

Quantitative Genomics and Genetics - Spring 2016

BTRY 4830/6830; PBSB 5201.01

Key to Homework 6 (Version 1 posted May 24)

Assigned April 8; Due 11:59PM April 15

Problem 1 (Easy)

Consider the (slightly idealized) genetics behind one of Mendel's famous experiments with pea plants (look Mendel up on wikipedia if you are new to genetics), where for two alleles A_1 and A_2 the phenotype of a pea is guaranteed to be 'yellow' if the genotype is either A_1A_1 or A_1A_2 and guaranteed to be 'green' if the genotype is A_2A_2 . If you were to code a random variable for this system $Y(\text{yellow}) = 1$ and $Y(\text{green}) = 0$ and used a (genetic) linear regression to model this case, what would be the true values of the parameters β_a , β_d , and σ_ϵ^2 ?

The linear regression model needs to precisely return a specific phenotype for each genotype, requiring $\sigma_\epsilon^2 = 0$ (or $\sigma_\epsilon^2 \rightarrow 0$). The β parameters need to be set such that $[Y(\text{yellow}) = 1|g = A_1A_1 \text{ or } A_1A_2]$ and $[Y(\text{yellow}) = 0|g = A_2A_2]$ such that:

$$Y = \beta_\mu - \beta_a - \beta_d = 1 \quad (1)$$

$$Y = \beta_\mu + \beta_d = 1 \quad (2)$$

$$Y = \beta_\mu + \beta_a - \beta_d = 0 \quad (3)$$

where solving these equations with the following steps:

$$\beta_\mu = 1 - \beta_d \quad (4)$$

$$1 - \beta_d - \beta_a - \beta_d = 1; \beta_a = -2\beta_d \quad (5)$$

$$1 - \beta_d - 2\beta_d - \beta_d = 0; \beta_d = 0.25 \quad (6)$$

$$\beta_\mu = 1 - 0.25 = 0.75 \quad (7)$$

$$0.75 - \beta_a - 0.25 = 1; \beta_a = -0.5 \quad (8)$$

leads to $\beta_\mu = 0.75$, $\beta_d = 0.25$ and $\beta_a = -0.5$.

Problem 2 (Medium)

For the ‘case / control’ phenotype data in the file ‘QG16_HW6_phenotypes.txt’ and the genotype data in the file ‘QG16_HW6_genotypes.txt’ (where each column is a genotype coded as -1 or 1 for homozygotes and 0 for heterozygotes, each row contains the genotypes for an individual, and assume these genotypes are in the same order along a chromosome as they are presented in the columns) answer the following questions concerning a GWAS analysis of these data. For parts of this question that require R code answers, provide your code in a text file that is easy for us to run to produce the requested output (NOTE THAT FOR FULL CREDIT = easy to run code in a txt file, plus figures, and written answers as appropriate).

See [html key file for coding and output answers](#).

- a. When using a logistic regression model applied independently to each of these N markers, state the null and alternative hypothesis that will be assumed in each case.

$$H_0 : \beta_a = 0 \cap \beta_d = 0 \quad (9)$$

$$H_A : \beta_a \neq 0 \cup \beta_d \neq 0 \quad (10)$$

- b. For each of the N markers, calculate the MLE of the β parameters under the null hypothesis using the IRLS algorithm. For your answer, provide R code that outputs the IRLS calculated estimates for each parameter of each marker AND for one marker of your choice, write down the estimates you obtained (make sure you note which marker you are reporting!).
- c. For each of the N markers, calculate the MLE of the β parameters under the alternative hypothesis using the IRLS algorithm. For your answer, provide R code that outputs the IRLS calculated estimates for each parameter AND for the same marker you selected in part ‘b’, write down the estimates you obtained.
- d. Why is it a reasonable assumption that the parameter estimates you obtained in part ‘c’ are not exactly the true $MLE(\hat{\beta})$? Also, using no more than one sentence, explain why this is not a problem in practice?

The true MLE are the parameter values that produce a maximum value of the likelihood function for the logistic regression model and since the IRLS algorithm is designed to stop when reaching parameter values that produce a likelihood value close to this maximum (i.e., a trade-off between accuracy and time) the result will (almost certainly) not be the true MLE. However, since the values of the parameters returned by the IRLS algorithm are extremely close to the parameter values of the true MLE, the impact on our estimation and p-value calculations are so small that they do not change the results compared to the case if we had used the true MLE.

- e. For each of the N markers, calculate the LRT. For your answer, provide R code that outputs the LRT for each marker AND for the same marker you selected in part ‘b’ and ‘c’, write down the value of the LRT you obtained.

- f. For each of the N markers, calculate the p-value using your LRT calculated in part ‘e’, when assuming this statistic is distributed as a $\chi_{df=2}^2$ under the null hypothesis. For your answer, provide R code that outputs the p-value for each marker AND write down the p-value for the MOST SIGNIFICANT MARKER (make sure you note which marker that is).
- g. Plot a Manhattan plot for these N markers (provide your R code as well as the plot).
- h. Plot a QQ plot for these N markers (provide your R code as well as the plot).
- i. What is an appropriate study-wide Type I error used to assess whether you have evidence of the location of causal polymorphism? Present your cutoff and justify why you selected this value.

One acceptable answer, others are possible: an appropriate study-wide Type I error α_B can be calculated using a Bonferroni correction $\alpha_B = \frac{\alpha}{N}$ for a target Type I error α , which we could set to $\alpha = 0.05$.

- j. In no more than three sentences, describe why your results in parts ‘g-i’ indicate that you have identified the location of a causal polymorphism.

The QQ plot, where the p-values are along the 45 degree line and have an tail for the extreme p-values, indicates that the statistical model fit to each marker is appropriate for the purposes of a GWAS, where the goal is to assign a p-value to the null hypothesis that each marker has no impact on the phenotype. Given that we have fit an appropriate statistical model, we can interpret the significant p-values as indicating a position in the genome that includes an unmeasured causal genotype, where the significant markers we are testing are in linkage disequilibrium with this causal polymorphism.

Problem 3 (Difficult)

Prove that for the (genetic) logistic regression model, if the null hypothesis $H_0 : \beta_a \cap \beta_d$ is true, this means that $Cov(Y, X_a) = 0 \cap Cov(Y, X_d) = 0$.

There are many possible approaches, one acceptable answer:

If H_0 is true, then $\beta_a = 0 \cap \beta_d = 0$ such that:

$$Y = \frac{e^{\beta_\mu + X_a \beta_a + X_d \beta_d}}{1 + e^{\beta_\mu + X_a \beta_a + X_d \beta_d}} + \epsilon = \frac{e^{\beta_\mu}}{1 + e^{\beta_\mu}} + \epsilon \quad (11)$$

where $\frac{e^{\beta_\mu}}{1 + e^{\beta_\mu}}$ is a constant (say c) and since

$$Pr(\epsilon) \sim bern(p); p = \frac{e^{\beta_\mu}}{1 + e^{\beta_\mu}} \quad (12)$$

we have $Y = c - \epsilon$ and $Pr(Y)$ is the same regardless of the genotype, and therefore the same regardless of the value of X_a or X_d . This means that

$$Pr(Y \cap X_a) = Pr(Y)(X_a) \rightarrow Cov(Y, X_a) = 0 \quad (13)$$

and similarly for X_d .