

Noise-augmented directional clustering of genetic association data identifies distinct mechanisms underlying obesity

Andrew J. Grant^{*1}, Dipender Gill^{2,3,4,5}, Paul D. W. Kirk^{1,6}, and Stephen Burgess^{1,7}

¹MRC Biostatistics Unit, University of Cambridge, Cambridge, UK

²Department of Epidemiology and Biostatistics, School of Public Health, St Mary's Hospital, Imperial College London, London, UK

³Clinical Pharmacology and Therapeutics Section, Institute of Medical and Biomedical Education and Institute for Infection and Immunity, St George's, University of London, London, UK

⁴Clinical Pharmacology Group, Pharmacy and Medicines Directorate, St George's University Hospitals NHS Foundation Trust, London, UK

⁵Novo Nordisk Research Centre Oxford, Old Road Campus, Oxford, United Kingdom

⁶Cambridge Institute of Therapeutic Immunology & Infectious Disease (CITIID), University of Cambridge, Cambridge, UK

⁷Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK

*Corresponding author. Email address: andrew.grant@mrc-bsu.cam.ac.uk

Abstract

Clustering genetic variants based on their associations with different traits can provide insight into their underlying biological mechanisms. Existing clustering approaches typically group variants based on the similarity of their association estimates for various traits. We present a new procedure for clustering variants based on their proportional associations with different traits, which is more reflective of the underlying mechanisms to which they relate. The method is based on a mixture model approach for directional clustering and includes a noise cluster that provides robustness to outliers. The procedure performs well across a range of simulation scenarios. In an applied setting, clustering genetic variants associated with body mass index generates groups reflective of distinct biological pathways. Mendelian randomization analyses support that the clusters vary in their effect on coronary heart disease, including one cluster that represents elevated body mass index with a favourable metabolic profile and reduced coronary heart disease risk. Analysis of the biological pathways underlying this cluster identifies inflammation as playing a key role in mediating the effects of increased body mass index on coronary heart disease.

Introduction

In recent years, the number of genome-wide association studies (GWAS) has grown enormously [1]. Such studies provide valuable information linking genetic variants across the human genome to a wide range of traits. What often remain less understood are the underlying mechanisms by which the associated genetic variants affect the traits. Insight into these mechanisms may be gained by investigating the pattern of associations with other related traits: genetic variants that share similar association patterns may be thought to act via similar mechanisms [2]. For example, some genetic variants associated with type 2 diabetes are also associated with obesity related traits such as body mass index (BMI), whereas others are instead associated with traits such as triglycerides, suggesting that the variants influence type 2 diabetes risk via different biological mechanisms [3].

A number of techniques have been implemented to cluster genetic variants based on their associations with traits that are believed to be relevant in informing biological pathways. The traits often include separate risk factors or potential mediators of some disease outcome(s) of interest. A common approach is to use hierarchical clustering, which groups observations based on their distance from each other [4, 5, 6, 7]. The number of clusters is then chosen heuristically. Other clustering approaches which have been applied to genetic variant-trait association estimates include fuzzy c-means [6] and Bayesian nonnegative matrix factorization [3]. A related approach which aims to determine distinct components of genetic variant-trait associations uses truncated singular value decomposition [8].

A key characteristic of previously implemented approaches is that they cluster based on the Euclidean distance between vectors of the genetic variant-trait association estimates, defined as the length of the line between the association estimates plotted as points on a graph. However, when trying to determine shared biological mechanisms, a more relevant clustering target is the proportional associations of each genetic variant with the set of traits. If two variants influence a set of related traits via a common mechanism, the genetic associations may differ considerably in magnitude due to one variant having a stronger effect than the other. However, their proportional associations across the traits will be similar for both variants. Equivalent to looking at proportional associations is to consider the direction of the association vector from the origin. That is, in order to distinguish between variants which act via different mechanisms, it is the direction of the association

vector rather than its location in space which is of most importance. This is illustrated graphically in Fig. 1.

In this paper we introduce a novel procedure for clustering genetic variants based on their associations with a given set of traits to identify groups with common biological mechanisms. We develop the NAvMix (Noise-Augmented von Mises–Fisher Mixture model) clustering method, which extends a directional clustering approach to include a noise cluster as well as a data-driven method for choosing the number of clusters. The method is shown in a simulation study to perform well in identifying true clusters and to outperform alternative approaches across a range of scenarios. We further apply the procedure to cluster genetic variants associated with body mass index (BMI). We study the downstream effects of the different components of BMI on coronary heart disease (CHD) using Mendelian randomization, which uses genetic variants as instrumental variables to study potential causal effects of a risk factor on an outcome [9, 10]. We identify a BMI increasing cluster of variants associated with a favourable cardiometabolic profile and lower CHD risk. Analysis of the biological pathways which underlie each group of variants suggests that a key difference of this cluster compared with the others is its distinct effect on systemic inflammation. The clustering method demonstrated in this work is thus able to identify distinct pathways underlying complex traits, in turn highlighting specific mechanisms for therapeutic intervention.

Results

Overview of the proposed clustering approach

We use a mixture model approach to clustering, which supposes that each observation is a realisation from one of a fixed number of probability distributions. Since we are interested in clustering based on direction of association, we fit a mixture of von Mises–Fisher (vMF) distributions, which is a distribution characterised by the mean direction of the observations from the origin and a dispersion parameter. A mixture model of vMF distributions has previously been described by Banerjee et al. [11]. We augment this approach by including a noise cluster, in recognition of the fact that not all observed vectors of genetic variant–trait association estimates are expected to fit well within the set of specified distributions. The noise cluster will contain outliers to the specified model, providing robustness to the identification of clusters. Our method of clustering is thus to fit a Noise-Augmented von Mises–Fisher Mixture model (NAvMix).

The NAvMix algorithm outputs a probability for each observation belonging to each cluster based on the given data. Each observation can then be assigned according to which cluster it has the highest probability of membership (referred to as hard clustering). The approach also provides the ability for soft clustering, which is where an observation is assigned to any cluster for which it has a probability of membership over a certain level, so that observations may belong to more than one cluster. Although the algorithm requires a fixed number of clusters to be specified, we repeat the procedure for varying numbers of clusters then chose the final number using the Bayesian Information Criterion (BIC). Full details of the procedure are given in the Methods section.

Let $\hat{\beta}_j$ be the vector of association estimates of genetic variant j with the set of traits under consideration, and let $\hat{\Sigma}_j$ be the covariance matrix of this vector. We assume that the genetic variants are independent of each other (that is, no linkage disequilibrium). We also note that the association estimates do not need to have been taken in the same sample, so we can consider sets of associations between genetic variants and any trait for which corresponding GWAS summary statistics are available. Although it is possible to input the raw association estimates into the

algorithm, we propose inputting the standardised association estimates, given by $\hat{\Sigma}_j^{-1/2}\hat{\beta}_j$, for the j th variant. The standardisation means that each element of the input vector is independent and has the same standard error. It thus is able to account for correlation between association estimates. Assuming all genetic associations are estimated with the same sample size for a given trait, this will not distort the direction vector. Otherwise, the direction vector will be weighted toward traits for which the associations are more precisely estimated. The first step in the algorithm is to transform each input vector to have a magnitude of one. This is done by dividing each vector by its Euclidean distance from the origin. We shall refer to this as normalisation. The normalised vectors represent the proportional association estimates.

The diagonal elements of the covariance matrices represent the variances of the genetic variant-trait association estimates. The off-diagonal elements represent the covariances between these estimates. If the genetic associations are estimated in separate samples for each trait, these covariances will be theoretically equal to zero. If the association estimates are taken from the same sample, the covariances will still be approximately zero if the traits are independent. If the traits are correlated, an estimate of this correlation is required to estimate the full covariance matrix in the one sample setting. This is easily computed using individual level data (Methods). If published GWAS summary statistics are being used, this information will not always be available. Nonetheless, the simulation study presented in the following section shows the clustering approach still performs well in the case where traits are truly correlated but the correlation estimates are set to zero.

Simulation results

We performed a simulation study in order to evaluate the performance of the proposed method and to compare it with alternative clustering approaches. We chose two methods for comparison. The first was to fit Gaussian mixture models to the standardised association estimates using the `mclust` algorithm in R [12]. The method was chosen for comparison because it is a model-based approach that is able to estimate the number of clusters by fitting multiple models and choosing between them using a principled model selection criterion. The second approach used for comparison was to fit Gaussian mixture models using the proportional association estimates. This is a case of model misspecification, since the association estimates after normalisation will not follow Gaussian distributions, even if the association estimates themselves do (see, for example, Fig. 1). It thus demonstrates the result of applying a method for clustering based on Euclidean distance to proportional associations. Note that other R packages which implement a form of directional clustering were not used for comparison because they either do not allow for estimation of the number of clusters (for example, `skmeans` [13], which uses the spherical k-means algorithm) or do not incorporate a noise cluster (for example, `movMF` [14]), and so performance cannot easily be compared.

