

# Local PCA Shows How the Effect of Population Structure Differs Along the Genome

Han Li, Peter Ralph

September 25, 2016

## Abstract

Population structure leads to systematic patterns in measures of mean relatedness between individuals in large genomic datasets, which are often discovered and visualized using dimension reduction techniques such as principal component analysis (PCA). Mean relatedness is an average of the relationships across locus-specific genealogical trees, which can be strongly affected on intermediate genomic scales by linked selection and other factors. We show how to use local principal components analysis to describe this meso-scale heterogeneity in patterns of relatedness, and apply the method to genomic data from three species, finding in each that the effect of population structure can vary substantially across only a few megabases. In a global human dataset, localized heterogeneity is likely explained by polymorphic chromosomal inversions. In a range-wide dataset of *Medicago truncatula*, factors that produce heterogeneity are shared between chromosomes, correlate with local gene density, and may be caused by background selection or local adaptation. In a dataset of primarily African *Drosophila melanogaster*, large-scale heterogeneity across each chromosome arm is explained by known chromosomal inversions thought to be under recent selection, and after removing samples carrying inversions, remaining heterogeneity is correlated with recombination rate and gene density, again suggesting a role for linked selection. The visualization method provides a flexible new way to discover biological drivers of genetic variation, and its application to data highlights the strong effects that linked selection and chromosomal inversions can have on observed patterns of genetic variation.

## 1 Introduction

Wright (1949) defined *population structure* to encompass “such matters as numbers, composition by age and sex, and state of subdivision”, where “subdivision” refers to restricted migration between subpopulations. The phrase is also commonly used to refer to the genetic patterns that result from this process, as for instance reduced mean relatedness

between individuals from distinct populations. However, it is not necessarily clear what aspects of demography should be included in the concept. For instance, Blair (1943) defines *population structure* to be the sum total of “such factors as size of breeding populations, periodic fluctuation of population size, sex ratio, activity range and *differential survival of progeny*” (emphasis added). The definition is similar to Wright’s, but differs in including the effects of natural selection. On closer examination, incorporating differential survival or fecundity makes the concept less clear: should a randomly mating population consisting of two types that are partially reproductively isolated from each other be said to show population structure or not? Whatever the definition, it is clear that due to natural selection, the effects of population structure – the *realized* patterns of genetic relatedness – differ depending on which portion of the genome is being considered. For instance, strongly locally adapted alleles of a gene will be selected against in migrants to different habitats, increasing genetic differentiation between populations near to this gene. Similarly, newly adaptive alleles spread first in local populations. These observations motivate many methods to search for genetic loci under selection, as for example in Huerta-Sánchez et al. (2013) and Martin et al. (2016).

These realized patterns of genetic relatedness summarize the shapes of the genealogical trees at each location along the genome. Since these trees vary along the genome, so does relatedness, but averaging over sufficiently many trees we hope to get a stable estimate that doesn’t depend much on the genetic markers chosen. This is not guaranteed: for instance, relatedness on sex chromosomes is expected to differ from the autosomes; and positive or negative selection on particular loci can dramatically distort shapes of nearby genealogies (Barton 2000; Charlesworth et al. 1993; Kim and Stephan 2002). Indeed, many species show chromosome-scale variation in diversity and divergence (e.g., (Langley et al. 2012)); species phylogenies can differ along the genome due to incomplete lineage sorting, adaptive introgression and/or local adaptation (e.g., (Ellegren et al. 2012; Nadeau et al. 2012; Pease and Hahn 2013; Pool 2015; Vernot and Akey 2014)); and theoretical expectations predict that geographic patterns of relatedness should depend on selection (Charlesworth et al. 2003).

Patterns in genome-wide relatedness are often summarized by applying principal components analysis (PCA, (Patterson et al. 2006)) to the genetic covariance matrix, as pioneered by Menozzi et al. (1978). The results of PCA can be related to the genealogical history of the samples, such as time to most recent common ancestor and migration rate between populations (McVean 2009; Novembre and Stephens 2008), and sometimes produce “maps” of population structure that reflect the samples’ geographic origin distorted by rates of gene flow (Novembre et al. 2008).

Modeling such “background” kinship between samples is essential to genome-wide association studies (GWAS, (Astle and Balding 2009; Price et al. 2006)), and so understanding variation in kinship along the genome could lead to more generally powerful methods, and may be essential for doing GWAS in species with substantial heterogeneity in realized patterns of mean relatedness along the genome.

A note on nomenclature: In this work we describe variation in patterns of relatedness using local PCA, where “local” refers to proximity along the genome. A number of general methods for dimensionality reduction also use a strategy of “local PCA” (e.g., (Kambhatla and Leen 1997; Manjón et al. 2013; Roweis and Saul 2000; Weingessel and Hornik 2000)), performing PCA not on the entire dataset but instead on subsets of observations, providing local pictures which are then stitched back together to give a global picture. At first sight, this differs from our method in that we restrict to subsets of *variables* instead of subsets of observations. However, if we flip perspectives and think of each genetic variant as an observation, our method shares common threads, although our method does not subsequently use adjacency along the genome, as we aim to identify similar regions that may be distant.

As reviewed above, it is common to describe variation along the genome of simple statistics such as  $F_{ST}$ , and also to use methods such as PCA to visualize large-scale patterns in genome-wide relatedness. In this paper we apply PCA locally along the genome, thus describing in an unbiased, descriptive way how patterns of mean relatedness vary systematically along the genome, in a way particularly suited to large samples from geographically distributed populations. Geographic population structure sets the stage by establishing similar patterns of relatedness across much of the genome; we then aim to describe how this structure is affected by selection and other factors. Our aim is not to identify outlier loci, but rather to describe larger-scale variation shared by many parts of the genome. By doing this with three taxa with diverse population histories, we can compare relative contributions of different sorts of variation across taxa; since the method is an visualization method, we allow ourselves to be suprised by unexpected signals in the data.

## 2 Results

As depicted in Figure 1, the general steps to the method are: (1) divide the genome into windows, (2) summarize the patterns of relatedness in each window, (3) measure dissimilarity in relatedness between each pair of windows, (4) visualize the resulting dissimilarity matrix using multidimensional scaling (MDS), and (5) combine similar windows to more accurately visualize local effects of population structure using PCA. Details of how we carried out each step are given in the Methods, and code to run the analyses is provided as an R package at [https://github.com/petrelharp/local\\_pca](https://github.com/petrelharp/local_pca).

We applied the method to genomic datasets with good geographic sampling: 380 African *Drosophila melanogaster* from the Drosophila Genome Nexus (Lack et al. 2015), a world-wide dataset of humans, 3,965 humans from several locations worldwide from the POPRES dataset (Nelson et al. 2008), and 263 *Medicago truncatula* from 24 countries around the Mediterranean basin a range-wide dataset of the partially selfing weedy annual plant from the *Medicago truncatula* Hapmap Project (Tang et al. 2014),

In all three species, PCA plots vary along the genome in a systematic way, showing

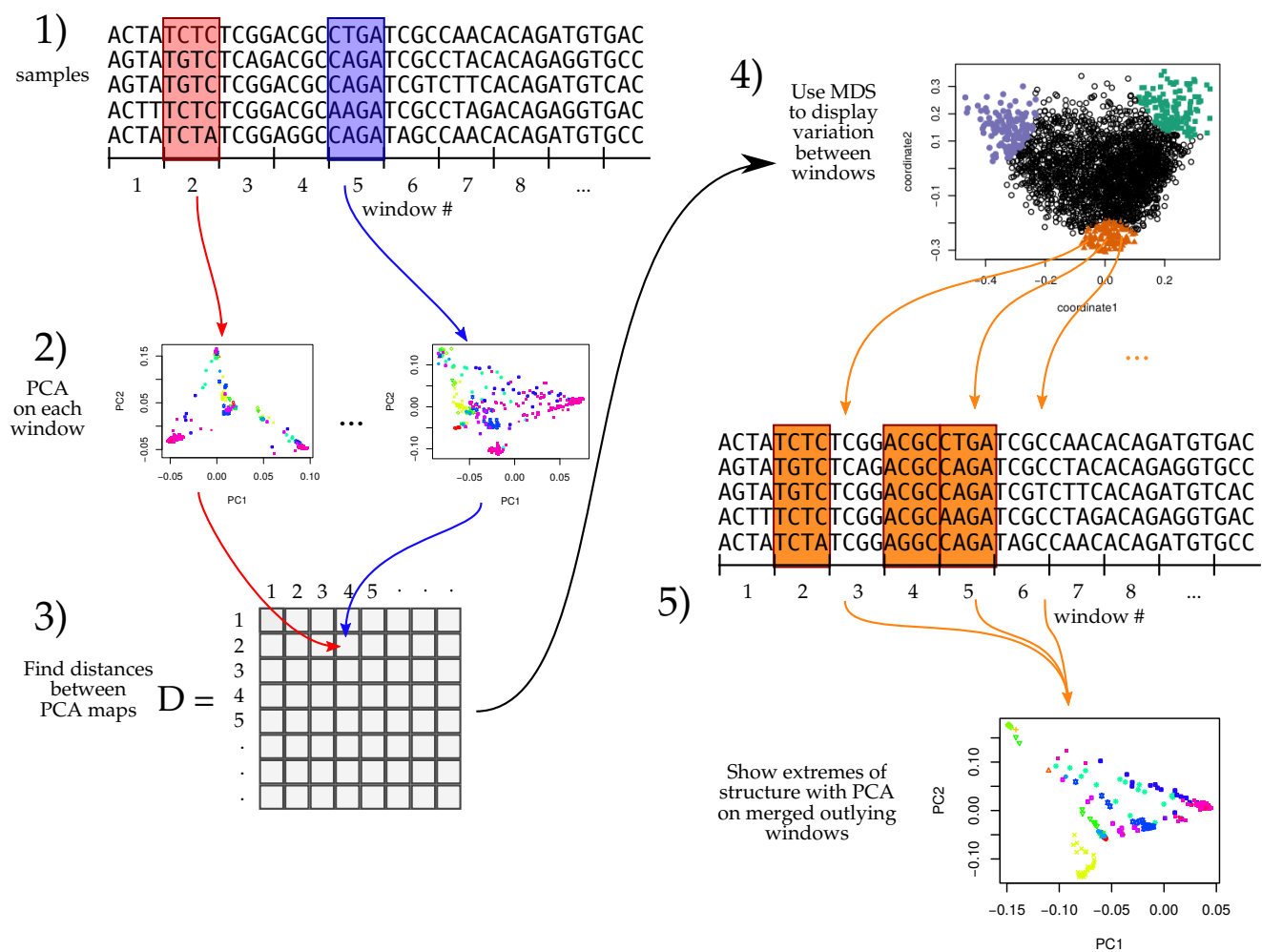


Figure 1: An illustration of the method; see Methods for details.

strong chromosome-scale correlations. This implies that variation is due to meaningful heterogeneity in a biological process, since noise due to randomness in choice of local genealogical trees is not expected to show long distance correlations. Below, we discuss the results and likely underlying causes.

## 2.1 *Drosophila melanogaster*

We applied the method to windows of average length 9 Kbp, across chromosome arms 2L, 2R, 3L, 3R and X separately. The first column of Figure 2 is a multidimensional scaling (MDS) visualization of the matrix of dissimilarities between genomic windows: in other words, genomic windows that are closer to each other in the MDS plot show more similar patterns of relatedness. For each chromosome arm, the MDS visualization roughly resembles a triangle, sometimes with additional points. Since the relative position of each window in this plot shows the similarity between windows, this suggests that there are at least three extreme manifestations of population structure typified by windows found in the “corners” of the figure, and that other windows’ patterns of relatedness may be a mixture of those extremes. The next two columns of Figure 2 respectively depict the two MDS coordinates of each window, plotted against the window’s position along the genome, to show how the plot of the first column is laid out along the genome.

To help visualize how clustered windows with similar patterns of relatedness are along each chromosome arm, we selected three “extreme” windows in the MDS plot and the 5% of windows that are closest to it in the MDS coordinates, then highlighted these windows’ positions along the genome, and created PCA plots for the windows, combined. Representative plots are shown for three groups of windows on each chromosome arm in Figure 2 (groups are shown in color), and in Supplemental Figure S1 (PCA plots). The latter plots are quite different, showing that genomic windows in different regions of the MDS plot indeed show quite different patterns of relatedness.

The most striking variation in patterns of relatedness turns out to be explained by several large inversions that are polymorphic in these samples, discussed in Corbett-Detig and Hartl (2012) and Langley et al. (2012). To depict this, Figure 3 shows the PCA plots in Figure S1 recolored by the orientation of the inversion for each sample. Taking chromosome arm 2L as an example, the two regions of similar, extreme patterns of relatedness shown in green in the first row of Figure 2 lie directly around the breakpoints of the inversion *In(2L)t*, and the PCA plots in the first rows of Figure 3 shows that patterns of relatedness here are mostly determined by inversion orientation. The regions shown in purple on chromosome 2L lie near the centromere, and have patterns of relatedness reflective of two axes of variation, seen in Figure S1, which correspond roughly to latitude within Africa and to degree of cosmopolitan admixture respectively (see Lack et al. (2015) for more about admixture in this sample). The regions shown in orange on chromosome 2L mostly lie inside the inversion, and show patterns of relatedness that are a mixture between the other two, as expected due to recombination within the (long) inversion (Guerrero et al. 2011).

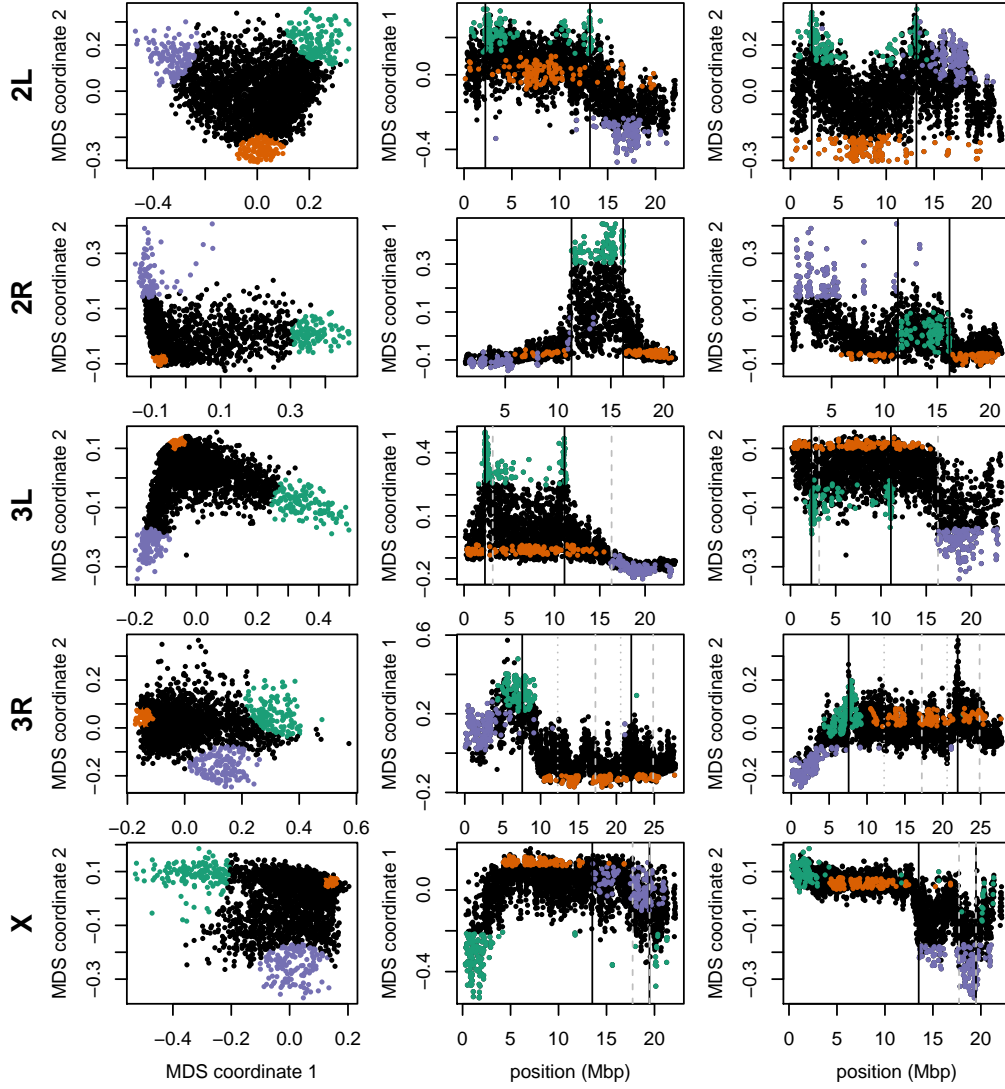


Figure 2: Variation in patterns of relatedness for windows across *Drosophila melanogaster* chromosome arms. In all plots, each point represents one window along the genome. The first column shows the MDS visualization of relationships between windows, and the second and third columns show the two MDS coordinates against the midpoint of each window; rows correspond to chromosome arms. Colors are consistent for plots in each row. Vertical lines show the breakpoints of known polymorphic inversions. Solid black lines are for the inversions we used in Figure 3, while dotted grey lines are for other known inversions.

