Xishuang Dong

Center for Computational Systems Biology
Prairie View A&M University
Prairie View, TX 77446, USA
xidong@pvamu.edu

Xiangfang Li

CCSB and Department of ECE
Prairie View A&M University
Prairie View, TX 77446, USA
xili@pvamu.edu

Abstract—Analyzing micro-array data faces many challenges such as high dimension, low sample size and sparse data. Feature selection is a technique to select more relevant features to implement dimension reduction to mitigate these challenges. In this paper, we propose a novel framework of feature selection based on recurrent neural network (RNN) to select a subset of features. Specifically, the proposed framework has been applied to select features from micro-array data for cell classification. We implement four feature selection models with different architectures of recurrent neural networks under the proposed framework, where these architectures include gated recurrent unit (GRU), long short-term memory (LSTM), RNN and bi-directional LSTM. The advantages of the proposed framework is demonstrated via real-world micro-array datasets.

Index Terms—recurrent neural network, feature selection, sparse data, gated recurrent unit, long-short term memory cell

I. INTRODUCTION

The analysis of gene expression data becomes a very important research topic for early disease diagnosis and drug development [1]. However, gene expression data processing is very challenging due to the rapid development of experimental technologies such as micro-array, next generation sequencing and mass spectrometer [2] as they generate high dimensional data. These types of data usually contain large number of features, however, typically only a portion of the features are relevant to the research problem at hand. If all the features are treated equally while performing machine learning (ML) such as classification on the data, it will degrade the performance of the ML model. Moreover, the model can suffer from the risk of overfitting through poor generalization ability [3]. These challenging issues of high dimensional data are posed as “the curse of dimensionality” [4]. A promising approach for the analysis of high dimensional biomedical data is to reduce the number of features, a technique known as feature selection. The goal of feature selection is to select an optimal subset of features so that the data can be presented in a more computationally feasible fashion. As a result, the classification performance will be improved through feature selection even though some features are dropped or ignored [5].

Traditional feature selection methods are categorized into four different types, namely, filter approach [6], wrapper approach [7], embedded approach [8], and hybrid approach [9]. Filter methods are simplest and computational efficient compared with other methods. It evaluates the value of features without any prior knowledge of learning algorithm. In this method, features dependencies and interaction between classifiers are ignored which leads the model to be misclassified. On the contrary, wrapper method considers the interaction in feature which guarantees better accuracy to classify the algorithm. The main disadvantage of the approach is its high complexity and poor generality. Wrapper method suffers from overfitting on small training set whereas filter method can be used as large number of features. On the other hand, embedded approach considers variable subset of selections to learn intensive feature dependencies. Hybrid approach is the combination of filter method and wrapper method. The goal of this approach is to gain best performance by intensive learning procedure.

Most recently, deep neural network has been achieved dramatic advancement in selecting feature from high dimensional of data [10]. Deep neural network is framed by multiple layers with non-linear activation functions which leads the model to mine more efficiently the pattern of complex feature format. It takes the advantages of its non-linear pattern recognition to dig deep inside of the data. However, it suffers from over fitting and high variance gradients for low sample size. Deep neural pursuit (DNP) [10] selects subset of features by overcoming the challenges. It incrementally selects and learns features and add them through multiple dropout technique to train the model for high dimension, low sample size data. However, for the micro-array data, DNP is not able to fully utilize relations between features to accomplish feature selection while genes, as the features, are correlated to each other.

In this paper, we propose recurrent neural network based feature selection model to extract features by directly using relations between feature on micro-array data. Firstly, we divide the features into two categories: selected features and candidate features. Then, we start with empty subset of feature

and consider the bias as selected features. In each step, the model chooses an individual feature from candidate feature and compute gradients through back propagation. Then the model calculates the average gradients and through this way, the model can include and exclude the features from select and candidate features accordingly. The main contribution of our work is to use recurrent model to select and update features from high dimension data, where the feature relations can be built by the recurrent connections of the neurons in the recurrent neural networks. Compared to DNP, the advantage of this proposed model is during each computation time, it is not only using its input feature information, but also using the information of neighbor feature information to enhance performance of feature selection.

