

CS 194 - 17.2 HW Q

Collaborators: Zaid Ahmad

8 Oct 4th
2021

a) Prove (1) $\frac{e^{\beta_0 + \beta_1 \cdot x}}{1 + e^{\beta_0 + \beta_1 \cdot x}} = (2) \frac{1}{1 + e^{-\beta_0 - \beta_1 \cdot x}}$ Stefan Bielmeier

→ take (1): $\frac{e^{\beta_0 + \beta_1 \cdot x} \cdot e^{-\beta_0 - \beta_1 \cdot x}}{(1 + e^{\beta_0 + \beta_1 \cdot x})(e^{-\beta_0 - \beta_1 \cdot x})} = \frac{e^{\beta_0 + \beta_1 \cdot x - \beta_0 - \beta_1 \cdot x}}{e^{-\beta_0 - \beta_1 \cdot x} + e^{\beta_0 + \beta_1 \cdot x - \beta_0 - \beta_1 \cdot x}}$

$= \frac{e^0}{e^{-\beta_0 - \beta_1 \cdot x} + e^0} = \frac{1}{1 + e^{-\beta_0 - \beta_1 \cdot x}} = (2) \quad \checkmark \quad | = 1 = 1$

b) $g(z) = \frac{1}{1 + e^{-z}} \Rightarrow g'(z) = \frac{0 \cdot (1 + e^{-z})' - (1 + e^{-z})' \cdot 1}{(1 + e^{-z})^2}$

$= \frac{- (0 - e^{-z})}{(1 + e^{-z})^2} = \frac{e^{-z}}{(1 + e^{-z})^2} = \underbrace{\frac{1}{(1 + e^{-z})}}_{g(z)} \cdot \underbrace{\frac{e^{-z}}{(1 + e^{-z})}}_{\text{get it to } 1 - g(z)}$

$= \frac{1}{(1 + e^{-z})} \cdot \left(\frac{1 + e^{-z} - 1}{1 + e^{-z}} \right) = \frac{1}{1 + e^{-z}} \cdot \left(\frac{1 + e^{-z}}{1 + e^{-z}} - \frac{1}{1 + e^{-z}} \right) = \frac{1}{1 + e^{-z}} \cdot \left(1 - \frac{1}{1 + e^{-z}} \right)$

$= g(z) \cdot (1 - g(z)) \quad \checkmark \quad = g'(z)$

number one,
sorry for
any
confusion.

$$1c) f(x; \beta_0, \beta_1) = \frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x)}} = \frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x)}}$$

\Rightarrow with $z = \beta_0 + \beta_1 \cdot x$

$$\Rightarrow f(x; \beta_0, \beta_1) = g(\beta_0 + \beta_1 \cdot x) \text{ with } g(z) = \frac{1}{1+e^{-z}}$$

1d) We will go from general to specific to answer this question.

Probability of observing a $Y=1$ given a specific datapoint $x \& \beta_0, \beta_1$ (model)

$$\Rightarrow P(Y=1 | x; \beta_0, \beta_1) = \frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x)}} \text{, using rule of total probability}$$

& knowing only 2 outcomes,

$$P(Y=0 | x; \beta_0, \beta_1) = 1 - \frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x)}} \Leftarrow 0 \& 1$$

$$\Rightarrow \text{for any given data point } (y_i, x_i): P(Y=y_i | x_i; \beta_0, \beta_1) = \left(\frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x_i)}} \right)^{y_i} \cdot \left(\frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x_i)}} \right)^{1-y_i}$$

(Bernoulli probability mass function)

\Rightarrow Now, for the equivalent of a sequence in Markov model, a dataset with n observations, the likelihood of a logistic regression is:

$$\Rightarrow L(\beta_0, \beta_1) = \prod_{i=1}^n P(Y=y_i | X=x_i; \beta_0, \beta_1), \text{ using } \frac{1}{1+e^{-(\beta_0 + \beta_1 \cdot x_i)}} = g(\beta_0 + \beta_1 \cdot x_i)$$

$$= \prod_{i=1}^n \left(g(\beta_0 + \beta_1 \cdot x_i) \right)^{y_i} \cdot \left(1 - g(\beta_0 + \beta_1 \cdot x_i) \right)^{1-y_i}$$

$$\Rightarrow \log L(\beta_0, \beta_1) = \sum_{i=1}^n \log \left(g(\beta_0 + \beta_1 \cdot x_i) \right)^{y_i} + \log \left(1 - g(\beta_0 + \beta_1 \cdot x_i) \right)^{1-y_i}$$

$$= \sum_{i=1}^n y_i \cdot \log[g(\beta_0 + \beta_1 \cdot x_i)] + (1-y_i) \cdot \log[1-g(\beta_0 + \beta_1 \cdot x_i)]$$

$$= \sum_{i=1}^n y_i \cdot \underline{\log[g(\beta_0 + \beta_1 \cdot x_i)]} + (1-y_i) \cdot \underline{\log[1-g(\beta_0 + \beta_1 \cdot x_i)]}$$