Similar results are found in other chromosome arms, albeit complicated by the coexistence of more than one polymorphic inversion; however, each breakpoint visibly affects patterns in the MDS coordinates (see vertical lines in Figure 2).

To see how patterns of relatedness vary in the absence of polymorphic inversions, we performed the same analyses after removing, for each chromosome arm, any samples carrying inversions on that arm. In the result, shown in Supplemental Figure S5, the striking peaks associated with inversion breakpoints are gone, and previously smaller-scale variation now dominates the MDS visualization. For instance, the majority of the variation along 3L in Figure 2 is on the left end of the arm, dominated by two large peaks around the inversion breakpoints; there is also a relatively small dip on the right end of the arm (near the centromere). In contrast, Supplemental Figure S5 shows that after removing polymorphic inversions, remaining structure is dominated by the dip near the centromere. Without inversions, variation in patterns of relatedness shown in the MDS plots follows similar patterns to that previously seen in *D. melanogaster* recombination rate and diversity (Langley et al. 2012; Mackay et al. 2012). Indeed, correlations between the recombination rate in each window and the position on the first MDS coordinate are highly significant (Spearman’s  $\rho = 0.54$ ,  $p < 2 \times 10^{-16}$ ; Figures 4 and S6). This is consistent with the hypothesis that variation is due to selection, since the strength of linked selection increases with local gene density, measured in units of recombination distance. The number of genes – measured as the number of transcription start and end sites within each window – was not significantly correlated with MDS coordinate ( $p = 0.22$ ).

## 2.2 Human

As we did for the *Drosophila* data, we applied our method separately to all 22 human autosomes. On each, variation in patterns of relatedness was dominated by a small number of windows having similar patterns of relatedness to each other that differed dramatically from the rest of the chromosome. These may be primarily inversions: outlying windows coincide with three of the six large polymorphic inversions described in Antonacci et al. (2009), notably a particularly large, polymorphic inversion on 8p23 (Figure 5). Similar plots for all chromosomes are shown in Supplementary Figures S7, S8, and S9. PCA plots of many outlying windows show a characteristic trimodal shape (shown for chromosome 8 in Figure S10), presumably distinguishing samples having each of the three diploid genotypes for each inversion orientation (although we do not have data on orientation status). This trimodal shape has been proposed as a method to identify inversions (Ma and Amos 2012), but distinguishing this hypothesis from others, such as regions of low recombination rate, would require additional data.

We also applied the method on all 22 autosomes together, and found that, remarkably, the inversion on chromosome 8 is still the most striking outlying signal (Figure S11). Further investigation with a denser set of SNPs, allowing a finer genomic resolution, may yield other patterns.

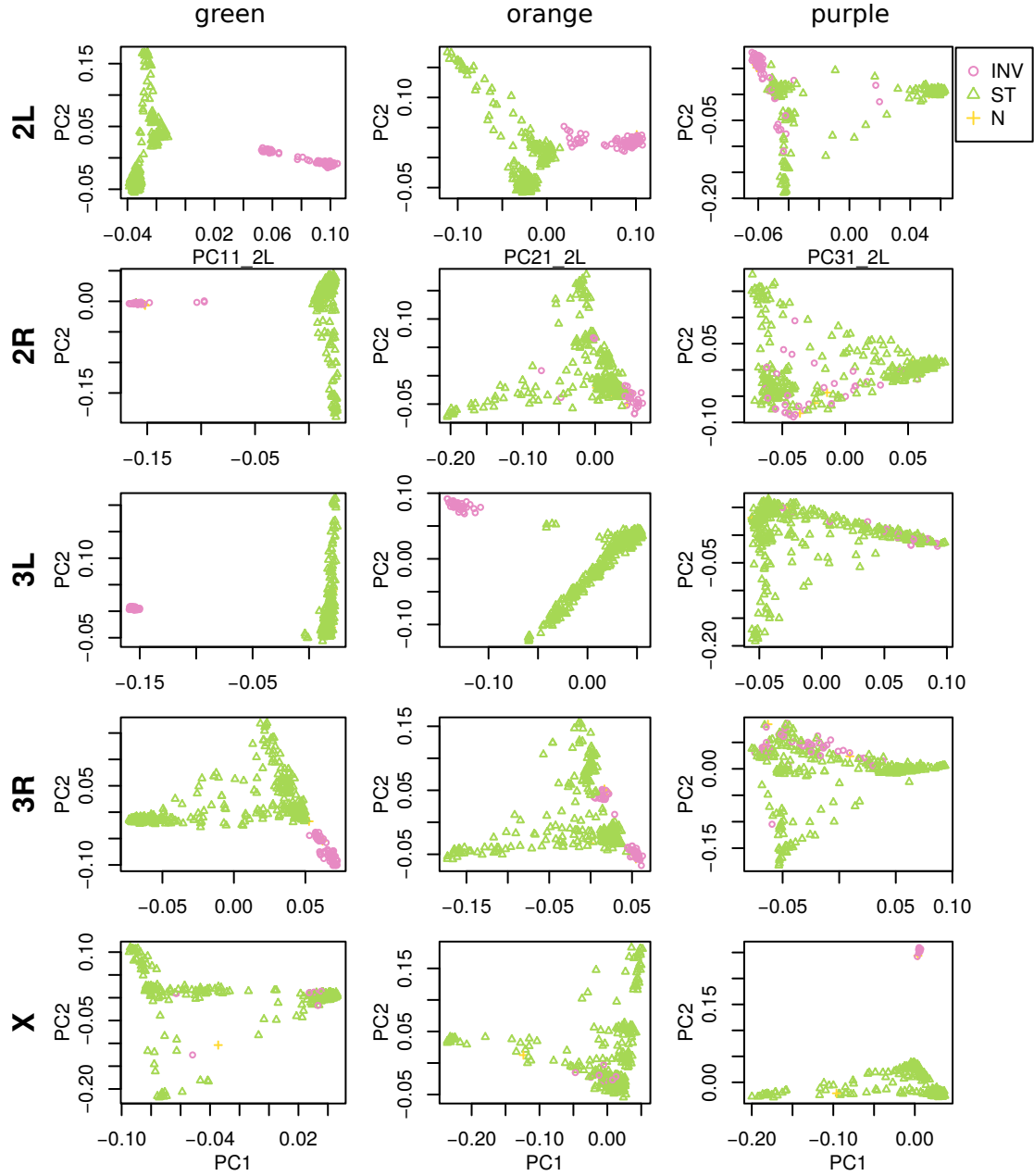


Figure 3: As in Figure 2, except that samples are colored by orientation of the polymorphic inversions In(2L)t, In(2R)NS, In(3L)OK, In(3R)K and In(1)A respectively. (data from (Lack et al. 2015)) In each “INV” denotes an inverted genotype, “ST” denotes the standard orientation, and “N” denotes unknown.



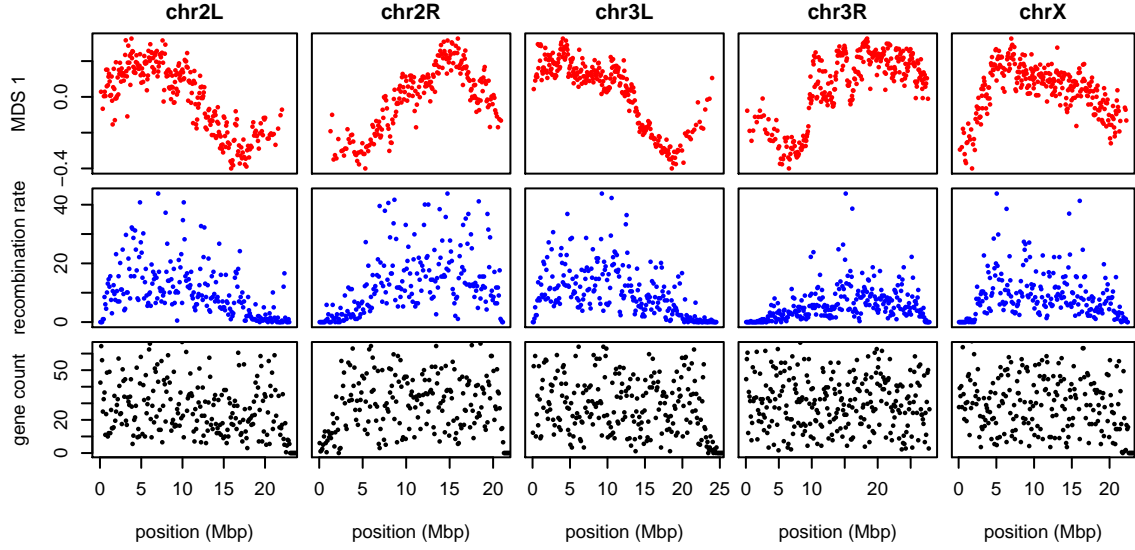


Figure 4: The effects of population structure without inversions is correlated to recombination rate in *Drosophila melanogaster*. The first plot (in red) shows the first MDS coordinate along the genome for windows of 10,000 SNPs, obtained after removing samples with inversions. (A plot analogous to Figure 2 is shown in Supplemental Figure S5.) The second plot (in blue) shows local average recombination rates in cM/Mbp, obtained as midpoint estimates for 100Kbp windows from the *Drosophila* recombination rate calculator (Fiston-Lavier et al. 2010) release 5, using rates from Comeron et al. (2012). The third plot (in black) shows the number of genes' transcription start and end sites within each 100Kbp window, divided by two. Transcription start and end sites were obtained from the RefGene table from the UCSC browser. The histone gene cluster on chromosome arm 2L is excluded.

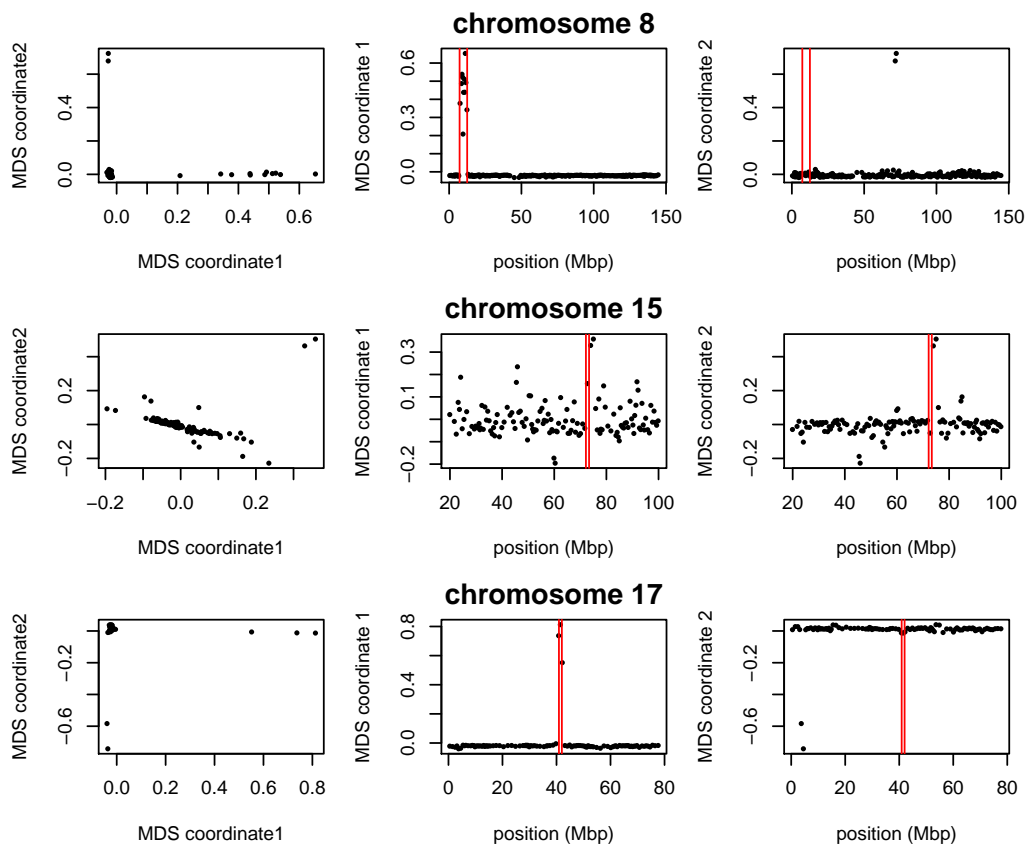


Figure 5: Variation in structure between windows on human chromosomes 8, 15, and 17. Each point in each plot represents a window. The first column shows the MDS visualization of relationships between windows; the second and third columns show the two MDS coordinates of each window against its position (midpoint) along the chromosome. Rows, from top to bottom show chromosomes 8, 15, and 17. The vertical red lines show the breakpoints of known inversions from Antonacci et al. (2009).

### 2.3 *Medicago truncatula*

Unlike the other two species, the method applied separately on all eight chromosomes of *Medicago truncatula* showed similar patterns of gradual change in patterns of relatedness across each chromosome, with no indications of chromosome-specific patterns. This consistency suggests that the factor affecting the population structure for each chromosome is the same, as might be caused by varying strengths of linked selection. To verify that variation in the effects of population structure is shared across chromosomes, we applied the method to all chromosomes together. Results for chromosome 3 are shown in Figures 6 and 6, and other chromosomes are similar: across chromosomes, the high values of the first MDS coordinate coincide with the position of the heterochromatic regions surrounding the centromere, which often have lower gene density and may therefore be less subject to linked selection. To verify that this is a possible explanation, we counted the number of genes found in each window using gene models in Mt4.0 from [jcvi.org](http://jcvi.org) (Tang et al. 2014), which are shown juxtaposed with the first MDS coordinate of each window in Figure 7, and are significantly correlated, as shown in Supplemental Figure S12. (Values shown are the number of start and end positions of each predicted mRNA transcript, divided by two, assigned to the nearest window.) However, other genomic features, such as distance to centromere show roughly the same patterns, so we cannot rule out alternative hypotheses. In particular, the recombination rates estimated by Paape et al. (2012) appear visually to be similar, but were not available in a form that can be compared to our MDS values.

We also found nearly identical results when choosing shorter windows of 1,000 SNPs; or choosing windows of equal length in base pairs rather than SNPs. Similarly, the results were not substantially changed when using weighted PCA to downweight the large group of Tunisian samples.

## 3 Discussion

Our investigations have found substantial variation in the patterns of relatedness formed by population structure across the genomes of three diverse species, revealing distinct biological processes driving this variation in each species. More investigation, particularly on more species and datasets, will help to uncover what aspects of species history can explain these differences. With growing appreciation of the heterogeneous effects of selection across the genome, especially the importance of adaptive introgression and hybrid speciation (Brandvain et al. 2014; Fitzpatrick et al. 2010; Hufford et al. 2013; Pool 2015; Staubach et al. 2012), local adaptation (Lenormand 2002; Wang and Bradburd 2014), and inversion polymorphisms (Kirkpatrick 2010; Kirkpatrick and Barrett 2015), local PCA may prove to be a useful exploratory tool to discover important genomic features.