In summary, the contributions of this research work are as follows:

- We propose a novel feature selection framework based on recurrent neural network to select features from high dimension data. We implement different feature selection models with various types of recurrent models, namely, gate recurrent unit (GRU), long-short term memory (LSTM) and bi-directional LSTM (Bi-lstm) to verify the proposed framework by testing on micro-array data.
- We validate our proposed model by experimenting on two types of high dimensional low sample size sparse biomedical micro-data namely Colon and Leukemia dataset and observe that proposed model performs better than deep neural pursuit (DNP).

The rest of the paper is organized as follows. Section II reviews some relevant works. In section III, we briefly describe our proposed approach. Dataset details, experimental set up and performance evaluation are presented in section IV. Finally, we conclude in section V.

II. RELATED WORK

In this section, we review some existing feature selection approaches namely linear and non-linear approach. In linear approach, data is mapped into a lower dimensional space. PCA (Principal Component Analysis) is popular linear model for feature selection. PCA maximizes the variance of data and employs orthogonal transformation to convert the data into lower dimensional space. The features associated with large eigenvalue contain huge amount of information. When the eigenvalue is small, the PCA fails to project the features in low dimension space [11]. However, PCA does not inherently capture the feature information from the data and does not work well for classification of data. Linear Discriminant Analysis (LDA) proposed by Ronald Fisher [12] is a popular linear model for feature selection and classification of the data. It maximizes the distance between the means, normalized by a measure of the sample-class variability. However, the model can not perform well on high dimensional low sample data. L1 regularized approach (Lasso) is very popular approach for dealing with high dimension low sample size data [13]. It minimizes the loss by L1 norms. The main limitation of this sparse linear model is it can not capture non-linear relationship

among input features. Therefore, HSIC lasso are used to handle non-linear relation of data. It employs Hilbert-Schmidt Independence Criterion (HSIC) to calculate the dependency between variables. It also uses L1 norms to a subset of features which results in convex optimization problem. Fast correlation Based Filter (FCBF) uses symmetric uncertainty to compute best subset of features with sequential search [15]. It selects features by establishing high correlation with the target variable and little correlation with other variables. Sparse additive model (SpAM) is used to select feature by back-fitting algorithm [16]. It is based on the combination of sparse linear modeling and additive regression. The main limitation of this model is it can not capture important interactions among features due to its additive manner.

For the non-linear, Minimum redundancy maximum relevancy (mRMR) is a prominent non-linear approach which ranks the features based on minimum redundancy maximum relevancy (mRMR) [14] criteria. Through mRMR criteria, it selects high relevance features. It computes the relevancy using F-statistics and mutual information for discrete and continuous features respectively. Redundancy is calculated through person correlation coefficient. However, this model selects high relevant features with a high correlation with the class (output) and a presents low correlation between themselves which results in the loss of temporal data information. Deep neural network is employed to extract features by using the non-linear relations between features. Deep feature section (DFS) [17] is used to select input features in deep structure for multi-class classification. DFS performs better comparatively than Lasso. The potential weakness of the DFS model is it fails to achieve the sparse connections for high dimensional low sample size data. DNP [10] proposed multiple dropout approaches to alleviate the limitation of feature selection on sparse data. To motivated by DNP, we propose recurrent neural network based feature selection approach and enhance the performance of high dimensional low sample sparse data.

III. METHODOLOGY

We propose a novel framework of feature selection to select relevant feature set among sparse high dimensional low sample size data in order to enhance the biomedical sequencing data analysis. Firstly, the features are divided into two types of feature set: selected set and candidate set. Initially, selected feature set starts from a bias. All weights including bias in the neural network are initialized as zero. The input weight comprises with selected weights and candidate weights. The input weights is initialized through Xavier Initializer. Neural Network utilizes multiple dropout technique to avoid high variance gradients. It randomly drops neurons multiple times, computes gradients based on neurons and connections and averages multiple gradients. Such multiple dropout technique obtains averaged gradients with low variance. Then whole neural network will convergence until all candidate weights are zero. In this way, features are added and removed. When number of selected features and maximum number of features will be equal, the proposed model will have a subset of feature

selection. Then the feature subset is used to train machine learning classifiers such as Decision tree, Naive bayes and Support vector machine to complete data analysis. The training dataset is also fit into the model. After training procedure, the machine learning models are used to evaluate the model performance through test data. The flow of feature selection is shown on Figure 1.

