A deep learning based approach for genetic risk prediction



Science Changing Life

Raquel Dias, PhD. Senior Staff Scientist Scripps Research Translational Institute raqueld@scripps.edu, @RaquelDiasSRTI Ali Torkamani, PhD. atorkama@scripps.edu, @ATorkamani

Whole Genome Sequencing vs. Genotype array

Full Data (whole genome sequencing)

Sparse Data (genotype array)

? ? ? 0 0 1 1 ? ? ? ? 1 0 1 1 ?	0 0 ? ? 0 0 1 ? 0 1 ? ? 1 0 1 1 ? 0 0 ? ? ? 0 0 1 1 ?
???0011??	00???0011??e 01???1011??e 00????0011??1
? ? ? 1 0 1 1 ? ?	0 1 ? ? ? 1 0 1 1 ? ?
? ? ? 1 0 1 1 ? ?	? ? ? 1 0 1 1 ? ?
	? ? ? 0 0 1 1 ?
	0???0011??



Whole Genome Sequencing vs. Genotype array

Full Data ~80M genetic variants

Sparse Data ~4 million genetic

0	0	1	0	0	0	0	1	1	1	1	0	1	1	0	1
0	0	1	0	0	0	0	1	1	1	1	0	1	1	0	1
)	1	0	1	0	1	0	1	1	1	1	0	1	1	0	1
	_	_	_	_	_		_		_	_	_	_	_		_
0	0	1	0	0	0	0	1	1	0	0	1	0	1	1	0
0	1	0	1	0	1	0	1	1	1	1	0	1	1	0	1
~	1	0	1	1	1	0	1	1	4	1	0	1	1	0	1



Genetic imputation problem



A typical imputation approach



A typical imputation approach



Polygenic Risk Score (PRS)





Polygenic Risk Calculation



100,000+ subjects

Millions of known variants

Polygenic Risk Score

Cumulative sum

***Trait** can often be heterogeneous

e.g. coronary artery = heart attack, stroke, bypass surgery, etc.



Objectives

- 1. More accurate and faster imputation
- 2. Find important genetic variants
- 3. Better polygenic risk score calculation



Our proposed approach





Denoising autoencoder for image restoration



Research

Translational Institute

Bigdeli, Siavash Arjomand, and Matthias Zwicker. "Image restoration using autoencoding priors." *arXiv preprint arXiv:1703.09964* (2017).

Wang, Ruxin, and Dacheng Tao. "Non-local auto-encoder with collaborative stabilization for image restoration." *IEEE Transactions on Image Processing* 25.5 (2016): 2117-2129.

Genotype imputation case study example

Ground truth (whole genome sequencing)

Masked input
(genotype array)

0	0	1	0	0	0	0	1	1	1	1	0	1	1	0	1
0	0	1	0	0	0	0	1	1	1	1	0	1	1	0	1
0	1	0	1	0	1	0	1	1	1	1	0	1	1	0	1
0	0	1	0	0	0	0	1	1	0	0	1	0	1	1	0
0	1	0	1	0	1	0	1	1	1	1	0	1	1	0	1
0	1	0	1	1	1	0	1	1	1	1	0	1	1	0	1





Case study: 9p21.3 region of the genome

- Length: 59846 bp
- 846 genetic variants in reference panel (whole genome data)
 - Approx. 200 common variants
 - Approx. 600 rare variants
- Only 17-47 variants in genotype array!!!
- Strong association to coronary artery disease (CAD)
- Genotyped and sequenced in many studies



Training on the reference panel: Data augmentation strategy





Customized Sparsity Loss Function

Sparsity loss with Kullback-Leibler (KL) / cross entropy element:

$$D_{KL}(\rho || \hat{\rho}) = \rho * \log\left(\frac{\rho}{\hat{\rho}}\right) + (1 - \rho) * \log\left(\frac{1 - \rho}{1 - \hat{\rho}}\right)$$

Customized loss adjusted for hidden activation sparsity:

$$loss = MSE + \beta * \sum_{i=1}^{n} D_{KL(i)}$$

Mean Squared Error:

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)$$



Hyper parameters to be optimized

- •β
- •ρ
- Activation functions
- L1/L2 regularizers
- Learning rate
- Batch size



Parallel Grid Search Hyperparameter optimization approach





Grid Search Results: training accuracy



Grid Search results: assessing best hyperparameter values













Optimizing batch size: training accuracy



Optimizing batch size: training run time



Translational Institute

Scripps

Testing on multiple case studies

Atherosclerosis Risk in Communities (ARIC)

- More than 3000 samples
- Whole genome sequencing (846 variants, 0% mask, ground truth)
- Affymetrix 6.0 genotype array (17 variants, 98% mask, input data)

• Framingham Heart Study (FHS)

- More than 500 samples
- Whole genome sequencing (846 variants, 0% mask, ground truth)
- Illumina 500K genotype array (47 variants, 95% mask, input data)
- Illumina 5M (93 variants, 89% mask, input data)



Accuracy in additional case studies: Proposed approach versus common statistic methodology



Accuracy in additional case studies: Proposed approach versus common statistic methodology



Run time: Proposed approach versus common statistic methodology





Linkage disequilibrium structure: ARIC



Linkage disequilibrium structure: FHS



Interpretability: identifying representative genetic variants

Maximal information criteria





Conclusions

- Grid search was able to find high accuracy models (>0.90)
- Hyperparameters played an important role in training performance
- Reconstruction of genetic variants from very sparse data with high accuracy (>0.80)
- Superior computational performance, faster predictions
- Fine parameter tuning may be necessary



Future steps

- Expand to other genomic regions, fine parameter tuning
- Use imputation autoencoder results as input for polygenic risk score calculation



Feed Forward Neural Network



Future steps

Focal loss to compensate for rare variants





Lin, Tsung-Yi, et al. **"Focal loss for dense object detection."** *Proceedings of the IEEE international conference on computer vision*. 2017.

Limitations: expanding the methodology to other genomic regions



Acknowledgements



Translational Institute

Ali Torkamani PhD Shang-Fu Chen Elias Salfati PhD Doug Evans Shuchen Liu Alex Lippman Nathan W. PhD Emily Spencer PhD Eric Topol MD



Scripps Research



Johnny Israeli Carla Leibowitz Fernanda Foertter Brian Welker David Nola



National Institutes of Health

Funding

NIH/NCRR flagship CTSA Grant

Contact: raqueld@scripps.edu, @RaquelDiasSRTI atorkama@scripps.edu, @Atorkamani

Thanks for your attention!!