We now discuss possible implications of this variation in the effects of population structure, the impact of various parameter choices in implementing the method, and possible additional applications.

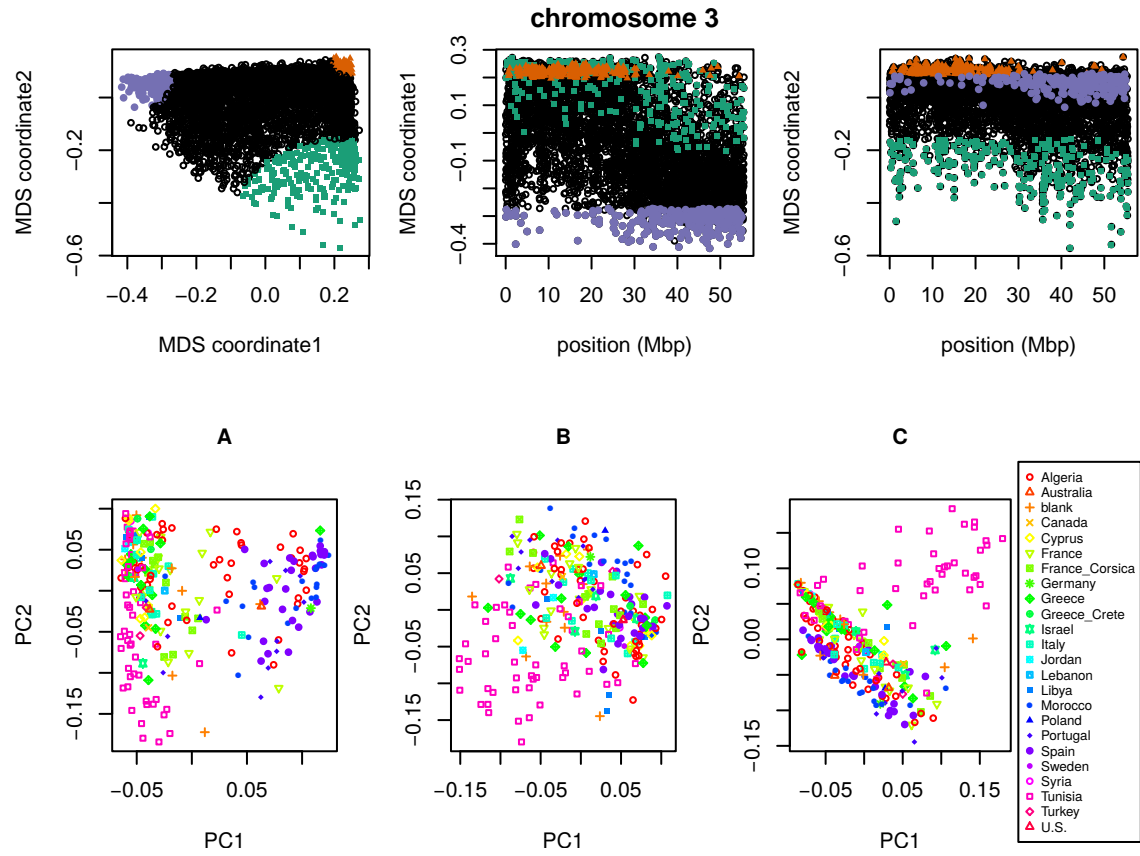


Figure 6: MDS visualization of patterns of relatedness on *M. truncatula* chromosome 3, with corresponding PCA plots. Each point in the plot represents a window; the structure revealed by the MDS plot is strongly clustered along the chromosome, with windows in the upper-right corner of the MDS plot (colored red) clustered around the centromere, windows in the upper-left corner (purple) furthest from the centromere, and the remaining corner (green) intermediate. Plots for remaining chromosomes are shown in Supplemental Figure S13. **(below)** PCA plots for the sets of genomic windows colored (A) green, (B) orange, and (C) purple in Figure 6. Each point corresponds to a sample, colored by country of origin. Plots for remaining chromosomes are shown in Supplemental Figure S14.

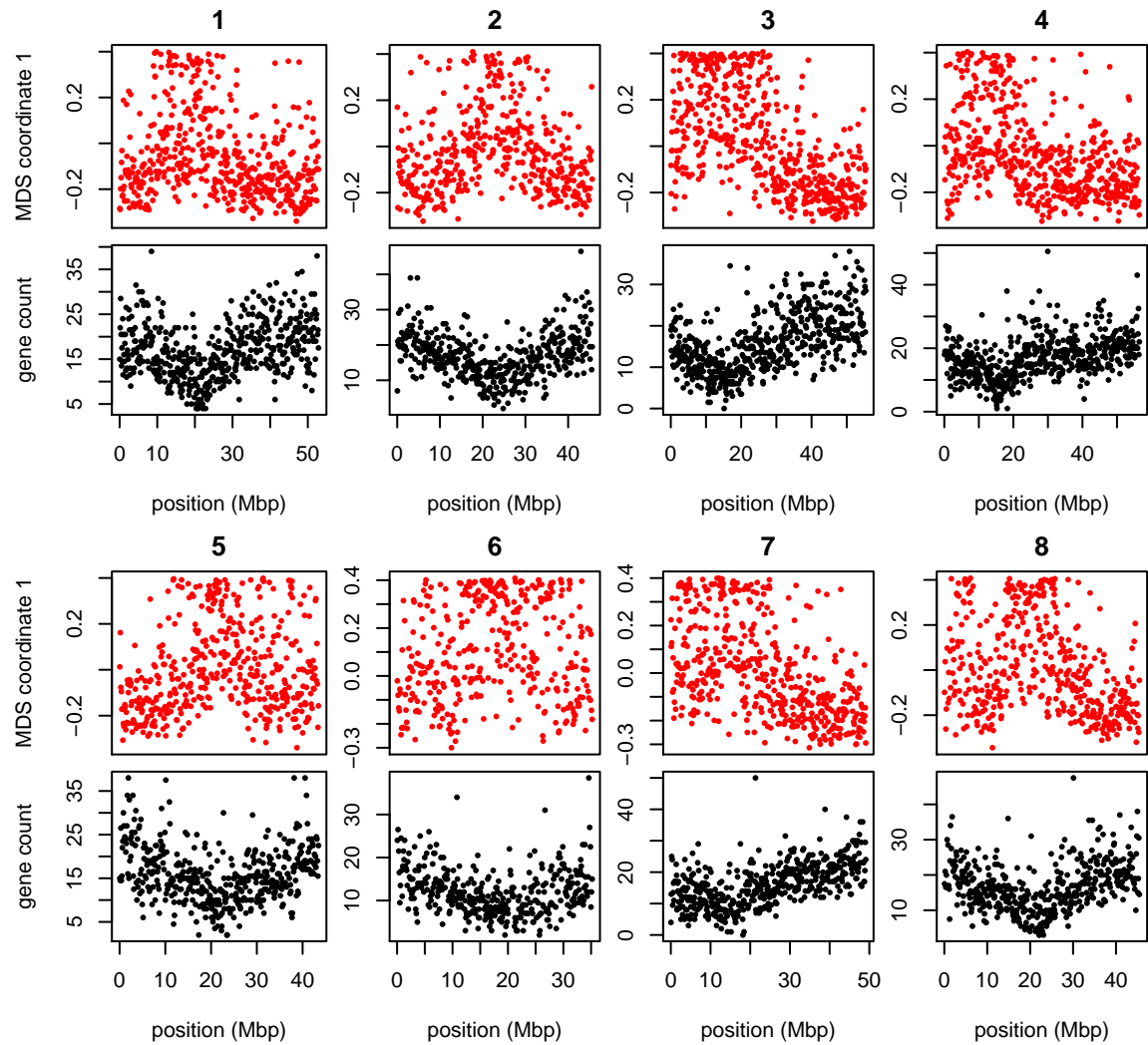


Figure 7: MDS coordinate and gene density for each window in the *Medicago* genome, for chromosomes 1–8 (numbered above each pair of figures). For each chromosome, the red plot above is first coordinate of MDS against the middle position of each window along each chromosome. The black plot below is gene count for each window against the position of each window.

**Chromosomal inversions** A major driver of variation in patterns of relatedness in two datasets we examined are inversions. This may be common, but the example of *Medicago truncatula* shows that polymorphic inversions are not ubiquitous. PCA has been proposed as a method for discovering inversions (Ma and Amos 2012); however, the signal left by inversions likely cannot be distinguished from long haplotypes under balancing selection or simply regions of reduced recombination without additional lines of evidence. Inversions show up in our method because across the inverted region, most gene trees share a common split that dates back to the origin of the inversion. However, in many applications, inversions are a nuisance. For instance, SMARTPCA (Patterson et al. 2006) reduces their effect on PCA plots by regressing out the effect of linked SNPs on each other. Removing samples with the less common orientation of each inversion reduced, but did not eliminate, the signal of inversions seen in the *Drosophila melanogaster* dataset, demonstrating that the genomic effects of transiently polymorphic inversions may outlast the inversions themselves.

**The effect of selection** It seems that the variation in patterns of relatedness we see in the *Medicago truncatula* and *Drosophila melanogaster* datasets must be explained somehow by linked selection. Furthermore, the selection must be affecting many targets across the genome, since we see similar effects across long distances (even distinct chromosomes). For this reason, the most likely candidate may be selection against linked deleterious mutations, known as “background selection” (Charlesworth et al. 1993; Charlesworth 2013). Informally, background selection reduces the number of potential contributors to the gene pool in regions of the genome with many possible deleterious mutations (Hudson and Kaplan 1995); for this reason, if it acts in a spatial context, it is expected to induce samples from nearby locations to cluster together more frequently. Therefore, regions of the genome harboring many targets of local adaptation may show similar patterns, since migrant alleles in these regions will be selected against, and so locally gene trees will more closely reflect spatial proximity.

A related possibility is that variation in patterns of relatedness is due to recent admixture between previously separated populations, the effects of which were not uniform across the genome due to selection. For instance, it has been hypothesized that large-scale variation in amount of introgressed Neanderthal DNA along the genome is due to selection against Neanderthal genes, leading to greater introgression in regions of lower gene density (Harris and Nielsen 2016; Juric et al. 2016). African *Drosophila melanogaster* are known to have a substantial amount of recently introgressed genome from “cosmopolitan” sources; if selection regularly favors genes from one origin, this could lead to substantial variation in patterns of relatedness correlated with local gene density.

There has been substantial debate over the relative impacts of different forms of selection. These have been difficult to disentangle in part because for the most part theory makes predictions which are only strictly valid in randomly mating (i.e., unstructured) populations, and it is unclear to what extent the spatial structure observed in most real

populations will affect these predictions. It may be possible to design more powerful statistics that make stronger use of spatial information.

**Parameter choices** There are several choices in the method that may in principle affect the results. As with whole-genome PCA, the choice of samples is important, as variation not strongly represented in the sample will not be discovered. The effects of strongly imbalanced sampling schemes are often corrected by dropping samples in overrepresented groups; but downweighting may be a better option that does not discard data (and here we present a method to do this). Next, the choice of window size may be important, although in our applications results were not sensitive to this, indicating that we can see variation on a sufficiently fine scale. Finally, which collections of genomic regions are compared to each other (steps 3 and 4 in Figure 1), along with the method used to discover common structure, will affect results. We used MDS, applied to either each chromosome separately or to the entire genome; for instance, human inversions are clearly visible as outliers when compared to the rest of their chromosome, but genome-wide, their signal is obscured by the numerous other signals of comparable strength.

Besides window length, there is also the question of how to choose windows. In these applications we have used nonoverlapping windows with equal numbers of polymorphic sites. Alternatively, windows could be chosen to have equal length in genetic distance, so that each would have roughly the same number of independent trees. However, we found little change in results when using different window sizes or when measuring windows in physical distance (in bp).

Finally, our software allows different choices for how many PCs to use in approximating structure of each window ( $k$  in equation 1), and how many MDS coordinates to use when describing the distance matrix between windows, but in our exploration, changing these has not produced dramatically different results. These are all part of more general techniques in dimension reduction and high-dimensional data visualization; we encourage the user to experiment.

**Applications** So-called cryptic relatedness between samples has been one of the major sources of confounding in genome-wide association studies (GWAS) and so methods must account for it by modeling population structure or kinship (Astle and Balding 2009; Yang et al. 2014). Since the effects of population structure is not constant along the genome, this could in principle lead to an inflation of false positives in parts of the genome with stronger population structure than the genome-wide average. A method such as ours might be used to provide a more sensitive correction. Fortunately, in our human dataset this does not seem likely to have a strong effect: most variation is due to small, independent regions, possibly primarily inversions, and so may not have a major effect on GWAS. In the other species we examined, particularly *Drosophila melanogaster*, treating population structure as a single quantity would entail a substantial loss of power, and could potentially be

misleading.

## 4 Methods

### 4.1 PCA in genomic windows

We first recoded sampled genotypes as numeric matrices in the usual manner, by recording the number of nonreference alleles seen at each locus for each sample. We then divided the genome into contiguous segments (“windows”) and applied principal component analysis (PCA) as described in McVean (2009) separately to the submatrices that corresponded to each window. The choice of window length entails a tradeoff between signal and noise, since shorter windows allow better resolution along the genome but provide less precise estimates of relatedness. A method for choosing a window length to balance these considerations is given in Appendix A. Precisely, denote by  $Z$  the  $L \times N$  recoded genotype matrix for a given window ( $L$  is the number of SNPs and  $N$  is the sample size), and by  $\overline{Z}_s$  the mean of non-missing entries for allele  $s$ , so that  $\overline{Z}_s = \frac{1}{n_s} \sum_j Z_{sj}$ , where the sum is over the  $n_s$  nonmissing genotypes. We first compute the mean-centered matrix  $X$ , as  $X_{si} = Z_{si} - \overline{Z}_s$ , and preserving missingness. (This mean-centering makes the result not depend on the choice of reference allele, exactly if there is no missing data, and approximately otherwise.) Next, we find the covariance matrix of  $X$ , denoted  $C$ , as  $C_{ij} = \frac{1}{m_{ij}-1} \sum_s X_{si} X_{sj} - \frac{1}{m_{ij}(m_{ij}-1)} (\sum_s X_{si}) (\sum_s X_{sj})$ , where all sums are over the  $m_{ij}$  sites where both sample  $i$  and sample  $j$  have nonmissing genotypes. The principal components are the eigenvectors of  $C$ , normalized to have Euclidean length equal to one, and ordered by magnitude of the eigenvalues.

The top few principal components are generally good summaries of population structure; we follow common practice and usually only use the first two (referred to as  $PC1$  and  $PC2$ ). The above procedure can be performed on any subset of the data; for future reference, denote by  $PC1_j$  and  $PC2_j$  the result after applying to all SNPs in the  $j^{\text{th}}$  window.

Since the definition of eigenvectors does not specify their sign, when comparing between windows we choose signs to best match each other: after choosing  $PC1_1$ , for instance, if  $u$  is the first eigenvector obtained from the covariance matrix for window  $j$ , then we next choose  $PC1_j = \pm u$ , where the sign is chosen according to which of  $\|PC1_1 - u\|$  or  $\|PC1_1 + u\|$  is smaller.

Several of the datasets we use have unbalanced representations of diverged populations, which can have a strong impact on the results of PCA. (The principal axes may describe variation *within* an overrepresented group rather than more significant variation between groups.) Therefore, to check that sampling patterns do not affect our results, we compared to a variant of PCA that gives roughly equal weight to each group of samples, rather than to each sample. The rationale and implementation of this method are described in Appendix B.