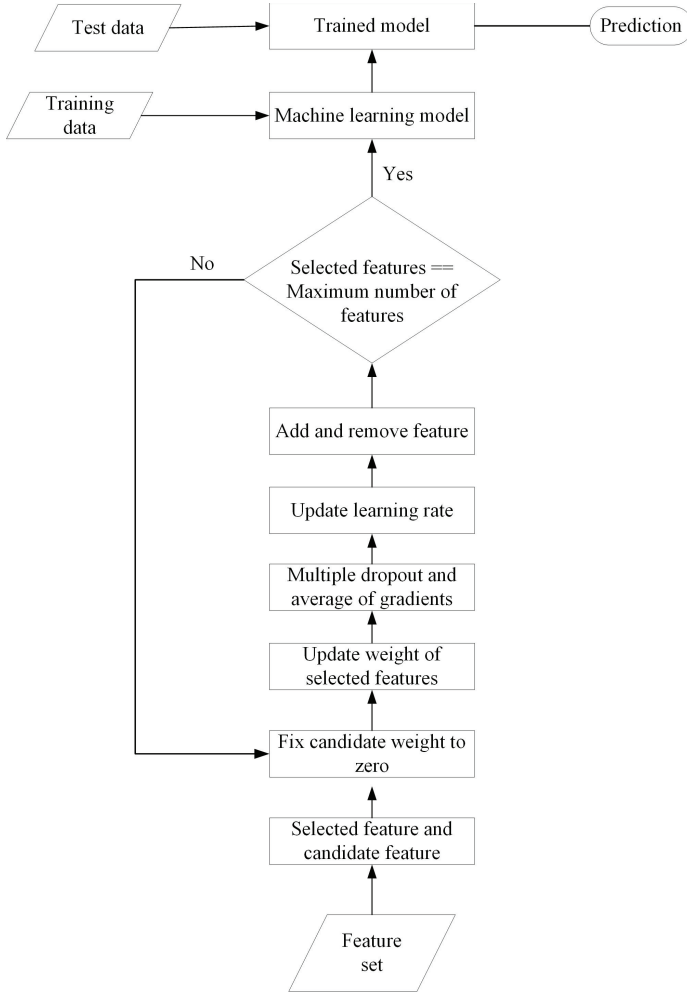


Fig. 1. Flow of feature selection under the proposed framework.

Specially, we used recurrent neural network (RNN) to select features for cell classification on micro-array data. We motivated for this framework through DNP model. In DNP model, the deep neural network is used. Fig. 2 represents how DNP works [10]. Firstly, DNP trains smaller sub-network and incrementally selects features to find local optima. It drops neurons multiple times and uses back-propagation technique to add and remove features. DNP uses the deep neural network architecture whereas in the proposed framework we take the advantage of recurrent neural network model which helps the model to store feature information in its memory and enhance the performance. It computes each feature of input sequence and transforms it into vector format using following equation (1) and (2) [18].

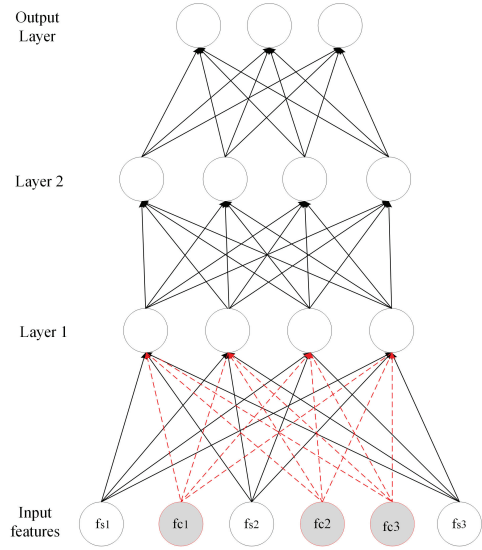


Fig. 2. The architecture of feature selection via deep neural network, where solid lines indicate selected feature and red dashed lines indicate candidate feature.

$$h_t = H(U_{xh}x_t + U_{hh}h_{t-1} + b_h). \quad (1)$$

$$y_t = U_{hy}h_t + b_y. \quad (2)$$

where U_{xh} , U_{hh} , U_{hy} denote the weight matrices of input-hidden, hidden-hidden and hidden-output processes, respectively. h_t is the vector of hidden states that derive the information from current input x_t and the previous hidden state h_{t-1} . Fig. 3 represents how RNN works features in each sequence.