We simulated data for genetic variants across six scenarios, where the number of traits (denoted by m) was either 2 or 9 and the number of clusters (K) was either 1, 2 or 4. In each scenario, 100 genetic variants were generated from these clusters and 20 additional noise genetic variants were generated. In the primary simulation study presented here, the genetic variant-trait associations were estimated in a single sample of 20 000 individuals. The traits were correlated but the off-diagonal entries of the covariance matrices were set to zero. This emulates the scenario in which genetic variant-trait associations are estimated in the same sample but where only GWAS summary data, with no trait correlation estimates, are available. The Supplementary Information presents the results of two further simulation studies. In the first, the estimated trait correlation from

individual level data is incorporated into the procedure, so the full estimated covariance matrices are used. In the second, the association estimates are taken from different samples of different sizes. Full details of the simulation parameters are given in the Methods section.

We evaluated the performance of each method using three measures: the Rand index [15]; the mean number of clusters estimated; and the mean number of observations assigned to the noise cluster. The Rand index is a similarity measure between the true and estimated cluster memberships, and shows how well each method allocated the observations. The closer to 1, the closer the estimated cluster membership is to the truth. Fig. 2 shows boxplots of the Rand index for each method and scenario. In calculating the Rand index we exclude genetic variants that truly belong to the noise cluster. Table 1 shows the mean number of clusters estimated and the mean size of the noise cluster for each method and scenario.

NAvMix performed very well in terms of allocating the observations to the correct clusters, with a median Rand index over 0.965 in all scenarios. It selected the correct number of clusters in almost all repetitions of all scenarios, with the exception that in the 9-dimensional scenario with 4 true clusters, it tended to slightly overestimate (estimating 4.45 on average). NAvMix also outperformed the other methods across each of the metrics considered. The median Rand index was higher, and the spread of the Rand indices was lower, in each scenario. The mclust algorithm, on average, underestimated the number of clusters, whereas mclust using proportional associations, on average, overestimated the number of clusters. Furthermore, mclust tended to allocate fewer observations to the noise cluster than NAvMix, and mclust using proportional associations, on average, allocated more observations to the noise cluster than NAvMix. A point of particular note is that the approaches which used mclust tended to find a number of clusters where there were no truly distinct clusters (that is, in the $K = 1$ scenarios), whereas NAvMix did not find spurious clusters in these null scenarios.

When incorporating trait correlation estimates, the median Rand index and its spread, as well as the mean number of clusters, were very similar to the results obtained without these estimates (Fig. S1 and Table S1). This suggests that the procedure is robust to missing trait correlation estimates. The results were also similar when estimating associations in separate samples (Fig. S2 and Table S2).

Clustering BMI associated genetic variants

We applied our procedure to cluster BMI associated genetic variants identified by the GWAS of Pulit et al. [16]. We considered genetic variants associated with BMI at a p-value $< 5 \times 10^{-8}$ and pruned at $r^2 < 0.001$. The clustering was performed in relation to the genetic associations with nine traits: body fat percentage; systolic blood pressure (SBP); triglycerides; high-density lipoprotein cholesterol (HDL); educational attainment; physical activity; lifetime smoking score; waist-to-hip ratio (WHR); and type 2 diabetes. These are lifestyle or cardiometabolic traits which have previously been shown to be related to BMI and which may offer insight into the pathways to downstream effects of BMI such as CHD [17, 18]. The genetic association estimates with these traits were all obtained from publicly available GWAS summary statistics (Methods). We clustered the 539 genetic variants that were available across all datasets. The full list of genetic variants and their allocated cluster, along with their probabilities of membership for each cluster, is given in Table S3.

Five clusters were identified, with 1 genetic variant allocated to the noise cluster. Fig. 3 shows a heat map of the proportional genetic association estimates with each trait by cluster and Fig. 4

plots the means of each fitted vMF distribution, representing the proportional associations for an observation at the centre of each cluster. The largest four clusters, labelled Clusters 1–4, contain genetic variants with very similar positive average proportional associations with fat percentage, WHR and type 2 diabetes. Variants in Cluster 3 have close to zero average association with SBP, whereas those in Clusters 1, 2, and 4 have positive average association with SBP. Variants in Cluster 2 have close to zero average association with smoking, whereas those in Clusters 1, 3 and 4 have positive average association with smoking. Variants in Cluster 4 have positive average association with HDL and negative average association with triglycerides, in contrast with those in Clusters 1–3.

Cluster 5 contains 20 genetic variants. These variants, on average, are positively associated with HDL and negatively associated with SBP, triglycerides, WHR and type 2 diabetes. These variants also have close to zero average association with smoking, physical activity and education, as well as weaker positive association with fat percentage compared with the other four clusters.

Mendelian randomization estimates of the effect of BMI on CHD

Mendelian randomization has previously suggested that BMI has a positive causal effect on CHD risk using as instruments 94 genetic variants identified by Locke et al. [19] [20]. We applied two-sample Mendelian randomization [21] using as instruments the set of BMI associated genetic variants which were used for clustering, as well as separately using the sets of variants for each cluster in turn (Methods). As well as applying the inverse-variance weighted (MR-IVW) method [22], we also performed as sensitivity analyses the MR-Median method [23], the Contamination Mixture (MR-ConMix) method [24] and the MR-PRESSO method [25]. Each of these methods provides a valid test for the causal null hypothesis under different sets of assumptions (Methods).

Fig. 5 shows scatterplots of the genetic association estimates with BMI against their association estimates with CHD risk for each set of instruments considered, as well the results of the Mendelian randomization analyses. When using the full set of genetic variants as instruments, the results suggest a positive effect of increased BMI on CHD risk, with an estimated odds ratio (OR) from MR-IVW of 1.50 (95 % confidence interval of 1.40–1.62) per 1 standard deviation increase in genetically predicted BMI. All sensitivity analyses gave similar estimates. This is in line with the results of Larsson et al. [20]. A similar result was obtained using the largest two clusters, with an estimated OR of 1.83 (1.68–2.00) using Cluster 1 and of 1.54 (1.38–1.72) using Cluster 2. When using the Cluster 3 genetic variants as instruments, the estimate attenuated toward the null, with an estimated OR of 1.22 (0.99–1.50). When using Cluster 4 genetic variants as instruments, there was no evidence that increased BMI is associated with CHD risk, with an estimated OR of 0.94 (0.69–1.29). When using Cluster 5 genetic variants as instruments, the results suggest a decrease in CHD risk from increased BMI, with an estimated OR of 0.34 (0.19–0.64). Note that the MR-Egger intercept test [26] did not show evidence of directional pleiotropy in any of these analyses (Table S4).

Exploring the biological pathways of clusters of BMI associated variants

We conducted gene set analysis on the BMI associated variants using the Functional Mapping and Annotation Platform [27] in order to examine the biological pathways relating to each cluster. The variants were mapped to genes based on positional and eQTL mappings, which were in turn tested for enrichment in gene sets from various pathway databases (Methods). A number of

distinct patterns emerge: Cluster 1 variants are associated with pathways related to cell division and differentiation; Cluster 3 variants with pathways related to cellular signalling; Cluster 4 variants with pathways related to lipid metabolism; and Cluster 5 variants with pathways related to inflammation. Cluster 2 variants were not found to be significantly enriched with any of the tested pathways. The full set of pathways associated with the mapped genes is given in Table S5.

The role of Cluster 5 variants in inflammation is of particular interest given its proposed relation to favourable adiposity. In order to confirm the role of these variants in inflammation, we conducted a Mendelian randomization analysis to examine the association of genetically predicted BMI, using all variants and each cluster separately, with C-reactive protein (CRP), a measure of systemic inflammation (Methods). The results from the MR-IVW method are shown in Fig. 6. When using all variants as instruments, MR-IVW estimated an increase in CRP of 0.44 standard deviations (95% confidence interval of 0.38–0.50) per standard deviation increase in genetically predicted BMI. The results when using Clusters 1–4 as instruments were in line with this. However, there was no evidence that the component of BMI predicted by Cluster 5 variants is associated with CRP (MR-IVW estimate of 0.01, 95% confidence interval of -0.24–0.27). These findings were supported in sensitivity analyses (see Fig. S3).

To further explore the pathways by which the various clusters affect inflammation, we performed separate Mendelian randomization analyses with the 41 cytokines and growth factors studied by Ahola-Olli et al. [28] and Kalaoja et al. [29] as outcomes (see Table S6 for the full list of cytokines and growth factors considered). Fig. 7 shows the MR-IVW estimates for each cluster and outcome. There was evidence of variation in the effects of BMI predicted by Cluster 5 variants on the cytokines compared with the effects of BMI predicted by the other clusters. For a number of inflammatory traits, such as hepatocyte growth factor (HGF) and TNF-related apoptosis inducing ligand (TRAIL), BMI predicted by Cluster 5 variants showed a weaker association than the other clusters. In some cases, such as for monocyte chemoattractant protein-1 (MCP1), the MR-IVW estimates using Cluster 5 variants were in the opposite direction to the other clusters.

