Faculty of Engineering and Information Technology

University of Technology Sydney

# Disease Gene Recognition and Editing Optimization Through Knowledge Learned from Domain Feature Spaces

A thesis submitted in partial fulfillment of

the requirements for the degree of

**Doctor of Philosophy**

by

## Hui Peng

May 2019

# CERTIFICATE OF AUTHORSHIP/ORIGINALITY

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Candidate

Production Note:
Signature removed
prior to publication.

i

# Acknowledgments

Foremost, I hope to express my sincere gratitude to my supervisor Prof. Jinyan Li for his continuous support during my PhD study and research, for his patience, motivation, and immense knowledge. His guidance helped me improve both my research skills and other necessary abilities such as scientific writing, academic communicating and presentation skills. The completion of this thesis and related researches would not be possible without his constructive advice for improving them and his valuable time and efforts to make them perfect.

I would like to appreciate Dr. Tao Liu and Prof. Gyorgy Hutvagner, two of my research partners, for their patience to discuss the problems of our research and for their insightful suggestions to finish the research tasks. Many thanks to my co-supervisor Prof. Dacheng Tao for his kind suggestions and help to a part of my research work.

I really appreciate our two team members Yi and Chaowang, who are also my friends and roommates, for their accompany during the past three and a half years, for their great help in not only my study and research but also my life abroad. I am grateful to the three former team members Dr. Jing Ren, Dr. Renhua Song and Dr. Shameek Ghosh, for their warm responses when I requested help from them during the first two years in UTS and even after their graduation. I also thank the other team members Zhixun, Yuansheng, Xiaocai, Xuan and Tao, who joined us in recent two years. I feel very happy to meet all of you in UTS and hold a lot of fantastic activities with you. Those happy moments will be kept in my mind forever.

# Contents

# List of Figures

# List of Tables

# List of Publications

Below is the list of journal and conference papers associated with my PhD research:

**Journal Papers Published**

- **Peng, H.**, Zheng, Y., Blumenstein, M., Tao, D., & Li, J. (2018). CRISPR/Cas9 cleavage efficiency regression through boosting algorithms and Markov sequence profiling. **Bioinformatics**, 34(18), pp.3069-3077.

- **Peng, H.**, Zheng, Y., Zhao, Z., Liu, T., & Li, J. (2018). Recognition of CRISPR/Cas9 off-target sites through ensemble learning of uneven mismatch distributions. **Bioinformatics**, 34(17), pp.i757-i765. (Oral Presentation at the 17th Europe Conference on Computational Biology (ECCB 2018))

- Zheng, Y., **Peng, H.**, Ghosh, S., Lan, C., & Li, J. (2018). Inverse Similarity and Reliable Negative Samples for Drug Side-effect Prediction. **BMC Bioinformatics**, 19(13), p.554.

- Zheng, Y., **Peng, H.**, Zhang, X., Zhao, Z., & Li, J. (2018). Predicting adverse drug reactions of combined medication from heterogeneous pharmacologic databases. **BMC Bioinformatics**, 19(19), p.517.

- Lan, C., **Peng, H.**, McGowan, G., Hutvagner, G., & Li, J. (2018). An isomiR expression panel based novel breast cancer classification

approach using improved mutual information. **BMC Medical Genomics**, 11(6), p.118.

- Ho, N., **Peng, H.**, Mayoh, C., Liu, P. Y., Atmadibrata, B., Marshall, G. M., ... & Liu, T. (2018). Delineation of the frequency and boundary of chromosomal copy number variations in paediatric neuroblastoma. **Cell Cycle**, 17(6), pp.749-758. **(co-first author)**

- Zhao, Z., **Peng, H.**, Lan, C., Zheng, Y., Fang, L., & Li, J. (2018). Imbalance learning for the prediction of N6-Methylation sites in mRNAs. **BMC Genomics**, 19(1), p.574.

- **Peng, H.**, Lan, C., Liu, Y., Liu, T., Blumenstein, M., & Li, J. (2017). Chromosome preference of disease genes and vectorization for the prediction of non-coding disease genes. **Oncotarget**, 8(45), p.78901.

- **Peng, H.**, Lan, C., Zheng, Y., Hutvagner, G., Tao, D., & Li, J. (2017). Cross disease analysis of co-functional microRNA pairs on a reconstructed network of disease-gene-microRNA tripartite. **BMC Bioinformatics**, 18(1), p.193.

- Liu, Y., **Peng, H.**, Wong, L., & Li, J. (2017). High-speed and high-ratio referential genome compression. **Bioinformatics**, 33(21), pp.3364-3372.

**Conference Papers**

- Zheng, Y., **Peng, H.**, Zhang, X., Gao, X., & Li, J. (2018). Predicting Drug Targets from Heterogeneous Spaces using Anchor Graph Hashing and Ensemble Learning. **International Joint Conference on Neural Networks**.