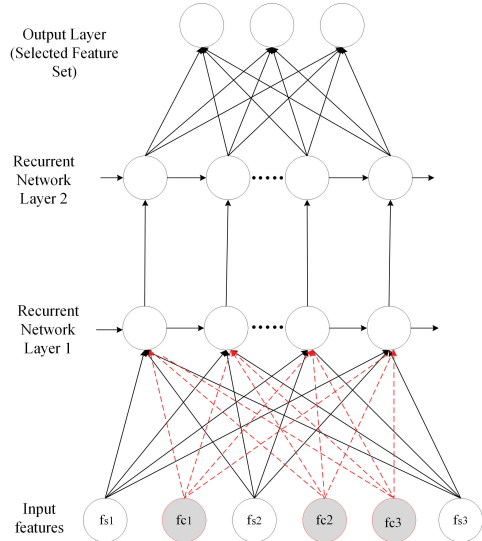


Fig. 3. RNN based feature selection model.

In our proposed model, the features are fed as input into the model. Initially, selected feature set starts as bias term. The input weights are initialized through Xavier Initializer. Each time the weights is updated through back propagation and compute

the dropout. As the weights are shared through all the layers in recurrent network, it contains the sequence information in its neuron. Hence, each time when it is computing its gradient, the gradient computation do not only based on current feature information but also based on sequence feature information which helps the model to enhance the performance.

We also used different types of RNN model namely Gated recurrent unit (GRU) [19], Long short term memory (LSTM) [20] and Bi-lstm to observe the model performance. The main limitation of recurrent model is sometimes it may suffer from vanishing gradients at the time of computation. GRU and LSTM are advanced type of RNN which structured in a special way so that they can deal with the limitation of RNN.

GRU consists of update gate and reset gate which helps the model to decide which feature should be passed through the output layer. It uses following equations (3)-(6) to compute how much information, GRU will carry forward through the network.

$$z_t = \sigma(W_z x_t + U_z h_{t-1}) \quad (3)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1}) \quad (4)$$

$$h_{ti} = \tanh(W x_t + r_t \odot U h_{t-1}) \quad (5)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot h_{ti} \quad (6)$$

Here, z_t , r_t and h_{ti} denote the equation for update gate, reset gate and current memory respectively. W and U represent the weight matrices of each gate. In the final memory content, element wise multiplication is applied to update the information in update gate and determines how much information will be hold through the network.

On the other hand, LSTM has three gates: input, output and forget gate to regulate dataflow in its memory. The computation equation that LSTM uses are presented as follows:

$$i_t = \sigma(U_{xi} x_t + U_{hi} h_{t-1} + U_{ci} c_{t-1} + b_i) \quad (7)$$

$$f_t = \sigma(U_{xf} x_t + U_{hf} h_{t-1} + U_{cf} c_{t-1} + b_f) \quad (8)$$

$$c_t = f_t c_{t-1} + i_t RELU(U_{xc} x_t + U_{hc} h_{t-1} + b_{cc}) \quad (9)$$

$$o_t = \sigma(U_{xo} x_t + U_{ho} h_{t-1} + U_{co} c_t + b) \quad (10)$$

$$h_t = o_t RELU(c_t) \quad (11)$$

where U and σ indicate weight matrices and logistic sigmoid function respectively. Different gates refer as different indices, like input gate as i , forget as f , cell as c and output as o . These gates and activation functions soothe LSTM to avoid the limitation of vanishing gradients by storing long term dependencies terms.

