# Non-linear Genetic Effects for Complex Traits

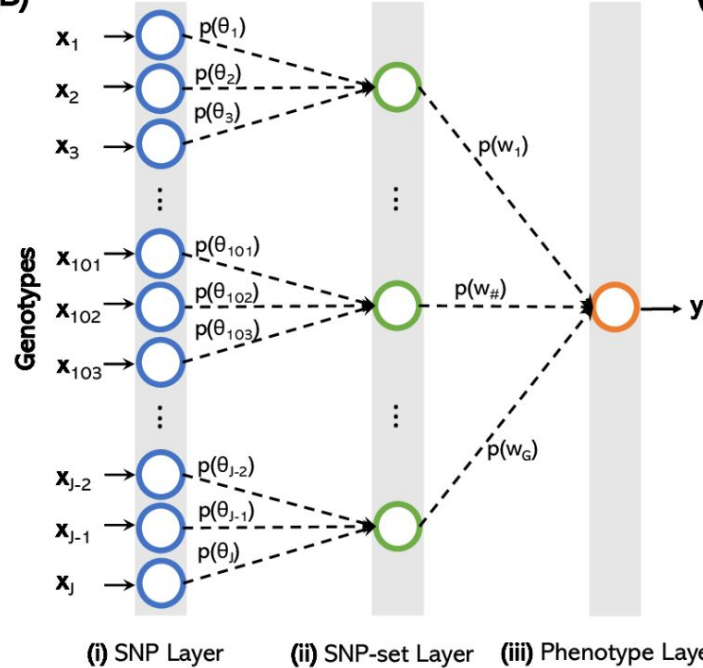Yu-Wei Chen, Zhengtong Liu, Jacob Wallin

# Introduction

- Polygenic Risk Score
  - Common method of evaluating association between genotypes and phenotypes
  - Captures linear relationship
- BANN (Biologically Annotated Neural Networks)
  - Machine learning model that can capture non-linear relationships
  - Interpretable statistical output
- VIPRS (Variational Inference of Polygenic Risk Scores)
  - Very efficient, uses VI instead of MCMC for posterior inference

# Biologically Annotated Neural Networks (BANNs)

# Kernel BANNs

**(A)**



**(i)** SNP Layer  **(ii)** SNP-set Layer (linear)  **(iii)** SNP-set Layer (nonlinear)  **(iv)** Phenotype Layer

**(B)**

**Modified Model:**

$$\boldsymbol{y} = \sum_{g=1}^{c}(\boldsymbol{X}_g\boldsymbol{\theta}_g + \mathbf{1}b_g^{(1)})w_g^{(1)} + \sum_{g=1}^{c} h(\boldsymbol{X}_g)\boldsymbol{w}_g^{(2)} + \mathbf{1}b^{(2)}$$

**SNP Layer and SNP-set Level (linear):**

$$\boldsymbol{y}_{linear} = \sum_{g=1}^{c}(\boldsymbol{X}_g\boldsymbol{\theta}_g + \mathbf{1}b_g^{(1)})w_g^{(1)} + \mathbf{1}b^{(2)}$$

**SNP-set Level (nonlinear):**

$$\boldsymbol{y}_{nonlinear} = \boldsymbol{y} - \hat{\boldsymbol{y}}_{linear}$$

$$\boldsymbol{y}_{nonlinear} = \sum_{g=1}^{c} h(\boldsymbol{X}_g)\boldsymbol{w}_g^{(2)}$$
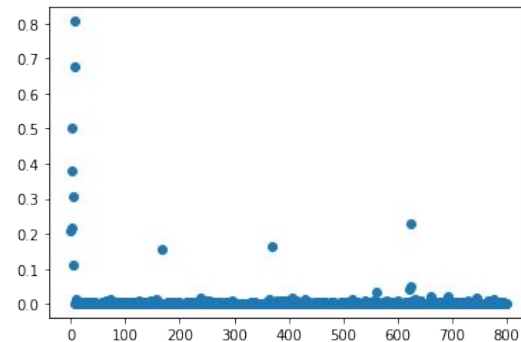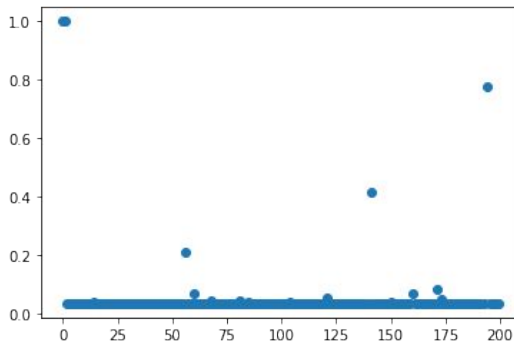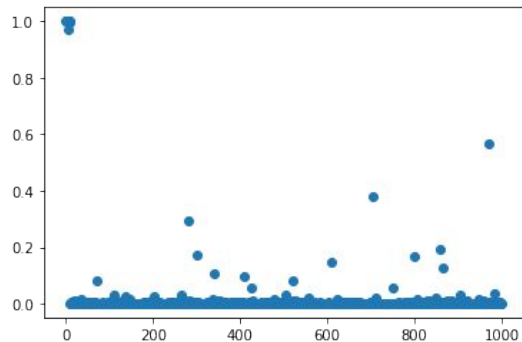
4

# Simulation Design

We assume

$$\mathbf{y} = \sum_{c \in C} \mathbf{x}_c \theta_c + \mathbf{W}\boldsymbol{\varphi} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \tau^2 \mathbf{I}),$$

where $\boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}_{c \times 1}, \mathbf{I}_{c \times c})$ is the additive effect size, $\mathbf{W} \sim \mathcal{N}(\mathbf{0}_{n \times e}, \mathbf{I}_{n \times n})$ represents the pairwise interactions between causal SNPs with corresponding effects $\boldsymbol{\varphi} \sim \mathcal{N}(\mathbf{0}_{e \times 1}, \mathbf{I}_{e \times e})$ ($e$ is the number of epistasis effects).