Discussion

In this paper we have presented a procedure for clustering genetic variants based on their associations with a given set of traits using the NAVMix method. The method uses a directional clustering algorithm to distinguish between genetic variants based on their proportional associations with the traits. Since it is a model-based clustering approach, it has many advantages over current methods that are employed for clustering genetic variants based on trait associations, such as a data-driven method for choosing the number of clusters and the ability to use soft clustering. The inclusion of a noise cluster provides robustness to outliers, offering greater confidence in the identified clusters. A simulation study showed the method performs well in a range of settings, and that it outperformed alternative clustering approaches in assigning observations based on proportional associations. Importantly, the method did not identify false positive clusters in the simulation setting when no true clusters existed in the data, in contrast to the other methods considered.

The application to clustering BMI associated genetic variants identified five clusters, suggesting that genetic predictors of BMI can be broken down into five separate mechanisms based on their associations with the traits considered. Mendelian randomization analyses provided evidence that the different pathways affecting BMI have different downstream effects on CHD risk. When using as instruments the set of genetic variants in Clusters 1 and 2, the Mendelian randomization estimate of BMI on CHD risk was positive, in line with the established overall effect of increased BMI. When

using as instruments the set of variants in Cluster 3, the estimate was still positive but attenuated to the null. The main difference between this cluster and Clusters 1 and 2 is that the variants do not, on average, associate with increased SBP. Previous evidence suggests that increased SBP is a downstream consequence of increased BMI [30], and has also been shown to have a causal effect on CHD [24]. Our results therefore support that the genetically predicted component of BMI that does not associate with increased SBP has a lower positive effect on CHD risk. However, there is still evidence of a positive causal effect, suggesting there are other mechanisms by which increased BMI may increase CHD risk [31].

When using as instruments the set of genetic variants in Cluster 4, which have average associations with increased HDL and decreased triglycerides, Mendelian randomization suggested there was no association with CHD risk. Furthermore, the Mendelian randomization estimate of the component of BMI predicted by the variants in Cluster 5 was negative. That is, in Cluster 5, we have identified genetic variants related to a BMI increasing pathway that is protective of CHD. Orientating to the BMI-increasing alleles, these genetic variants are associated with a favourable metabolic profile, namely increased HDL and decreased SBP, triglycerides, WHR and type 2 diabetes liability.

By analysing the biological pathways underpinning the different clusters, we found evidence supporting that inflammation plays a key role in mediating the effect of obesity on cardiovascular risk. Furthermore, our findings identify possible inflammatory pathways related to elevated BMI that represent therapeutic targets for preventing CHD. Specifically, the estimated effects of Cluster 5 variants, in contrast to the BMI increasing variants more generally, are consistent with lower levels of key inflammatory cytokines implicated in CHD pathogenesis, including HGF [32], MCP1 [33] and TRAIL [34]. By ameliorating the increased inflammation attributable to elevated BMI, its detrimental effects on CHD risk may also be mitigated.

A number of studies have previously sought to identify genetic variants associated with metabolically favourable adiposity. Huang et al. [35] conducted pairwise significance tests between adiposity traits and various other cardiometabolic traits to identify genetic variants which, for at least one such pairing, associate with an increase in the adiposity trait and a decrease in the cardiometabolic trait. A similar approach to identifying genetic variants associated with favourable adiposity has also been performed by Yaghootkar et al. [36]. Our approach differs to these in that our clusters are formed without using genetic associations with the risk factor or outcome of interest, in this case BMI and CHD, but rather in relation to the chosen traits. Therefore, any difference between clusters in their associations with CHD risk is a meaningful statistical test, rather than a difference driven by the clustering algorithm.

The proposed approach has some limitations. It uses as input the full covariance matrix of the genetic variant-trait associations. An estimate of the full covariance matrix relies on estimates of the trait correlations, either from individual level data or from a reference dataset. Without this, the full covariance matrix requires either the assumption that the traits are uncorrelated or that the genetic variant-trait associations are estimated in separate samples. In practice, it is unlikely that the entire set of traits will be uncorrelated, since they would typically be related at least via common association with the primary trait of interest. However, the simulation study suggested the method is robust to ignoring the genetic variant-trait association correlations. This also suggests that the approach is robust to some participant overlap in the samples.

Another limitation is that the results are dependent on the choice of traits used to cluster on. Domain knowledge should be used to select a set of traits which are believed to be informative of potential mechanisms of the genetic variants under consideration. Future research will look to extend the method to include feature selection [37], so that the inclusion of a moderate to large

number of traits, many of which may not distinguish between clusters, is possible. It should be noted that adding highly correlated traits does not add much extra information, and may impact the results if correlation estimates are not incorporated. Thus, if there are a number of traits of interest which are highly correlated, it is better to chose just one of them.

In conclusion, we have presented a procedure for clustering genetic variants based on their direction of association with relevant traits, in order to gain insight into their underlying biological mechanisms and pathways. We have demonstrated the utility of clustering genetic variants in this way by applying the method to BMI associated genetic variants and performing Mendelian randomization analyses to infer the differential effects of distinct BMI increasing pathways on CHD risk.

Methods

The von Mises–Fisher distribution

The m -dimensional von Mises–Fisher (vMF) distribution has probability density function

$$f(x | \mu, \kappa) = C_m(\kappa) e^{\kappa \mu' x},$$

where $\|x\| = \|\mu\| = 1$ and $C_m(\kappa)$ is a normalising constant given by

$$C_\nu(x) = \frac{x^{\nu/2-1}}{(2\pi)^{\nu/2} I_{\nu/2-1}(x)},$$

where $I_\nu(x)$ is the modified Bessel function of the first kind and order ν [38, 11]. The mean parameter μ is a unit vector which represents the direction from the origin in m -dimensional space. The concentration parameter κ represents the spread of observations around the mean. When $\kappa = 0$, the distribution is the uniform distribution on the $(m - 1)$ -dimensional unit sphere. As κ increases, the distribution becomes increasingly focused around the point on the unit sphere given by μ .

The noise-augmented von Mises–Fisher mixture model

Suppose we have m -dimensional observations $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ where $\|\mathbf{x}_j\| = 1$ for all i (if the observations are not normalised to have magnitude 1, then this normalisation is the first step in the procedure). Here, x_j represents the vector of proportional association estimates for genetic variant j with the m traits. Further suppose that each observation either belongs to one of K clusters, each cluster containing observations from a vMF distribution, or else belongs to none of these clusters and is therefore considered noise. We can represent this with the $K + 1$ component vMF mixture model given by

$$p(\mathbf{x}_j | \Theta) = \sum_{k=1}^{K+1} p(\mathbf{x}_j, z_j = k | \boldsymbol{\mu}_k, \kappa_k) = \sum_{k=1}^{K+1} \pi_k f(\mathbf{x}_j | \boldsymbol{\mu}_k, \kappa_k)$$

for the j th observation, where:

- $\Theta = \{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \kappa_1, \dots, \kappa_K, \pi_1, \dots, \pi_{K+1}\};$

- $\mathbf{z} = \{z_1, \dots, z_n\}$ denotes cluster membership (that is, $z_j = k$ if \mathbf{x}_j belongs to cluster k);
- π_k is the mixing proportion of cluster k , with $\sum_{k=1}^{K+1} \pi_k = 1$;
- $f(\mathbf{x} | \boldsymbol{\mu}, \kappa)$ is the density function of the m -dimensional vMF distribution;
- $\boldsymbol{\mu}_{K+1}$ is the unit vector which is fixed according to the global sample mean direction, given by

$$\boldsymbol{\mu}_{K+1} = \frac{\sum_{j=1}^n \mathbf{x}_j}{\left\| \sum_{j=1}^n \mathbf{x}_j \right\|};$$

- κ_{K+1} is fixed at a number close to zero (for example 0.0001).

In this model, cluster $K + 1$ is referred to as the noise cluster. With κ close to zero, the distribution function represents the uniform distribution on the $(m - 1)$ -dimensional unit sphere, and so observations which do not fit well to the other K clusters will tend to be assigned here.

The log-likelihood function is

$$l_K(\boldsymbol{\Theta}) = \sum_{j=1}^n \log \left\{ \sum_{k=1}^{K+1} \pi_k f(\mathbf{x}_j | \boldsymbol{\mu}_k, \kappa_k) \right\}.$$

In order to maximise the likelihood function to obtain estimates of the parameters $\boldsymbol{\Theta}$, we would require knowledge of the latent variables \mathbf{z} . Mixture models of this sort are thus fitted using the EM algorithm [39].

