

SOM: Bayesian analysis of genetic association
across tree-structured routine healthcare data in
the UK Biobank

January 11, 2017

Contents

1	Model description	2
2	Model fitting and Bayes factor calculation	4
3	Conditional analysis	6
4	Supplementary Figures	8

1 Model description

Consider a data set of N individuals, each of which is annotated with a series of categorical observations, which are themselves organised in a hierarchical structure reflecting increasing levels of resolution; *i.e.*, annotations are associated with nodes within a classification tree. Observations may be made at both terminal and internal nodes depending on resolution. We define an indicator, Z_{ij} , for the presence of at least one annotation $j \in T$ for individual i ; where T is the set of all annotations (organised as a tree). We model the distribution of Z_{ij} , conditional on the genotype of individual i at variant s , $G_{is} \in \{0, 1, 2\}$, using a logistic model, with an intercept (β_j^0) and separate coefficients for the heterozygous (β_j^1) and homozygous (β_j^2) states:

$$Y_{ijs} = \beta_j^0 + \beta_j^1 * I(G_{is} == 1) + \beta_j^2 * I(G_{is} == 2), \quad (1)$$

$$P(Z_{ij} = 1 | Y_{ijs}) = \frac{e^{Y_{ijs}}}{(1 + e^{Y_{ijs}})}. \quad (2)$$

To model the correlation structure of the genetic coefficients across categories, we allow the coefficient pair $\{\beta^1, \beta^2\}$ to evolve down the tree in a Markovian fashion. The coefficients attached to a parent node x can either be inherited by a child node y , with probability $e^{-\theta}$, or can transition to a new pair of values, with probability $1 - e^{-\theta}$. With probability $1 - \pi_1$ the new values are $\{0, 0\}$, and with probability π_1 they are drawn from a joint prior on β^1 and β^2 , $f(\beta^1, \beta^2)$. The state of the ancestral node in the tree is drawn from the stationary distribution of this process; *i.e.*, $\{0, 0\}$ with probability $1 - \pi_1$ or from $f(\beta^1, \beta^2)$ with probability π_1 . We use a non-local prior for $f(\beta^1, \beta^2)$, such that:

$$f(\boldsymbol{\beta}) = N_2(\mathbf{0}, \Sigma) * |\boldsymbol{\beta}|^k * e, \quad (3)$$

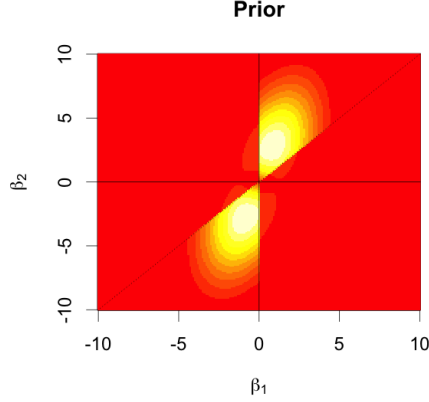
with

$$\Sigma = \begin{bmatrix} \sigma_1^2 & r\sigma_1\sigma_2 \\ r\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \quad (4)$$

and,

$$e = \begin{cases} 0.10, & \text{if } \beta_1 * \beta_2 < 0 \\ 0.10, & \text{if } \beta_1 > \beta_2 \\ 1, & \text{otherwise.} \end{cases} \quad (5)$$

The density of the prior is illustrated in Figure 1. The mixture prior on the coefficients (including the point mass at 0) is referred to as $f^*(\beta^1, \beta^2)$. Unless stated otherwise we use parameter values $\pi_1 = 0.001$, $\theta = 1/3$, $\sigma_1 = 2$, $\sigma_2 = 4$, $k = 1/2$ and $r = 0.5$, throughout. The unknown intercept term β_j^0 is chosen, for each value of $\{\beta_j^1, \beta_j^2\}$, to maximise the likelihood. That is,



Supplementary Figure 1: Prior on effect sizes for the full genetic model. Respecitvely, β_1 and β_2 are the log-odds coefficients for the heterozygotes and homozygotes. The heatmap indicates the relative density of the prior.

$$L(\beta_j^1, \beta_j^2 | \mathbf{Z}_j) = \max_{\beta_j^0} L(\beta_j^0, \beta_j^1, \beta_j^2 | \mathbf{Z}_j), \quad (6)$$

where \mathbf{Z}_j is $\{Z_{1j}, Z_{2j}, \dots, Z_{N_j}\}$.