✓

1e) The log function is a monotonically increasing function, as $f(x) = \log x$, and $f'(x) = \frac{1}{x}$ if log is to the base e. ~~Because~~

SB

Oct 4th

2021

With $x > 0$ as per definition of the log, $\frac{1}{x} > 0$! CS194-172
That means if we maximize x^* , $\log_e x$ will also maximize. \Rightarrow OK to use, and more convenient to compute.

*which we do,
 $\max L(\beta_0, \beta_1)$

1f) There are 2 ways to calculate the gradients - one extensive, one simpler by substitution. I'll do one each. $\log = \ln$

$$\begin{aligned} \frac{\partial \log L(\beta_0, \beta_1)}{\partial \beta_0} &= \sum_{i=1}^n y_i \cdot \frac{1}{g(\beta_0 + \beta_1 \cdot x_i)} \cdot \underbrace{g'(\beta_0 + \beta_1 \cdot x_i) \cdot (1 - g(\beta_0 + \beta_1 \cdot x_i))}_{\frac{\partial \log(g)}{\partial g}} \cdot \underbrace{\frac{1}{\beta_0 + \beta_1 \cdot x_i}}_{\frac{\partial}{\partial \beta_0}} \\ &\quad + (1 - y_i) \cdot \frac{1}{1 - g(\beta_0 + \beta_1 \cdot x_i)} \cdot (-1) \cdot g'(\beta_0 + \beta_1 \cdot x_i) \cdot (1 - g(\beta_0 + \beta_1 \cdot x_i)) \cdot 1 \\ &= \sum_{i=1}^n y_i \cdot \cancel{g'(\beta_0 + \beta_1 \cdot x_i)} - \cancel{g(\beta_0 + \beta_1 \cdot x_i)} - y_i \cdot \cancel{g'(\beta_0 + \beta_1 \cdot x_i)} \\ &= \sum_{i=1}^n y_i - g(\beta_0 + \beta_1 \cdot x_i) - \sum_{i=1}^n y_i - \frac{1}{1 + e^{-(\beta_0 + \beta_1 \cdot x_i)}} \end{aligned}$$

for $\frac{\partial \log L(\beta_0, \beta_1)}{\partial \beta_1}$ we use: $f = \log(g(h)) \& h = \beta_0 + \beta_1 \cdot x_i$ with $f' = \frac{\partial f}{\partial g(h)}$

$$\begin{aligned} &= \sum_{i=1}^n y_i \cdot f'(g(h)) \cdot g'(h) \cdot h' + (1 - y_i) \cdot f'(1 - g(h)) \cdot (1 - g(h))' \& h' = \frac{\partial h}{\partial \beta_1} \\ &= \sum_{i=1}^n y_i \cdot \frac{1}{g(h)} \cdot g'(h) (1 - g(h)) \cdot \frac{\partial h}{\partial \beta_1} + (1 - y_i) \frac{1}{1 - g(h)} \cdot (-1) \cdot g'(h) (1 - g(h)) \cdot x_i \\ &\quad - \sum_{i=1}^n (y_i - y_i \cdot g(h) - g(h) + y_i \cdot g(h)) \cdot x_i = \sum_{i=1}^n (y_i - g(h)) \cdot x_i \\ &= \sum_{i=1}^n \left(y_i - \frac{1}{1 + e^{-(\beta_0 + \beta_1 \cdot x_i)}} \right) \cdot x_i \quad \text{✓} \end{aligned}$$

Ig)

Initialize $\beta_0 = 0$, $\beta_1 = 0$;

Do :

$$\beta_0^{\text{new}} = \beta_0^{\text{old}} + \alpha \cdot \sum_{i=1}^n y_i - g(\beta_0^{\text{old}} + \beta_1^{\text{old}} \cdot x_i)$$

$$\beta_1^{\text{new}} = \beta_1^{\text{old}} + \alpha \cdot \sum_{i=1}^n (y_i - g(\beta_0^{\text{old}} + \beta_1^{\text{old}} \cdot x_i)) \cdot x_i$$

Until: $\|\beta_0^{\text{new}} - \beta_0^{\text{old}}\| < \text{threshold}$ and $\|\beta_1^{\text{new}} - \beta_1^{\text{old}}\| < \text{threshold}$

α being step size of gradient ascent.

