

Constrained ordination method as a tool for performing genome scan and environmental genomic studies

Eric Bazin,^{*,1} Keurcien Luu,² Michael G. B. Blum,²

¹LECA, Université de Grenoble

²TIMC, Université de Grenoble

***Corresponding author:** E-mail: eric.bazin@univ-grenoble-alpes.fr

Associate Editor:

Abstract

Key words:

1 Introduction

2 Performing genome scan in order to detect
3 genomic region of interest is a common task in
4 population genomic area (????). Some methods
5 aim at detecting genes that has suffered from
6 a loss of genetic diversity and increase of
7 linkage disequilibrium following the appearance
8 of a beneficial allele and its spread by the
9 mean of selective sweep. Others aim at picking
10 up alleles with strong correlation with some
11 environmental variable (e.g. Temperature,
12 drought) with the idea that these alleles
13 may confer a selective advantage to the
14 individuals (??). Finally, other methods aim
15 at detecting genomic region involved in local
16 adaptation process. These region should have
17 an increased differentiation between population
18 because different alleles tend to be beneficial
19 in each environment. Differentiation between

20 population is excepted under the hypothesis
21 of geographical isolation. Therefore, this region
22 can be detected by quantifying the level of
23 differentiation using some statistics and detecting
24 the regions with unexpectedly high values. A
25 common statistic and very easily comprehensible
26 in population genetic is Fst. Many methods use
27 this parameter as a basis in many different
28 implementation of genome scans (???). These
29 are model based method where parameters such
30 as Fst are usually inferred using likelihood or
31 Bayesian methods. This mean that users must
32 have some a priori on their parameter value and
33 the best model that fits their data in order to
34 expect the best from their analysis. However, it is
35 often difficult to get a satisfactory a priori picture
36 of the demographic and population structure of
37 the species one is interested in. Indeed many
38 species are not clearly structured in different
39 populations but more or less show a pattern of

40 isolation by distance without clear geographical

© The Author 2013. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution. All rights reserved.
For permissions, please email: journals.permissions@oup.com

41 barrier to gene flow. One solution would be to 73 of interpretation but are not involved in the
42 use more complex model that better reflect reality 74 inference process. It should be recalled that
43 but these are difficult to implement in Bayesian 75 natural selection is the result of a complex set of
44 framework. Additionally, these latter methods are 76 environmental pressures and that it most often
45 very time consuming and the increase of both 77 acts on several characters simultaneously and
46 model complexity and the amount of data to 78 that these characters are encoded by several
47 analyze in terms of the number of individuals and 79 genes which generally have weak effects. In
48 loci makes them more and more difficult to use. 80 order to extract the maximum of all available
49 A new path has opened recently with the use 81 information, it seems therefore necessary to use
50 of multivariate methods. The idea is to capture 82 approaches that are able to compile all kind of
51 the whole genome geographic structure using an 83 variable (e.g. alleles, phenotypic measurement,
52 ordination method such as ACP. Following this 84 biotic and abiotic variables). One natural way
53 analysis, outliers loci are detected if they have 85 to overcome this limitation would be to use
54 extremely high correlation with one or more 86 more sophisticated ordination method than ACP
55 ordination axis (??). These are very efficient 87 like methods. Constrained ordination methods
56 methods and simulations have shown that while 88 (i.e. Redundancy Analysis, RDA, Canonical
57 they are very fast, they show similar efficiency 89 Correspondence Analysis, CCA) are well-known
58 than classical Bayesian method and sometimes 90 set of approaches in Ecology for instance to
59 perform better when the simulated demographic 91 explain the species distribution pattern by
60 model drift from the model implemented in 92 the mean of environmental data. They have
61 bayesian method, usually the island model. 93 specifically been designed in order to deal with
62 For instance ? have shown their method to 94 biological complexity. In the population genomic
63 be better when population are structured in 95 era, it seems that data amount, complexity
64 hierarchical set or in isolation by distance pattern. 96 and heterogeneity is often a limitation to the
65 Nevertheless, one conundrum of such approaches 97 use of inference methods based on classical
66 is the difficulty to interpret ordination axis in term 98 population genetic models. Although they are
67 of ecological meanings. These are usually tight to 99 more difficult to interpret, such approaches
68 geographical axis (latitudinal or longitudinal) but 100 would be complementary to the model based
69 they are not necessary linked to an environmental 101 method because of their long-term use in ecology
70 variable such as Temperature, drought, diet habit, 102 and their efficiency on complex and large
71 etc. Therefore, when this information exists, 103 datasets. These method have sometimes been
72 it has to be a posteriori used as a mean 104 used in population genomic studies, not as a