## 4.2 Similarity of patterns of relatedness between windows

We think of the local effects of population structure as being summarized by *relative* position of the samples in the space defined by the top principal components. However, we do not compare patterns of relatedness of different genomic regions by directly comparing the PCs, since rotations or reflections of these imply identical patterns of relatedness. Instead, we compare the low-dimensional approximations of the local covariance matrices obtained using the top  $k$  PCs, which is invariant under reflections and rotations and yet contains all other information about the PCs. (For results shown here, we use  $k = 2$ ; results using larger numbers of PCs were nearly identical.) Furthermore, to remove the effect of artifacts such as mutation rate variation, we also rescale each approximate covariance matrix so that the underlying data matrix has trace norm equal to one.

To do this, define the  $N \times k$  matrix  $V(i)$  so that  $V(i)_{\cdot\ell}$ , the  $\ell^{\text{th}}$  column of  $V(i)$ , is equal to the  $\ell^{\text{th}}$  principal component of the  $i^{\text{th}}$  window, multiplied by  $(\lambda_{\ell i} / \sum_{m=1}^k \lambda_{mi})^{1/2}$ , where  $\lambda_{\ell i}$  is the  $\ell^{\text{th}}$  eigenvalue of the genetic covariance matrix. Then, the rescaled, rank  $k$  approximate covariance matrix for the  $i^{\text{th}}$  window is

$$M(i) = \sum_{\ell=1}^k V(i)_{\cdot\ell} V(i)_{\cdot\ell}^T. \quad (1)$$

To measure the similarity of patterns of relatedness for the  $i^{\text{th}}$  window and  $j^{\text{th}}$  window, we then use Euclidean distance  $D_{ij}$  between the matrices  $M(i)$  and  $M(j)$ :  $D_{ij}^2 = \sum_{k,\ell} (M(i)_{k,\ell} - M(j)_{k,\ell})^2$ .

The goal of comparing PC plots up to rotation and reflection turned out to be equivalent to comparing rank- $k$  approximations to local covariance matrices. This suggests instead directly comparing entire local covariance matrices. However, with thousands of samples and tens of thousands of windows, computing the distance matrix would take months of CPU time, while as defined above,  $D$  can be computed in minutes using the following method. Since for square matrices  $A$  and  $B$ ,  $\sum_{ij} (A_{ij} - B_{ij})^2 = \sum_{ij} (A_{ij}^2 + B_{ij}^2) - 2 \text{tr}(A^T B)$ , then due to the orthogonality of eigenvectors and the cyclic invariance of trace,  $D_{ij}$  can be computed efficiently as

$$D_{ij} = \frac{\sum_{\ell=1}^k \lambda_{\ell i}^2}{(\sum_{\ell=1}^k \lambda_{\ell i})^2} + \frac{\sum_{\ell=1}^k \lambda_{\ell j}^2}{(\sum_{\ell=1}^k \lambda_{\ell j})^2} - 2 \sum_{\ell,m=1}^k (V(i)^T V(j))_{\ell m}^2. \quad (2)$$

## 4.3 Visualization of results

We use multidimensional scaling (MDS) to visualize relationships between windows as summarized by the dissimilarity matrix  $D$ . MDS produces a set of  $m$  coordinates for each window that give the arrangement in  $m$ -dimensional space that best recapitulates the original distance matrix. For results here, we use  $m = 2$  to produce one- or two-dimensional visualizations of relationships between windows' patterns of relatedness.

We then locate variation in patterns of relatedness along the genome by choosing collections of windows that are nearby in MDS coordinates, and map their positions along the genome. A visualization of the effects of population structure across the entire collection is formed by extracting the corresponding genomic regions and performing PCA on all, aggregated, regions.

## 4.4 Datasets

The three publicly available datasets we used are summarized in Table 1. We converted each to a numeric matrix (with one row per polymorphic variant and one column per sample) by replacing each genotype with the number of nonreference alleles (or NA for missing data).

***Drosophila melanogaster*:** We used whole-genome sequencing data from the *Drosophila* Genome Nexus (<http://www.johnpool.net/genomes.html>, (Lack et al. 2015)), consisting of the *Drosophila* Population Genomics Project phases 1–3 (Langley et al. 2012; Pool et al. 2012), and additional African genomes (Lack et al. 2015). After removing 20 genomes with more than 8% missing data, we were left with 380 samples from 16 countries across Africa and Europe. Since the *Drosophila* samples are from inbred lines or haploid embryos, we treat the samples as haploid when recoding; regions with residual heterozygosity were marked as missing in the original dataset; we also removed positions with more than 20% missing data. Each chromosome arm we investigated (X, 2L, 2R, 3L, and 3R) has 2–3 million SNPs; PCA plots for each arm are shown in Figure S2.

**Human:** We also used genomic data from the entire POPRES dataset (Nelson et al. 2008), which has array-derived genotype information for 447,267 SNPs across the 22 autosomes of 3,965 samples in total: 346 African-Americans, 73 Asians, 3,187 Europeans and 359 Indian Asians. Since these data derive from genotyping arrays, the SNP density is much lower than the other datasets, which are each derived from whole genome sequencing. We excluded the sex chromosomes and the mitochondria. PCA plots for each chromosome, separately, are shown in Figure S3.

***Medicago truncatula*:** Finally, we used whole-genome sequencing data from the *Medicago truncatula* Hapmap Project (Tang et al. 2014), which has 263 samples from 24 countries, primarily distributed around the Mediterranean basin. Each of the 8 chromosomes has 3–5 million SNPs; PCA plots for these are shown in Figure S4. We did not use the mitochondria or chloroplasts.

species	# SNPs per window	mean window length (bp)	mean # windows per chromosome	mean % variance explained by top 2 PCs
<i>Drosophila melanogaster</i>	1,000	9,019	2,674	0.53
Human	100	636,494	203	0.55
<i>Medicago truncatula</i>	10,000	102,580	467	0.50

Table 1: Descriptive statistics for each dataset used. Columns 2–4 give statistics describing the window sizes used for most analyses for each dataset (windows had a fixed number of SNPs). The final column gives the percent variance explained by the top two principal components, averaged across independent PCA of each window.

## 5 Data access

The methods described here are implemented in an open-source R package available at [https://github.com/petrelharp/local\\_pca](https://github.com/petrelharp/local_pca), as well as scripts to perform all analyses from VCF files at various parameter settings.

Datasets are available as follows: human (POPRES) at dbGaP with accession number phs000145.v4.p2, *Medicago* at the Medicago Hapmap <http://www.medicagohapmap.org/>, and *Drosophila* at the Drosophila Genome Nexus, <http://www.johnpool.net/genomes.html>.

## Acknowledgements

We are indebted to John Pool, Russ Corbett-Detig, Matilde Cordeiro, and Peter Chang for assistance with obtaining data and interpreting results (especially inversion status of *D. melanogaster* samples). Thanks also go to Yaniv Brandvain, Barbara Engelhardt, Charles Langley, Graham Coop, and Jeremy Berg for helpful comments and for encouraging the project.

## Disclosure declaration

The authors declare no conflicts of interest.

## References

Francesca Antonacci, Jeffrey M Kidd, Tomas Marques-Bonet, Mario Ventura, Priscillia Siswara, Zhaoshi Jiang, and Evan E Eichler. Characterization of six human disease-associated inversion polymorphisms. *Human molecular genetics*, 18(14):2555–2566, 2009.

- William Astle and David J. Balding. Population structure and cryptic relatedness in genetic association studies. *Statistical Science*, 24(4):451–471, 11 2009. doi: 10.1214/09-STS307. URL <http://dx.doi.org/10.1214/09-STS307>.
- Nicholas H. Barton. Genetic hitchhiking. *Philos Trans R Soc Lond B Biol Sci*, 355(1403): 1553–1562, November 2000. doi: 10.1098/rstb.2000.0716. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1692896/>.
- Albert P. Blair. Population structure in toads. *The American Naturalist*, 77(773):563–568, 1943. ISSN 00030147, 15375323. URL <http://www.jstor.org/stable/2457848>.
- Yaniv Brandvain, Amanda M. Kenney, Lex Flagel, Graham Coop, and Andrea L. Sweigart. Speciation and introgression between *Mimulus nasutus* and *Mimulus guttatus*. *PLoS Genet*, 10(6):e1004410, 06 2014. doi: 10.1371/journal.pgen.1004410. URL <http://dx.doi.org/10.1371/journal.pgen.1004410>.
- Frank M.T.A. Busing, Erik Meijer, and Rien Van Der Leeden. Delete-m jackknife for unequal m. *Statistics and Computing*, 9(1):3–8, 1999. ISSN 0960-3174. doi: 10.1023/A:1008800423698. URL <http://dx.doi.org/10.1023/A:1008800423698>.
- B Charlesworth, M T Morgan, and D Charlesworth. The effect of deleterious mutations on neutral molecular variation. *Genetics*, 134(4):1289–1303, August 1993. URL <http://www.genetics.org/content/134/4/1289>.
- Brian Charlesworth. Background selection 20 years on: The Wilhelmine E. Key 2012 invitational lecture. *Journal of Heredity*, 104(2):161–171, 2013. doi: 10.1093/jhered/ess136. URL <http://jhered.oxfordjournals.org/content/104/2/161.abstract>.
- Brian Charlesworth, Deborah Charlesworth, and Nicholas H. Barton. The effects of genetic and geographic structure on neutral variation. *Annual Review of Ecology, Evolution, and Systematics*, 34(1):99–125, 2003. doi: 10.1146/annurev.ecolsys.34.011802.132359. URL <http://arjournals.annualreviews.org/doi/abs/10.1146/annurev.ecolsys.34.011802.132359>.
- J M Comeron, R Ratnappan, and S Bailin. The many landscapes of recombination in *Drosophila melanogaster*. *PLoS Genet*, 8(10), 2012. doi: 10.1371/journal.pgen.1002905. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3469467/>.
- Russell B Corbett-Detig and Daniel L Hartl. Population genomics of inversion polymorphisms in *drosophila melanogaster*. *PLoS Genet*, 8(12):e1003056, 2012.
- B. Efron. *The Jackknife, the Bootstrap and Other Resampling Plans*. Society for Industrial and Applied Mathematics, 1982. doi: 10.1137/1.9781611970319. URL <http://epubs.siam.org/doi/abs/10.1137/1.9781611970319>.

- Hans Ellegren, Linnea Smeds, Reto Burri, Pall I. Olason, Niclas Backstrom, Takeshi Kawakami, Axel Kunstner, Hannu Makinen, Krystyna Nadachowska-Brzyska, Anna Qvarnstrom, Severin Uebbing, and Jochen B. W. Wolf. The genomic landscape of species divergence in *Ficedula* flycatchers. *Nature*, 491(7426):756–760, November 2012. ISSN 00280836. doi: 10.1038/nature11584. URL <http://dx.doi.org/10.1038/nature11584>.
- Barbara E. Engelhardt and Matthew Stephens. Analysis of population structure: a unifying framework and novel methods based on sparse factor analysis. *PLoS Genet*, 6(9), September 2010. doi: 10.1371/journal.pgen.1001117. URL <http://www.ncbi.nlm.nih.gov/pubmed/20862358>.
- Anna-Sophie Fiston-Lavier, Nadia D. Singh, Mikhail Lipatov, and Dmitri A. Petrov. *Drosophila melanogaster* recombination rate calculator. *Gene*, 463(1–2):18 – 20, 2010. ISSN 0378-1119. doi: <http://dx.doi.org/10.1016/j.gene.2010.04.015>. URL <http://www.sciencedirect.com/science/article/pii/S0378111910001769>.
- B M Fitzpatrick, J R Johnson, D K Kump, J J Smith, S R Voss, and H B Shaffer. Rapid spread of invasive genes into a threatened native species. *Proc Natl Acad Sci U S A*, 107(8):3606–3610, February 2010. doi: 10.1073/pnas.0911802107. URL <http://www.ncbi.nlm.nih.gov/pubmed/20133596>.
- Rafael F. Guerrero, François Rousset, and Mark Kirkpatrick. Coalescent patterns for chromosomal inversions in divergent populations. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 367(1587):430–438, 2011. ISSN 0962-8436. doi: 10.1098/rstb.2011.0246. URL <http://rstb.royalsocietypublishing.org/content/367/1587/430>.
- Kelley Harris and Rasmus Nielsen. The genetic cost of Neanderthal introgression. *Genetics*, 203(2):881–891, June 2016. URL <http://www.genetics.org/content/203/2/881>.
- R R Hudson and N L Kaplan. Deleterious background selection with recombination. *Genetics*, 141(4):1605–1617, December 1995. URL <http://www.genetics.org/content/141/4/1605>.
- Emilia Huerta-Sánchez, Michael DeGiorgio, Luca Pagani, Ayele Tarekegn, Rosemary Ekong, Tiago Antao, Alexia Cardona, Hugh E. Montgomery, Gianpiero L. Cavalleri, Peter A. Robbins, Michael E. Weale, Neil Bradman, Endashaw Bekele, Toomas Kivisild, Chris Tyler-Smith, and Rasmus Nielsen. Genetic signatures reveal high-altitude adaptation in a set of Ethiopian populations. *Molecular Biology and Evolution*, 30(8):1877–1888, 2013. doi: 10.1093/molbev/mst089. URL <http://mbe.oxfordjournals.org/content/30/8/1877.abstract>.
- Matthew B. Hufford, Pesach Lubinsky, Tanja Pyhäjärvi, Michael T. Devengenzo, Norman C. Ellstrand, and Jeffrey Ross-Ibarra. The genomic signature of crop-wild introgression.

- sion in maize. *PLoS Genet*, 9(5):e1003477, 05 2013. doi: 10.1371/journal.pgen.1003477. URL <http://dx.doi.org/10.1371%2Fjournal.pgen.1003477>.
- Ivan Juric, Simon Aeschbacher, and Graham Coop. The strength of selection against Neanderthal introgression. *bioRxiv*, 2016. doi: 10.1101/030148. URL <http://biorxiv.org/content/early/2016/07/22/030148>.
- N. Kambhatla and T. K. Leen. Dimension reduction by local principal component analysis. *Neural Computation*, 9(7):1493–1516, July 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.7.1493. URL [http://ieeexplore.ieee.org/xpl/freeabs\\_all.jsp?arnumber=6795533](http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?arnumber=6795533).
- Yuseob Kim and Wolfgang Stephan. Detecting a local signature of genetic hitchhiking along a recombining chromosome. *Genetics*, 160(2):765–777, 2002. URL <http://www.genetics.org/cgi/content/abstract/160/2/765>.
- Mark Kirkpatrick. How and why chromosome inversions evolve. *PLoS Biol*, 8(9), 2010. doi: 10.1371/journal.pbio.1000501. URL <http://www.ncbi.nlm.nih.gov/pubmed/20927412>.
- Mark Kirkpatrick and Brian Barrett. Chromosome inversions, adaptive cassettes and the evolution of species’ ranges. *Molecular Ecology*, 2015. ISSN 1365-294X. doi: 10.1111/mec.13074. URL <http://dx.doi.org/10.1111/mec.13074>.
- Justin B Lack, Charis M Cardeno, Marc W Crepeau, William Taylor, Russell B Corbett-Detig, Kristian A Stevens, Charles H Langley, and John E Pool. The *Drosophila* genome nexus: a population genomic resource of 623 *Drosophila melanogaster* genomes, including 197 from a single ancestral range population. *Genetics*, 199(4):1229–1241, 2015.
- C H Langley, K Stevens, C Cardeno, Y C Lee, D R Schrider, J E Pool, S A Langley, C Suarez, R B Corbett-Detig, B Kolaczowski, S Fang, P M Nista, A K Holloway, A D Kern, C N Dewey, Y S Song, M W Hahn, and D J Begun. Genomic variation in natural populations of *Drosophila melanogaster*. *Genetics*, 192(2):533–598, October 2012. doi: 10.1534/genetics.112.142018. URL <http://www.ncbi.nlm.nih.gov/pubmed/22673804>.
- Thomas Lenormand. Gene flow and the limits to natural selection. *Trends in Ecology & Evolution*, 17(4):183 – 189, 2002. ISSN 0169-5347. doi: DOI:10.1016/S0169-5347(02)02497-7. URL <http://www.sciencedirect.com/science/article/pii/S0169534702024977>.
- J Ma and C I Amos. Investigation of inversion polymorphisms in the human genome using principal components analysis. *PLoS One*, 7(7), 2012. doi: 10.1371/journal.pone.0040224. URL <http://www.ncbi.nlm.nih.gov/pubmed/22808122>.