The joint distribution of different annotations across individuals has substantial non-independence. For example, the same individual might be recorded as having different subtypes of a disorder on separate visits to a hospital, the recording of a specific disease subtype will mean that other subtypes are less likely to be recorded for the same individual and a disease may have multiple diagnostic features. However, rather than attempt to capture such structure, we make the approximation that annotations are independent conditional on an individual's genotype (and evaluate the impact of this approximation). Hence, the likelihood for a given vector of $\{\beta^1, \beta^2\}$ values across annotations, β , is given by the product over all nodes in the tree T :

$$L(\beta | \mathbf{Z}) = \prod_{j \in T} L(\beta_j | \mathbf{Z}_j), \quad (7)$$

where $\beta_j = \{\beta_j^1, \beta_j^2\}$. The prior density for β can be calculated by considering the state of the ancestral node, A , and all transitions between parent and child nodes:

$$P(\beta) = p(\beta_A) \prod_{p,c} q(\beta_p, \beta_c), \quad (8)$$

where $q(\beta_p, \beta_c)$ is the transition probability between the coefficients of the parent and child nodes. Because of the structure of the model, it is possible to

sum the likelihood over all possible values of β using dynamic programming. To achieve this, for each node j we calculate an integrated likelihood

$$L_j = \int P_j(D|\beta) f^*(\beta) d\beta, \quad (9)$$

where $P_j(D|\beta)$ is given by the likelihood function in Equation 6 when j is a terminal node, or by

$$P_j(D|\beta) = \prod_{i \in \gamma(j)} [e^{-\theta} P_i(D|\beta) + (1 - e^{-\theta}) L_i], \quad (10)$$

when j is an intermediate node. Here, $e^{-\theta}$ is the stay transition probability in β and $(1 - e^{-\theta})$ is the switch transition probability in β , which results in uncorrelated genetic coefficients between nodes. Note that in practice we evaluate the functions over a grid of values for β .

The full likelihood (i.e. by summing over all possible coefficients) is given by summing the values at the ancestral node A :

$$L_{full} = L_A. \quad (11)$$

The likelihood under the model of no genetic association across all nodes in the tree, L_\emptyset , is calculated by summing the likelihood over all nodes with $\beta = \mathbf{0}$, and the prior on this:

$$L_j(\beta = \mathbf{0}) = \prod_{i \in \gamma(j)} p_{00} L_i(\beta = \mathbf{0}), \quad (12)$$

where $p_{00} = e^{-\theta} + (1 - e^{-\theta})(1 - \pi_1)$. For terminal nodes, $L_j(\beta = \mathbf{0})$ is calculated directly from the likelihood function by evaluating Equation 6 at $\beta = \mathbf{0}$. It follows that the null likelihood L_\emptyset is given by

$$L_\emptyset = (1 - \pi_1) L_A(\beta = \mathbf{0}). \quad (13)$$

2 Model fitting and Bayes factor calculation

There are two objectives to the analysis. First, to calculate the evidence for association between a genetic variant and any of the annotations, thus identifying variants that have association to at least one annotation. Second, for variants with some association, to identify those annotations with non-zero coefficients.

