

FeCO3, constructing the network biomarkers using the inter-feature correlation coefficients

Shenggen Lin^{1,2,*}, Yuqi Lin^{1,*}, Kexin Wu^{1,*}, Yueying Wang³, Zixuan Feng¹, Meiyu Duan¹, Shuai Liu¹, Jingxuan Huang⁴, Zhao Wang⁴, Yusi Fan⁴, Lan Huang¹, Fengfeng Zhou^{1,#}.

¹ College of Computer Science and Technology, and Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun, Jilin 130012, China.

² State Key Laboratory of Microbial Metabolism, and School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, Shanghai 200240, China.

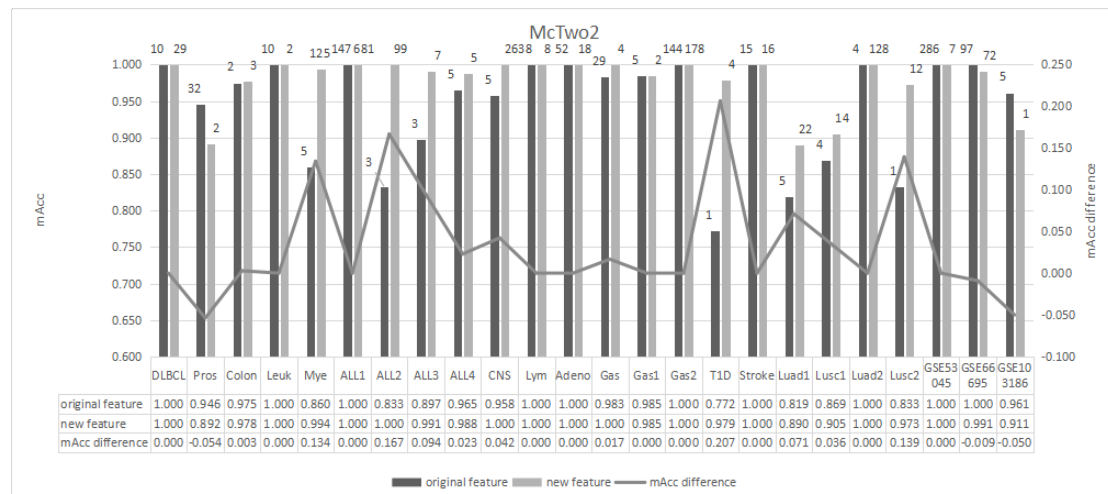
³ Department of Epidemiology and Biostatistics, School of Public Health, Jilin University, Changchun, Jilin Province, China.

⁴ College of Software, and Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, Changchun, Jilin 130012, China.

* Contributed equally to this study.

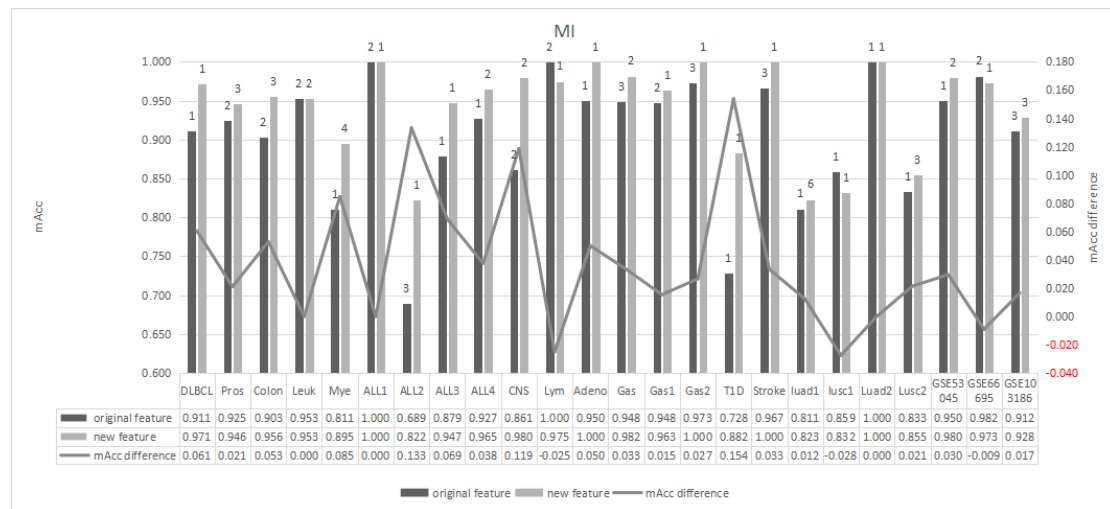
Correspondence may be addressed to Fengfeng Zhou: FengfengZhou@gmail.com or ffzhou@jlu.edu.cn . Lab web site: <http://www.healthinformatics.org/> . Phone: +86-431-8516-6024. Fax: +86-431-8516-6024.