- Trudy F. C. Mackay, Stephen Richards, Eric A. Stone, Antonio Barbadilla, Julien F. Ayroles, Dianhui Zhu, Sonia Casillas, Yi Han, Michael M. Magwire, Julie M. Cridland, Mark F. Richardson, Robert R. H. Anholt, Maite Barron, Crystal Bess, Kerstin Petra Blankenburg, Mary Anna Carbone, David Castellano, Lesley Chaboub, Laura Duncan, Zeke Harris, Mehwish Javaid, Joy Christina Jayaseelan, Shalini N. Jhangiani, Katherine W. Jordan, Fremiet Lara, Faye Lawrence, Sandra L. Lee, Pablo Librado, Raquel S. Linheiro, Richard F. Lyman, Aaron J. Mackey, Mala Munidasa, Donna Marie Muzny, Lynne Nazareth, Irene Newsham, Lora Perales, Ling-Ling Pu, Carson Qu, Miquel Ramia, Jeffrey G. Reid, Stephanie M. Rollmann, Julio Rozas, Nehad Saada, Lavanya Turlapati, Kim C. Worley, Yuan-Qing Wu, Akihiko Yamamoto, Yiming Zhu, Casey M. Bergman, Kevin R. Thornton, David Mittelman, and Richard A. Gibbs. The *Drosophila melanogaster* genetic reference panel. *Nature*, 482(7384):173–178, February 2012. ISSN 00280836. URL <http://dx.doi.org/10.1038/nature10811>.
- José V Manjón, Pierrick Coupé, Luis Concha, Antonio Buades, D Louis Collins, and Montserrat Robles. Diffusion weighted image denoising using overcomplete local PCA. *PloS one*, 8(9):e73021, 2013.
- Simon Henry Martin, Markus Moest, Wiliam J Palmer, Camilo Salazar, W. Owen McMillan, Francis M Jiggins, and Chris D Jiggins. Natural selection and genetic diversity in the butterfly *Heliconius melpomene*. *Genetics*, 203(1):525–541, May 2016. doi: 10.1101/042796. URL <http://www.genetics.org/content/203/1/525>.
- Gil McVean. A genealogical interpretation of principal components analysis. *PLoS Genet*, 5(10):e1000686, 2009.
- P Menozzi, A Piazza, and L Cavalli-Sforza. Synthetic maps of human gene frequencies in Europeans. *Science*, 201(4358):786–792, September 1978. URL <http://www.ncbi.nlm.nih.gov/pubmed/356262>.
- N J Nadeau, A Whibley, R T Jones, J W Davey, K K Dasmahapatra, S W Baxter, M A Quail, M Joron, R H French Constant, M L Blaxter, J Mallet, and C D Jiggins. Genomic islands of divergence in hybridizing *Heliconius* butterflies identified by large-scale targeted sequencing. *Philos Trans R Soc Lond B Biol Sci*, 367(1587):343–353, February 2012. doi: 10.1098/rstb.2011.0198. URL <http://www.ncbi.nlm.nih.gov/pubmed/22201164>.
- M R Nelson, K Bryc, K S King, A Indap, A R Boyko, J Novembre, L P Briley, Y Maruyama, D M Waterworth, G Waeber, P Vollenweider, J R Oksenberg, S L Hauser, H A Stirnadel, J S Kooner, J C Chambers, B Jones, V Mooser, C D Bustamante, A D Roses, D K Burns, M G Ehm, and E H Lai. The Population Reference Sample, POPRES: a resource for population, disease, and pharmacological genetics research. *Am J Hum Genet*, 83(3):

- 347–358, September 2008. doi: 10.1016/j.ajhg.2008.08.005. URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2556436/?tool=pubmed>.
- John Novembre and Matthew Stephens. Interpreting principal component analyses of spatial population genetic variation. *Nature genetics*, 40(5):646–649, 2008.
- John Novembre, Toby Johnson, Katarzyna Bryc, Zoltán Kutalik, Adam R Boyko, Adam Auton, Amit Indap, Karen S King, Sven Bergmann, Matthew R Nelson, et al. Genes mirror geography within europe. *Nature*, 456(7218):98–101, 2008.
- Timothy Paape, Peng Zhou, Antoine Branca, Roman Briskine, Nevin Young, and Peter Tiffin. Fine-scale population recombination rates, hotspots, and correlates of recombination in the *Medicago truncatula* genome. *Genome Biology and Evolution*, 4(5):726–737, 2012. doi: 10.1093/gbe/evs046. URL <http://gbe.oxfordjournals.org/content/4/5/726.abstract>.
- Nick Patterson, Alkes L Price, and David Reich. Population structure and eigenanalysis. *PLoS Genetics*, 2(12):e190, 12 2006. doi: 10.1371/journal.pgen.0020190. URL <http://dx.plos.org/10.1371%2Fjournal.pgen.0020190>.
- J B Pease and M W Hahn. More accurate phylogenies inferred from low-recombination regions in the presence of incomplete lineage sorting. *Evolution*, 67(8):2376–2384, August 2013. doi: 10.1111/evo.12118. URL <http://www.ncbi.nlm.nih.gov/pubmed/23888858>.
- John E Pool. The mosaic ancestry of the *Drosophila* Genetic Reference Panel and the *D. melanogaster* reference genome reveals a network of epistatic fitness interactions. *Molecular Biology and Evolution*, 32(12):3236–3251, 2015. doi: 10.1101/014837. URL <http://mbe.oxfordjournals.org/content/32/12/3236.abstract>.
- John E. Pool, Russell B. Corbett-Detig, Ryuichi P. Sugino, Kristian A. Stevens, Charis M. Cardeno, Marc W. Crepeau, Pablo Duchon, J. J. Emerson, Perot Saelao, David J. Begun, and Charles H. Langley. Population genomics of sub-Saharan *Drosophila melanogaster*: African diversity and non-African admixture. *PLoS Genet*, 8(12):1–24, 12 2012. doi: 10.1371/journal.pgen.1003080. URL <http://dx.doi.org/10.1371%2Fjournal.pgen.1003080>.
- Alkes L Price, Nick J Patterson, Robert M Plenge, Michael E Weinblatt, Nancy A Shadick, and David Reich. Principal components analysis corrects for stratification in genome-wide association studies. *Nature genetics*, 38(8):904–909, 2006.
- Yixuan Qiu and Jiali Mei. *RSpectra: Solvers for Large Scale Eigenvalue and SVD Problems*, 2016. URL <https://CRAN.R-project.org/package=RSpectra>. R package version 0.11-0.



- Sam T. Roweis and Lawrence K. Saul. Nonlinear dimensionality reduction by locally linear embedding. *Science*, 290(5500):2323–2326, 2000. ISSN 0036-8075. doi: 10.1126/science.290.5500.2323. URL <http://science.sciencemag.org/content/290/5500/2323>.
- Fabian Staubach, Anna Lorenc, Philipp W. Messer, Kun Tang, Dmitri A. Petrov, and Diethard Tautz. Genome patterns of selection and introgression of haplotypes in natural populations of the house mouse (*Mus musculus*). *PLoS Genet*, 8(8):e1002891, 08 2012. doi: 10.1371/journal.pgen.1002891. URL <http://dx.doi.org/10.1371/journal.pgen.1002891>.
- Haibao Tang, Vivek Krishnakumar, Shelby Bidwell, Benjamin Rosen, Agnes Chan, Shiguo Zhou, Laurent Gentzbittel, Kevin L Childs, Mark Yandell, Heidrun Gundlach, et al. An improved genome release (version mt4. 0) for the model legume *Medicago truncatula*. *BMC genomics*, 15(1):1, 2014.
- Benjamin Vernot and Joshua M. Akey. Resurrecting surviving neandertal lineages from modern human genomes. *Science*, 2014. doi: 10.1126/science.1245938. URL <http://www.sciencemag.org/content/early/2014/01/28/science.1245938.abstract>.
- Ian J. Wang and Gideon S. Bradburd. Isolation by environment. *Molecular Ecology*, 23(23):5649–5662, 2014. ISSN 1365-294X. doi: 10.1111/mec.12938. URL <http://dx.doi.org/10.1111/mec.12938>.
- Andreas Weingessel and Kurt Hornik. Local PCA algorithms. *Neural Networks, IEEE Transactions on*, 11(6):1242–1250, 2000.
- Sewall Wright. The genetical structure of populations. *Annals of Eugenics*, 15(1):323–354, 1949. ISSN 2050-1439. doi: 10.1111/j.1469-1809.1949.tb02451.x. URL <http://dx.doi.org/10.1111/j.1469-1809.1949.tb02451.x>.
- Jian Yang, Noah A Zaitlen, Michael E Goddard, Peter M Visscher, and Alkes L Price. Advantages and pitfalls in the application of mixed-model association methods. *Nat Genet*, 46(2):100–106, February 2014. ISSN 10614036. URL <http://dx.doi.org/10.1038/ng.2876>.

## A Choosing window length

The choice of window length entails a balance between signal and noise. In very short windows, genealogies of the samples will only be represented by a few trees, so variation between windows represents demographic noise rather than meaningful variation in patterns of relatedness. Longer windows generally have more distinct trees (and SNPs), allowing for less noisy estimation of local patterns of relatedness. However, to better resolve

meaningful signal, i.e., differences in patterns of relatedness along the genome, we would like reasonably short windows.

Since we summarize patterns of relatedness using relative positions in the principal component maps, we quantify “noise” as the standard error of a sample’s position on PC1 in a particular window, averaged across windows and samples, and “signal” as the standard deviation of the sample’s position on PC1 over all windows, averaged over samples. (Recall that the signs for PCs are chosen to match each other.) Then, the mean variance across windows is

$$\sigma_{\text{signal}}^2 = \frac{1}{N} \sum_{j=1}^N \frac{1}{L} \sum_{i=1}^L (PC1_{ij} - \overline{PC1}_j)^2,$$

where  $PC1_{ij}$  is the position of the  $i^{\text{th}}$  individual on PC1 in window  $j$ , and  $\overline{PC1}_j = (1/N) \sum_{i=1}^N PC1_{ij}$ . We estimate the standard error for each  $PC1_{ij}$  using the block jack-knife (Busing et al. 1999; Efron 1982): we divide the  $j^{\text{th}}$  window into 10 equal-sized pieces, and let  $PC1_{ij,k}$  denote the first principal component of this region found after removing the  $k^{\text{th}}$  piece; then the estimate of the squared standard error is  $\sigma_{ij}^2 = \frac{9}{10} \sum_{k=1}^{10} (PC1_{ij,k} - \frac{1}{10} \sum_{\ell=1}^{10} PC1_{ij,\ell})^2$ . Averaging over samples and windows,

$$\sigma_{\text{noise}}^2 = \frac{1}{N} \sum_{j=1}^N \frac{1}{L} \sum_{i=1}^L \sigma_{ij}^2.$$

For the main analysis, we defined windows to each consist of the same number of neighboring SNPs, and calculated  $\sigma_{\text{signal}}^2$  and  $\sigma_{\text{noise}}^2$  for a range of window sizes (i.e., numbers of SNPs). For our main results we chose the smallest window for which  $\sigma_{\text{signal}}^2$  was consistently larger than  $\sigma_{\text{noise}}^2$  (but checked other sizes); the values for various window sizes across *Drosophila* chromosomes are shown in Table S1. In the cases we examined, we found nearly identical results after varying window size, and choosing windows to be of the same physical length (in bp) rather than in numbers of SNPs.

## B Weighted PCA

Principal components analysis can be thought of as finding a good low-dimensional matrix factorization (Engelhardt and Stephens 2010) that well-approximates the original data in the least-squares sense: if  $C$  is the  $N \times N$  genetic covariance matrix, then to find the top  $k$  principal components, we find an orthogonal  $N \times k$  matrix  $U$ , and a  $k \times k$  diagonal matrix  $\Lambda$  with diagonal entries  $\Lambda_{ii} = \lambda_i$  to minimize

$$\|C - U\Lambda U^T\|^2 = \sum_{ij} \left( C_{ij} - \sum_m \lambda_m U_{im} U_{jm} \right)^2. \quad (3)$$

chrom. arm		window length (SNPs)				
		100	500	1,000	10,000	100,000
2L	$\sigma_{\text{noise}}^2$	2.05	1.64	1.18	0.17	0.04
	$\sigma_{\text{signal}}^2$	2.76	2.69	2.23	0.68	0.31
2R	$\sigma_{\text{noise}}^2$	2.18	1.92	1.63	0.58	0.13
	$\sigma_{\text{signal}}^2$	2.78	2.70	2.65	2.31	1.82
3L	$\sigma_{\text{noise}}^2$	2.08	2.00	1.64	0.73	0.25
	$\sigma_{\text{signal}}^2$	2.60	2.52	2.40	1.68	1.89
3R	$\sigma_{\text{noise}}^2$	1.95	1.76	1.44	0.59	0.20
	$\sigma_{\text{signal}}^2$	2.58	2.51	2.44	1.96	1.40
X	$\sigma_{\text{noise}}^2$	2.48	2.04	1.54	1.62	0.17
	$\sigma_{\text{signal}}^2$	2.61	2.43	2.30	0.32	1.14

Table S1: Measures of signal and noise, computed separately for each chromosome arm in the *Drosophila* dataset, at different window sizes. All values are multiplied by 1,000 (so typical variation is of order of 50% of the actual values). Starting at windows of 1,000 SNPs, the signal (variation of PC1 between windows) starts to be substantially larger than the noise (standard error of PC1 for each window).

The columns of  $U$ , known as the principal components, are the eigenvectors of  $C$ , the entries of  $\lambda$  are the eigenvalues of  $C$ , and the proportion of variance explained by the  $m^{\text{th}}$  component is

$$\frac{\lambda_m^2}{\sum_{\ell} \lambda_{\ell}^2} = \frac{\sum_{ij} (\lambda_m U_{im} U_{jm})^2}{\sum_{ij} C_{ij}^2}.$$

Thinking about the problem as a least-squares approximation problem makes it clear why unbalanced sample sizes can result in undesirable outcomes. If we want to describe variation *between* populations, but 80% of the samples are from a single population, then unless populations are highly differentiated, a better approximation to  $C$  may be obtained by using the columns of  $U$  to describe variation *within* the overrepresented population rather than between the populations. A common workaround is to remove samples, but a more elegant solution can be found by reweighting the objective function in (3). Let  $w_i$  be a weight associated with sample  $i$ ,  $W$  the diagonal matrix with  $w$  along the diagonal, and instead seek to minimize

$$\|W^{1/2}(C - U\Lambda U^T)W^{1/2}\|^2 = \sum_{ij} w_i w_j \left( G_{ij} - \sum_m \lambda_m U_{im} U_{jm} \right)^2, \quad (4)$$

and now for convenience we require  $U$  to be orthogonal in  $\ell_2(w)$ , i.e., that  $U^T W U = I$ . We then would choose  $w$  to give roughly equal weight to each *population*, instead of each

individual. We have used with good results the weightings  $w_i = 1/\max(10, n_i)$ , where  $n_i$  is, if there are discrete populations, the number of samples in the same population as sample  $i$ ; or, for continuously sampled individuals, the number of samples within a certain distance of sample  $i$ .

To solve (4), let  $\lambda$  and  $V$  denote the top  $k$  eigenvalues and eigenvectors of  $W^{1/2}CW^{1/2}$ , so that  $V\Lambda V^T$  is the rank  $k$  matrix closest in least squares to  $W^{1/2}CW^{1/2}$ ; so if we define  $U = W^{-1/2}V$  then  $U^T W U = V^T V = I$ , and

$$W^{-1/2}V\Lambda V^T W^{-1/2} = U\Lambda U^T$$

is the low-dimensional approximation to  $C$ . The proportion of variance explained is calculated from eigenvalues as before, but has the interpretation

$$\frac{\lambda_m^2}{\sum_{\ell} \lambda_{\ell}^2} = \frac{\sum_{ij} w_i w_j (\lambda_m U_{im} U_{jm})^2}{\sum_{ij} w_i w_j C_{ij}^2}.$$

In our R implementation we use the Spectra library (Qiu and Mei 2016) to find only the top  $k$  eigenvectors.