Our first objective can be met by calculating a Bayes factor that compares the likelihood integrated over all possible values of β in which at least one node is active, L^+ , to the likelihood under which all nodes are inactive. By noting that there is only one way in which all nodes can be inactive and that it is easy to calculate both the prior, π_\emptyset , and likelihood, L_\emptyset , for this state, we can obtain the Bayes factor as follows. First, note that we can rewrite the full likelihood function L_{full} in Equation 11 as:

$$L_{full} = \pi_{\emptyset} L_{\emptyset} + \sum_{p \in \emptyset'} \pi_p L_p, \quad (14)$$

which sums over the path where all nodes are inactive and all possible path with at least one active node ($p \in \emptyset'$). Then, we can solve for the likelihood L^+ :

$$L^+ = \frac{L_{full} - \pi_{\emptyset} L_{\emptyset}}{(1 - \pi_{\emptyset})}. \quad (15)$$

The desired Bayes factor is then calculated by taking the ratio of the two likelihoods:

$$\text{BF}_{\text{tree}} = \frac{L^+}{L_{\emptyset}} = \frac{L_{full} - \pi_{\emptyset} L_{\emptyset}}{(1 - \pi_{\emptyset}) L_{\emptyset}}. \quad (16)$$

Using the same framework, it is also possible to compute Bayes factors for the cases where there is no correlation in state between parent and child nodes (*i.e.*, $\theta \rightarrow \infty$), and where all states are active and either share a single set of coefficients (*i.e.*, $\pi_1 \rightarrow 1$, $\theta \rightarrow 0$) or are independent (*i.e.*, $\pi_1 \rightarrow 1$, $\theta \rightarrow \infty$). In theory it would be possible either to estimate π_1 and θ or to integrate over a hyper-prior.

For those variants where there is evidence for association within the annotation tree, it is possible to identify active nodes and estimate coefficients of association for each node by using the forward and backward algorithms, also known as the inside and outside algorithms when applied to tree-like Markov models. The forward (inside) algorithm has been described above, though for completeness and consistency of notation, it is repeated below.

In the forward (inside) algorithm we are iterating up from the terminal nodes towards the root of the tree calculating the joint likelihood of the subtree each node subtends. To initialise, let j be a terminal node, so $F_j(\boldsymbol{\beta})$ is the probability of the observed data at node j for a given value of $\boldsymbol{\beta}$,

$$F_j(\boldsymbol{\beta}) = P_j(D|\boldsymbol{\beta}). \quad (17)$$

We can then integrate over the values of $\boldsymbol{\beta}$ to calculate the integrated likelihood at node j as in Equation 9,

$$L_j = \int F_j(\boldsymbol{\beta}) f^*(\boldsymbol{\beta}) d\boldsymbol{\beta}. \quad (18)$$

For intermediate nodes we calculate F_j recursively up the tree. First, let j be an intermediate node and $\gamma(j)$ the set of child nodes of j . For each $i \in \gamma(j)$ we define,

$$G_i(\boldsymbol{\beta}) = e^{-\theta} F_i(\boldsymbol{\beta}) + (1 - e^{-\theta}) L_i, \quad (19)$$

It follows that for an internal node

$$F_j(\boldsymbol{\beta}) = \prod_{i \in \gamma(j)} G_i(\boldsymbol{\beta}). \quad (20)$$

We can then calculate the integrated likelihood at node j using Equation 18 as for the terminal nodes and continue the algorithm up the tree until the ancestral node, A , is reached.

In the backward (outside) algorithm we calculate the probability density of β starting from the root of the tree and moving recursively down the tree. The quantity we are aiming to calculate is the likelihood for the data not subtended by the node of interest.

To initialise, let A be the ancestral node, so $B_A(\beta)$ is given by the prior on β ,

$$B_A(\beta) = f^*(\beta). \quad (21)$$

We then iterate down the tree. Let i and j be such that $j \in \gamma(i)$ (that is i is the parent of j), then

$$B_j(\beta) = \int B_i(\beta')q(\beta, \beta') \frac{F_i(\beta')}{G_j(\beta)} d\beta', \quad (22)$$

where $q(\beta, \beta')$ is the (transition) probability of state β in the daughter node given state β' in the parent node. Note that because of the structure of the model there are only two types of transition, which enables efficient calculation. The posterior density for β in node j can then be calculated from

$$\pi_j(\beta|D) = \frac{F_j(\beta) \times B_j(\beta)}{L_{full}}, \quad (23)$$

and from this distribution we can integrate to estimate the probability of $\beta \neq \mathbf{0}$ and the 95% credible sets for β .

