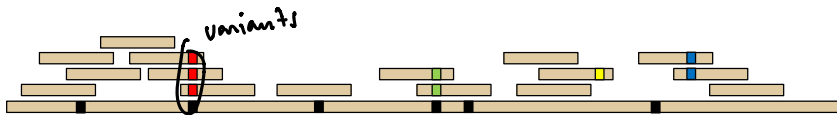# Lecture 5: GWAS (cont.)



ECE 365 - Data Science and Genomics

# Announcements:

- ☐ Lab 2 (Sequence alignment) due on Thursday
- ☐ Lab 3 (GWAS) released tomorrow
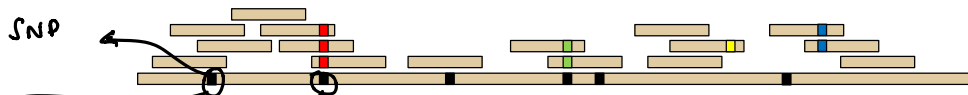
# Genotype data

- Focus on a set of common variants in the genome

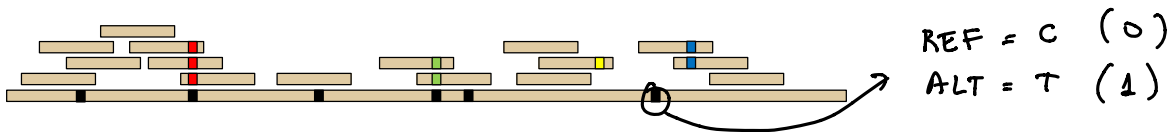# Genotype data

☐ Focus on a set of common variants in the genome

SNP

☐ Genotype data (VCF file) looks like this:

person 1   paternal   maternal     person n

$$
\begin{array}{c}
\text{SNP 1} \\
\text{SNP 2} \\
\\
\\
\\
\\
\text{SNP m}
\end{array}
\begin{bmatrix}
0|0 & 1|0 & 1|0 & 0|0 & 0|0 & 0|1 \\
0|1 & 0|0 & 1|1 & 1|0 & 0|0 & 1|1 \\
0|1 & 0|1 & 0|1 & 0|0 & 0|0 & 0|0 \\
0|0 & 0|0 & 0|0 & 1|1 & 0|1 & 1|0 \\
1|0 & 0|1 & 1|1 & 0|0 & 0|0 & 0|0 \\
0|0 & 1|1 & 0|0 & 0|0 & 0|1 & 0|0 \\
0|1 & 0|1 & 0|0 & 1|0 & 0|0 & 1|1
\end{bmatrix}
$$

# Genotype data

- Focus on a set of common variants in the genome



REF = C (0)
ALT = T (1)

- Genotype data (VCF file) looks like this:

# Genome-Wide Association Studies (GWAS)

|          | SNP 1 |   |   |   |   | SNP m |
|----------|-------|---|---|---|---|-------|
| person 1 | 0 | 1 | 1 | 0 | 0 | 1 |
|          | 1 | 0 | 2 | 1 | 0 | 2 |
|          | 1 | 1 | 1 | 0 | 0 | 0 |
|          | 0 | 0 | 0 | 2 | 1 | 1 |
|          | 1 | 1 | 1 | 0 | 0 | 0 |
|          | 0 | 2 | 0 | 0 | 1 | 0 |
| person n | 1 | 1 | 0 | 1 | 0 | 2 |

phenotype

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

genome-wide

e.g. has Diabetes

(doesn't need to be binary. e.g., height)

# Genome-Wide Association Studies (GWAS)

associated

|  | SNP 1 | | SNP 200 | | SNP m | | | phenotype |
|---|---|---|---|---|---|---|---|---|

person 1
$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 2 & 1 & 0 & 2 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 2 \end{bmatrix}$$
person n

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

new person     1   0   1   1   2   0        ? (predict)

- ☐ Which SNPs are associated with a given phenotype?
- ☐ Given a new individual's genotype, can you predict their phenotype?

# Revisiting Logistic Regression

□ Predict binary variable from real-valued features

$$p(1 \mid \underline{x}) = \frac{e^{\beta_0 + \underline{\beta}^T \underline{x}}}{1 + e^{\beta_0 + \underline{\beta}^T \underline{x}}} = \frac{1}{1 + e^{-(\beta_0 + \underline{\beta}^T \underline{x})}}$$

$$\overset{''}{p}$$

$$1 + e^{-(\beta_0 + \underline{\beta}^T \underline{x})} = \frac{1}{p} \implies e^{-(\beta_0 + \underline{\beta}^T \underline{x})} = \frac{1}{p} - 1 = \frac{1-p}{p}$$

$$\longrightarrow \ln\left(\frac{p}{1-p}\right) = \underbrace{\beta_0 + \underline{\beta}^T \underline{x}}_{\text{linear model}}$$

$$\frac{p}{1-p} : \text{odds ratio} \in (0, \infty)$$

$$\ln \frac{p}{1-p} : \text{log odds ratio} \in (-\infty, \infty)$$

# Revisiting Logistic Regression

□ Can we use Logistic Regression for GWAS?

million

|  | SNP 1 |  |  |  |  | SNP m | phenotype |
|---|---|---|---|---|---|---|---|
| person 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
|  | 1 | 0 | 2 | 1 | 0 | 2 | 0 |
| $\underline{x}$ | 1 | 1 | 1 | 0 | 0 | 0 | 1 |
|  | 0 | 0 | 0 | 2 | 1 | 1 | 1 |
|  | 1 | 1 | 1 | 0 | 0 | 0 | 1 |
|  | 0 | 2 | 0 | 0 | 1 | 0 | 0 |
| person n | 1 | 1 | 0 | 1 | 0 | 2 | 1 |

$$p(1 \mid \underline{x}) = \frac{1}{1 + e^{-(\beta_0 + \underline{\beta}^T \underline{x})}}$$

risk of having a disease

# Revisiting Logistic Regression

☐ Can we use Logistic Regression for GWAS?