## C Supplementary Figures

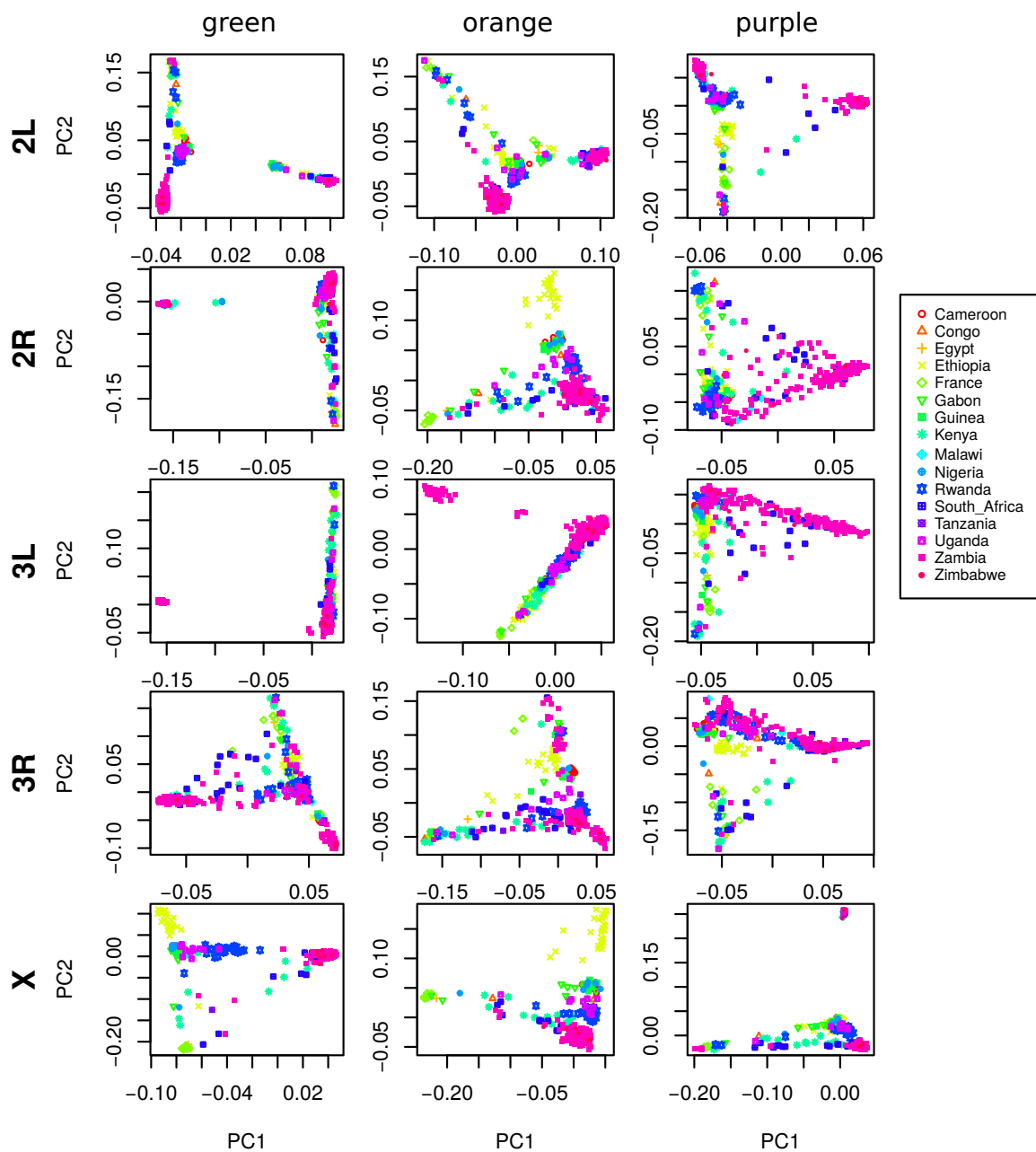


Figure S1: PCA plots for the three sets of genomic windows colored in Figure 2, on each chromosome arm of *Drosophila melanogaster*. In all plots, each point represents a sample. The first column shows the combined PCA plot for windows whose points are colored green in Figure 2; the second is for orange windows; and third is for purple windows.

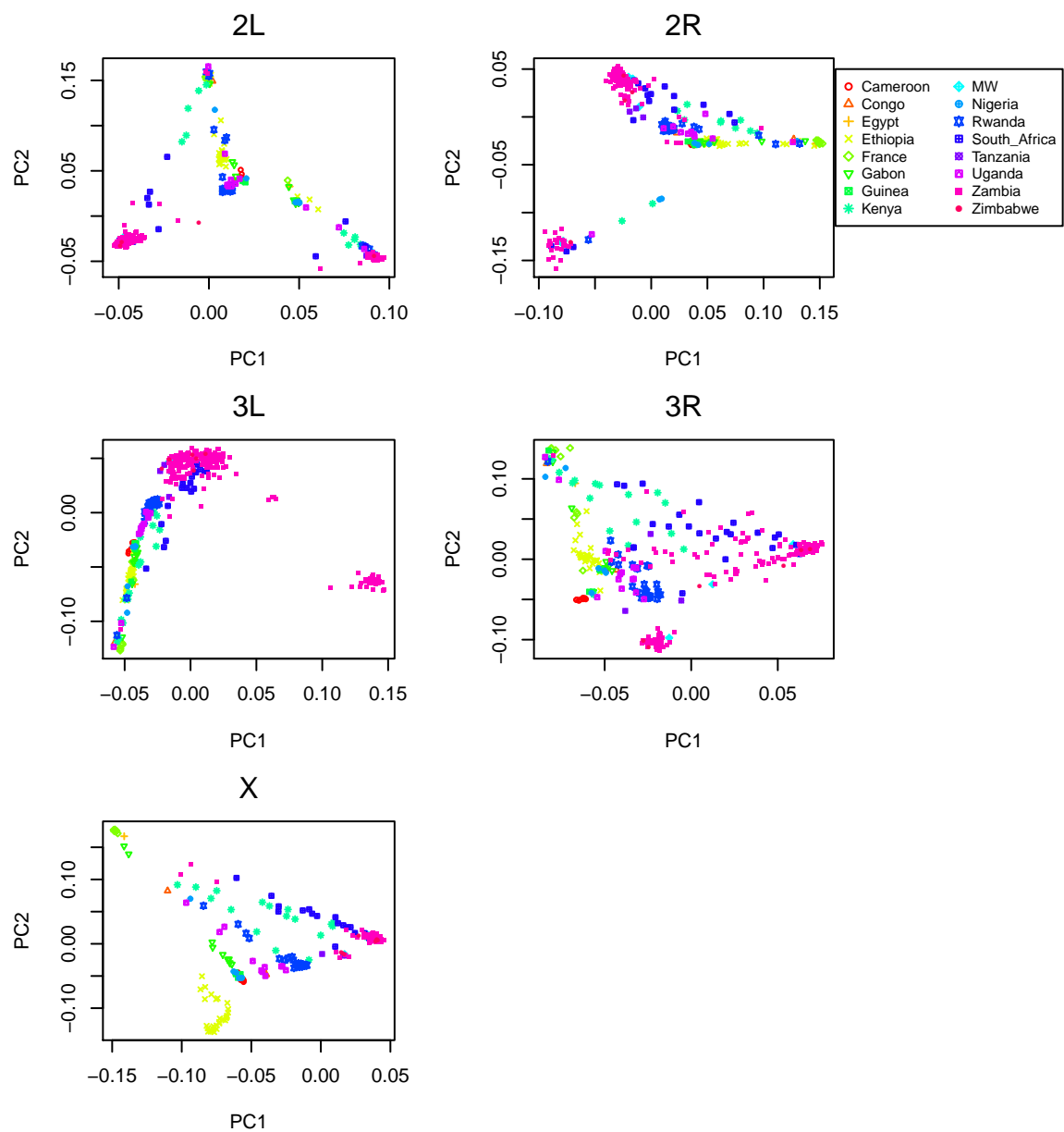


Figure S2: PCA plots for chromosome arms 2L, 2R, 3L, 3R and X of the *Drosophila melanogaster* dataset.

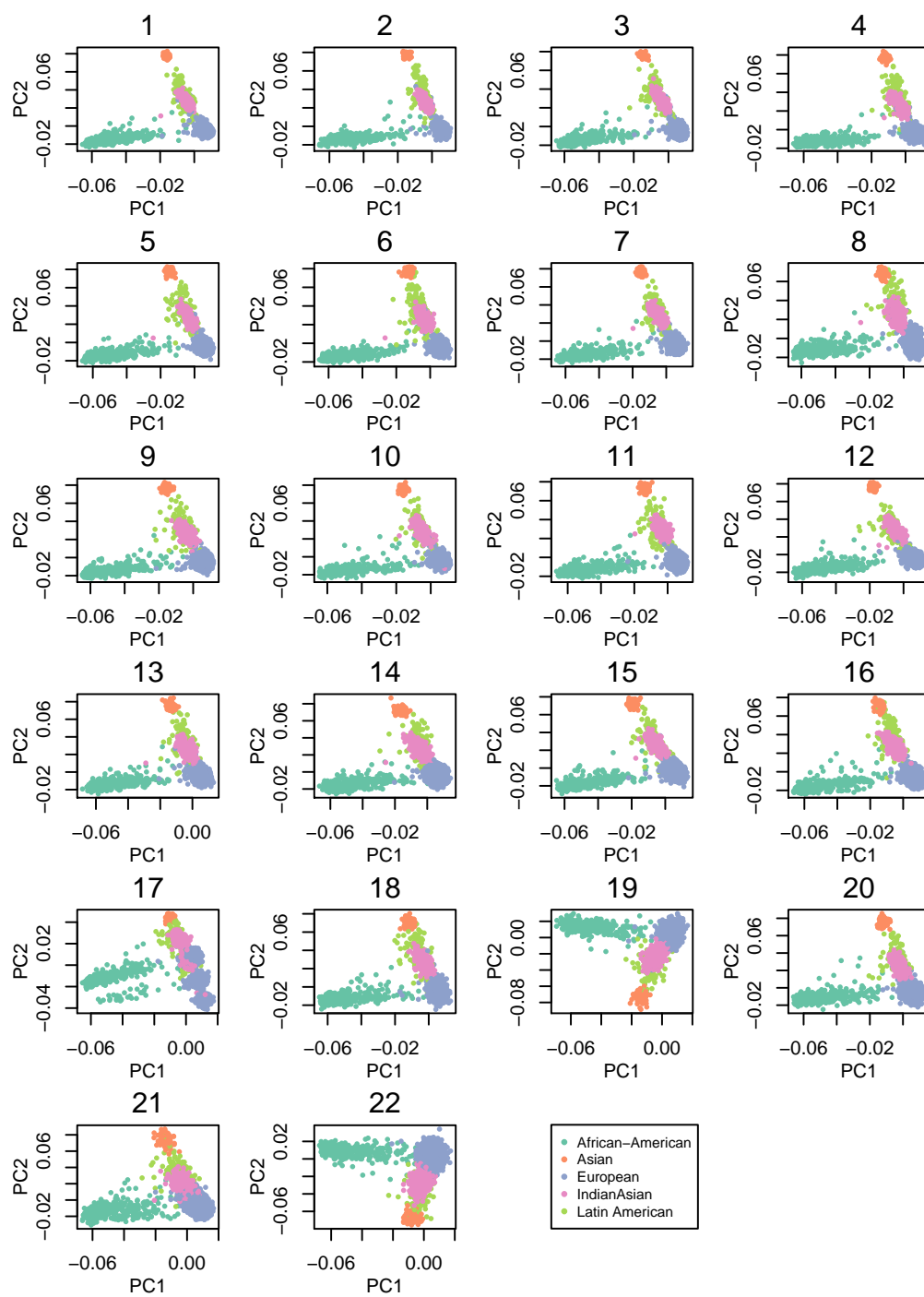


Figure S3: PCA plots for all 22 huan autosomes from the POPRES data.

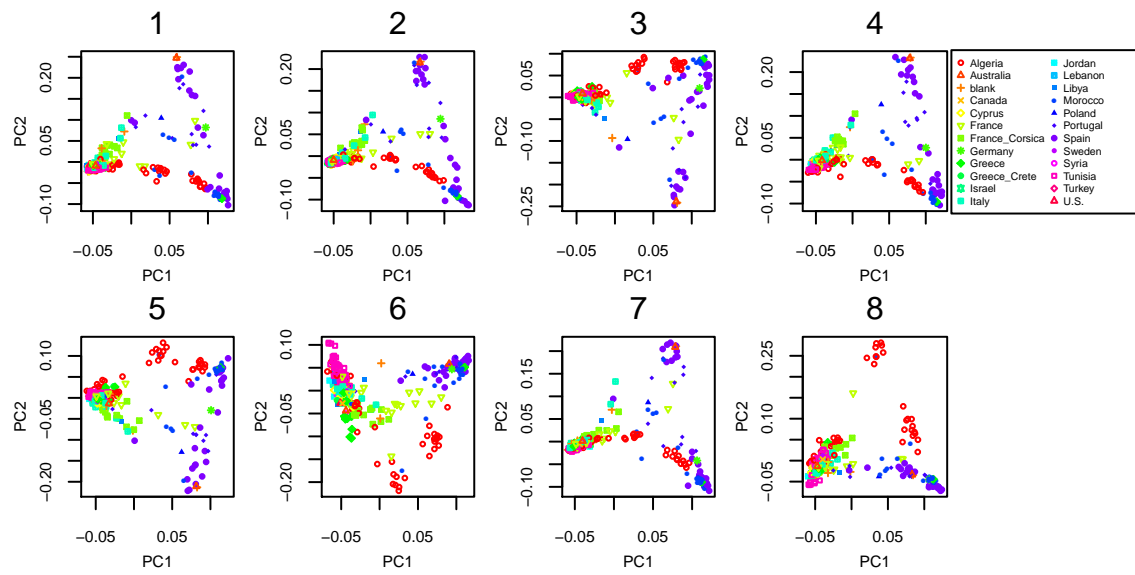


Figure S4: PCA plots for all 8 chromosomes in the *Medicago truncatula* dataset.