3 Conditional analysis

To account for linkage disequilibrium in the MHC and to identify independent associations with HLA alleles we performed conditional analysis. For each of the datasets (SR and HES) we first analysed each imputed HLA and identified the allele with the strongest evidence of association, as measured by the BF_{tree} statistic. We then continue to analyse the remaining HLA alleles in an iterative approach, where at each iteration we controlled for previous identified HLA alleles, through conditional analysis. To account for these covariates in the analysis we use an approximation to the likelihood function. Let Δ_{ij} quantify the aggregated risk effects due to covariates in individual i in annotation j ,

$$\Delta_{ij} = \sum_k [\hat{\beta}_{jk}^1 \times I(G_{ik} == 1) + \hat{\beta}_{jk}^2 \times I(G_{ik} == 2)], \quad (24)$$

where the genetic coefficients $\{\hat{\beta}_{jk}^1, \hat{\beta}_{jk}^2\}$ are the MAP estimates inferred for the HLA allele with the largest BF_{tree} in round k for annotation j , and $G_{ik} \in \{0, 1, 2\}$ are the genotypes for individual i in the HLA allele identified in round k .

To model the distribution of Z_{ij} we modified the logistic model in Equation 1 and 2 to account for the aggregate risk effect due to HLA alleles identified in previous rounds:

$$Y_{ijs}^c = \beta^0 + \beta_j^1 * I(G_{is} == 1) + \beta_j^2 * I(G_{is} == 2) + \Delta_{ij}, \quad (25)$$

and,

$$P(Z_{ij} = 1 | Y_{ijs}^c) = \frac{e^{Y_{ijs}^c}}{1 + e^{Y_{ijs}^c}}. \quad (26)$$

The conditional likelihood function is then given by the binomial distribution,

$$L_j^c(\boldsymbol{\beta} | \mathbf{Z}_j) = \prod_{i=1}^N p_{ij}^{c^{Z_{ij}}} (1 - p_{ij}^c)^{1 - Z_{ij}}, \quad (27)$$

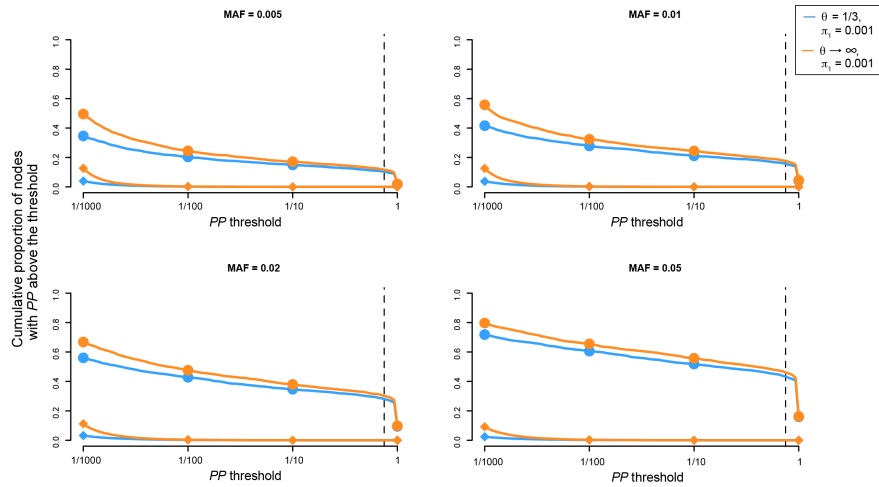
where we let $p_{ij}^c = P(Z_{ij} = 1 | Y_{ijs}^c)$.

