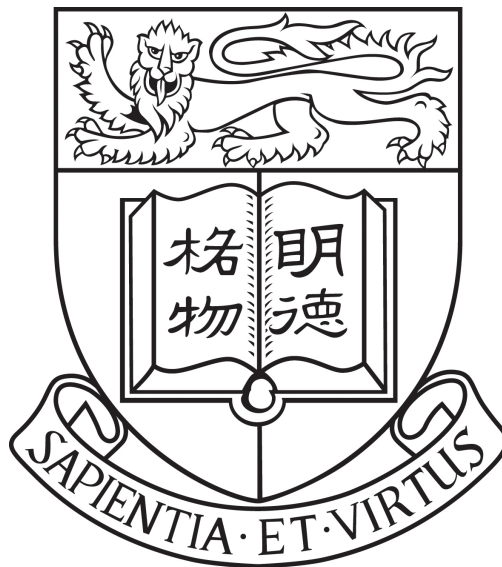


# Heritability Estimation and Risk Prediction in Schizophrenia

Choi Shing Wan

A thesis submitted in partial fulfillment of the requirements for  
the Degree of Doctor of Philosophy



Department of Psychiatry  
University of Hong Kong  
Hong Kong  
September 2, 2015



# Declaration



# Acknowledgements



# Abbreviations

**GWAS** Genome Wide Association Study. 7

**LD** Linkage Disequilibrium. 8

**SCZ** Schizophrenia. 13

**SNP** Single Nucleotide Polymorphism. 7–9





# Contents

<b>Declaration</b>	<b>i</b>
<b>Acknowledgments</b>	<b>iii</b>
<b>Abbreviations</b>	<b>v</b>
<b>Contents</b>	<b>vii</b>
<b>Introduction</b>	<b>1</b>
<b>1 Literature Review</b>	<b>5</b>
1.1 Twin Studies . . . . .	5
1.2 Searching for Genetic Variants . . . . .	5
1.2.1 Role of Common Variants . . . . .	5
1.2.2 Role of Rare Variants . . . . .	5
1.3 Narrow Sense Heritability . . . . .	5
1.4 Risk Prediction . . . . .	5
1.5 Summary . . . . .	5
<b>2 Heritability Estimation</b>	<b>7</b>
2.1 Introduction . . . . .	7
2.2 Methodology . . . . .	7
2.2.1 Heritability Estimation . . . . .	7
2.2.2 Inverse of the Linkage Disequilibrium matrix . . . . .	10
2.2.3 Quantitative Trait . . . . .	11
2.2.4 Case Control Studies . . . . .	11
2.2.5 Extreme Phenotype Selections . . . . .	11
2.3 Simulation . . . . .	11
2.3.1 Quantitative Trait . . . . .	11
2.3.2 Case Control Studies . . . . .	11
2.3.3 Extreme Phenotype Selections . . . . .	11
2.4 Result . . . . .	11
2.5 Discussion . . . . .	11
<b>3 Heritability of Schizophrenia</b>	<b>13</b>
3.1 Introduction . . . . .	13
3.2 Heritability Estimation . . . . .	13
3.2.1 Methodology . . . . .	13
3.2.2 Result . . . . .	13
3.3 Brain development and Schizophrenia . . . . .	13
3.3.1 Methodology . . . . .	13
3.3.2 Result . . . . .	13
3.4 Discussion . . . . .	13

<b>4</b>	<b>Heritability of Response to antipsychotic treatment</b>	<b>15</b>
4.1	Introduction . . . . .	15
4.2	Methodology . . . . .	15
4.3	Result . . . . .	15
4.4	Discussion . . . . .	15
<b>5</b>	<b>Risk Prediction</b>	<b>17</b>
5.1	Methodology . . . . .	17
5.1.1	Simulation . . . . .	17
5.2	Result . . . . .	17
5.3	Discussion . . . . .	17
<b>6</b>	<b>Conclusion</b>	<b>19</b>

# Introduction

---

# Some considerations

1. PRSice requires the phenotype to aid its selection (More information= stronger)
2. It seems like LDSC doesn't necessary perform badly in oligogenic situation. Rather, it is that when the trait is oligogenic, it is more likely for LDSC to behaviour in a strange way.
3. For each condition: extreme phenotype, quantitative trait, case control, we can have a separated review. Discuss on the benefits and challenges of each condition and the method we deal with them. So we can have two chapters (case control, quantitative trait) where extreme phenotype can be a big subsection within quantitative trait.
4. For each chapter, there will be this introduction (review on the method), our methodology (Calculation, implementation and also simulation), result (the simulation result). Then we can have the application (PGC, network)