Supplementary Figure S1



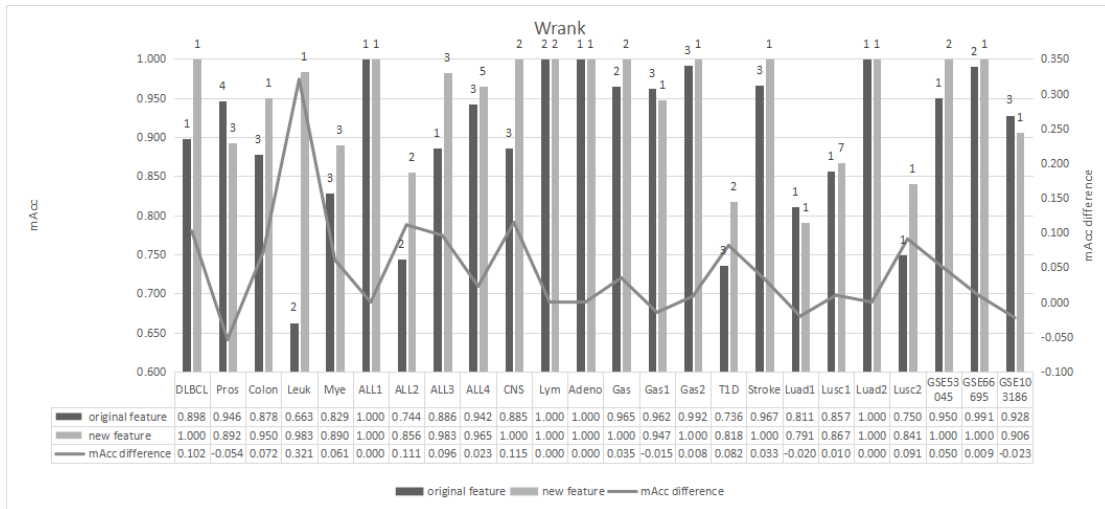
Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm McTwo2.