For simplicity, we assume $\mathbb{V}[\mathbf{y}] = 1$, so the broad-sense heritability $H^2 = \mathbb{V}[\sum \mathbf{x}_c \theta_c] + \mathbb{V}[\mathbf{W}\boldsymbol{\varphi}]$. The additive effect makes up $\rho$ while the pairwise interaction makes up $(1 - \rho)$ of the genetic variance, so the data is rescaled such that $\mathbb{V}[\sum \mathbf{x}_c \theta_c] = \rho H^2$ and $\mathbb{V}[\mathbf{W}\boldsymbol{\varphi}] = (1 - \rho)H^2$. With the constraint that $\mathbb{V}[\mathbf{y}] = 1$, the noise is rescaled so that $\mathbb{V}[\boldsymbol{\epsilon}] = 1 - H^2$.

# Preliminary Results

- Phenotypic Variance Explained (PVE)
  - Ground truth: $H^2 = 0.6$
  - SNP Layer, SNP-set Layer (linear), SNP-set Layer (nonlinear): 0.48887597222414664, 0.4525421533848896, 0.1233089950726313

- Posterior Inclusion Probability (PIP)
  - Ground truth: First 10 SNPs are causal ones

# VIPRS

PRS methods that operate on GWAS summary statistics and sparse Bayesian methods have shown competitive predictive ability.

1. **Markov Chain Monte Carlo (MCMC)** for posterior inference: computationally inefficient
2. **Variational Inference**:
   a. Propose a simple parametric distribution $q(\beta, s)$
   b. Optimize the parameters to match the true posterior as closely as possible
   c. Closeness between the true posterior and the proposed distribution is measured by the Kullback-Leibler (KL) divergence and maximizing ELBO (Evidence Lower BOund)

# VIPRS Design

- GWAS **summary statistics** for trait of interest
- Linkage-Disequil-ibrium **(LD) matrices**

**Standard Linear Model:**
$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$
**With inferred posterior distribution:**
$$p(\boldsymbol{\beta} \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) = \frac{p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\theta})p(\boldsymbol{\beta} \mid \boldsymbol{\theta})}{\int p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\theta})p(\boldsymbol{\beta} \mid \boldsymbol{\theta})d\boldsymbol{\beta}}$$

**Assign Spike and Slab Prior:**
$$\beta_j \sim \pi \mathcal{N}(\beta_j; 0, \sigma_{\boldsymbol{\beta}}^2) + (1 - \pi)\delta_0$$

**Variational Inference:**
$$q(\boldsymbol{\beta}, \mathbf{s}) = \prod_j^M q(\beta_j, s_j) = \prod_j^M \mathcal{N}(\beta_j; \mu_j, \sigma_j^2)Bern(s_j; \gamma_j)$$

Shadi Zabad, Simon Gravel, and Yue Li. (2022) "Fast and Accurate Bayesian Polygenic Risk Modeling with Variational Inference," bioRxiv 2022.05.10.491396

# Example Results

- Examine posterior estimates for the model parameters

|  | CHR | SNP | A1 | A2 | BETA | PIP | VAR_BETA |
|---|---|---|---|---|---|---|---|
| 0 | 22 | rs131538 | A | G | -1.061113e-05 | 0.014341 | 3.336974e-08 |
| 1 | 22 | rs9605903 | C | T | 1.118776e-05 | 0.014497 | 3.423175e-08 |
| 2 | 22 | rs5746647 | G | T | -2.928591e-07 | 0.012277 | 2.175343e-08 |
| 3 | 22 | rs16980739 | T | C | -2.519748e-05 | 0.019512 | 6.657493e-08 |
| 4 | 22 | rs9605923 | A | T | 1.276606e-05 | 0.015085 | 3.797483e-08 |
| ... | ... | ... | ... | ... | =posterior mean ... | ... | =variance. |
| 15930 | 22 | rs8137951 | A | G | -1.235467e-05 | 0.014874 | 3.645719e-08 |
| 15931 | 22 | rs2301584 | A | G | -8.020855e-04 | 0.237438 | 2.486760e-06 |
| 15932 | 22 | rs3810648 | G | A | 2.263451e-05 | 0.018560 | 5.996763e-08 |
| 15933 | 22 | rs2285395 | A | G | 5.651122e-06 | 0.012978 | 2.556586e-08 |
| 15934 | 22 | rs28729663 | A | G | -5.423681e-06 | 0.013021 | 2.585062e-08 |

15935 rows × 7 columns

- Examine inferred hyperparameters

|  | Parameter | Value |
|---|---|---|
| 0 | Residual_variance | 9.942623e-01 |
| 1 | Heritability | 5.618044e-03 |
| 2 | Proportion_causal | 2.282775e-02 |
| 3 | Average_effect_variance | 2.789654e-07 |
| 4 | sigma_beta | 1.222045e-05 |

# Comparisons

| | BANNs | VIPRS |
|---|---|---|
| **Similarity** | Bayesian paradigm | |
| | Same prior assumption | |
| | Variational approximation | |
| **Differences** | Two layers, interpretable statistical output | One layer, computation efficiency |
| | Relies on gene range information | GWAS summary statistics Scalable and flexible |
| | Worse performance on sparse data | Able to perform train-test split on data |

# Thanks!