genome scan but in order to quantify multilocus adaptation to an environmental gradient (????). These studies whereby relationships between environmental data and large multilocus data is explored are becoming more and more popular and are often coined as Ecological Genomics or Landscape Genomics studies. However the concept of using constrained ordination methods to analyse genomic data has never been tested on simulated datasets. This paper aims at filling this gap. First, we show how one can make use of a constrained ordination method namely Redundancy Analysis (RDA) as an efficient and robust genome scan method. We discarded the other constraint ordination methods such as CCA since they are very similar in their principles. RDA has already been used for instance by ? to perform genome scan in order to detect loci involved in the adaptation to climate in *Arabidopsis thaliana*. Outliers were identified as SNPs with the greatest squared scores along the first RDA axis (i.e. those in the 0.5 % tail). We build on this idea to develop a comprehensive and robust statistical test that allows to search for outliers on an arbitrary number of RDA axis simultaneously and allows to control precisely for the false discovery rate. Using simulations, we show that it has better results than PCA-based method. Second, thanks to these simulations, we show that RDA can indeed help to identify important environmental gradient that better explain the adaptive variation in the data. It is therefore a proof of concept of the idea of using constrained ordination method as an environmental genomic tool to identify relevant selective gradient in the environmental data. Finally, to give a concrete illustration of RDA approach in population genomics, we apply this method to the detection of outliers on a real data set.

Material and method

Genome scan

Redundancy analysis (RDA) was first introduced by (?) and is clearly described in (?) section 11.1. It is the direct extension of multiple regression to the modeling of multivariate response data. Typically the data to be analysed are separated in two sets, a response matrix Y of variable to be explained (e.g. species abundance in a set of sites; m sites and n species) and an explanatory matrix X (e.g. a set of environmental variable within each site; m sites and p environment). In the following analysis, species are replaced by loci and sites by individuals. In other word, we wish to project on a reduced space the proportion of variance in genetic difference between individuals which is better explained by environmental data. After this ordination, we follow the ? methodology to compute pvalues. First we compute the test statistic by regressing each of the p SNPs by the K ordination axis X_1, \dots, X_K .

$$G_j = \sum_{k=1}^K \beta_{jk} X_k + \epsilon_j, j=1, \dots, p$$

where β_{jk} is the regression coefficient corresponding to the j-th SNP regressed by the k-th ordination axis, and ϵ_j is the residuals

vector. To summarize the result of the regression analysis for the j -th SNP, we return a vector of z -scores $z_j = (z_{j1}, \dots, z_{jK})$ where z_{jk} corresponds to the z -score obtained when regressing the j -th SNP by the k -th ordination axis. The test statistic is a robust Mahalanobis distance D computed using `covRob` function of the `robustR` package. D should be Khi2 distributed after a correction with inflation factor (Luu et al., 2016). P -values are computed using K degree of freedom. We use the FDR approach to control for false positives. Q -value are computed with `qvalueR` package and a loci is considered as an outlier if its q -value is less than 10%. For the analysis of simulated dataset (see below), we retain the first four ordination axis to compute Mahalanobis distances as they seem to explain most of the variance in the data. To perform the ordination, we use the 10th environmental variables as input in the explanatory matrix. In the following example, we don't use phenotypic informations since these informations are often lacking in environmental genomics. Neither we use geographical coordinates (i, j) which is sometimes added to control for the geographical covariation in the differentiation pattern (?).

To emphasize the utility of RDA, we compared to `pcadapt` from which the idea of using multivariate method for genome scan is based. On the simulated dataset, we retain $K=3$ axis to compute Mahalanobis distances as it seems to explain the main amount of variance in the data using scatter plots. To control for false positive, we used the same q -value threshold (i.e. $q=10\%$).

Environmental genomic

Once outliers have been identified, we isolate them in a separate matrix A defining an "adaptively enriched genetic space" as coined by ?. Following their methodology, we perform a second constrained ordination (RDA) on matrix A against environmental data. The rationale of this analysis is to remove neutral variation before performing ordination in order to have a better picture of which environmental gradients have the strongest association with the adaptive genetic space. On the simulated dataset, we report the R^2 statistics between `env1`, `env2` and `env3` and the first three ordination axis to have an idea of which they are better associated with and if the ordination space succeed in separating the environmental effect on different axis.