SB

Homework 2 - ct'd

SB

Comp Gen

Oct 8th 2021

- 2a) Prove that $FWER \leq \alpha$ for rejecting all H_0^j at $p_j \leq \frac{\alpha}{m}$ (after Bonferroni correction)

$$\begin{aligned} FWER &= P(\text{rejecting } \geq 1 H_0^j \text{ falsely}) \quad \text{for all } H_1, \dots, H_m \\ &= P\left(\bigcup_{j=1}^m \text{rejecting } H_0^j \text{ falsely}\right) \end{aligned}$$

knowing that $P(A \cup B) \leq P(A) + P(B)$, no matter if A and B are independent

$$\Rightarrow FWER = P\left(\bigcup_{j=1}^m \text{rejecting } H_0^j \text{ falsely}\right) \leq \sum_{j=1}^m P(\text{rejecting } H_0^j \text{ falsely})$$

$$\Rightarrow FWER(\alpha) \leq \sum_{j=1}^m \alpha, \quad \text{with Bonferroni correction follows}$$

$$\Rightarrow FWER(\alpha) \leq \sum_{j=1}^m \frac{\alpha}{M} = m \cdot \frac{\alpha}{m} = \alpha$$

$$\Rightarrow FWER(\alpha) \leq \alpha$$

Bonferroni correction controls $FWER(\alpha)$ at $\leq \alpha$! for a family of hypothesis. The p-values need not be independent because of beer cese of inequality - we do not require to have j events / Hypotheses are independent. (which is probably helpful in variant analysis because of LD)

Oct 15th 2021

CS 194-172 – Computational Genomics – HW 2

Stefan Bielmeier

Problem 1h)

With perfectly linearly separable data, why might gradient ascent still not converge to stable values of β_0 and β_1 ?

=> one reason could be that the maximum of the log likelihood function does not exist. That can be the case with perfectly linearly separable data, where the data points can be completely separated by a linear function into two same-size categories, where $X^*\beta_1 < \text{(threshold of separation)}$ corresponds to $Y = 0$, and $X^*\beta_1 > \text{(threshold of separation)}$ correspond to $Y = 1$.

In that case, if we initialize β_0 – the intercept – with 0, it would converge quickly to an appropriate value: the threshold. β_1 however could be made artificially large, yielding the same, linear separation at that certain threshold X-value, and it still would not converge to a global maximum of the log likelihood function. It could be increased infinitely, with which it is not converging.

Problem 1i)

This is hard to say, multiple options are mentioned in [literature](#): Bayesian inference (estimate β_1 via a probability prior), or exact logistic regression.

Intuitively I would say add one noise data point to the perfectly linearly separable data so that it's not perfectly linearly separable anymore – basically some value $X < X\text{-threshold}$ with $Y = 1$, and calculate a β_1 to see if we can reject it or not. It works because we can no longer separate our data by a linear function with an intercept, but need to actually run the regression.