The EM algorithm

Suppose we have an estimate of $\boldsymbol{\Theta}$, denoted by $\hat{\boldsymbol{\Theta}}$. Let $Q(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}) = E_{\mathbf{z} | X, \hat{\boldsymbol{\Theta}}} l_K(\boldsymbol{\Theta})$. Then

$$Q(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}}) = \sum_{j=1}^n \sum_{k=1}^{K+1} \gamma_{jk} \log \{ \pi_k f(\mathbf{x}_j | \boldsymbol{\mu}_k, \kappa_k) \},$$

where

$$\gamma_{jk} = \Pr(z_j = k | \mathbf{x}_j, \hat{\boldsymbol{\Theta}}) = \frac{\pi_k f(\mathbf{x}_j | \boldsymbol{\mu}_k, \kappa_k)}{\sum_{l=1}^{K+1} \pi_l f(\mathbf{x}_j | \boldsymbol{\mu}_l, \kappa_l)}, \quad k = 1, \dots, K + 1.$$

Computing the γ_{jk} for a given $\hat{\boldsymbol{\Theta}}$ is the E step in the EM algorithm.

Given the γ_{jk} , we can estimate $\boldsymbol{\Theta}$ by maximising $Q(\boldsymbol{\Theta} | \hat{\boldsymbol{\Theta}})$. Following Banerjee et al. [11], the parameter estimates are obtained from

$$\hat{\boldsymbol{\mu}}_k = \frac{\sum_{j=1}^n \gamma_{jk} \mathbf{x}_j}{\left\| \sum_{j=1}^n \gamma_{jk} \mathbf{x}_j \right\|}, \quad k = 1, \dots, K,$$

$$\frac{I_{m/2}(\hat{\kappa}_k)}{I_{m/2-1}(\hat{\kappa}_k)} = \frac{\left\| \sum_{j=1}^n \gamma_{jk} \mathbf{x}_j \right\|}{\left\| \sum_{j=1}^n \gamma_{jk} \right\|}, \quad k = 1, \dots, K \quad (1)$$

$$\hat{\pi}_k = \frac{1}{n} \sum_{j=1}^n \gamma_{jk}, \quad k = 1, \dots, K + 1.$$

This is the M step of the EM algorithm. Note that we do not update the noise cluster parameters, $\boldsymbol{\mu}_{K+1}$ and κ_{K+1} , but we do update the proportion of observations which are assigned to the noise cluster, $\hat{\pi}_{K+1}$. Now, (1) does not give a closed form solution for computing $\hat{\kappa}_k$. However, a number of methods for approximating these solutions have been proposed which allow the concentration parameter estimates to be easily updated. Banerjee et al. [11] proposed the approximation

$$\hat{\kappa}_k = \frac{\bar{r}_k m - \bar{r}_k^3}{1 - \bar{r}_k^2},$$

where

$$\bar{r}_k = \frac{\left\| \sum_{j=1}^n \gamma_{jk} \mathbf{x}_j \right\|}{\left\| \sum_{j=1}^n \gamma_{jk} \right\|}.$$

Hornik and Grün [14] summarise several other approximation methods and provide software for implementing each of them. Note that, in practice, values of \bar{r} very close to 1 can cause numerical problems, (due to the fact that this relates to the case where the observations are almost all at the same point, and the precision is thus close to infinity). To get around this, we cap the value that $\hat{\kappa}_k$ can take at 500.

The EM algorithm can be started at either the E step, given an initial estimate of $\boldsymbol{\Theta}$, or at the M step, given initial values of the γ_{jk} . The algorithm is iterated until the absolute value of the difference between successive values of $l_K(\hat{\boldsymbol{\Theta}})$ is less than some predefined convergence threshold. In our simulation study and applied example, we used 10^{-4} as the convergence threshold.

Initialisation of the algorithm

In order to initialise the algorithm, we must first set an initial proportion of observations which belong in the noise cluster, which we will denote by $0 < \hat{\pi}_{K+1}^{(0)} < 1$. We then perform the spherical k-means procedure [13], which clusters observations based on similarity of their direction from the origin, analogous to the k-means procedure which clusters observations based on Euclidean distance. We take as initial values, for $i = 1, \dots, n$,

$$\gamma_{ik} = \begin{cases} 1 - \hat{\pi}_{K+1}^{(0)}, & \text{if observation } i \text{ is assigned to cluster } k \\ 0, & \text{otherwise} \end{cases}, \quad k = 1, \dots, K$$

$$\gamma_{i(K+1)} = \hat{\pi}_{K+1}^{(0)}.$$

We then begin the EM algorithm at the M step. Note that the spherical k-means procedure relies on an initial random set of cluster means, and thus its results are sensitive to this randomisation. There is a possibility that certain initial values from the procedure will result in the EM algorithm converging to a local, rather than global, maximum. We therefore run the algorithm a number of times in practice, each time beginning with different initial values. We take as final parameter estimates those which result in the EM algorithm converging to the greatest maximum. In our simulation study and applied example, we ran the algorithm with 5 different initialisations.

Choosing the number of clusters

In practice, we will not know the number of clusters to fit to the data. The number of clusters can be determined using an information criterion, for example BIC [40, 41]. For successive values of K , we perform the algorithm above and compute

$$\phi_m(K) = -2l_K(\hat{\Theta}) + r_m(K) \log(n),$$

where $r_m(K) = (m+2)K + m$ is the number of parameters estimated. We continue until $\phi_m(K)$ increases for successive iterations. The final number of clusters is then taken to be $\arg \min_K \phi_m(K)$.

Assigning cluster membership

Output from the procedure for fitting the mixture model is a set of probabilities for each observation belonging to each cluster (that is, the γ_{ik} parameters). The simplest approach for assigning cluster membership is to assign each observation to the cluster for which it has the greatest probability of membership (that is, $\hat{z}_i = \arg \max_k \gamma_{ik}$). This is the approach used in both the simulation study and the applied example presented in this paper.

Mixture model approaches to clustering allow for flexibility in the way that cluster membership is assigned. For increased confidence in the clusters, a threshold could be set such that an observation is only assigned to a cluster if the probability of membership is greater than a certain level. Those which do not meet the threshold for any cluster remain unassigned. Finally, soft clustering is possible, whereby observations are assigned to any cluster for which its probability of membership is greater than a certain level. Under the soft clustering approach, an observation may be assigned to more than one cluster.

Genetic variant-trait association covariance matrix

For variant j , the (k, l) th element of $\hat{\Sigma}_j$ is given by

$$\text{se}(\hat{\beta}_{jk}) \text{se}(\hat{\beta}_{jl}) \text{cor}(\hat{\beta}_{jk}, \hat{\beta}_{jl}),$$

where $\text{se}(\hat{\beta}_{jk})$ is the standard error of $\hat{\beta}_{jk}$. If the genetic variant-trait associations are estimated in separate, non-overlapping, samples, then $\text{cor}(\hat{\beta}_{jk}, \hat{\beta}_{jl}) = 0$ and $\hat{\Sigma}_j$ can be taken to be the diagonal matrix with k th diagonal entry equal to $\text{se}^2(\hat{\beta}_{jk})$. If the traits are estimated in the same sample, then the off-diagonal entries of $\hat{\Sigma}_j$ will be non-zero. Although the correlation between $\hat{\beta}_{jk}$ and $\hat{\beta}_{jl}$ is not easily estimated, provided the j th genetic variant explains only a small proportion of the variance in the k th and l th traits, then $\text{cor}(\hat{\beta}_{jk}, \hat{\beta}_{jl}) \approx \text{cor}(X_k, X_l)$ [42]. We can therefore compute the (k, l) th entry of $\hat{\Sigma}_j$, $i \neq j$, by

$$\text{se}(\hat{\beta}_{jk}) \text{se}(\hat{\beta}_{jl}) \widehat{\text{cor}}(X_k, X_l),$$

where $\widehat{\text{cor}}(X_k, X_l)$ is an estimate of the correlation between the k th and l th traits. As a result of this, if the traits are assumed to be independent, then the off-diagonal entries of $\hat{\Sigma}_j$ can be approximated by zeros, and the covariance matrix taken to be diagonal as in the separate samples case.

Simulation study

We simulated $n = 120$ independent genetic variants for N_l individuals, denoted G_{ij} for individual i and genetic variant j , and m traits, denoted X_{ik} for individual i and trait l , from the following model

$$\begin{aligned} G_{ij} &\sim \text{Binomial}(2, \text{maf}_j) \\ U_i, \varepsilon_{i1}, \dots, \varepsilon_{im} &\sim N(0, 1), \text{ independently} \\ X_{il} &= \sum_{j=1}^n \beta_{jl} G_{ij} + \gamma_l U_i + \sqrt{1 - \gamma_l^2} \varepsilon_{il}, \end{aligned}$$

for $i = 1, \dots, N_k$ and $l = 1, \dots, m$. The common variable U_i is included to induce correlation between the errors, while maintaining the same amount of variation explained by the genetic variants. The γ_l values, which determine the magnitude and direction of correlation between the traits, were generated from the Uniform $(-0.8, 0.8)$ distribution. The minor allele frequencies (denoted in the model by maf_j) were generated from the Uniform $(0.01, 0.5)$ distribution. The number of traits was either $m = 2$ or 9 .