Simulations

To test for the efficiency of RDA in population genomic, we performed simulations using `simuPop` python library (?). We compared our approach to `PCAdapt` method to perform genome scans. Both approach are equivalent except their ordination method. Finally we use these simulations to evaluate RDA approach as a mean to detect selective environmental gradient. A lattice of 8×8 populations is simulated (i.e. 64 populations in total). Each population is initialized with 200 diploid individual with random genotypes. Migration is set to 0.5 so that population structure

must be very smooth and genetic differentiation must show an isolation by distance pattern over the 64 populations. This is where pcadapt is best designed for. Loci are biallelic (0 or 1) like SNPs. Allele frequency of the whole population is initialized at 0.5. 1000 loci are defined. They are separated in 200 chunks of 5 SNPs in physical linkage with recombination rate between adjacent loci fixed at 0.1. 3 different Traits are coded by a group of 10 different loci. The first trait is coded by loci 1, 11, 21, ..., 91. The trait value is simply the sum of genotype value and therefore can take value between 0 and 20. For the sake of realism, we add to each trait a random noise (non heritable variation) drawn from a normal distribution $N(0,2)$. The second trait is coded by loci 101, 111, ..., 191 and the third is coded by loci 201, 211, ..., 291. Each trait is therefore coded by free recombining SNP loci. In other words, there are 30 coding SNPs among 1000. Selection can have an effect on linked loci, for instance, loci 2, 3, 4 and 5 can be impacted by selection on locus 1. However, recombination is high enough (0.1) to expect a limited linkage effect. We have defined 10 different environmental variables. The first one determines the selective pressure on trait 1, the second one on trait 2 and the third one on trait 3. The first environment variable is a quadratic gradient coded by function $env1 = -(\cos(\theta) * (i - 3.5))^2 - (\sin(\theta) * (j - 3.5))^2 + 18$, $\theta = \pi/2$, i and j being the population indicator on the 8x8 lattice. The second one is a linear plan gradient coded by function $env2 = h * \cos(\theta) * (i - 1) + h * \sin(\theta) * (j - 1) + k$ with $h = 2$, $\theta = \pi/4$ and $k = 3$. The third environment variable simulates a coarse environment with value $env3 = 2$ for all populations except population $(i,j) = (2,2)$, $(2,3)$, $(3,2)$, $(3,3)$, $(6,2)$, $(6,3)$, $(7,2)$, $(7,3)$, $(2,6)$, $(2,7)$, $(3,6)$, $(3,7)$, $(6,6)$, $(6,7)$, $(7,6)$, $(7,7)$ for which $env3 = 18$. Env4, env5 and env6 have exactly the same equation than env1, env2 and env3 respectively. The remaining 4 environment variable are similar to env2 but with different value of h and θ . Env7 has $h = 2$, $\theta = 0$ and $k = 3$. Env8 has $h = 2$, $\theta = \pi/4$ and $k = 0$. Env9 has $h = 1$, $\theta = \pi/4$ and $k = 4$. Env10 has $h = 0.5$, $\theta = \pi/4$ and $k = 8$. Graphical representation of mean environmental value for environment 1, 2 and 3 is given in Fig. ???. Environment 4, 5 and 6 have respectively the same mean value spatial distribution. For a graphical representation of environment 7 to 10, see supplementary material. Environmental equation gives a mean value of the environment variable. To avoid colinearity between environments variable, we added noise by drawing an environment value within a normal distribution $N(\mu = env, \sigma = 1)$. Fitness for each trait is set to be $-e^{((x - env)^2 / (2 * \omega^2))}$, x being the quantitative trait value, env the environmental value and ω is defining selection strength and has been set to 10 which in our experience seems sufficient for loci to be often detected. To get the overall fitness for a given individual, fitness associated to each trait are multiplied. Fitness

are relative and selection arises on parents and determine their number of offsprings. Simulations are made across 500 generations. At the end of simulation, we sample 10 individuals per population. Therefore, we have a sample of 640 individuals with 1000 SNP-like loci.

Real dataset

The Loblolly pine dataset is a sample of 682 individuals genotyped on 1,730 SNPs selected in ESTs (?). 60 climatic variables were available and summarized by the authors in the five first axis of a PCA. The first axis, PC1 is mainly linked to latitude, longitude, temperature, and winter aridity. PC2 is linked to longitude, spring-fall aridity, and precipitation. We inputed the missing data using a very simple algorithm implement in function sing.im of the R package linkim (?). It imputes the missing value based on the observed data proportions. We used $K=4$ axis to compute Malahanobis distances.