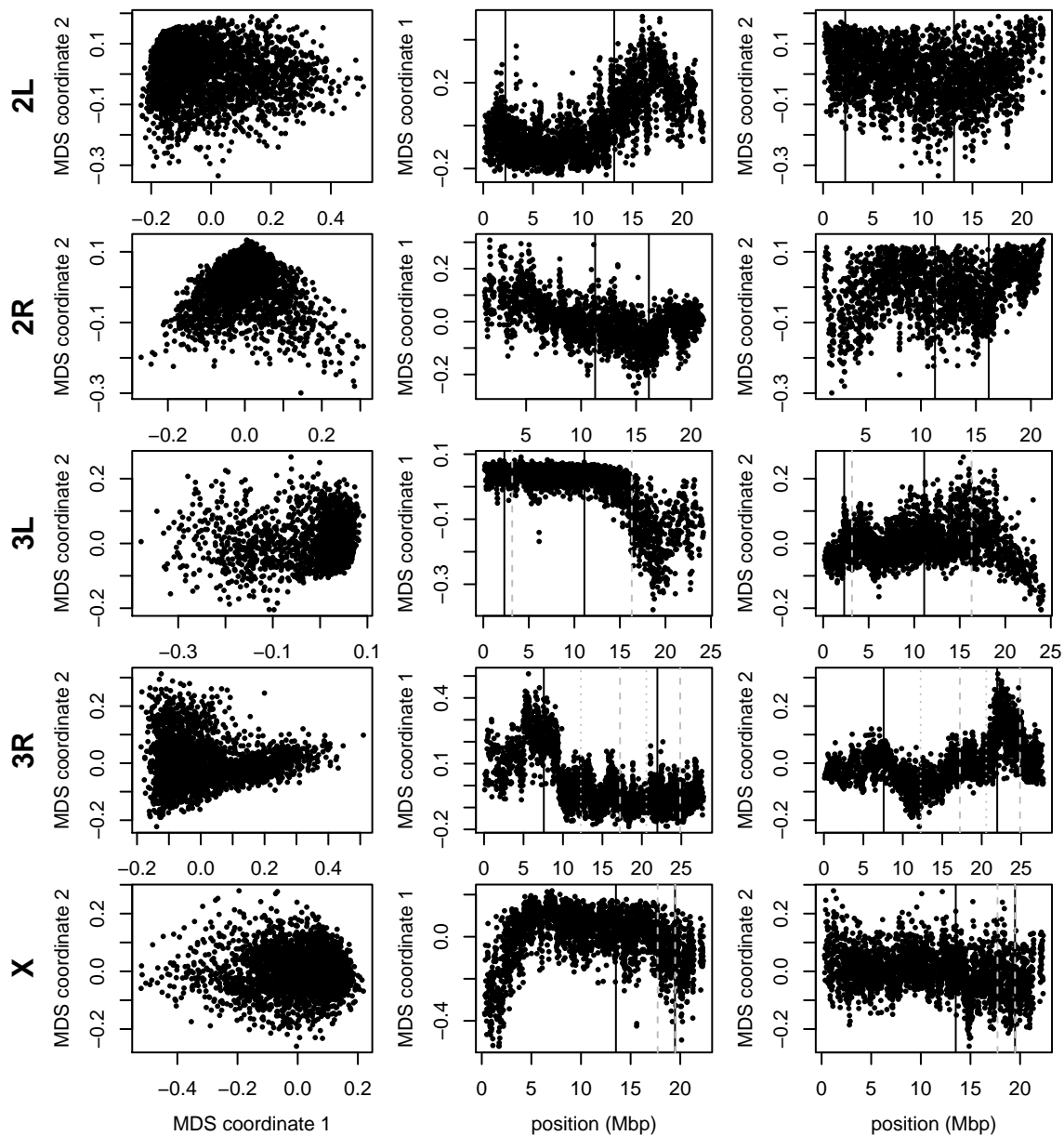


Figure S5: Variation in structure for windows of 1,000 SNPs across *Drosophila melanogaster* chromosome arms: without inversions. As in Figure 2, but after omitting for each chromosome arm individuals carrying the less frequent orientation of any inversions on that chromosome arm. The values differ from those in 4 in the window size used and that some MDS values were inverted (but relative orientation is meaningless as chromosome arms were run separately, unlike for *Medicago*). In all plots, each point represents one window along the genome. The first column shows the MDS visualization of relationships between windows, and the second and third columns show the midpoint of each window against the two MDS coordinates; rows correspond to chromosome arms. Colors are consistent for plots in each row. Vertical lines show the breakpoints of known polymorphic inversions.

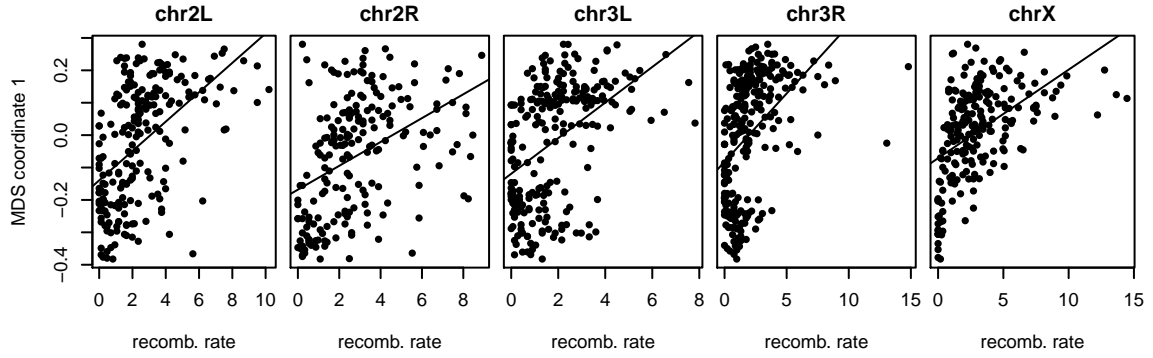


Figure S6: Recombination rate, and the effects of population structure for *Drosophila melanogaster*: this shows the first MDS coordinate and recombination rate (in cM/Mbp), as in Supplementary Figure 4, against each other. Since the windows underlying estimates of Supplementary Figure 4 do not coincide, to obtain correlations we divided the genome into 100Kbp bins, and for each variable (recombination rate and MDS coordinate 1) averaged the values of each overlapping bin with weight proportional to the proportion of overlap. The correlation coefficient and  $p$ -values for each linear regression are as follows: 2L: correlation = 0.52,  $r^2 = 0.27$ ; 2R: correlation = 0.43,  $r^2 = 0.18$ ; 3L: correlation = 0.47,  $r^2 = 0.21$ ; 3R: correlation = 0.46,  $r^2 = 0.21$ ; X: correlation = 0.50,  $r^2 = 0.24$ .

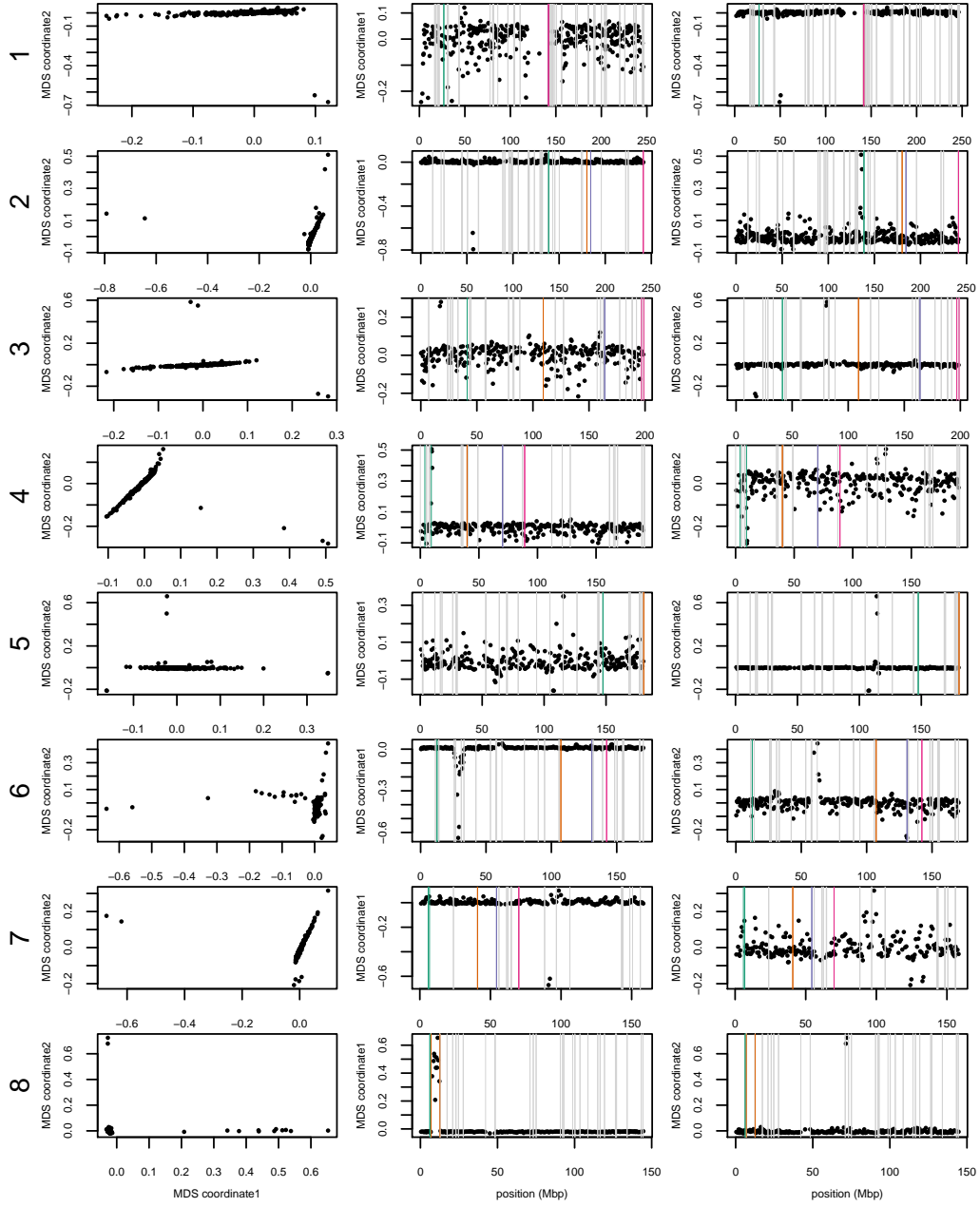


Figure S7: MDS plots for human chromosome 1-8. The first column shows the MDS visualization of relationships between windows, and the second and third columns show the midpoint of each window against the two MDS coordinates; rows correspond to chromosomes. Colorful vertical lines show the breakpoints of known valid inversions, while grey vertical lines show the breakpoints of predicted inversions.

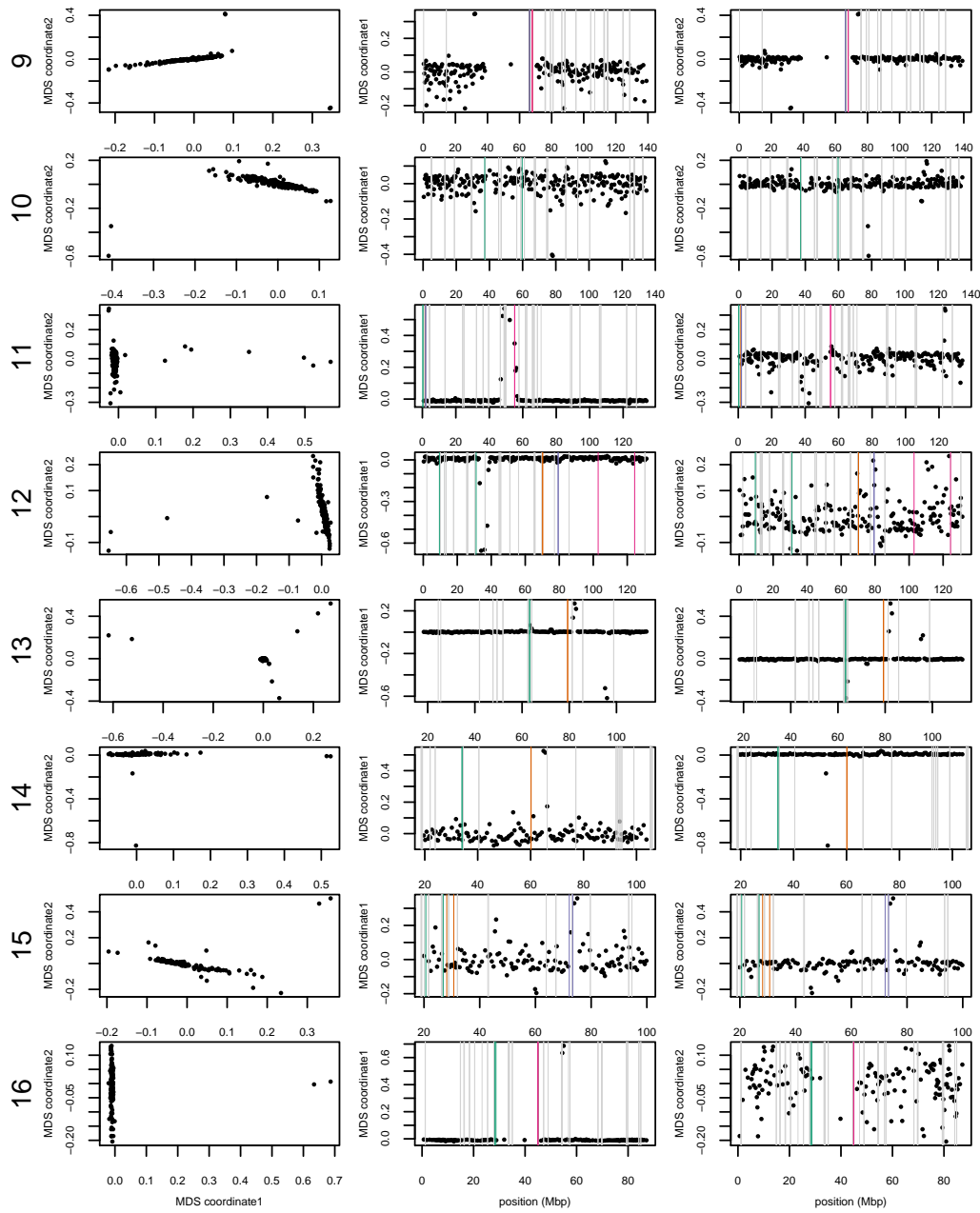


Figure S8: MDS plots for human chromosome 9-16.

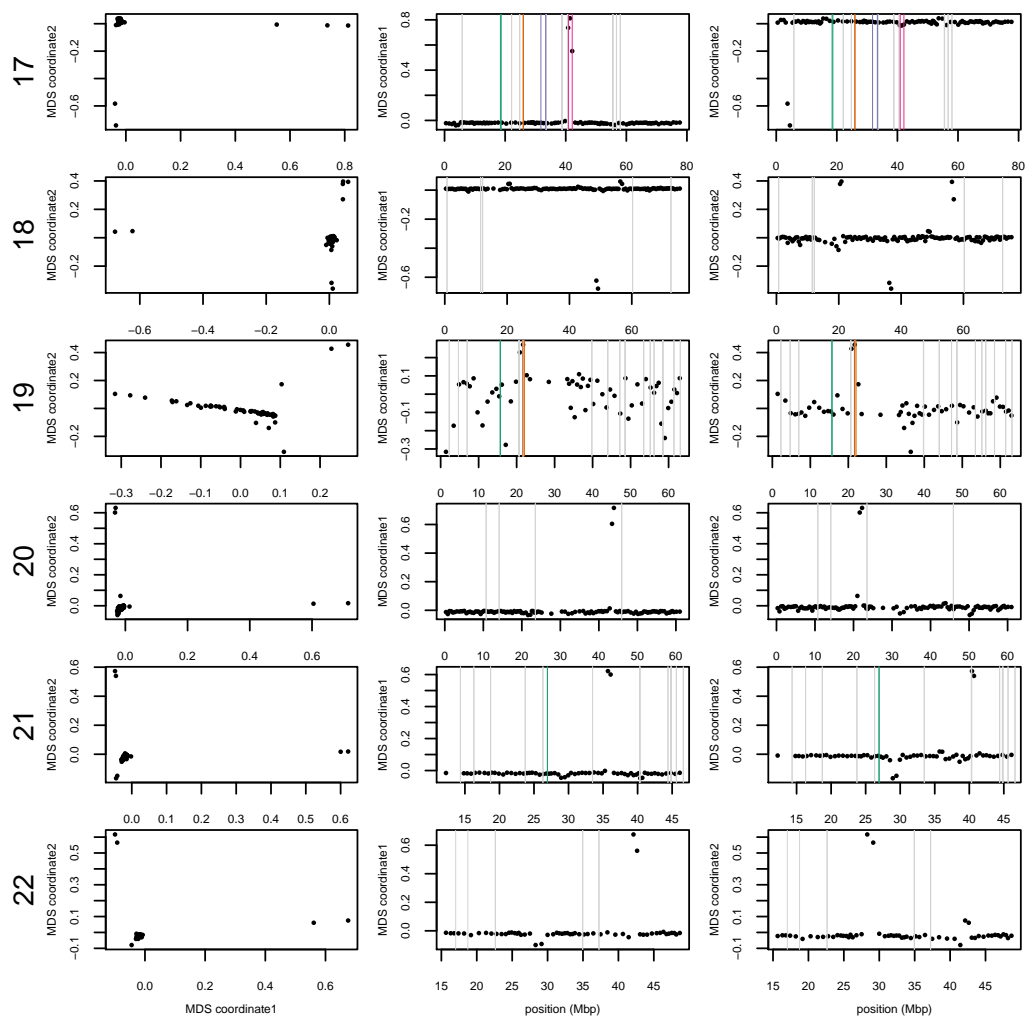


Figure S9: MDS plots for human chromosome 17-22.