To compute the above conditional likelihood we use an approximation by taking the 2nd order Taylor expansion around $\Delta = 0$. After evaluation of the first and second derivatives of $\log(L_j^c(\boldsymbol{\beta} | \mathbf{Z}_j))$ at $\Delta = 0$ and simplifying terms we obtain:

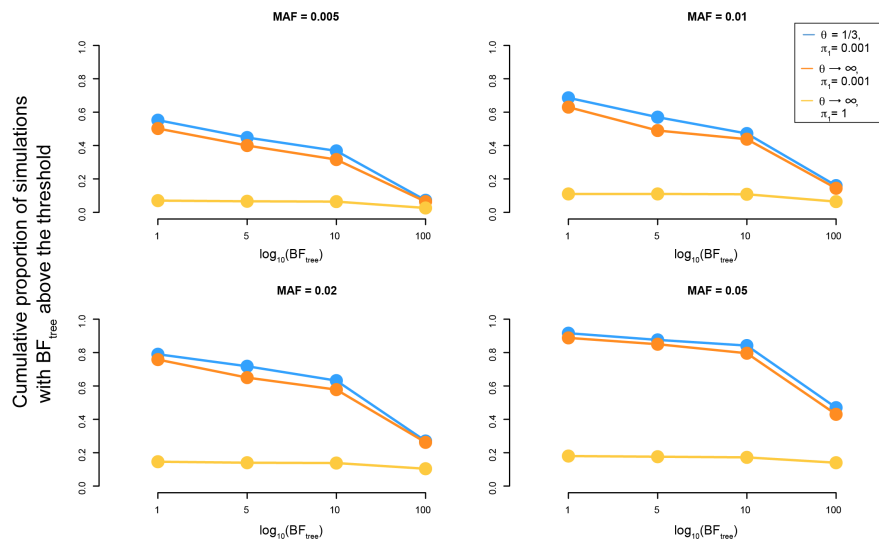
$$\log(L_j^c(\boldsymbol{\beta} | \mathbf{Z}_j)) \approx \log(L_j(\boldsymbol{\beta} | \mathbf{Z}_j)) + \sum_{i=1}^N [\Delta_{ij}(Z_{ij} - p_{ij}) - \frac{\Delta_{ij}^2}{2} p_{ij}(1 - p_{ij})], \quad (28)$$

where $L_j(\boldsymbol{\beta} | \mathbf{Z}_j)$ and p_{ij} are given by the equivalent functions when we don't account for covariates. We note that while the approximation works well for early rounds, its accuracy is likely to decrease after multiple rounds of conditioning. Extensions that enable re-estimation at later steps will be explored in subsequent work.