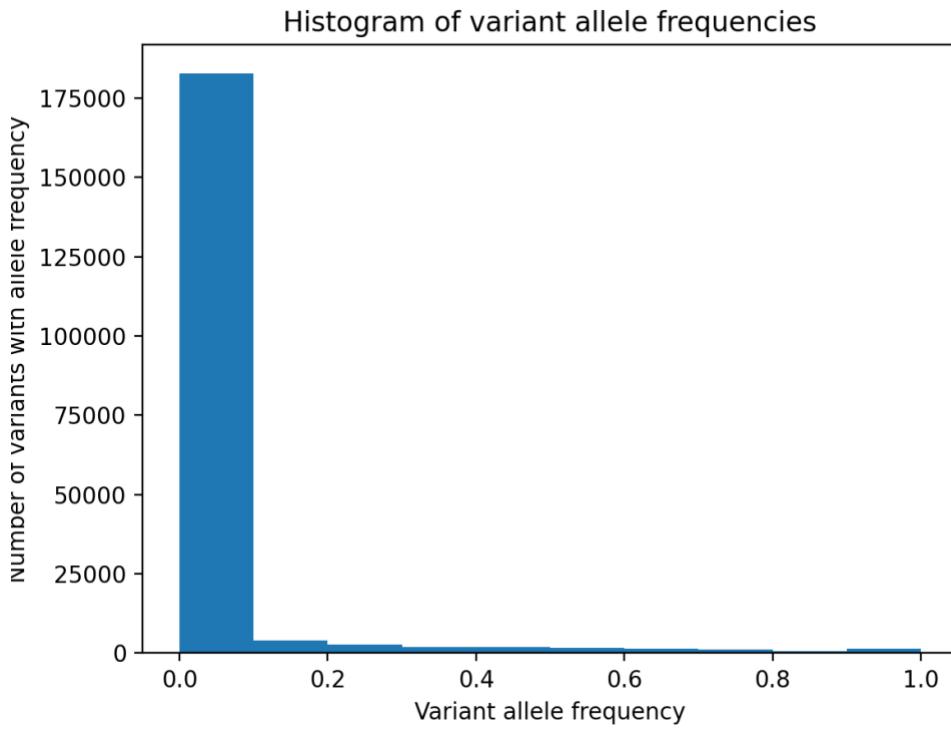
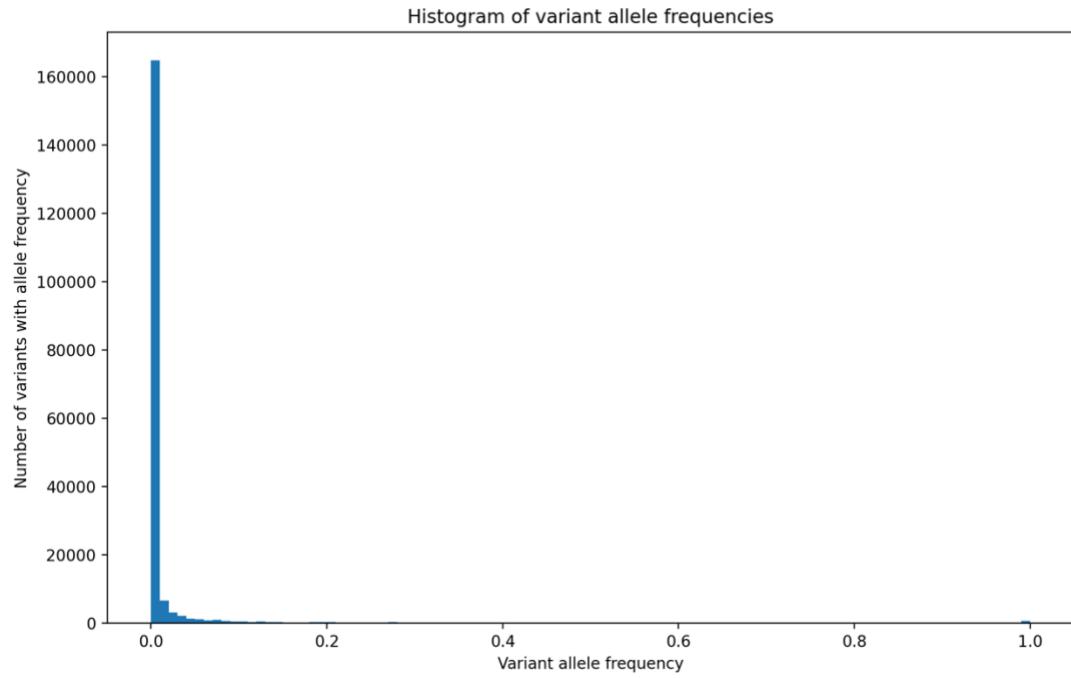
This may work for large datasets where noise is tolerable.

Problem 3a)

There are 190896 SNPs listed in the file, out of a total of 198102 variants.

Problem 3b)

Provided histogram in 2 different x-axis scales...



The number of variants with allele frequencies $< 1\%$ is 164.785, which is most of them.

Problem 3c)

There are 2504 samples (individuals).

The average individual contains approx. 11663.579872204473 many distinct variants.

To check if this is correct, sum of all alternative alleles divided by the number of samples should be bigger than the average distinct variant frequency (above): 16042.289137380192

True!

Problem 3d)

A dbSNP is a known SNP / genetic variation in the free public archive for known genetic variations for different species by National Center for Biotechnology Information (NCBI).

It has an rs number, e.g. rs328 which is without clinical significance (C => (A, G))

Going through all SNPs (rows) in the file, not a single one has an rs number in their ID column.