The first 100 genetic variant-trait associations were split into 1, 2 or 4 clusters. For the 2 cluster scenarios, each cluster contained 50 genetic variants. For the 4 cluster scenarios, the first cluster contained 40 genetic variants, and the other three clusters contained 20 genetic variants. The β_{j1} values were generated from the Uniform $(0.05, 0.4)$ distribution for the first two clusters, and from the Uniform $(-0.4, 0.05)$ distribution for the second two clusters. For each $l = 2, \dots, m$, the β_{jl} values were generated by $\tan(v_j) \beta_{j1}$. The v_j values were generated from the truncated normal distribution with mean δ_k , variance 0.2^2 and truncation points $(\delta_k - 0.2, \delta_k + 0.2)$, where the δ_k values are as follows. For the 1 cluster scenarios, $\delta_1 = \pi/4$, $m = 2$, and

$$\delta_1 = (\pi/4, \pi/4, \pi/4, \pi/4, \pi/4, \pi/4, \pi/4, \pi/4),$$

$m = 9$. For the 2 cluster scenarios, $\delta_1 = \pi/4$ and $\delta_2 = -\pi/4$, $m = 2$, and

$$\begin{aligned} \delta_1 &= (\pi/4, \pi/4, \pi/8, \pi/8, \pi/8, \pi/8, \pi/8, \pi/8) \\ \delta_2 &= (-\pi/4, -\pi/4, \pi/8, \pi/8, \pi/8, \pi/8, -\pi/8, -\pi/8), \end{aligned}$$

$m = 9$. For the 4 cluster scenarios, $\delta_1 = \delta_3 = \pi/4$ and $\delta_2 = \delta_4 = -\pi/4$, $m = 2$, and

$$\begin{aligned} \delta_1 &= (\pi/4, \pi/4, \pi/8, \pi/8, \pi/8, \pi/8, \pi/8, \pi/8) \\ \delta_2 &= (\pi/4, \pi/4, \pi/4, \pi/4, 0, 0, 0, 0) \\ \delta_3 &= (-\pi/4, -\pi/4, \pi/8, \pi/8, \pi/8, \pi/8, -\pi/8, -\pi/8), \\ \delta_4 &= (-\pi/4, -\pi/4, -\pi/4, -\pi/4, 0, 0, 0, 0), \end{aligned}$$

$m = 9$. The l th element of δ_k represents the mean angle from the origin, in radians, of the values (β_{j1}, β_{jl}) , $j = 1, \dots, 100$. The effect of this setup is that the angle from the origin of the pair (β_{j1}, β_{jl}) is centered around a cluster specific mean and bounded according to the truncation points. Each cluster can thus be considered to contain genetic variants acting in similar directions on the given traits. The final 20 genetic variants were simulated to represent noise. For these genetic variants, the β_{jl} values, $l = 1, \dots, m$, were generated from a normal distribution with mean 0 and variance 0.2^2 .

The estimated genetic variant-trait associations were computed using simple linear regression of each trait on each genetic variant in turn. For the primary simulation study presented in the Results section, the associations were estimated in the same sample of size 20 000. The Supplementary Information also presents results for scenarios where the associations were estimated in separate samples with sizes varying between 20 000 and 44 000.

The resulting datasets were clustered using NAvMix with an initial proportion of genetic variants in the noise cluster of 0.05, and using mclust with an initial noise cluster of 5 randomly selected genetic variants.

Clustering BMI associated genetic variants

Genetic variant association estimates with BMI were taken from the GWAS of Pulit et al. [16]. Variants with p-value $< 5 \times 10^{-8}$ were pruned using the TwoSampleMR package in R [43] with $r^2 = 0.001$.

Genetic variant association estimates with body fat percentage, SBP, triglycerides and HDL were taken from results from the Neale Lab which are based on the UK Biobank dataset (<http://www.nealelab.is/uk-biobank/>). Genetic variant associations for educational attainment were taken from the GWAS of Okbay et al. [44]; for physical activity, the GWAS of Doherty et al. [45]; for lifetime smoking score, the GWAS of Wootton et al. [46]; for WHR the GWAS of Pulit et al. [16]; and for type 2 diabetes, the GWAS of Mahajan et al. [6]. Note that for the educational attainment dataset, one BMI associated genetic variant (rs10761785) was replaced with a proxy (rs2163188) with $r^2 = 0.9842$ (identified using PhenoScanner [47, 48]). All studies used were performed on samples of individuals of European ancestry or predominantly European ancestry. All genetic variant trait-association estimates were orientated with respect to the alleles such that the associations with BMI were positive.

Clustering was performed using NAvMix with an initial proportion of genetic variants in the noise cluster of 0.05, and 5 separate initialisations of the algorithm was used. The probability of membership of each genetic variant to each cluster produced by the algorithm is shown in Table S3.

Mendelian randomization analyses

A genetic variant is a valid instrumental variable for a Mendelian randomization analysis if it is: associated with the risk factor; independent of any confounders of the risk factor-outcome relationship; and has no causal pathway to the outcome other than via the risk factor [49]. Under the two-sample framework, the genetic variant-risk factor and genetic variant-outcome associations are estimated in separate samples [21]. Under the assumption that all variants in the analysis are valid instruments, MR-IVW produces a statistically consistent estimator of the causal effect and a test for the causal null hypothesis [22]. The three methods used for sensitivity analyses were chosen since they each produce a valid estimate of the causal effect of BMI on CHD under different assumptions [50]: MR-Median (a majority of the genetic variants are valid instrument); the Contamination Mixture method (a plurality of the genetic variants are valid instruments); and the MR-PRESSO method (the InSIDE assumption is met). The intercept test from the MR-Egger method was used to test for the presence of unmeasured directional pleiotropy. Analyses were carried out using the MendelianRandomization [51, 52] and MRPRESSO [25] packages.

Genetic variant association estimates with CHD were taken from the CARDIoGRAMplusC4D

dataset of Nikpay et al. [53] and accessed using PhenoScanner [47, 48]. Genetic variant associations with CRP were taken from results from the Neale Lab which are based on the UK Biobank dataset (<http://www.nealelab.is/uk-biobank/>). Genetic variant association estimates with the 41 cytokines and growth factors were taken from the data supporting Ahola-Olli et al. [28] and Kalaoja et al. [29]. Table S7 gives a list of the BMI associated genetic variants which were not available in each of the outcome datasets and were therefore excluded from the relevant Mendelian randomization analyses.

Gene mapping and gene set analysis

The 539 BMI associated genetic variants were mapped to genes using the SNP2GENE function in FUMA [27]. Summary statistics for each cluster of variants were uploaded separately, and were identified as pre-defined lead SNPs. Both positional and eQTL mapping was performed. For the eQTL mapping, tissue types were selected as all those from the following sources: eQTL catalogue; PsychENCODE; van der Wijst et al. scRNA eQTLs; DICE; eQTLGen; Blood eQTLs; MuTHER; xQTLServer; ComminMind Consortium; BRAINEAC; and GTEx v8. All other default settings were used. Gene set analysis was performed using the GENE2FUNC function. The results presented in Table S5 include all canonical pathways from MsigDB which associate with the mapped genes using hypergeometric tests (with multiple test correction applied per cluster).

Data availability

All the data used in this paper are publicly available and can be accessed via the references given.

Code availability

R code for performing the NAvMix clustering algorithm, and for reproducing the simulation results and applied analysis, can be found at <https://github.com/aj-grant/navmix>.

Acknowledgments

A.J.G. and S.B. are supported by a Sir Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal Society (grant number 204623/Z/16/Z). D.G. is supported by the British Heart Foundation Research Centre of Excellence (RE/18/4/34215) at Imperial College London and a National Institute for Health Research Clinical Lectureship (CL-2020-16-001) at St. George's, University of London. P.D.W.K. is supported by the UK Medical Research Council (MC_UU_00002/13). This research was funded by the NIHR Cambridge Biomedical Research Centre. The views expressed are those of the authors and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care. For the purpose of open access, the author has applied a CC-BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

Competing interests statement

Dipender Gill is employed part-time by Novo Nordisk. The other authors declare no competing interests.

