

Superior Machine Learning Method for breast cancer cell lines identification

Sajid Farooq

Center for Lasers and Applications(CLA)
Instituto de Pesquisas Energéticas e Nucleares (IPEN)
Sao Paulo, Brazil
sajiddahar@gmail.com

Amanda Caramel-Juvino

Center for Lasers and Applications(CLA)
Instituto de Pesquisas Energéticas e Nucleares (IPEN)
Sao Paulo, Brazil
caramel@usp.br

Matheus Del-Valle

Center for Lasers and Applications(CLA)
Instituto de Pesquisas Energéticas e Nucleares (IPEN)
Sao Paulo, Brazil
matheus.valle@usp.br

Sofia Santos

Center for Radiopharmaceutics
Instituto de Pesquisas Energéticas e Nucleares (IPEN)
Sao Paulo, Brazil
snsantos85@gmail.com

Emerson Soares Bernardes

Center for Radiopharmaceutics
Instituto de Pesquisas Energéticas e Nucleares(IPEN)
Sao Paulo, Brazil
emerson.bernardes@gmail.com

Denise Maria Zzell

Center for Lasers and Applications(CLA)
Instituto de Pesquisas Energéticas e Nucleares (IPEN)
Sao Paulo, Brazil
zezell@usp.br

Abstract—We propose an artificial intelligence platform based on machine learning (ML) algorithm using Neighborhood Component analysis and K-Nearest Neighbors for breast cancer cell lines recognition. Our model presents up to 97% accuracy for identification of breast cancer cell lines.

Index Terms—Machine learning, Cell lines, FTIR, Accuracy, Breast Cancer

I. INTRODUCTION

Breast cancer (BC) is the leading cause of cancer-related fatality globally among women [1]. Due to complex tumour micro-environment system, this cancer is a big challenge of clinical decision making procedures [2]. Therefore, the breast cancer stratification based on molecular cell lines plays a remarkable role during treatment. Moreover, appropriate classification also assists not only to choose an accurate therapy methodology but also sort patients with prognostic diagnosis [5]. Further, breast cancer cell lines are up to 92 and are evolved based on the expression of receptors. Although, these cell lines can be classified into distinct subtypes: as describe by luminal A, luminal B, HER2, and triple-negative [3]. Using the erratic stratification of these cell lines, we are overwhelmed with BC cell lines nonexistent with a lot of characteristics documentation and consistent BC cell lines [5].

Therefore, models of identifying BC cell lines are urgently required to avoid the under-diagnosis of tumours before metastasis and to reduce the over-treatment of low risk disease,

that could assist to reduce the demand of aggressive systemic therapy.

To cope with the challenge, we aim to develop a ML algorithm for the recognition the BC cell lines accurately and rapidly. A ML algorithm is designed to use a semi-supervised and supervised classifiers i.e. Neighborhood Component Analysis (NCA) and K-Nearest Neighbors (KNN) [6]. The performance to achieve accurate prognostication perfectly with NCA-KNN algorithm that evidences de novo learning algorithm, is an important factor to identify BC cell lines.

II. METHODS AND MATERIALS

We choose breast cancer cell lines such as BT-474 i.e. a luminal B (ER/PR/HER2 positive), and SKBR-3, a HER2 (ER/PR negative and HER2 positive). To compute the model, we obtained an input data using Fourier transform Infrared Spectroscopy (FTIR). The data is collected in the shape of spectral image acquisitions and number of spectrums are >10k. The aforementioned data is preprocessed using Savitzky-Golay(SG) and then, extended multiplicative signal correction (EMSC). Later on, it is normalized to perform a digital de-waxing [7]. To apply NCA-KNN method on the BC cell lines, the NCA-KNN method is employed with cross-validation (k fold = 10). The receiver operating characteristics (ROC) curve is explored to evaluate accuracy. Others different individual classifiers were used as Support Vector Machine (SVM), AdaBoost and KNN for comparison purpose.