**Papers to be Submitted/Under Review/Accepted**

- Liu, P., Tee, A., Milazzo, G., Hannan, K., Maag, J., Mondal, S., Atmadibrata, B., Bartonicek, N., **Peng, H.**, Ho, N., Mayoh, C., Sun, Y., Welham, Z., Hulme, A., Henderson, M., Wong, M., Lan, Q., Cheung, B., Wang, J., Simon, T., Fischer, M., Zhang, X., Marshall, G., Norris, M., Haber, M., Vandesompele, J., Li, J., Mattick, J., Mestdagh, P., Hannan, R., Dinger, M., Perini, G., & Liu, T. (2018). The novel long noncoding RNA lncNB1 promotes tumorigenesis by interacting with ribosomal protein RPL35. **Nature Communications**.

- Lan, C., **Peng, H.**, Hutvagner, G., & Li, J. (2018). Construction of Competing Endogenous RNA Networks from Paired RNA-seq Data Sets by Pointwise Mutual Information Theories. **Bioinformatics**. (Major revision)

# Abstract

This thesis presents computational methods used for the recognition of disease genes and for the optimal design of disease gene CRISPR/Cas9 editing systems. The key innovation in these computational methods is the feature space and characteristics captured from the biology domain knowledge through machine learning algorithms.

The disease-gene association prediction problems are studied in Chapters 3-5. Disease gene recognition is a hot topic in various fields, especially in biology, medicine and pharmacology. Non-coding genes, a type of genes without protein products, have been proved to play important roles in disease development. Particularly, the two kinds of non-coding gene products such as microRNA (miRNA) and long non-coding RNA (lncRNA) have caught much attention as they are abundantly expressed in various tissues and frequently interact with other biomolecules, e.g. DNA, RNA and protein. The disease-ncRNA relationships remain largely unknown. Computational methods can immensely help replenish this kind of knowledge. To overcome existing computational methods' limitations such as significantly relying on network structures and similarity measurements, or lacking reliable negative samples, this thesis presents two novel methods.

One is the precomputed kernel matrix support vector machine (SVM) method to predict disease related miRNAs in Chapter 3. The precomputed kernel matrix was built by integrating several kinds of similarities computed with effective characteristics for miRNAs and diseases. The reliable negative samples were collected through analyzing the published array and sequencing

data. This binary classification method accurately predicts disease-miRNA associations, which outperforms those state-of-the-art methods. In Chapter 4, the predicted novel disease-miRNA associations were combined with known relationships of diseases, miRNAs and genes to reconstruct a disease-gene-miRNA (DGR) tripartite network. Reliable multi-disease associated co-functional miRNA pairs were extracted from this DGR for cross-disease analysis by defining the co-function score. This not only proves the proposed method's effectiveness but also contributes to the study of multi-purpose miRNA therapeutics.

Another is the bagging SVM-based positive-unlabeled learning method for disease-lncRNA prioritizing that is described in Chapter 5. It creatively characterized a disease with its related genes' chromosome distribution and pathway enrichment properties. The disease-lncRNA pairs were represented as novel feature vectors to train the bagging SVM for predicting disease-lncRNA associations. This novel representation contributes to the superior performance of the proposed method in disease-lncRNA prediction even when a given disease has no currently recognized lncRNA genes.

After confirming the relationships between genes and diseases, one of the most difficult tasks is to investigate the molecular mechanism and treatment of the diseases considering their related genes. The CRISPR/Cas9 system is a promising gene editing tool for operating the genes to achieve the goals of disease-gene function clarification and genetic disease curing. Designing an optimal CRISPR/Cas9 system can not only improve its editing efficiency but also reduce its side effect, i.e. off-target editing. Furthermore, the off-target site detection problem involves genome-wide sequence observing which makes it a more challenging job. The CRISPR/Cas9 system on-target cutting efficiency prediction and off-target site detection questions are discussed in Chapters 6 and 7 respectively.

To accurately measure the CRISPR/Cas9 system's cutting efficiency, the profiled Markov properties and some cutting position related features were merged into the feature space for representing the single-guide RNAs

(sgRNAs). These features were learned by a two-step averaging method where an XGBoost's predictions and an SVM's predictions were averaged as the final results. Later performance evaluations and comparisons demonstrate that this method can predict a sgRNA's cutting efficiency with consistently good performance no matter it is expressed from a U6 promoter in cells or from a T7 promoter in vitro.

In the off-target site detection, a sample was defined as an on-target-off-target site sequence pair to turn this problem into a classification issue. Each sample was numerically depicted with the nucleotide composition change features and the mismatch distribution properties. An ensemble classifier was constructed to distinguish real off-target sites and no-editing sites of a given sgRNA. Its excellent performance was confirmed with different test scenarios and case studies.