References

- [1] Visscher, P. M., et al. 10 years of GWAS discovery: Biology, function, and translation. *Am J Hum Genet* **101**, 5 – 22 (2017).
- [2] Winkler, T. W., et al. A joint view on genetic variants for adiposity differentiates subtypes with distinct metabolic implications. *Nat Commun* **9**, 1946 (2018).
- [3] Udler, M. S., et al. Type 2 diabetes genetic loci informed by multi-trait associations point to disease mechanisms and subtypes: A soft clustering analysis. *PLoS Med* **15**, 1–23 (2018).
- [4] Dimas, A. S., et al. Impact of type 2 diabetes susceptibility variants on quantitative glycemic traits reveals mechanistic heterogeneity. *Diabetes* **63**, 2158–2171 (2014).
- [5] Scott, R. A., et al. An expanded genome-wide association study of type 2 diabetes in Europeans. *Diabetes* **66**, 2888–2902 (2017).
- [6] Mahajan, A., et al. Refining the accuracy of validated target identification through coding variant fine-mapping in type 2 diabetes. *Nat Genet* **50**, 559–571 (2018).
- [7] Ruth, K. S., et al. Using human genetics to understand the disease impacts of testosterone in men and women. *Nat Med* **26**, 252–258 (2020).
- [8] Tanigawa, Y., et al. Components of genetic associations across 2,138 phenotypes in the UK Biobank highlight adipocyte biology. *Nat Commun* **10**, 4064 (2019).
- [9] Davey Smith, G. and Ebrahim, S. ‘Mendelian randomization’: can genetic epidemiology contribute to understanding environmental determinants of disease? *Int J Epidemiol* **32**, 1–22 (2003).
- [10] Lawlor, D. A., Harbord, R. M., Sterne, J. A. C., Timpson, N., and Davey Smith, G. Mendelian randomization: Using genes as instruments for making causal inferences in epidemiology. *Stat Med* **27**, 1133–1163 (2008).
- [11] Banerjee, A., Dhillon, I. S., Ghosh, J., and Sra, S. Clustering on the unit hypersphere using von Mises-Fisher distributions. *J Mach Learn Res* **6**, 1345–1382 (2005).
- [12] Scrucca, L., Fop, M., Murphy, T. B., and Raftery, A. E. mclust 5: Clustering, classification and density estimation using Gaussian finite mixture models. *R J* **8**, 289–317 (2016).
- [13] Dhillon, I. S. and Modha, D. S. Concept decompositions for large sparse text data using clustering. *Mach Learn* **42**, 143–175 (2001).
- [14] Hornik, K. and Grün, B. movmf: An R package for fitting mixtures of von Mises-Fisher distributions. *J Stat Softw* **58**, 1–31 (2014).

- [15] Rand, W. M. Objective criteria for the evaluation of clustering methods. *J Am Stat Assoc* **66**, 846–850 (1971).
- [16] Pulit, S. L., et al. Meta-analysis of genome-wide association studies for body fat distribution in 694 649 individuals of European ancestry. *Hum Mol Genet* **28**, 166–174 (2019).
- [17] Van Gaal, L. F., Mertens, I. L., and De Block, C. E. Mechanisms linking obesity with cardiovascular disease. *Nature* **444**, 875–880 (2006).
- [18] Davies, N. M., Dickson, M., Davey Smith, G., van den Berg, G. J., and Windmeijer, F. The causal effects of education on health outcomes in the UK Biobank. *Nat Hum Behav* **2**, 117–125 (2018).
- [19] Locke, A. E., et al. Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197–206 (2015).
- [20] Larsson, S. C., Bäck, M., Rees, J. M. B., Mason, A. M., and Burgess, S. Body mass index and body composition in relation to 14 cardiovascular conditions in UK Biobank: a Mendelian randomization study. *Eur Heart J* **41**, 221–226 (2019).
- [21] Burgess, S., et al. Using published data in Mendelian randomization: a blueprint for efficient identification of causal risk factors. *Eur J Epidemiol* **30**, 543–552 (2015).
- [22] Burgess, S., Butterworth, A., and Thompson, S. G. Mendelian randomization analysis with multiple genetic variants using summarized data. *Genet Epidemiol* **37**, 658–665 (2013).
- [23] Bowden, J., Davey Smith, G., Haycock, P. C., and Burgess, S. Consistent estimation in Mendelian randomization with some invalid instruments using a weighted median estimator. *Genet Epidemiol* **40**, 304–314 (2016).
- [24] Burgess, S., Foley, C. N., Allara, E., Staley, J. R., and Howson, J. M. M. A robust and efficient method for Mendelian randomization with hundreds of genetic variants. *Nat Commun* **11**, 376 (2020).
- [25] Verbanck, M., Chen, C.-Y., Neale, B., and Do, R. Detection of widespread horizontal pleiotropy in causal relationships inferred from Mendelian randomization between complex traits and diseases. *Nat Genet* **50**, 693–698 (2018).
- [26] Bowden, J., Davey Smith, G., and Burgess, S. Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *Int J Epidemiol* **44**, 512–525 (2015).
- [27] Watanabe, K., Taskesen, E., van Bochoven, A., and Posthuma, D. Functional mapping and annotation of genetic associations with FUMA. *Nat Commun* **8**, 1826 (2017).
- [28] Ahola-Olli, A. V., et al. Genome-wide association study identifies 27 loci influencing concentrations of circulating cytokines and growth factors. *Am J Hum Genet* **100**, 40–50 (2017).
- [29] Kalaoja, M., et al. The role of inflammatory cytokines as intermediates in the pathway from increased adiposity to disease. *Obesity* **29**, 428–437 (2021).

- [30] Marini, S., et al. Mendelian randomization study of obesity and cerebrovascular disease. *Ann Neurol* **87**, 516–524 (2020).
- [31] Gill, D., et al. Risk factors mediating the effect of body-mass index and waist-to-hip ratio on cardiovascular outcomes: Mendelian randomization analysis. Preprint at <https://www.medrxiv.org/content/10.1101/2020.07.15.20154096v1> (2020).
- [32] Morishita, R., Aoki, M., Yo, Y., and Ogihara, T. Hepatocyte growth factor as cardiovascular hormone: Role of HGF in the pathogenesis of cardiovascular disease. *Endocr J* **49**, 273–284 (2002).
- [33] Georgakakis, M. K., et al. Genetically determined levels of circulating cytokines and risk of stroke. *Circulation* **139**, 256–268 (2019).
- [34] Bernardi, S., Bossi, F., Toffoli, B., and Fabris, B. Roles and clinical applications of OPG and TRAIL as biomarkers in cardiovascular disease. *BioMed Res Int* **2016**, 1752854 (2016).
- [35] Huang, L. O., et al. Genome-wide discovery of genetic loci that uncouple excess adiposity from its comorbidities. *Nat Metab* **3**, 228–243 (2021).
- [36] Yaghootkar, H., et al. Genetic evidence for a link between favorable adiposity and lower risk of type 2 diabetes, hypertension, and heart disease. *Diabetes* **65**, 2448–2460 (2016).
- [37] Law, M. H., Jain, A. K., and Figueiredo, M. A. T. Feature selection in mixture-based clustering. In *Adv Neural Inf Process Syst* **15**, 641–648 (2003).
- [38] Mardia, K. V. and Jupp, P. *Directional Statistics*. John Wiley & Sons Chichester (2000).
- [39] Dempster, A. P., Laird, N. M., and Rubin, D. B. Maximum likelihood from incomplete data via the EM algorithm. *J R Stat Soc Series B Stat Methodol* **39**, 1–22 (1977).
- [40] Schwarz, G. Estimating the dimension of a model. *Ann Stat* **6**, 461–464 (1978).
- [41] Banfield, J. D. and Raftery, A. E. Model-based Gaussian and non-Gaussian clustering. *Biometrics* **49**, 803–821 (1993).
- [42] Sanderson, E., Spiller, W., and Bowden, J. Testing and correcting for weak and pleiotropic instruments in two-sample multivariable Mendelian randomisation. Preprint at <https://www.biorxiv.org/content/10.1101/2020.04.02.021980v1> (2020).
- [43] Hemani, G., et al. The MR-Base platform supports systematic causal inference across the human phenome. *eLife* **7**, e34408 (2018).
- [44] Okbay, A., et al. Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* **533**, 539–542 (2016).
- [45] Doherty, A., et al. GWAS identifies 14 loci for device-measured physical activity and sleep duration. *Nat Commun* **9**, 5257 (2018).
- [46] Wootton, R. E., et al. Evidence for causal effects of lifetime smoking on risk for depression and schizophrenia: a Mendelian randomisation study. *Psychol Med* **50**, 2435–2443 (2020).