Both GRU and LSTM has similar types of architecture, GRU is much simpler and trains faster than LSTM. Besides, GRU performs better when the model does not require long term dependencies information and trained on less training data and. On the other hand, There are two different states in bi-directional LSTM. In forward states, they compute future sequence information whereas backward states compute past sequence information and finally generate the output o_t by integrating two hidden states computation results.

IV. EXPERIMENT

A. Datasets

We used public available bio-medical micro-array data to evaluate the performance of our proposed framework. We focused on two different kinds of bio-medical data namely, Colon and Leukemia [21]. Each of the dataset consists of large amount of sparse data. Colon cancer dataset contains 62 samples and 2000 genes. The dataset is classified as tumor and normal tissues. Leukemia dataset has total 7070 number of genes and 72 number of observation sets. All of the samples are collected from Leukemia patients either they have acute lymphoblastic leukemia (ALL) or acute myelogenous leukemia (AML). The details of both datasets are represented in table I.

TABLE I
DATASET DETAILS

Data	Colon	Leukemia
Sample Size	62	72
Feature Size	2000	7070
Zero	51554	222326
Non-zero Value	72446	286715
Sparsity	41.58%	43.67%

B. Experimental setup

The key parameters for the proposed methodology are: Learning rate: 0.1, Dropout rate: 0.5, Dropout iteration: 50, Maximum iteration: 25. We used DNP as baseline model as it outperforms the traditional feature selection models [10]. In addition, we implemented other types of recurrent model using proposed framework to observe the performance. We used same experimental set up every time. The details of the experimental set up is shown represented in table II:

TABLE II
PARAMETERS OF PROPOSED METHODOLOGY

Name	Description
Input	Feature set
Number of feature	25
Recurrent model	RNN, GRU, LSTM, Bi-LSTM
Number of layers	2
Number of neurons in each layer	[30, 20]

C. Evaluation metrics

We perform 10-fold cross validation to select feature from high dimensional data. After extracting features from high dimensional sparse data using RNN model, the data is fit into traditional machine learning classifiers to evaluate the performances. We employ confusion matrices namely Precision, Recall and F-score to demonstrate the evaluation. Precision [22] defines how accurately and exactly a model can recognize correct category whereas Recall [23] indicates the percentage of total relevant results correctly classified through a model. F-score [24] is the average of precision and recall value.

$$Precision = \frac{TP}{TP + FP} \quad (12)$$

$$Recall = \frac{TP}{TP + FN} \quad (13)$$

$$Fscore = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (14)$$

whereas TP (True Positive) counts total number of predicted class matched with actual class. FP (False positive) counts total number of predicted class does not match with the actual class. FN (False negative) measures actual labels not present in the predicted labels.

D. Result and Discussion

The experimental results are summarized in table III and IV using different evaluation matrices and employed different classifiers to evaluate the model performance. In table III, the comparison result shows that the proposed model outperforms DNP model in some cases. For instance, Fscore is improved around 20% for decision tree and SVM classifiers using RNN with the comparison of using DNP model. On the other hand, the improvement of Fscore in terms of NB is around 2%.

The comparison results for leukemia dataset are presented in table IV. It is observed that the proposed model shows the effectiveness for all the classifiers. The Fscore is improved by 11% and 15% for decision tree and NB classifiers in terms of RNN based feature selection. For RNN based feature selection for SVM classifier, the model performance degrades by 15%. The reason behind the lower performance of the classifier may be the data is too sparse to be suitable for the classifier. Hence, we are observing lower performance for SVM classifier on different RNN model for leukemia dataset.

Based on comparing the performance of the proposed model on different classifiers, it is observed that RNN based feature selection model outperforms the DNP model. Comparing two tables, it is found that GRU performs best whereas LSTM model shows poor performance among all the recurrent models. The reason behind poor performance of lstm is due to independency of data pattern. LSTM performs best when it arises to restore information in long dependencies of features in sequence data. Moreover, LSTM performs better for large training data. Comparatively, GRU performs best for low sample size data. Here, in our experiment, the dataset is low sample size and does not require long term dependencies among feature. Hence, we obtain poor performance for LSTM model.