The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms

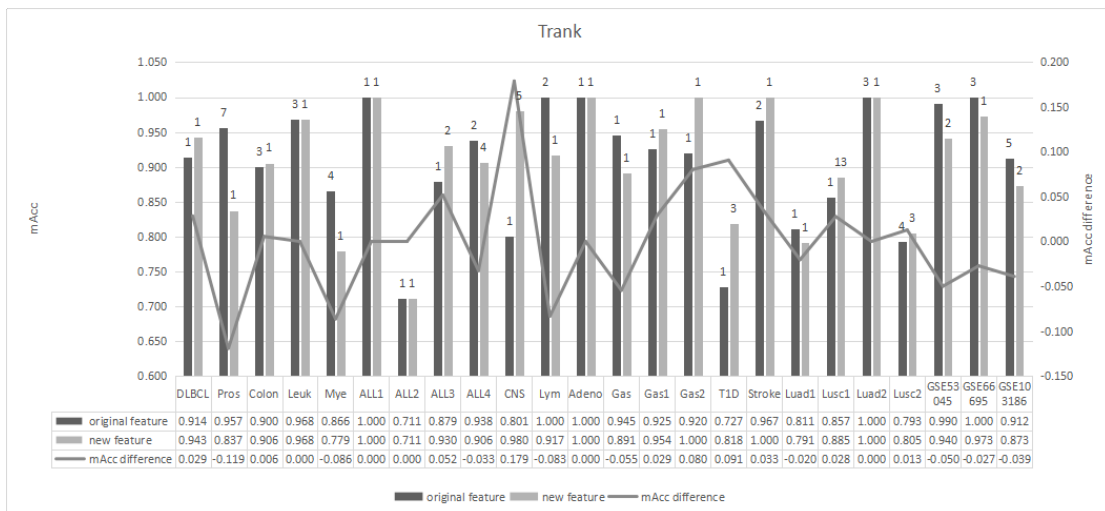


Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm MI.

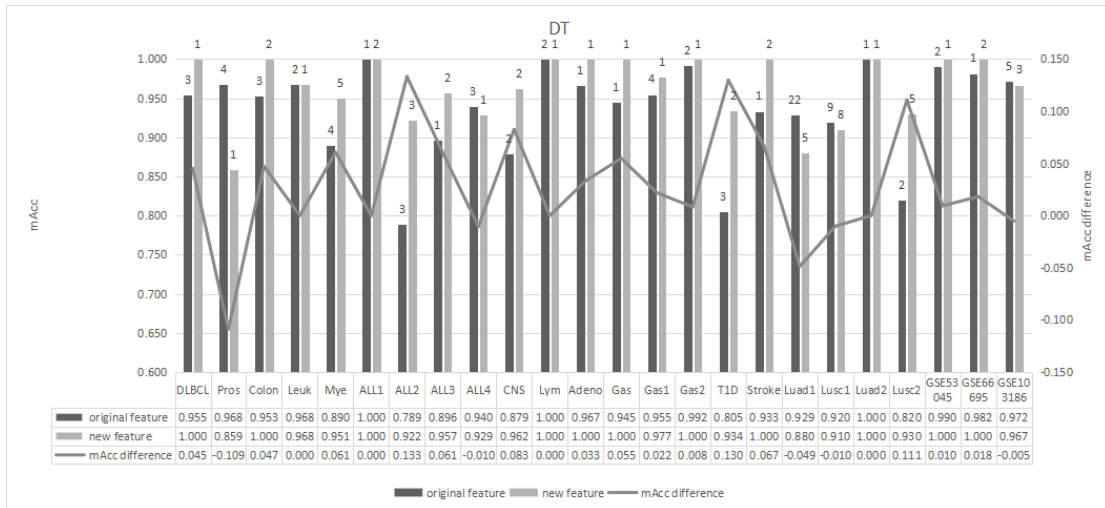
The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms



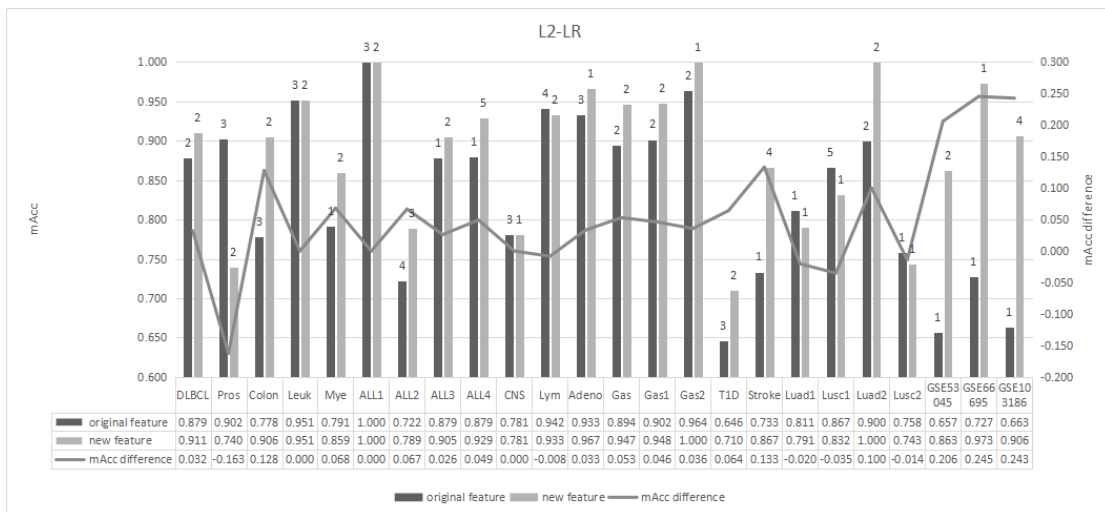
Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm Wrnk. The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms



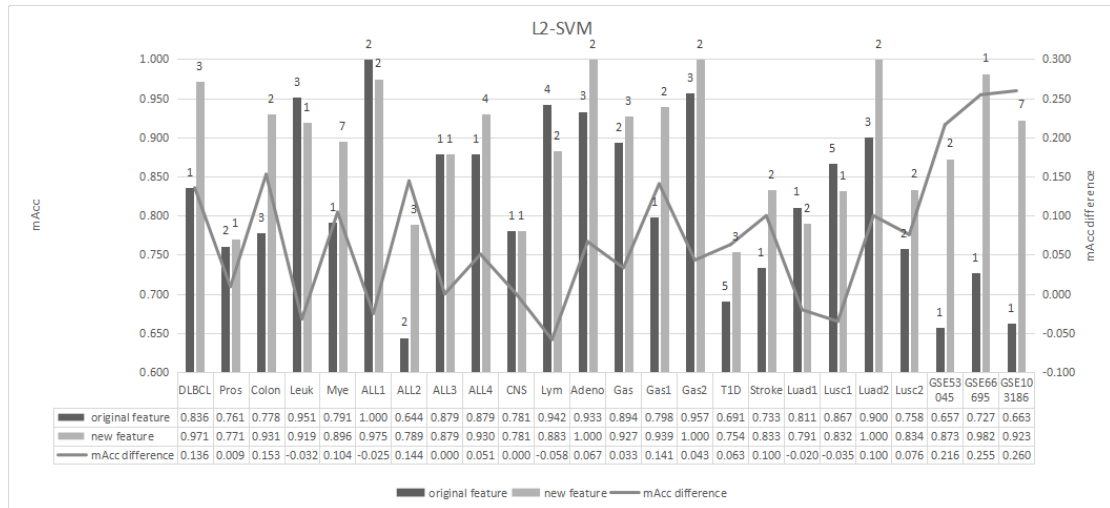
Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm Trank. The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms



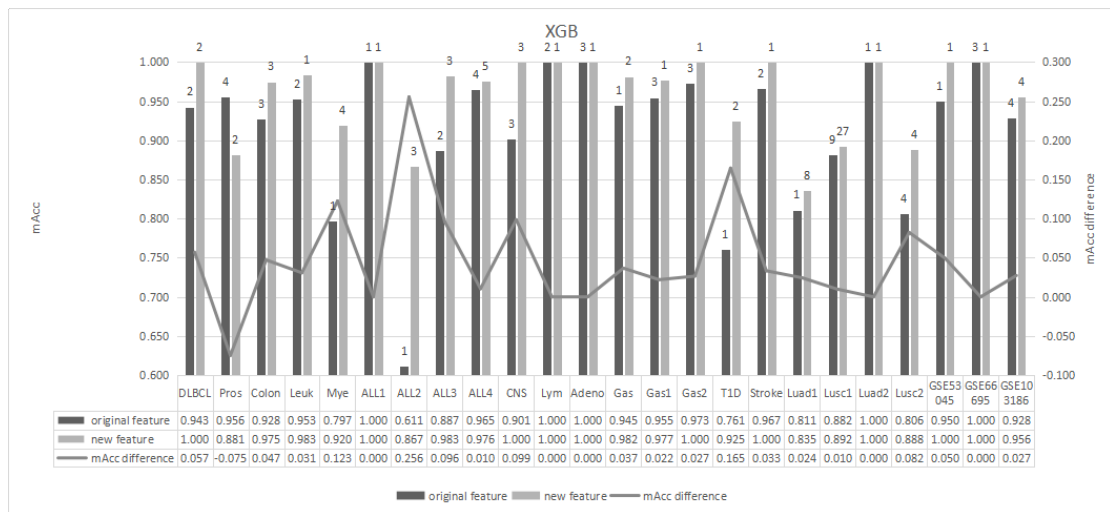
Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm DT. The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms



Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm L2-LR. The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms



Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm L2-SVM. The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms



Comparison of the classification metric mAcc and the number of features between the original features and the FeCO3 features using the feature selection algorithm XGB. The series “original” gives the data of the original features, and the series “FeCO3” illustrates the status of the FeCO3 features. The series “Improvement” is the mAcc value of the FeCO3 features minus the mAcc value of the original features. The classification accuracy is calculated using the 5-fold cross validation strategy. The data was calculated using the feature selection algorithms

References

- Cherkassky, V. The nature of statistical learning theory~. *IEEE Trans Neural Netw* 1997;8(6):1564.
- Ge, R., *et al.* McTwo: a two-step feature selection algorithm based on maximal information coefficient. *BMC Bioinformatics* 2016;17:142.
- Geva, S. and Sitte, J. Adaptive nearest neighbor pattern classification. *IEEE Trans Neural Netw* 1991;2(2):318-322.
- Guo, P., *et al.* Gene expression profile based classification models of psoriasis. *Genomics* 2014;103(1):48-55.
- Haridas, V., *et al.* TRANK, a novel cytokine that activates NF-kappa B and c-Jun N-terminal kinase. *J Immunol* 1998;161(1):1-6.
- Li, Y., *et al.* Association of miR-155 and Angiotensin Receptor Type 1 Polymorphisms with the Risk of Ischemic Stroke in a Chinese Population. *DNA Cell Biol* 2019.
- Liu, W.M., *et al.* Analysis of high density expression microarrays with signed-rank call algorithms. *Bioinformatics* 2002;18(12):1593-1599.
- Mortazavi, A. and Moattar, M.H. Robust Feature Selection from Microarray Data Based on Cooperative Game Theory and Qualitative Mutual Information. *Adv Bioinformatics* 2016;2016:1058305.
- Nguyen, D.H. and Patrick, J.D. Supervised machine learning and active learning in classification of radiology reports. *J Am Med Inform Assoc* 2014;21(5):893-901.
- Pal, M. Random forest classifier for remote sensing classification. *International journal of remote sensing* 2005;26(1):217-222.
- Pregibon, D. Logistic regression diagnostics. *The Annals of Statistics* 1981;9(4):705-724.
- Senders, J.T., *et al.* Natural Language Processing for Automated Quantification of Brain Metastases Reported in Free-Text Radiology Reports. *JCO Clin Cancer Inform* 2019;3:1-9.
- Waldmann, P. On the Use of the Pearson Correlation Coefficient for Model Evaluation in Genome-Wide Prediction. *Front Genet* 2019;10:899.
- Wang, Y., *et al.* In Silico Prediction of Human Intravenous Pharmacokinetic Parameters with Improved Accuracy. *J Chem Inf Model* 2019;59(9):3968-3980.
- Ye, Y., *et al.* RIFS: a randomly restarted incremental feature selection algorithm. *Sci Rep* 2017;7(1):13013.
- Yu, X., *et al.* Individual-specific edge-network analysis for disease prediction. *Nucleic Acids Res* 2017;45(20):e170.