The Chinook salmon consists of 19 703 SNP loci genotyped on 1956 total individuals pooled in 46 collections. Five variables (MigDistKM, StreamOrder, bio03, bio17 and bio18) have been used among 24 different climate and environmental variables because they have been tested as significantly associated with the SNP variation rangewide citepHecht2015. MigDistKM stands for Migration distance from collection site to ocean (km), StreamOrder for Stream Order of collection site using Strahler method, bio03 for Isothermality, bio17 for Precipitation of Driest

Quarter (mm) and bio18 for Precipitation of Warmest Quarter (mm). We could have tested more variable but this is just an illustration and is by no mean an extensive study of this species. Since data are pooled, we have randomly created a sample of 100 individuals for each collection based on the allele frequencies to be able to analyze the data following our individual based pipeline. We used $K=4$ axis to compute Malahanobis distances.

Results

Genome scan

When looking at the analysis on one simulation, the pcadapt method seems successful at detecting QTL2 SNPs (Fig. ??) but fails at detecting QTL1 and QTL3 SNPs. On the other hand, RDA succeeds at detecting QTL2 SNPs and also some of the QTL1 and QTL2 SNPs (Fig. ??). The ordination seems to correctly detect environmental variable 1 and 3 as drivers of genetical variance in the data. Over the 100 simulations, we have measured the average FDR and power for both pcadapt and RDA (Fig ??).

Environmental genomics

We then performed a second RDA on the ”adaptively enriched genetic space” as performed by ? on the same simulated dataset as in Fig. ?? and ?? and display its results on Fig. ???. We did the same analyis and measured the mean R^2 between env1, env2 and env3 and each of the first three ordination axis. This is summarized in Fig. ??.

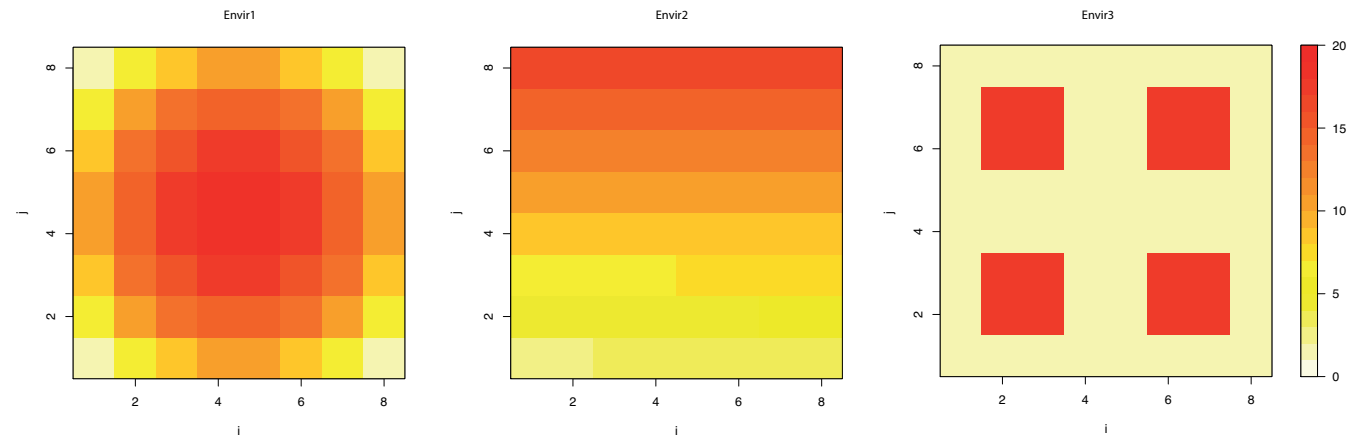
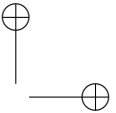
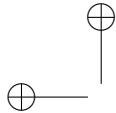


FIG. 1. Graphical representation of mean environmental value for environment 1, 2 and 3

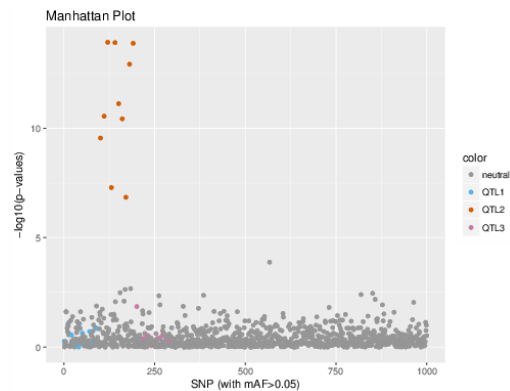


FIG. 2. Manhattan plot of the result of pcadapt on a simulated data set.