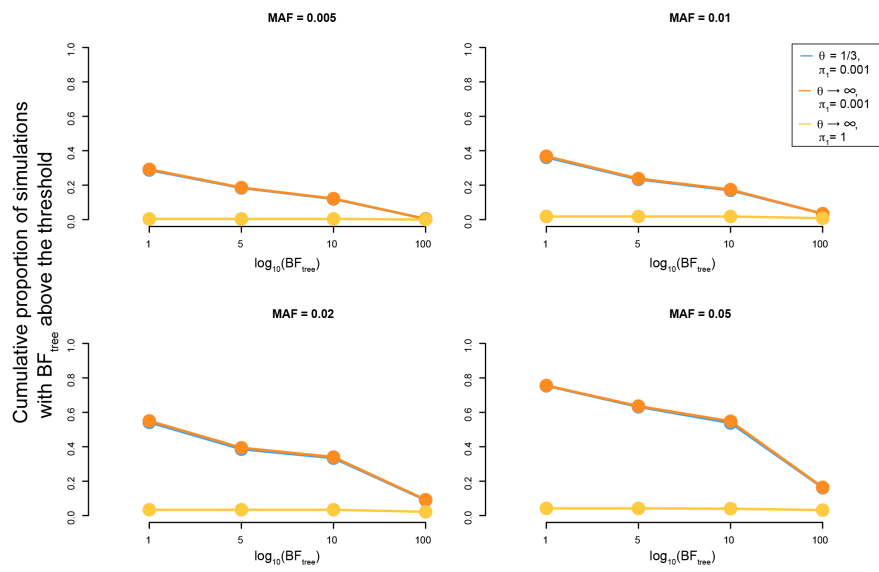
4 Supplementary Figures



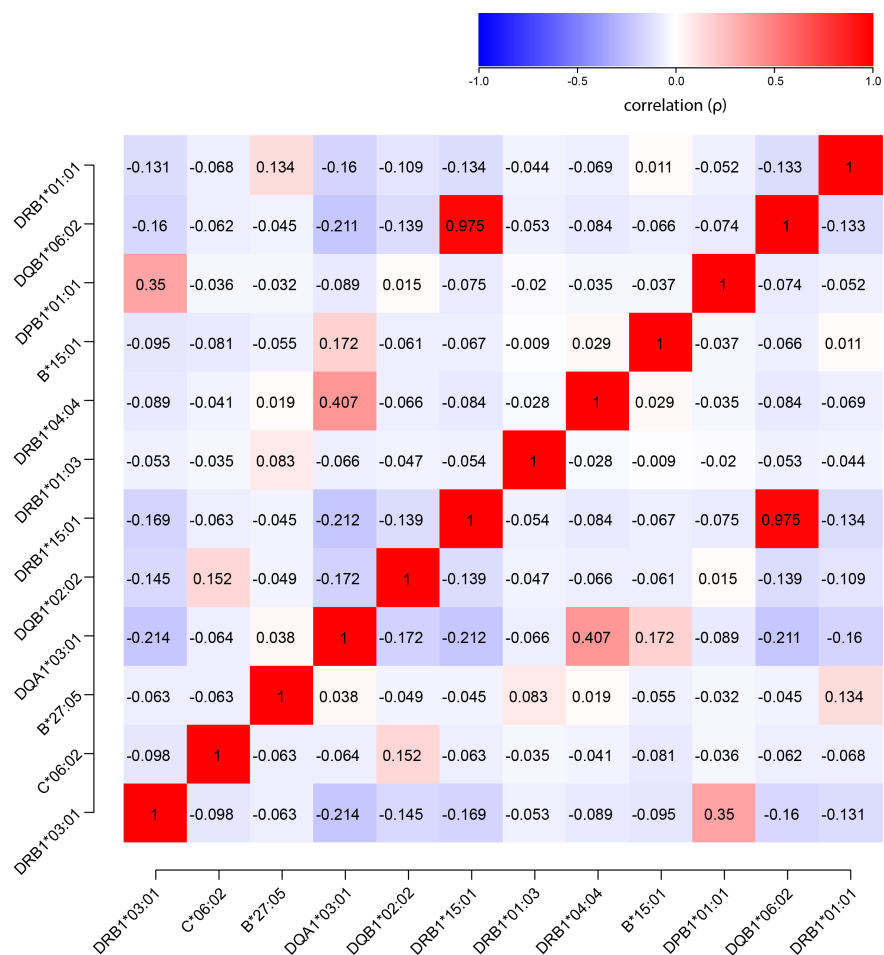
Supplementary Figure 2: Rate of active node identification at increasing posterior probability (PP) thresholds and different simulated allele frequencies of the causal genetic variant, for the TreeWAS method ($\theta = 1/3$ and $\pi_1 = 0.001$) and a model assuming complete independence among phenotypes ($\theta \rightarrow \infty$ and $\pi_1 = 0.001$). We simulated data for 500 replicates where the genetic variant affects clinical annotations found distributed in the tree. Rate of active node identification was calculated for the five affected clinical annotations (\bullet) and for the rest of the annotations is the tree with zero genetic coefficients (\diamond).