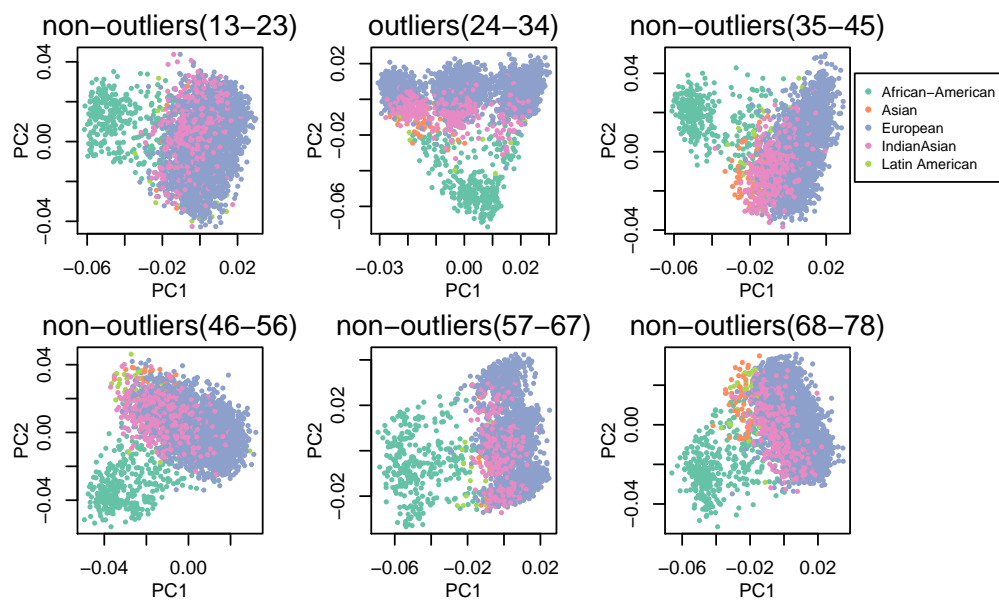


Figure S10: PCA plots comparison for outlier windows and non-outlier windows.

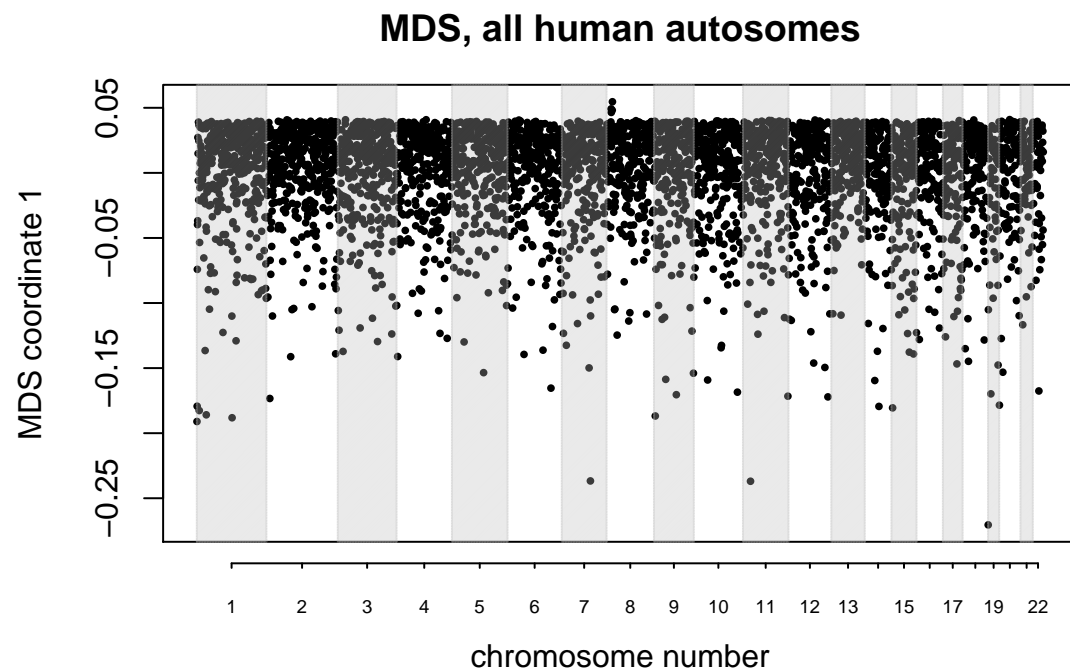


Figure S11: MDS visualization of variation in the effects of population structure amongst windows across *all* human autosomes simultaneously. The small group of windows with positive outlying MDS values lie around the inversion at 8p23.

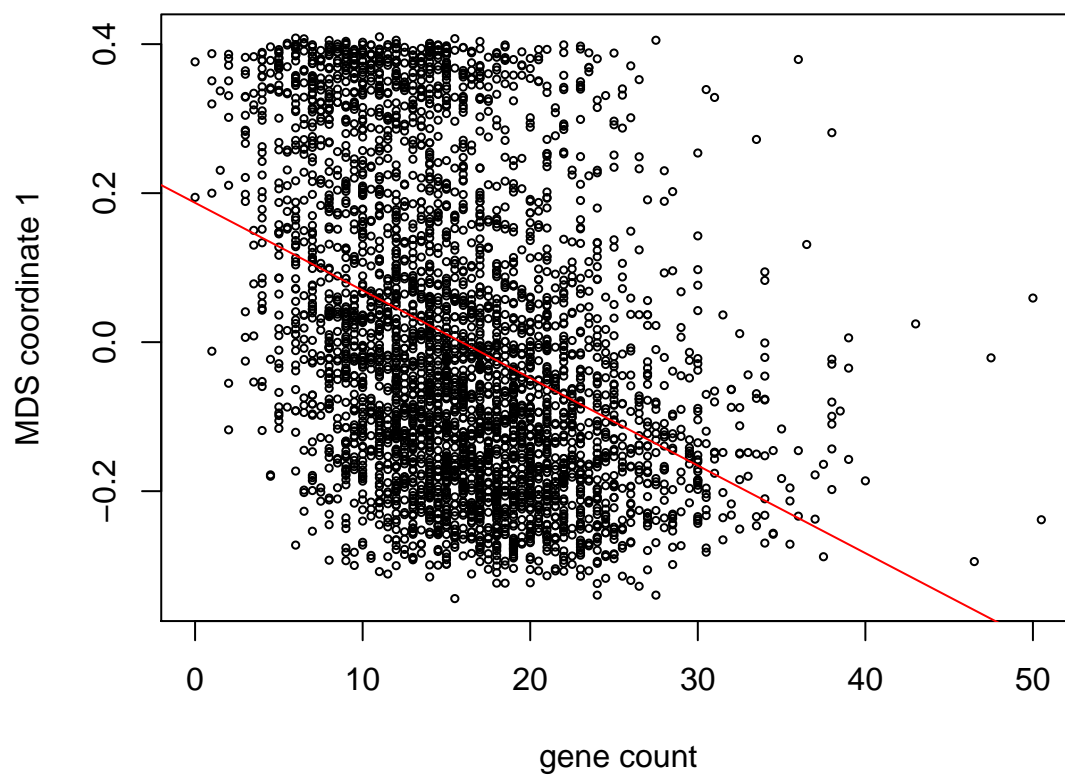


Figure S12: First MDS coordinate against gene density for all 8 chromosomes of *M. truncatula*. The first MDS coordinate is significantly correlated with gene count ( $r = 0.149$ ,  $p = 2.2 \times 10^{-16}$ ).



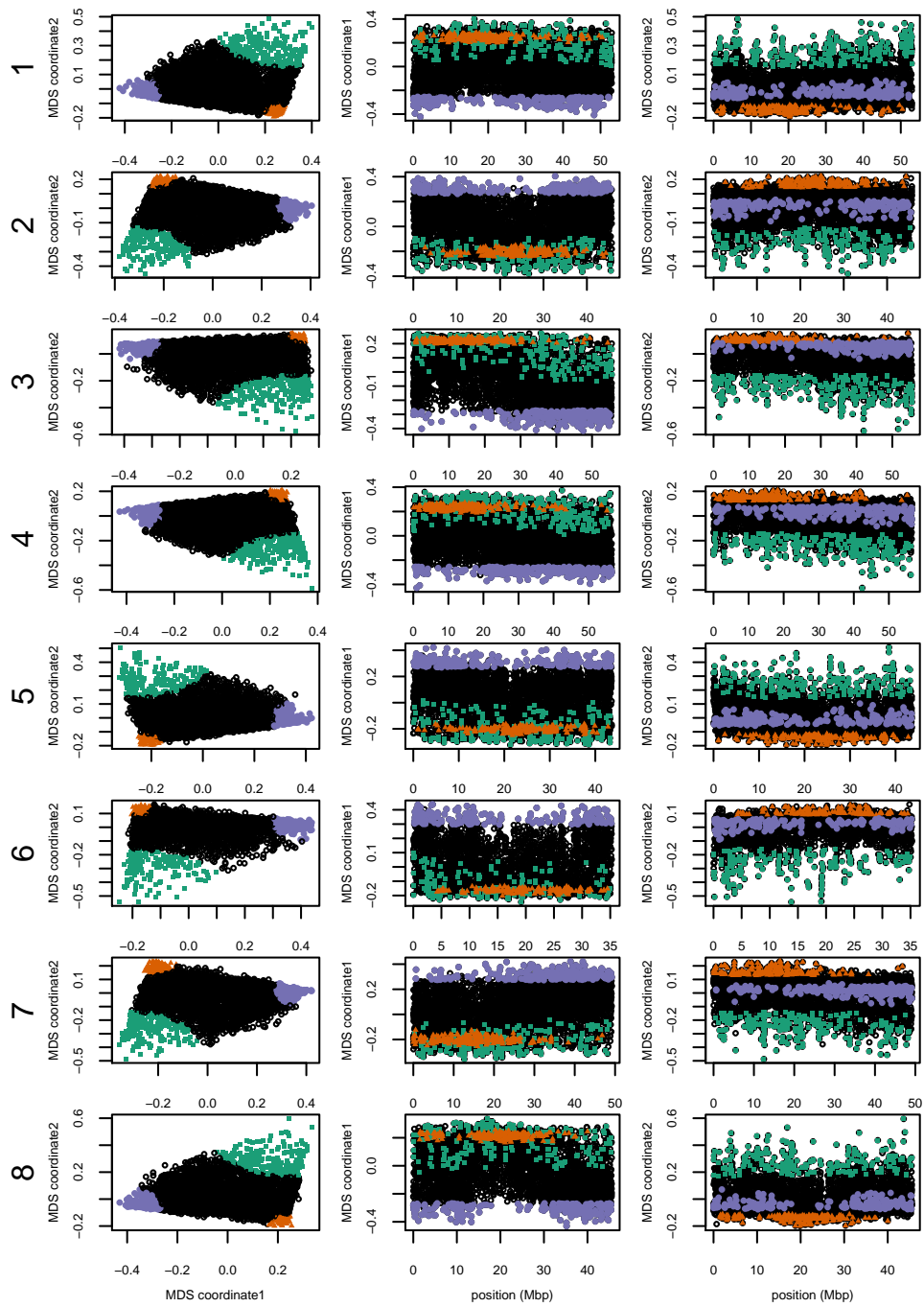


Figure S13: MDS visualizations of the effects of population structure for all 8 chromosomes of the *Medicago truncatula* data.

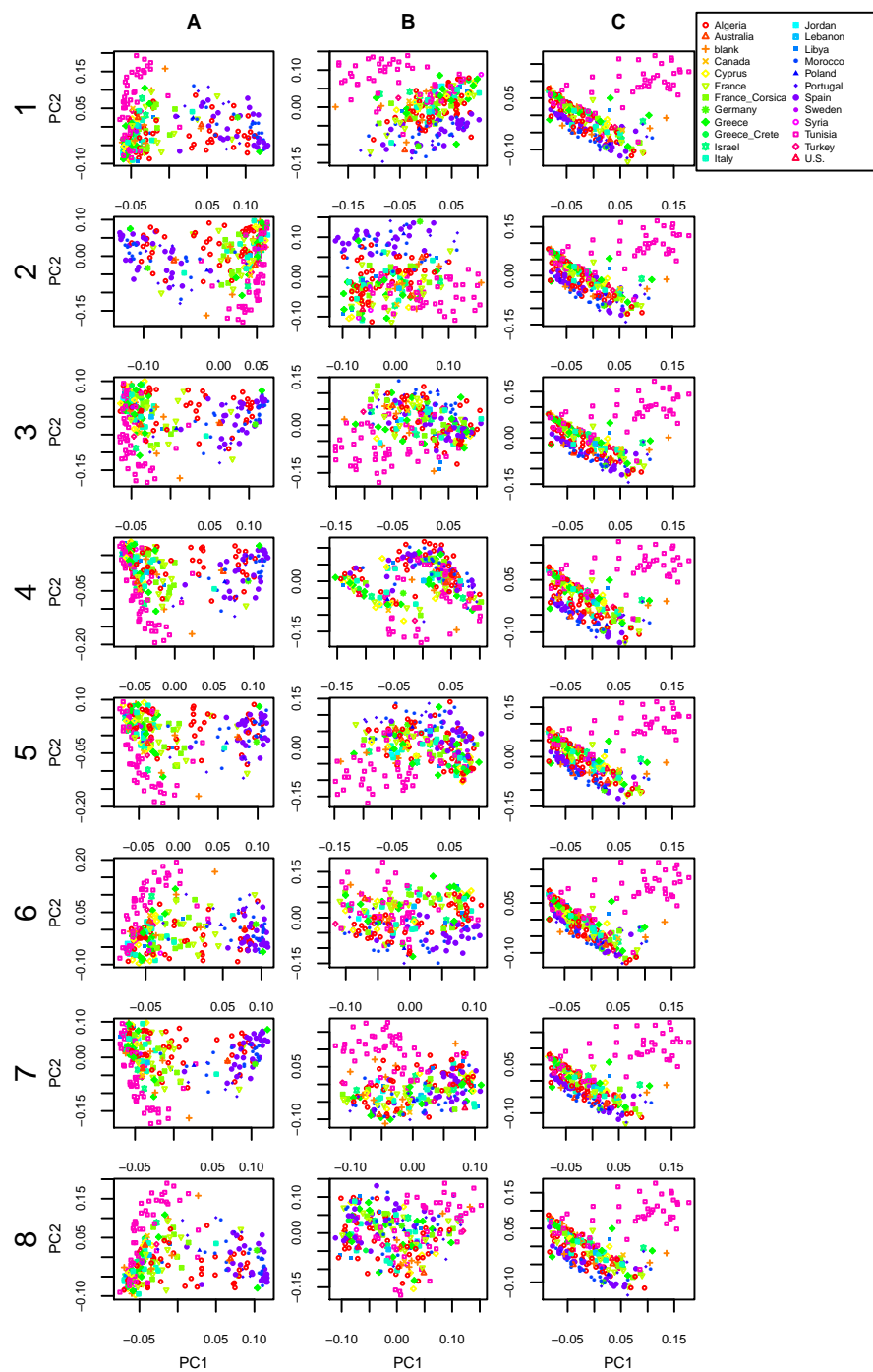


Figure S14: PCA plots for regions colored in Figure S13 on all 8 chromosomes of *Medicago truncatula*: (A) green, (B) orange, and (C) purple.