Analyzing the performance of Colon and Leukemia data on different classifiers, we can find that Naive Bayes (NB) performs better than other two machine learning classifiers (Decision tree and Support vector machine). The reason is the sparsity which seriously affect model performance. Leukemia data is more sparse than Colon data. Among other classifiers, NB is less affected through sparsity [25]. Hence, we have better performance for NB than other classifiers.

V. CONCLUSION

In this paper, we proposed the recurrent neural network based feature selection to improve the cell classification on high dimension low sample size data. Initially, we start with empty subset of features as bias. Then, we incrementally add features by averaging gradients through multiple dropout technique. We fit the selected features in machine learning classifiers to evaluate the model's performance. Experimental results show that our model has better performance of selecting subset of features in some cases. In the future, we plan to extend the proposed model to implement feature selection on on sparse high dimension large sample data like the single cell sequencing data for improving single cell classification.

ACKNOWLEDGMENT

This research work is supported in part by the Texas A&M Chancellor's Research Initiative (CRI), the U.S. National Science Foundation (NSF) award 1464387 and 1736196, and by the U.S. Office of the Under Secretary of Defense for Research and Engineering (OUSD(R&E)) under agreement number FA8750-15-2-0119. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. National Science Foundation (NSF) or the U.S. Office of the Under Secretary of Defense for Research and Engineering (OUSD(R&E)) or the U.S. Government.

REFERENCES

- [1] M. Yamada, J. Tang, J. Lugo-Martinez, E. Hodzic, R. Shrestha, A. Saha, H. Ouyang, D. Yin, H. Mamitsuka, C. Sahinalp, and others, "Ultra high-dimensional nonlinear feature selection for big biological data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 30, pp. 1352–1365, 2018.
- [2] Y. Li and L. Chen, "Big biological data: Challenges and opportunities," *Genomics, Proteomics & Bioinformatics.*, vol. 12, no.5, pp.187–189, 2014.
- [3] V. Pappu and P. M. Pardalos, "High-Dimensional Data Classification," in *Clusters, Orders, and Trees: Methods and Applications*, pp. 119–150, Springer, New York, 2014.
- [4] R. Bellman, *Dynamic Programming*, 1st ed. Princeton, NJ, USA: Princeton University Press, 1957.
- [5] D. Donoho and J. Jin, "Higher criticism thresholding: Optimal feature selection when useful features are rare and weak," in *Proceedings of the National Academy of Sciences*, vol. 105, no. 39, pp. 14790–14795, 2008.
- [6] X. He, D. Cai, and P. Niyogi, "Laplacian score for feature selection," in *Proceedings of the 18th International Conference on Neural Information Processing Systems*, pp. 507–514, 2005.
- [7] S. Tabakhi, P. Moradi, and F. Akhlaghian, "An unsupervised feature selection algorithm based on ant colony optimization," *Engineering Applications of Artificial Intelligence*, vol. 32, no. 6, pp. 112–123, 2014.
- [8] S. Wang, J. Tang, and H. Liu, "Embedded unsupervised feature selection," in *Proceedings of the Twenty-Ninth MAI Conference on Artificial Intelligence*, pp. 470–476, 2015.
- [9] H. Lu, J. Chen, K. Yan, Q. Jin, Y. Xue, and Z. Gao, "A hybrid feature selection algorithm for gene expression data classification," *Elsevier*, vol. 256, pp. 56–62, 2017.
- [10] B. Liu, Y. Wei, Y. Zhang and Q. Yang, "Deep Neural Networks for High Dimension, Low Sample Size Data," in *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 2287–2293, 2017.