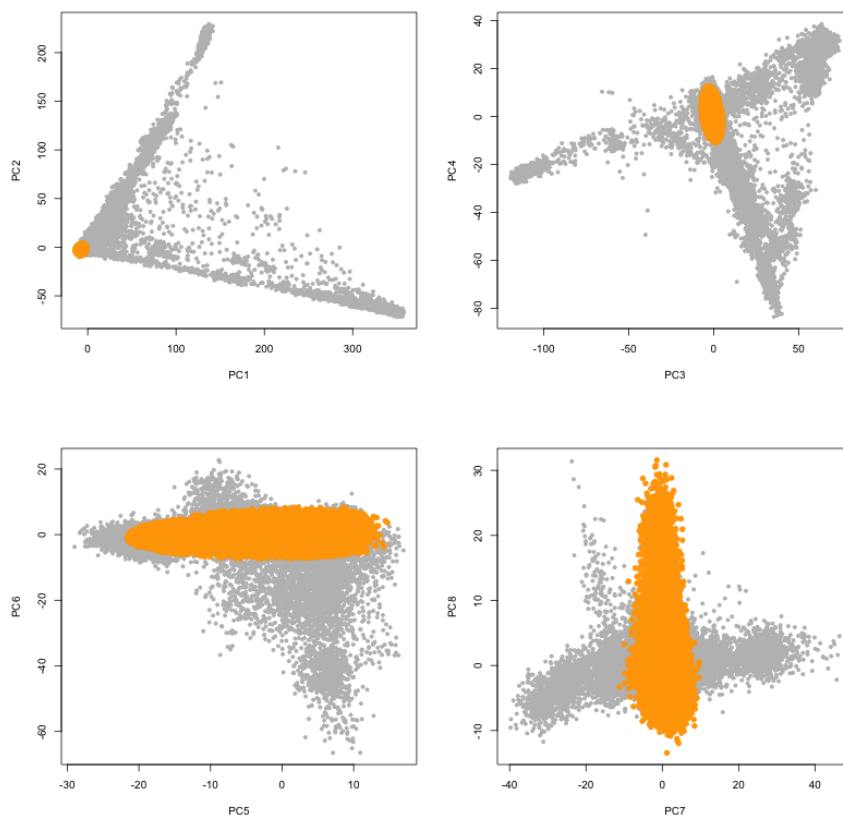
Supplementary Figure 3: Sensitivity analysis in the detection of genetic association at the tree level as measured by the BF_{tree} statistic. We simulated data where the causal variant affected clustered nodes in the tree and fitted the TreeWAS method (blue) and the models where we assume complete independence among phenotypes (orange) and where we assume complete independence among phenotypes and all nodes to be active (yellow).



Supplementary Figure 4: Sensitivity analysis in the detection of genetic association at the tree level as measured by the BF_{tree} statistic. We simulated data where the causal variant affected distributed nodes in the tree and fitted the TreeWAS method (blue) and the models where we assume complete independence among phenotypes (orange) and where we assume complete independence among phenotypes and all nodes to be active (yellow).



Supplementary Figure 5: Linkage disequilibrium between independent HLA associations found in the analyses for the SR and HES datasets. Each allele shown was found in at least one of the analyses, seven of which were found in both. With the exception of the alleles *HLA-DRB1*15:01* and *HLA-DQB1*06:02*, all identified associations were not in linkage disequilibrium ($r^2 < 0.02$). The alleles *HLA-DRB1*15:01* and *HLA-DQB1*06:02* were identified in the SR and HES analyses, respectively, and both are in high linkage disequilibrium ($\rho = 0.98$) and both were fine-mapped to the same phenotypes.



Supplementary Figure 6: Ancestry analysis of UK Biobank individuals using principal component analysis. 120,286 individuals plotted in orange were retained in the analysis and these co-cluster with European ancestry populations.