# Chapter 1

## Literature Review

### 1.1 Twin Studies

Should briefly talk about how Twin modeling was used for finding the GE contribution. Should also mention the ACE model. At the end, we can talk about the heritability estimates of SCZ and AD

### 1.2 Searching for Genetic Variants

#### 1.2.1 Role of Common Variants

##### Genome Wide Association Study

Should talk about what is GWAS and how it is used. Should also talk about the current GWAS studies in SCZ and AD

#### 1.2.2 Role of Rare Variants

##### Exome Sequencing

Similar to the GWAS. Talk about the Pros and Cons. Need to briefly mention the Denovo paper and Shaun's paper.

##### Whole Genome Sequencing

Very very brief description of WGS and the current status.

### 1.3 Narrow Sense Heritability

### 1.4 Risk Prediction

### 1.5 Summary





## Chapter 2

# Heritability Estimation

This chapter should be used in similar way as the general method section in Clara's thesis. Considering that the subsequent chapters all rely on this implementation.

### 2.1 Introduction

### 2.2 Methodology

The work in this chapter were done in collaboration with my colleagues who have kindly provide their support and knowledges to make this piece of work possible.

#### 2.2.1 Heritability Estimation

The narrow-sense heritability is defined as

$$h^2 = \frac{\text{Var}(X)}{\text{Var}(Y)}$$

where  $\text{Var}(X)$  is the variance of the genotype and  $\text{Var}(Y)$  is the variance of the phenotype. In a Genome Wide Association Study (GWAS), regression were performed between the Single Nucleotide Polymorphisms (SNPs) and the phenotypes, giving

$$Y = \beta X + \epsilon \quad (2.1)$$

where  $Y$  and  $X$  are the standardized phenotype and genotype respectively.  $\epsilon$  is then the error term, accounting for the non-genetic elements contributing to the phenotype (e.g. Environment factors). Based on equation 2.1, one can then have

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(\beta X) + \text{Var}(\epsilon) \\ \text{Var}(Y) &= \beta^2 \text{Var}(X) \\ \beta^2 \frac{\text{Var}(X)}{\text{Var}(Y)} &= 1 \end{aligned} \quad (2.2)$$

$\beta^2$  is then considered as the portion of phenotype variance explained by the variance of genotype, which can also be considered as the narrow-sense heritability of the phenotype.

A challenge in calculating the heritability from GWAS data is that usually only the test-statistic or p-value were provided and one will not be able to directly calculate the heritability based on equation 2.2. In order to estimation the heritability of a trait from the GWAS test-statistic, we first observed that when both  $X$

and  $Y$  are standardized,  $\beta^2$  will be equal to the coefficient of determination ( $r^2$ ). Then, based on properties of the Pearson product-moment correlation coefficient:

$$r = \frac{t}{\sqrt{n-2+t^2}} \quad (2.3)$$

where  $t$  follows the student-t distribution and  $n$  is the number of samples. One can then obtain the  $r^2$  by taking the square of 2.3

$$r^2 = \frac{t^2}{n-2+t^2} \quad (2.4)$$

It is observed that  $t^2$  will follow the F-distribution and when  $n$  is big,  $t^2$  will converge into  $\chi^2$  distribution.

When the effect size is small and  $n$  is big,  $r^2$  will be approximately  $\chi^2$  distributed with mean  $\sim 1$ . We can then approximate equation 2.4 as

$$r^2 = \frac{\chi^2}{n} \quad (2.5)$$

and define the *observed* effect size of each SNP to be

$$f = \frac{\chi^2 - 1}{n} \quad (2.6)$$

When there are Linkage Disequilibrium (LD) between each individual SNPs, the situation will become more complicated as each SNPs' observed effect will contains effect coming from other SNPs in LD with it.

$$f_{observed} = f_{true} + f_{LD} \quad (2.7)$$

To account for the LD structure, we first assume our phenotype  $\mathbf{Y}$  and genotype  $\mathbf{X} = (X_1, X_2, \dots, X_m)^t$  are standardized and that

$$\begin{aligned} \mathbf{Y} &\sim f(0, 1) \\ \mathbf{X} &\sim f(0, \mathbf{R}) \end{aligned}$$

Where  $\mathbf{R}$  is the LD matrix between SNPs.