- [47] Staley, J. R., Blackshaw, J., Kamat, M. A., et al. PhenoScanner: a database of human genotype–phenotype associations. *Bioinformatics* **32**, 3207–3209 (2016).
- [48] Kamat, M. A., et al. PhenoScanner V2: an expanded tool for searching human genotype–phenotype associations. *Bioinformatics* **35**, 4851–4853 (2019).
- [49] Greenland, S. An introduction to instrumental variables for epidemiologists. *Int J Epidemiol* **29**, 722–729 (2000).
- [50] Slob, E. A. W. and Burgess, S. A comparison of robust Mendelian randomization methods using summary data. *Genet Epidemiol* **44**, 313–329 (2020).
- [51] Yavorska, O. O. and Burgess, S. MendelianRandomization: an R package for performing Mendelian randomization analyses using summarized data. *Int J Epidemiol* **46**, 1734–1739 (2017).
- [52] Broadbent, J. R., et al. MendelianRandomization v0.5.0: updates to an R package for performing Mendelian randomization analyses using summarized data [version 2; peer review: 1 approved, 2 approved with reservations]. *Wellcome Open Res* **5** (2020).
- [53] Nikpay, M., et al. A comprehensive 1000 genomes–based genome-wide association meta-analysis of coronary artery disease. *Nat Genet* **47**, 1121–1130 (2015).

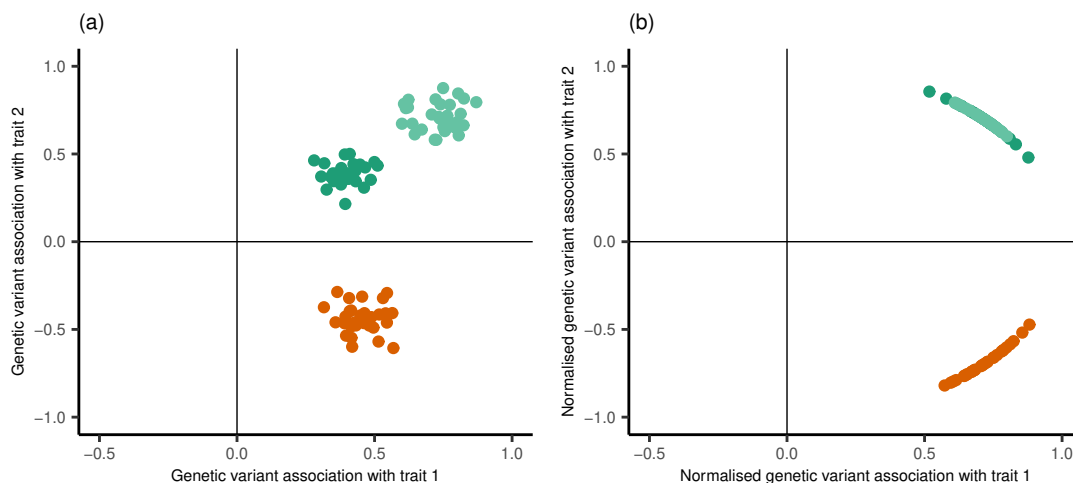


Figure 1: Illustrative figure showing the difference between clustering based on Euclidean distance compared with direction. Panel (a) plots 90 simulated points representing genetic associations with two traits. Each point was generated from one of three bivariate normal distributions. Panel (b) plots the normalised genetic associations, representing the proportional association of each genetic variant with respect to the two traits. All points sit on the unit circle. The green points represent genetic variants which are positively associated with each trait by similar magnitudes. The orange points represent genetic variants which are positively associated with trait 1 and negatively associated with trait 2, again by similar magnitudes. Methods based on Euclidean distance such as Gaussian mixture models and hierarchical clustering would consider there to be three clusters, distinguishing between the light and dark green points, as shown in Panel (a). Directional clustering approaches would consider there to be two clusters, grouping the green points in the same cluster. This is shown in Panel (b), where the points are clearly grouped in two separate clusters.

Table 1: Mean number of clusters estimated and mean number of observations allocated to the noise cluster for each simulated scenario using NAvMix, mclust, and mclust using proportional associations (pr).

Number of traits (m)	Number of clusters (K)	Number of clusters			Number of noise variants		
		NAvMix	mclust	mclust (pr)	NAvMix	mclust	mclust (pr)
2	1	1.00	1.93	3.73	18.99	15.05	19.00
	2	2.00	2.10	4.62	17.36	12.26	17.60
	4	4.00	2.15	6.13	14.24	14.49	15.11
9	1	1.01	2.18	1.04	21.72	18.06	23.01
	2	2.00	2.82	2.00	23.42	17.67	24.83
	4	4.45	4.21	4.93	21.20	18.21	23.39

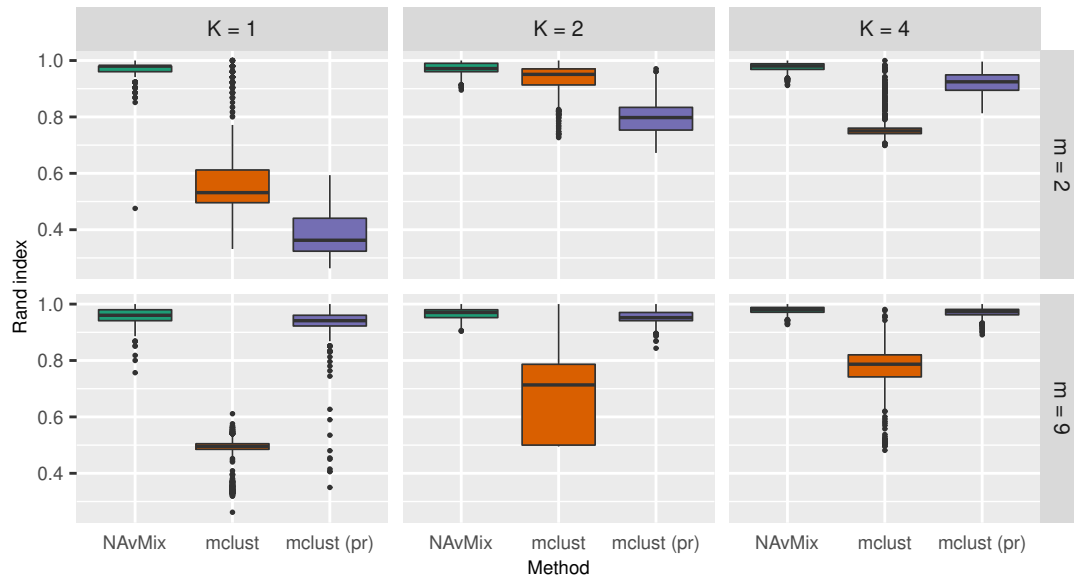


Figure 2: Boxplots of the Rand index for each scenario using NAvMix, mclust, and mclust using proportional associations (pr).

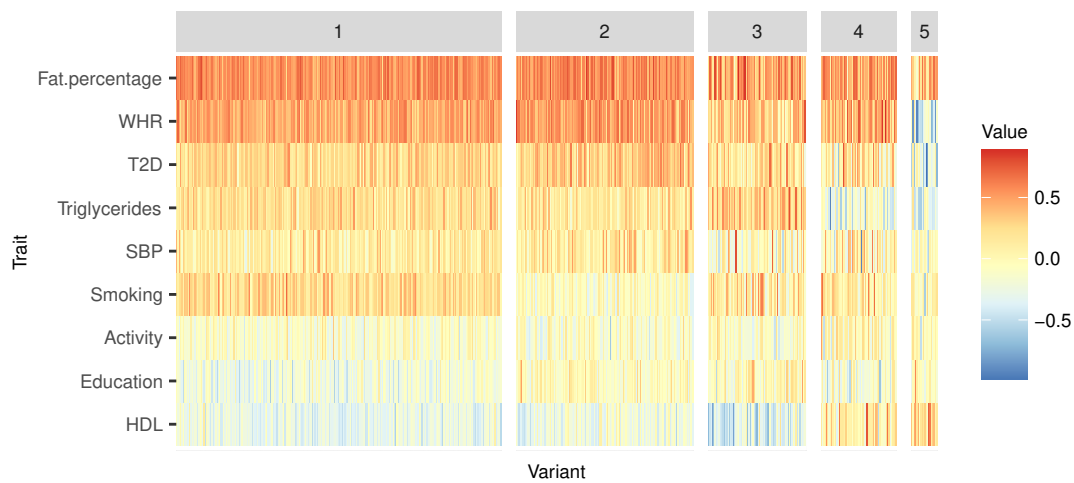


Figure 3: Heat map showing the association estimates of the genetic variants with each trait by cluster, excluding the noise cluster. The association estimates were first standardised by dividing by their standard errors, then normalised so that the vectors of association estimates for each variant have magnitude one. Thus, the values shown represent the proportional association estimates for each genetic variant on the set of traits.