That means that there are 0 dbSNPs in the VCF file.

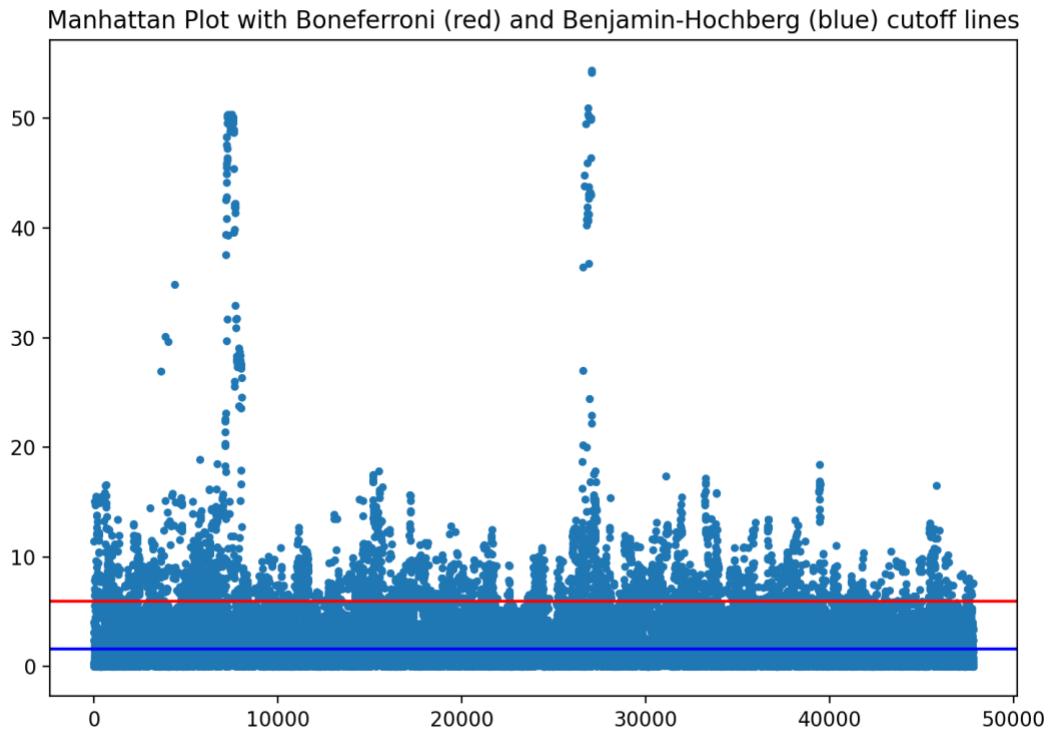
Problem 3e)

A Phred quality score is a measure of the quality of the identification of the nucleobases generated by automated DNA sequencing.

With a Phred score of 45, what is the accuracy of the variant call?

- 1) what is the error probability? $P = 10^{(-45/10)} = 0.0000316227766$.
- 2) The accuracy is thus $1 - P$, which is 99.99683772

Problem 4b)



Problem 4c)

Out of 47808 SNPs present in the dataset for chromosome 22,

4013 SNPs are reported as significant after the Bonferroni correction

and 24390 SNPs are reported as significant after the Benjamin-Hochberg correction.