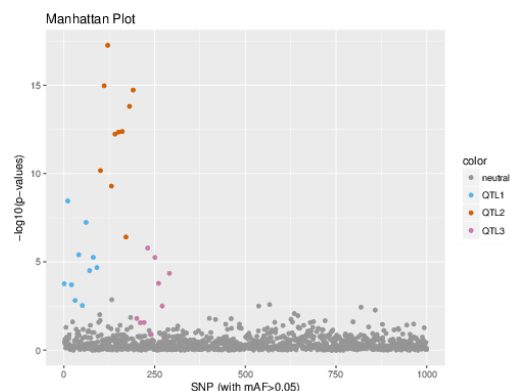


FIG. 3. Manhattan plot of the result of genome scan using RDA on a simulated data set.

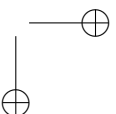
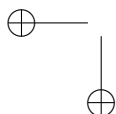
Loblolly Pine

Discussion

Fig ?? shows that pcadapt approach works well when the environmental gradient and the selective

pressures are acting in the same direction than the geographical pattern of isolation by distance. Whereas when the environmental gradient is quadratic on the geographical range (QTL1) or when it is a coarse environment (QTL3). Indeed, we can hypothesize that the PCA ordination fails at orienting the genetic space differentiation into the direction of environment 1 and environment 3 therefore leaving no chance to detect any outliers on the QTL influenced by these environmental variables. Fig ?? shows that RDA has a much better behavior than pcadapt by taking advantage of using informations of environmental local conditions.

Results summarized on Fig ?? is confirming that both methods have a good control of false discovery rate (8.36×10^{-2} for pcadapt and 8.51×10^{-2} for RDA) and that overall RDA shows better performance at detecting true outliers since it succeeds to detect quite often QTL1 and QTL3 SNPs. It seems however less efficient at detecting



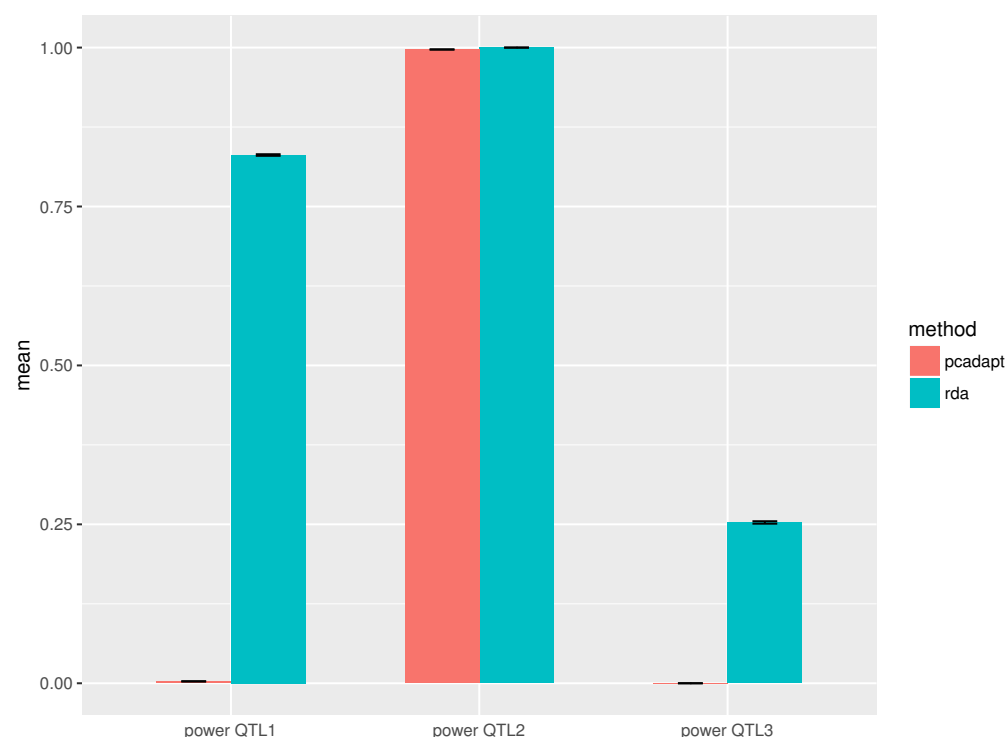


FIG. 4. Performance results of rda and pcadapt methods. Each performance value is averaged over 100 simulated dataset (error bars are displayed but hardly visible since they are very scarce). Power is given separately for loci coding for quantitative trait 1, 2 and 3.