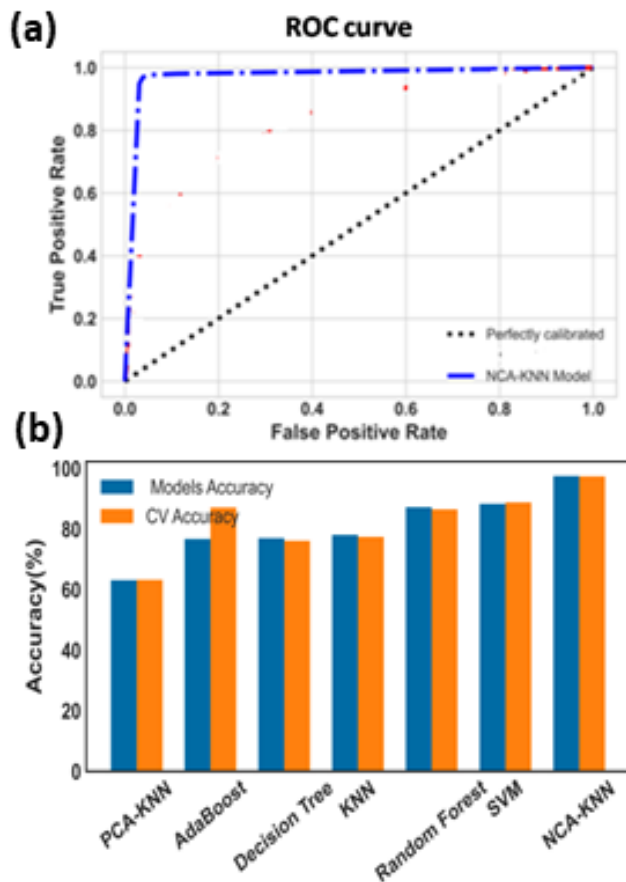


Fig. 3. The receiver operating characteristics (ROC) curve to identify the accuracy (a) and different model performance (b) (KNN–NCA).

IV. CONCLUSION

We presented a ML algorithm employing two supervised classifiers to evaluate the stratification of breast cancer cell lines. Our proposed method shows the higher performance up to 97% , which is around 9% higher than that of second best SVM method, presenting a hidden potential for the cancer cell lines identification with superior performance for clinical uptake.

ACKNOWLEDGMENT

This work was supported by FAPESP [17/50332-0] and CNPq [INCT-465763/2014-6, PQ-31451712021-9, 142229/2019-9].

V. CONFLICT OF INTEREST

The authors declare no conflict of interest.

REFERENCES

- [1] Sammut et al., “Multi-omic machine learning predictor of breast cancer therapy response,” *Nature* 601 **601**, 623–629, (2022).
- [2] Amiri Souri, E., Alicia Chenoweth, Anthony Cheung, Sophia N. Karagiannis, and Sophia Tsoka, “Cancer Grade Model: a multi-gene machine learning-based risk classification for improving prognosis in breast cancer,” *British Journal of Cancer* **125**, 748–758, no. 5 (2021).

- [3] Dai, Xiaofeng, Hongye Cheng, Zhonghu Bai, and Jia Li, “Breast cancer cell line classification and its relevance with breast tumor subtyping,” *Journal of Cancer* **8**, no. 16, 3131–(2017)
- [4] Horr, Christina, and Steven A. Buechler. ”Breast Cancer Consensus Subtypes: A system for subtyping breast cancer tumors based on gene expression.” *NPJ breast cancer* 7.1 (2021): 1-13.
- [5] Horr, Christina, and Steven A. Buechler, “Breast Cancer Consensus Subtypes: A system for subtyping breast cancer tumors based on gene expression,” *NPJ breast cancer* 7, 1-13 (2021)
- [6] Goldberger, Jacob, Geoffrey E. Hinton, Sam Roweis, and Russ R. Salakhutdinov, “Neighbourhood components analysis,” *Advances in neural information processing systems* **17**, (2004)
- [7] del-Valle, et al., “The impact of scan number and its preprocessing in micro-FTIR imaging when applying machine learning for breast cancer subtypes classification,” *Vibrational Spectroscopy* **117**,103309– (2021)