TABLE III
COMPARISON RESULTS FOR COLON DATASET

Feature Selection	Decision Tree			Naive Bayes			Support Vector Machine		
	Precision	Recall	F-score	Precision	Recall	F-score	Precision	Recall	F-score
DNP	0.47	0.47	0.46	0.60	0.62	0.62	0.38	0.62	0.47
RNN	0.68	0.64	0.65	0.65	0.64	0.64	0.68	0.70	0.66
GRU	0.72	0.68	0.69	0.60	0.60	0.60	0.66	0.66	0.66
LSTM	0.60	0.62	0.61	0.65	0.64	0.64	0.8	0.72	0.65
Bi-LSTM	0.46	0.54	0.67	0.36	0.60	0.45	0.36	0.60	0.45

TABLE IV
COMPARISON RESULTS FOR LEUKEMIA DATASET

Feature Selection	Decision Tree			Naive Bayes			Support Vector Machine		
	Precision	Recall	F-score	Precision	Recall	F-score	Precision	Recall	F-score
DNP	0.64	0.53	0.55	0.58	0.60	0.59	0.76	0.66	0.67
RNN	0.77	0.66	0.66	0.74	0.74	0.74	0.43	0.66	0.52
GRU	0.63	0.59	0.60	0.77	0.72	0.73	0.79	0.79	0.79
LSTM	0.58	0.52	0.53	0.57	0.57	0.57	0.78	0.67	0.56
Bi-LSTM	0.56	0.59	0.57	0.68	0.69	0.66	0.41	0.64	0.50

- [11] A. Malhi and R. X. Gao, "PCA-based feature selection scheme for machine defect classification," in IEEE Transactions on Instrumentation and Measurement, vol. 53, no. 6, pp. 1517–1525, 2004.
- [12] M. Masaeli, J. G. Dy and G. M. Fung, "From transformation-based dimensionality reduction to feature selection," in Proceedings of the 27th International Conference on Machine Learning (ICML-10), pp. 751 – 758, 2010.
- [13] M. Yamada, W. Jitkrittum, L. Sigal, E. P. Xing, and M. Sugiyama, "High-dimensional feature selection by feature-wise kernelized lasso", Neural Computation, vol. 26, no. 1, pp. 185-207, 2014.
- [14] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf, "Measuring statistical dependence with Hilbert-Schmidt norms", in Proceedings of 16th International Conference on Algorithmic Learning Theory, pp. 63-77, 2005.
- [15] L. Yu, and H. Liu, "Feature selection for high-dimensional data: A fast correlation-based filter solution", in Proceedings of the 20th International Conference on Machine Learning, pp. 856-863, 2003.
- [16] P. Ravikumar, J. Lafferty, H. Liu, and L. Wasserman, "Sparse additive models", Journal of the Royal Statistical Society: Series B, vol. 71, no. 5, pp. 1009-1030, 2009.
- [17] Y. Li, C. Chen, and W. W. Wasserman, "Deep feature selection: Theory and application to identify enhancers and promoters," in International Conference on Research in Computational Molecular Biology, pp 205–217, Springer, 2015.
- [18] "Recurrent Neural Networks Tutorial, Part 1 – Introduction to RNNs," [Online]. Available : <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns>. [Accessed: 14-Sep-2019]
- [19] "Understanding GRU Networks," [Online]. Available : <https://towardsdatascience.com/understanding-gru-networks-2ef37df6c9be>. [Accessed: 14-Sep-2019]
- [20] "Long short-term memory," [Online]. Available : https://en.wikipedia.org/wiki/Long_short-term_memory. [Accessed: 15-Sep-2019]
- [21] "The details on the six datasets used in this project can be found" [Online]. Available: <https://www.ntu.edu.sg/home/elhchen/data.htm>. [Accessed: 17-Sep-2019]
- [22] "Classification: Precision and Recall — Machine Learning Crash Course," [Online]. Available: <https://developers.google.com/machine-learning/crash-course/classification/precision-and-recall>. [Accessed: 17-Sep-2019].
- [23] "Evaluation Metrics for Machine Learning - Accuracy, Precision, Recall, and F1 Defined". [Online]. Available: <https://skymind.ai/wiki/accuracy-precision-recall-f1>. [Accessed: 17-Sep-2019].
- [24] "Classification Accuracy is Not Enough: More Performance Measures You Can Use" . [Online]. Available : <https://machinelearningmastery.com/classification-accuracy-is-not-enough-more-performance-measures-you-can-use/>. [Accessed: 17-Sep-2019].
- [25] J. Bissmark and O. Wärnling, "The Sparse Data Problem Within Classification Algorithms: The Effect of Sparse Data on the Naïve Bayes Algorithm," 2017.