millions

|  | SNP 1 | | | | | SNP m | phenotype |
|---|---|---|---|---|---|---|---|

$$
\text{person 1} \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 2 & 1 & 0 & 2 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 2 \end{bmatrix} \text{person n} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}
$$

1000s

☐ Problem: Number of SNPs can be ~$10^6$

# Revisiting Logistic Regression

☐ Can we use Logistic Regression for GWAS?

|            | SNP 1 |   |   |   |   | SNP m | phenotype |
|------------|-------|---|---|---|---|-------|-----------|

person 1
$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 2 & 1 & 0 & 2 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 2 \end{bmatrix} \quad \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$
person n

☐ Problem: Number of SNPs can be $\sim 10^6$

# GWAS via *univariate* logistic regression

- Run a separate logistic regression for each SNP

SNP 1     $x_i$     SNP m    phenotype

person 1
$$
\begin{bmatrix}
0 & 1 & 1 & 0 & 0 & 1 \\
1 & 0 & 2 & 1 & 0 & 2 \\
1 & 1 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 2 & 1 & 1 \\
1 & 1 & 1 & 0 & 0 & 0 \\
0 & 2 & 0 & 0 & 1 & 0 \\
1 & 1 & 0 & 1 & 0 & 2
\end{bmatrix}
\begin{bmatrix}
0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1
\end{bmatrix}
$$
person n

$$ p(1 \mid x_i) = \frac{1}{1 + e^{-(\beta_0 + \beta_i x_i)}} $$

Captures association
between $i$ th SNP
and phenotype

- Use this to identify small subset of SNPs associated with phenotype
- Let's look at some examples on a Jupyter notebook

# GWAS via *univariate* logistic regression

- We will get a $\beta$ for each SNP

|  | SNP 1 |  |  |  |  | SNP m |
|---|---|---|---|---|---|---|
| person 1 | 0 | 1 | 1 | 0 | 0 | 1 |
|  | 1 | 0 | 2 | 1 | 0 | 2 |
|  | 1 | 1 | 1 | 0 | 0 | 0 |
|  | 0 | 0 | 0 | 2 | 1 | 1 |
|  | 1 | 1 | 1 | 0 | 0 | 0 |
|  | 0 | 2 | 0 | 0 | 1 | 0 |
| person n | 1 | 1 | 0 | 1 | 0 | 2 |

phenotype

$$\begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

$\beta_1 \quad \beta_2 \quad \cdots \quad \beta_m$

model for ith SNP:

$$\ln\left(\frac{p(1|x_i)}{1 - p(1|x_i)}\right) = \beta_0 + \beta_i x_i$$

LOR for reference genome (no ALT SNPs)

# GWAS via *univariate* logistic regression

☐ Idea: combine all beta coefficients into a single model:

$$\ln\left(\frac{p(1|x)}{1 - p(1|x)}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m$$

# GWAS via *univariate* logistic regression

☐ Idea: combine all beta coefficients into a single model:

$$\ln\left(\frac{p(1|x)}{1-p(1|x)}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m$$

☐ Problems with this approach:

# GWAS via *univariate* logistic regression

□ Idea: combine all beta coefficients into a single model:

$$\ln\left(\frac{p(1|x)}{1 - p(1|x)}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m$$

□ Problems with this approach:

   □ Since $m$ is very large, some $\beta_i$s will be large **by chance**

# GWAS via *univariate* logistic regression

☐ Idea: combine all beta coefficients into a single model:

$$\ln\left(\frac{p(1|x)}{1 - p(1|x)}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m$$

☐ Problems with this approach:

 ☐ Since $m$ is very large, some $\beta_i$s will be large **by chance**

 ☐ Some $x_i$s are correlated

$$\left(e.g., \; anyone \; with \; x_3 = 1 \; has \; x_4 = 1\right)$$

# Identifying statistically significant SNPs

☐ To measure the significance of the association, we use the *p*-value
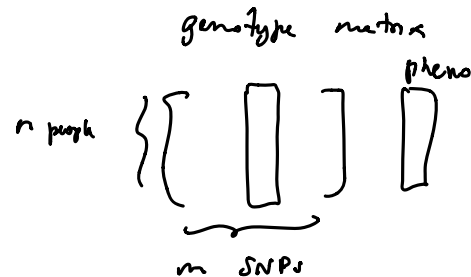
# Identifying statistically significant SNPs

□ To measure the significance of the association, we use the *p*-value

   ◘ Probability that the coefficient $\beta_i$ would be obtained by chance if there was **no** association

   $$P\left(|\hat{\beta_i}| > |\beta_i| \text{ found} \mid \text{no association}\right)$$

□ We can use the statsmodels Python package to perform the logistic regressions and compute the *p*-values

□ Let's return to the jupyter notebook

# Manhattan plots

- Allow us to see the significance of all SNPs in the genome
- We plot $-\log_{10}(p-\text{value})$

*genotype matrix*

*pheno*

*n people* $\{$ $[$ $|$ $]$ $]$ $[$

*m SNPs*

Manhattan plot of SNP association with Irritable Bowel Disease



$10^{-12}$

$-\log_{10}$ P

$-\log_{10}$ p-value

threshold