We can then express equation 2.1 in matrix form:

$$\mathbf{Y} = \boldsymbol{\beta}^t \mathbf{X} + \epsilon \quad (2.8)$$

Definition of heritability will then become

$$\begin{aligned} \text{Heritability} &= \frac{\text{Var}(\boldsymbol{\beta}^t \mathbf{X})}{\text{Var}(\mathbf{Y})} \\ &= \text{Var}(\boldsymbol{\beta}^t \mathbf{X}) \end{aligned} \quad (2.9)$$

If we then assume now that  $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_m)^t$  has distribution

$$\begin{aligned} \boldsymbol{\beta} &\sim f(0, \mathbf{H}) \\ \mathbf{H} &= \text{diag}(\mathbf{h}) \\ \mathbf{h} &= (h_1^2, h_2^2, \dots, h_m^2)^t \end{aligned}$$

where  $\mathbf{H}$  is the variance of the true effect. It is shown that heritability can be expressed as

$$\begin{aligned}
\text{Var}(\boldsymbol{\beta}^t \mathbf{X}) &= \text{E}_X \text{Var}_{\beta|X}(\mathbf{X}^t \boldsymbol{\beta}) + \text{Var}_X \text{E}_{(\beta|X)}(\boldsymbol{\beta}^2 \mathbf{X}) \\
&= \text{E}_X(\mathbf{X}^t \boldsymbol{\beta} \boldsymbol{\beta}^T \mathbf{X}) \\
&= \text{E}_X(\mathbf{X}^t \mathbf{H} \mathbf{X}) \\
&= \text{E}(\mathbf{X})^t \mathbf{H} \text{E}(\mathbf{X}) + \text{Tr}(\text{Var}(\mathbf{X} \mathbf{H})) \\
&= \text{Tr}(\text{Var}(\mathbf{X} \mathbf{H})) \\
&= \sum_i h_i^2
\end{aligned} \tag{2.10}$$

Now if we consider the covariance between SNP  $i$  ( $X_i$ ) and  $Y$ , we have

$$\begin{aligned}
\text{Cov}(\mathbf{X}_i, \mathbf{Y}) &= \text{Cov}(\mathbf{X}_i, \boldsymbol{\beta}^t \mathbf{X} + \epsilon) \\
&= \text{Cov}(\mathbf{X}_i, \boldsymbol{\beta}^t \mathbf{X}) \\
&= \sum_j \text{Cov}(\mathbf{X}_i, \mathbf{X}_j) \boldsymbol{\beta}_j \\
&= \mathbf{R}_i \boldsymbol{\beta}_j
\end{aligned} \tag{2.11}$$

As both  $X$  and  $Y$  are standardized, the covariance will equal to the correlation and we can define the correlation between SNP  $i$  and  $Y$  as

$$\rho_i = \mathbf{R}_i \boldsymbol{\beta}_j \tag{2.12}$$

In reality, the *observed* correlation usually contains error. Therefore we define the *observed* correlation to be

$$\hat{\rho}_i = \rho_i + \frac{\epsilon_i}{\sqrt{n}} \tag{2.13}$$

for some error  $\epsilon_i$ . The distribution of the correlation coefficient about the true correlation  $\rho$  is approximately

$$\hat{\rho}_i \sim f(\rho_i, \frac{(1 - \rho^2)^2}{n})$$

By making the assumption that  $\rho_i$  is close to 0 for all  $i$ , we have

$$\begin{aligned}
\text{E}(\epsilon_i | \rho_i) &\sim 0 \\
\text{Var}(\epsilon_i | \rho_i) &\sim 1
\end{aligned}$$

We then define our  $z$ -statistic and  $\chi^2$ -statistic as

$$\begin{aligned}
z_i &= \hat{\rho}_i \sqrt{n} \\
\chi^2 &= z_i^2 \\
&= \hat{\rho}_i^2 n
\end{aligned}$$

From equation 2.13 and equation 2.12,  $\chi^2$  can then be expressed as

$$\begin{aligned}
\chi^2 &= \hat{\rho}^2 n \\
&= n(\mathbf{R}_i \boldsymbol{\beta}_j + \frac{\epsilon_i}{\sqrt{n}})^2
\end{aligned}$$

