Towards reliable experimental recommendations of gene-gene interactions

Tiffany Tang

Applied and Computational Mathematics and Statistics University of Notre Dame

<u>ttang4@nd.edu</u>

Our Interdisciplinary Team





Bin Yu

Ben Brown

Students/Postdocs





Ana Kenney





Tiffany Tang

Omer Ronen

Abhi Agarwal















Euan Ashley

James Priest Victoria Parikh











Chad Weldy

Weston Hughes









Rima Arnaout Atul Butte 2

Outline

1

Problem Background

2 Low-signal Iterative Random Forest (lo-siRF) For recommending genes and gene-gene interactions



Problem Background

A Public Health Crisis: Cardiovascular Disease

Cardiovascular disease (CVD) is the leading cause of death globally and in the US

17.9 million people die each year from CVDs

32% of all deaths worldwide

2x more Americans died from CVD than COVID-19 in 2020

Risk factors for CVD: poor diet, physical inactivity, tobacco use, genetic factors

Traditional tools for detecting genetic effects

Historically, research has focused on identifying genetic variants that have **marginally additive effects** on the phenotype



Common tools: genome-wide association studies (GWAS), polygenic risk scores [Khera et al. (2018), Bycroft et al. (2018), Shah et al. (2020), Pirruccello et al. (2020), Meyer et al. (2020), Harper et al. (2021), Khurshid et al. (2022), ...]

Our biological system is far more complex than this \rightarrow "missing heritability"

Beyond the marginally additive curtain: Epistasis

Epistasis is the non-additive interaction of genetic effects*



* Many technical definitions: see Bateson (1909), Fisher (1919), Wade et al. (2001), Cordell (2002), Ritchie and Van Steen (2018)

Computational Challenges: exhaustive search is computationally expensive

Ex. 100,000 genetic variants \rightarrow ~5 billion pairwise interactions!

Statistical Challenges: unclear what is an *appropriate* statistical model What about higher-order interactions? (Not just pairwise) What about nonlinear interactions? (e.g., $f(x_A)f(x_B)$, not just $x_A x_B$)

Experimental Challenges: difficult to experimentally assess small effect sizes; requires very precise, high-throughput measurements

Our aim (very broadly)

To develop an **end-to-end pipeline** for identifying **genes** and **gene-gene interactions** that affect **cardiovascular disease.**

Gene / interaction recommendation system

Wet-lab experimental validation

Computationally-tractable interaction search engine for **higher-order**, **nonlinear** interactions

Precise, high-throughput phenotyping via microfluidics-enabled gene silencing experiments

Which cardiovascular phenotype?

Major roadblock with **HCM:**

- ~50% balanced classification accuracy
- Severe under-diagnosis \rightarrow noisy labels

Attempt #1:

Hypertrophic Cardiomyopathy (HCM)

- High prevalence (~1 in 500)
- Team's clinical expertise
- Experimental capabilities for measuring cell size

Left Ventricular Hypertrophy

Attempt #1: Hypertrophic Cardi

- High prevalence
- Team's clinical e
- Experimental car measuring cell s

Normal Heart Left Ventricular Hypertrophy



Left Ventricular Hypertrophy (LVH)

Carries significant independent risk for incident heart failure, atrial arrhythmia, and sudden cardiac death

A distinguishing clinical feature of HCM



Left Ventricular Hypertrophy (LVH)



Carries significant independent risk for incident heart failure, atrial arrhythmia, and sudden cardiac death

A distinguishing clinical feature of HCM

A quantitative proxy for LVH, **left ventricular mass indexed by body surface area (LVMi)**, can be extracted from cardiac MRIs using deep learning [Bai et al. (2018)]



UK Biobank Data

n = 30K patients from white British unrelated population with cardiac MRIs

p = 15 million imputed SNPs



14

Epistasis...

Our contribution

We develop an **end-to-end pipeline** for identifying **genes** and **gene-gene interactions** that affect **left ventricular mass.**

Gene / interaction recommendation system



Wet-lab experimental validation



Computationally-tractable interaction search engine for higher-order, nonlinear interactions

Tailored for **low-signal** phenotypes

Suitable for **high-dimensional** data

Precise, high-throughput phenotyping via microfluidics-enabled knockdown experiments

Low-signal iterative random forest (lo-siRF)



Challenges

Low Signal iRF (lo-siRF)

High-dimensionality

Finding interactions

Very low signal

Domain-inspired dimension reduction via GWAS

Iterative Random Forest (iRF)

A **computationally-efficient** search engine to find stable, higher-order, nonlinear interactions [Basu, Kumbier, Brown, Yu (2018)]



Binarize LVMi phenotype

"Simplifying" the problem



A new RF feature importance score

Leverages SNP correlations to aggregate weak SNP-level importances into more stable, stronger gene-level importances

lo-siRF for gene (interaction) recommendations



Dimension reduction

Fit iRF on binarized LVMi

Rank genes / interactions

lo-siRF for gene (interaction) recommendations



```
Dimension reduction
```

Fit iRF on binarized LVMi

Rank genes / interactions

Dimension Reduction via GWAS

Run GWAS using BOLT-LMM [Loh et al. (2015)] and PLINK [Purcell et al. (2015)]