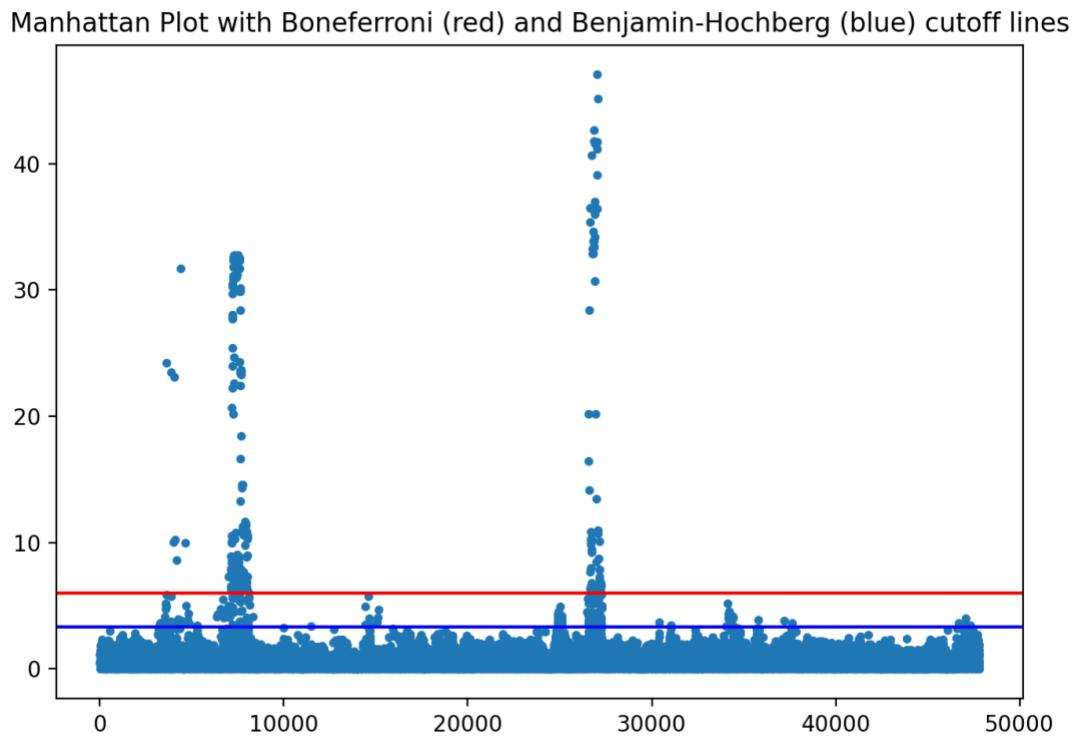
Problem 4d)

One reason why I see so many significant SNPs even after controlling for FWER and FDR at alpha = 0.05 is population stratification. Due to linkage disequilibrium, many non-causal SNPs in the same chromosome are inherited together with the causal SNP. That is especially the case when the different subpopulations are observed. That would mean that the proband who has the causal mutation (SNP) would show allele mutations for the other, non-causal SNPs, as well, simply because they are part of a certain subpopulation. When our model encounters such a non-causal SNP, it will still come to the conclusion that that non-causal SNP is significantly associated with the disease (phenotype == 1), even though the non-causal SNP is only significantly associated with the causal SNP through ancestry.

As I study only chromosome 22, it makes sense that my analysis reports high significance levels for many non-causal SNPs, because these non-causal SNPs show similar significant association to the phenotype (according to our model) but are in fact only associated with the 5 functional SNPs due to inheritance. Correcting for FWER and FDR doesn't make a difference, as the significance

levels for the non-causal, LD SNPs are just as high as the causal-SNPs because of very similar alternative allele counts for SNPs that are inherited together.

Problem 4e)



After including the Top 3 principal Components as covariates in the logistic regression analysis of 47808 SNPs present in the dataset for chromosome 22,

253 SNPs are now reported as significant after the Bonferroni correction,
and **481** SNPs are now reported as significant after the Benjamin-Hochberg correction.

Problem 4f)

Including the top (3) principal components for each individual in the logistic regression models reduces the number of significant SNPs because they reduce the multi-collinearity (mutual dependence) between our “features” for each individual. For example, the presence of alternative alleles in non-causal SNPs and causal SNPs for a certain individual caused by ancestry and joint inheritance in a subgroup of a population (population stratification).

PCA can do that by representing all 2504 SNPs for each individual in only 3 top principal components. The 3 principle components represent variance of alternative alleles that occur in one

individual. For example, when we use the principle components as covariates in one logistic regression model to estimate the effect of one single SNP on the phenotype, the principle components account for most of the variance (per definition) for each sample/individual in the model that is caused by population stratification. If the principle components capture most of the variance caused by alternative alleles, the SNP under consideration must itself show greater variance over all samples so that the model determines the its (**Beta₁**) to be significantly associated with the phenotype.

Therefore, less SNPs reach **Beta₁s** with p-values low enough to pass the threshold of, e.g. $0.05 / \text{number of observations}$ after a Bonferroni correction, leading to less rejected H_0 Hypothesis (β_{11}). This makes it easier to detect causal SNPs.

Problem 4e)

The model may not perform as well on individuals of African ancestry because they show generally different alternative allele frequencies across their genome, compared to individuals of European ancestry: population stratification. The model may interpret the presence of knowingly non-phenotype causing alternative alleles in individuals of African ancestry as phenotype-causing alternative alleles, simply because it doesn't recognize these alternative alleles to be commonly present, as the model was trained on genomes of people of European ancestry, for whom the absence of such alternative alleles is common.