The expectation of  $\chi^2$  is then

$$\begin{aligned} E(\chi^2) &= n(\mathbf{R}_i \boldsymbol{\beta} \boldsymbol{\beta}^t \mathbf{R}_i + 2\mathbf{R}_i \boldsymbol{\beta} \frac{\epsilon_i}{\sqrt{n}} + \frac{\epsilon_i^2}{n}) \\ &= n\mathbf{R}_i \mathbf{H} \mathbf{R}_i + 1 \end{aligned}$$

To derive least square estimates of  $h_i^2$ , we need to find  $\hat{h}_i^2$  which minimizes

$$\begin{aligned} \sum_i (\chi_i^2 - E(\chi_i^2))^2 &= \sum_i (\chi_i^2 - (n\mathbf{R}_i \mathbf{H} \mathbf{R}_i + 1))^2 \\ &= \sum_i (\chi_i^2 - 1 - n\mathbf{R}_i \mathbf{H} \mathbf{R}_i)^2 \end{aligned}$$

If we define

$$f_i = \frac{\chi_i^2 - 1}{n} \quad (2.14)$$

we got

$$\begin{aligned} \sum_i (\chi_i^2 - E(\chi_i^2))^2 &= \sum_i (f_i - \mathbf{R}_i \mathbf{H} \mathbf{R}_i)^2 \\ &= \mathbf{f} \mathbf{f}^t - 2\mathbf{f}^t \mathbf{R}_{sq} \hat{\mathbf{h}} + \hat{\mathbf{h}}^t \mathbf{R}_{sq}^t \mathbf{R}_{sq} \hat{\mathbf{h}} \end{aligned} \quad (2.15)$$

where  $\mathbf{R}_{sq} = \mathbf{R} \circ \mathbf{R}$ . By differentiating equation 2.15 w.r.t  $\hat{\mathbf{h}}$  and set to 0, we get

$$\begin{aligned} 2\mathbf{R}_{sq}^t \mathbf{R}_{sq} \hat{\mathbf{h}} - 2\mathbf{R}_{sq} \mathbf{f} &= 0 \\ \mathbf{R}_{sq} \hat{\mathbf{h}} &= \mathbf{f} \end{aligned} \quad (2.16)$$

And the heritability is then defined as

$$\text{Heritability} = \mathbf{1}^t \mathbf{R}_{sq}^{-1} \mathbf{f} \quad (2.17)$$

### 2.2.2 Inverse of the Linkage Disequilibrium matrix

In order to obtain the heritability estimation, we will require to solve equation 2.17. If  $\mathbf{R}_{sq}$  is of full rank and positive semi-definite, it will be straight-forward to solve the matrix equation. However, the LD matrix almost always suffer from the problem of being multicollinear. This

**2.2.3 Quantitative Trait**

**2.2.4 Case Control Studies**

**2.2.5 Extreme Phenotype Selections**

**2.3 Simulation**

**2.3.1 Quantitative Trait**

**2.3.2 Case Control Studies**

**2.3.3 Extreme Phenotype Selections**

**2.4 Result**

**2.5 Discussion**



## Chapter 3

# Heritability of Schizophrenia

### 3.1 Introduction

### 3.2 Heritability Estimation

This will be a very simple section, focused on how to perform the heritability estimation on Schizophrenia (SCZ). Should also tokenize the heritability into subcategories (e.g. immune, neuron, etc)

#### 3.2.1 Methodology

#### 3.2.2 Result

### 3.3 Brain development and Schizophrenia

Here we will perform the WGCNA and brain development network. Seeing how the whether if any brain development network were enriched with SNPs that explain the variance of phenotype

#### 3.3.1 Methodology

#### 3.3.2 Result

### 3.4 Discussion





## Chapter 4

# Heritability of Response to antipsychotic treatment

### 4.1 Introduction

Here we try to use Beatrice's data and estimate the heritability explained in drug response. Should also repeat the region-wise heritability

### 4.2 Methodology

### 4.3 Result

### 4.4 Discussion



## Chapter 5

# Risk Prediction

### 5.1 Methodology

#### 5.1.1 Simulation

### 5.2 Result

### 5.3 Discussion



## Chapter 6

## Conclusion



# Appendix