Select union of top 1000 SNPs from each GWAS run \rightarrow ~1400 SNPs



lo-siRF for gene (interaction) recommendations



Dimension reduction

Fit iRF on binarized LVMi

Rank genes / interactions

- Run **GWAS** using two methods: BOLT-LMM and PLINK
- Select union of top 1000 SNPs from each GWAS method

LVMi Binarization

Binarize LVMi phenotype into high and low groups to **"simplify"** the problem (using multiple thresholds: 15%, 20%, 25%)



Fit iRF on the binarized iLVM

For each binarization threshold (15%, 20%, 25%): Fit iRF using GWAS-filtered SNP data to predict binarized LVMi



Decision Tree



Random Forest

A collection of decision trees, where

- each tree is fitted on a different **bootstrap** version of the data
- features are subsampled at each node



Iterative Random Forest (iRF)

• Core idea: iRF induces **stability** in the RF to improve **interpretability**



Iterative Random Forest (iRF)

Core idea: iRF induces **stability** in the RF to improve **interpretability** *without sacrificing prediction accuracy

iRF \rightarrow ~55% classification accuracy for LVMi and better than other ML methods



Iterative Random Forest (iRF)

Core idea: iRF induces **stability** in the RF to improve **interpretability** *without sacrificing prediction accuracy

iRF \rightarrow ~55% classification accuracy for LVMi and better than other ML methods



Limitations of iRF

1. iRF identifies candidate interactions based upon their stability within the RF

Problem: Low-signal phenotype + highly-correlated features \rightarrow SNP-SNP interactions are highly unstable

Solution: Aggregate SNPs at the gene level \rightarrow gene-gene interactions



Limitations of iRF

1. iRF identifies candidate interactions based upon their stability within the RF

Problem: Low-signal phenotype + highly-correlated features \rightarrow SNP-SNP interactions are highly unstable

Solution: Aggregate SNPs at the gene level \rightarrow gene-gene interactions

iRF ranks candidate interactions based upon their frequency within the RF
Problem: Longer genes will naturally be more frequent in the RF
Solution: A new gene-level (or group) feature importance score

lo-siRF for gene (interaction) recommendations



Dimension reduction

Fit iRF on binarized LVMi

Rank genes / interactions

- Run GWAS using two methods: BOLT-LMM and PLINK
- Select union of top 1000 SNPs from each GWAS method
- **Binarize** iLVM phenotype into high and low groups to *"denoise"* (using multiple thresholds: 15%, 20%, 25%)
- Fit **iRF** on SNP data to extract candidate gene interactions
- Using a new stability-based importance score to aggregate SNP-level importances from iRF into a gene-level score

A paradigm shift for feature importances



A new stability-based feature importance score for RF



Extending to gene-gene interactions



lo-siRF for gene (interaction) recommendations



Dimension reduction

Fit iRF on binarized LVMi

Rank genes / interactions

- Run **GWAS** using two methods: BOLT-LMM and PLINK
- Select union of top 1000 SNPs from each GWAS method
- **Binarize** iLVM phenotype into high and low groups to *"denoise"* (using multiple thresholds: 15%, 20%, 25%)
- Fit **iRF** on SNP data to extract candidate gene interactions
- Using a new stability-based importance score to aggregate SNP-level importances from iRF into a gene-level score



"Domain expert opinion solicitation with negative controls"

We presented three lists to our cardiology experts:

- 1. Top-ranked genes/interactions
- 2. Mid-ranked genes/interactions
- 3. Random genes/interactions

Collaborators (Chad and Euan) chose list #1:)

LVMi gene & gene-gene interaction recommendations

Considering only those genes and gene-gene interactions that were stably important across all three binarization thresholds, **lo-siRF identifies**

• **genes** that are **well-known** to impact cell size

TTN IGF1R

- plausible candidate genes that are known to be associated with the heart CCDC141 RSPO3 LSP1
- stable gene-gene interactions

CCDC141–IGF1R CCDC141–TTN CCDC141–TNKS

Experimental Validation





Qianru Wang

Nate Yo

How do the size of heart cells change when we silence a gene or pair of genes?

Silenced genes/gene pairs:

- 1. CCDC141
- 2. *IGF1R*
- 3. TTN
- 4. CCDC141 and IGF1R (interaction)
- 5. CCDC141 and TTN (interaction)

Across two cell lines:

- 1. Healthy cell line
- 2. HCM cell line



39

High-throughput microfluidics + image processing











Qianru Wang

Ana Kenney On

Omer Ronen

40

Knocking down genes led to decrease in heart cell sizes



Epistatic effect sizes



Size differences are most pronounced for large heart cells



Key Takeaways



Scientific Discovery

Experimentally validated epistatic regulation of cardiac hypertrophy



lo-siRF

A gene and gene-gene interaction search engine, tailored for low-signal data



Only possible through collaboration

Iterative cycle of feedback is crucial for a successful and impactful interdisciplinary collaboration

Thank you!

Q. Wang*, TT*, ..., B. Yu, E. Ashley. "<u>Epistasis regulates genetic</u> <u>control of cardiac hypertrophy</u>." Nature Cardiovascular Research (2025).