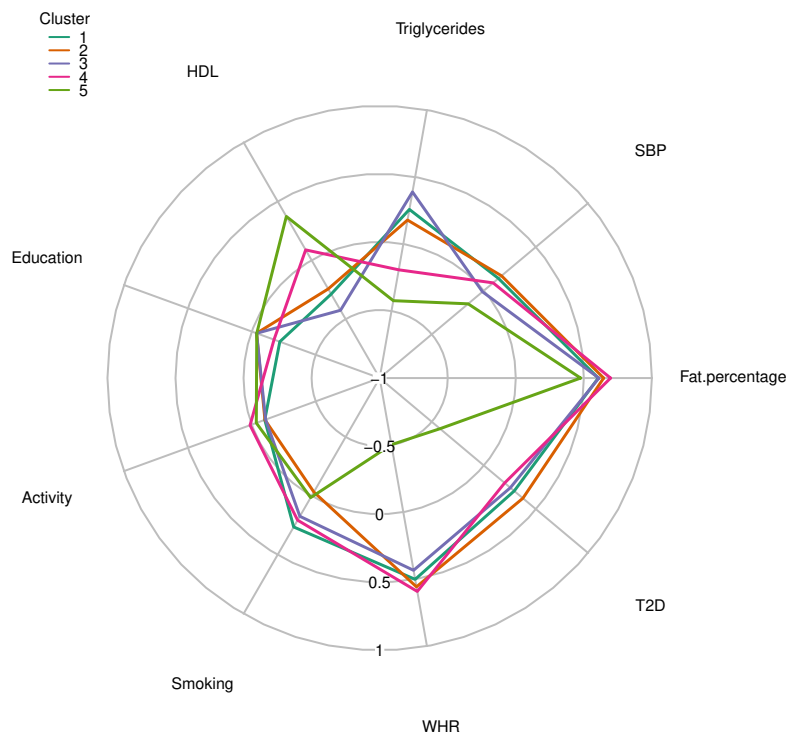


Figure 4: Radial plot of the mean vector of the fitted von Mises–Fisher distribution for each cluster. The plotted points represent the standardised proportional association with each trait for an observation at the centre of each cluster.

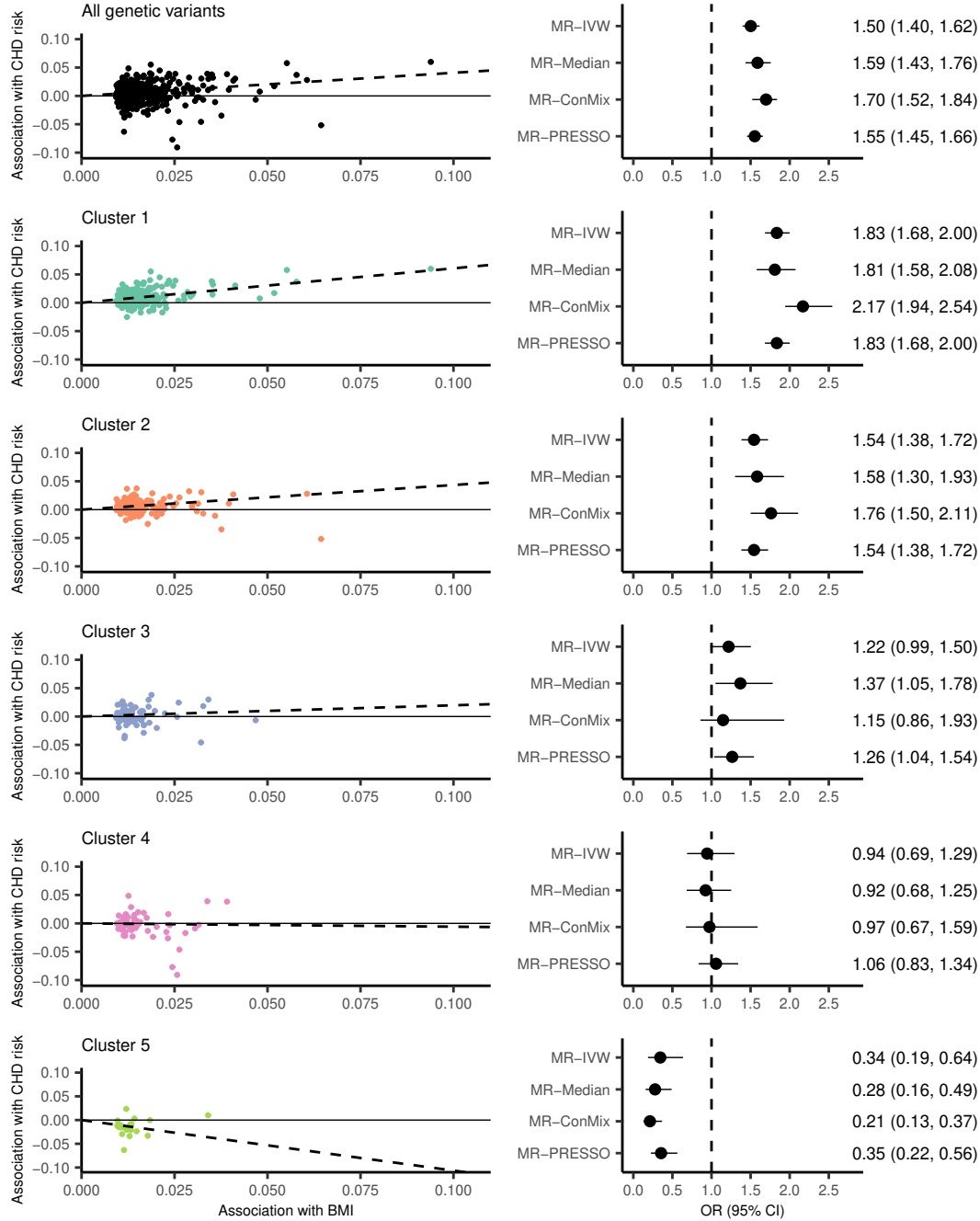


Figure 5: Scatterplots of the associations of each genetic variant with BMI (standard deviation units) and the log odds ratio of CHD risk, and forest plots showing estimates and 95% confidence intervals from Mendelian randomization, for all genetic variants and for each cluster. Mendelian randomization estimates represent the change in odds ratio of CHD risk per 1 standard deviation increase in genetically predicted BMI.

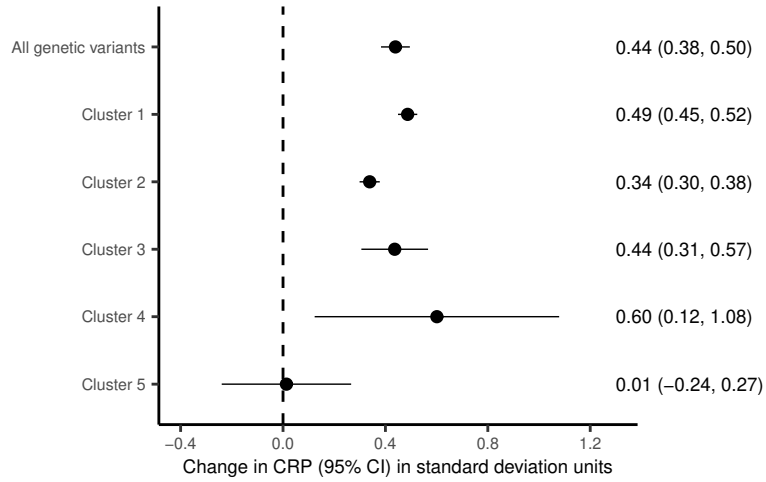


Figure 6: MR-IVW estimates and 95% confidence intervals of the association of genetically predicted BMI with CRP, for all genetic variants and for each cluster. The estimates represent the change in CRP in standard deviation units per 1 standard deviation increase in genetically predicted BMI.

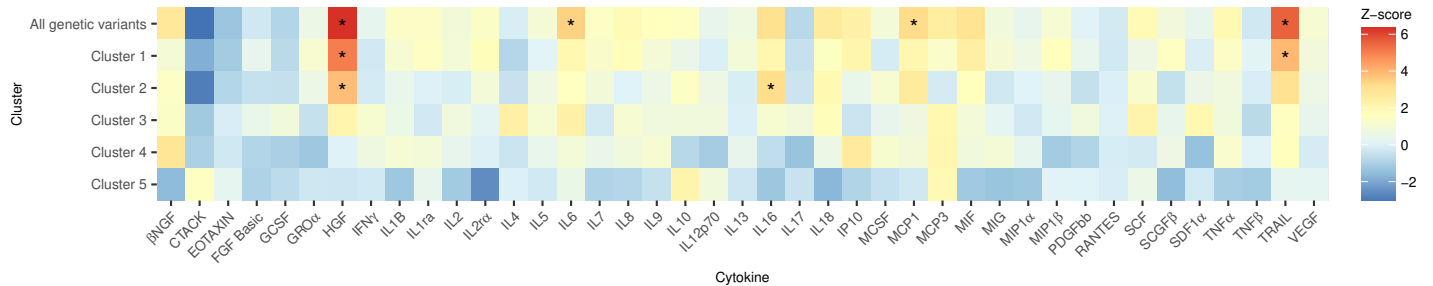


Figure 7: MR-IVW estimates (expressed as Z-scores, i.e. estimate divided by its standard error) for the association of genetically predicted BMI with 41 cytokines and growth factors. Values denoted with * have a p-value less than 0.05/41.