QTL3 outliers but this might be due to the fact that local adaptation on a coarse environment is more difficult than adaptation on a smooth environmental gradient as environment 1 and 2. These simulations plead in favor of using constrained ordination method instead of PCA when non genetic data such as environmental variable are available in order to orientate the axis in the direction of informative gradients.

When performing an RDA on the “adaptively enriched genetic space”, Fig. ?? and ?? show that the method succeed at detecting the relevant selective gradient and separating them on different axis at least on our simulations. This therefore serves as a proof of concept of ?’s approach to represent multilocus selective gradient and the possibility to use the ordination axis it to devise

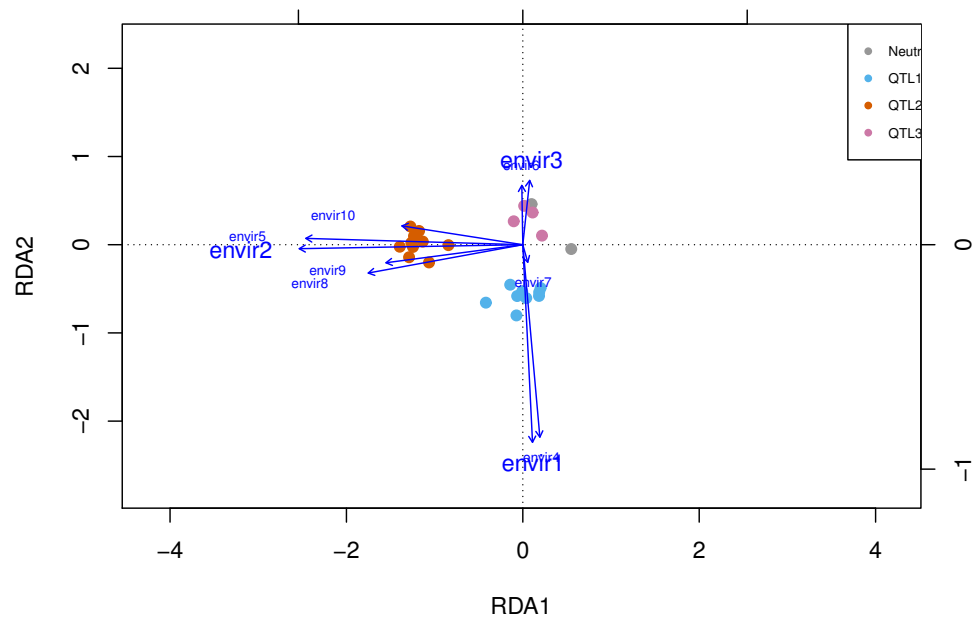


FIG. 5. RDA on the adaptively enriched genetic space. We discarded the individual points for readability. Dots represents outliers SNPs. R^2 of env1 with the first, second and third axis is (0.02%, 77.5%, 14.5%), env2 is (99.3%, 0.003%, 0001%) and env3 is (0.009%, 0.82%, 64.7%)

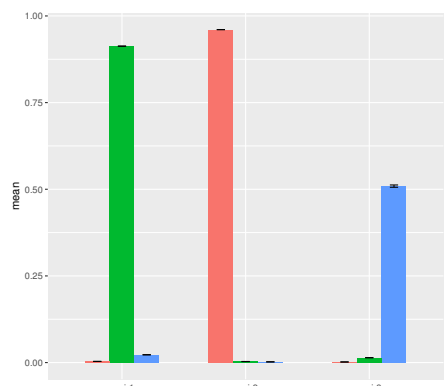


FIG. 6. R^2 between env1, env2 and env3 and each of the first three ordination axis. Values are averaged across the 100 simulated datasets.

Supplementary Material

Acknowledgments

References

Bazin, E., Dawson, K. J., and Beaumont, M. A. 2010. Likelihood-Free Inference of Population Structure and Local Adaptation in a Bayesian Hierarchical Model. *Genetics*, 185(2): 587–602.

Coop, G., Witonsky, D., Di Rienzo, A., and Pritchard, J. K. 2010. Using environmental correlations to identify loci underlying local adaptation. *Genetics*, 185(4): 1411–23.

De Kort, H., Vandepitte, K., Bruun, H. H., Closset-Kopp, D., Honnay, O., and Mergeay, J. 2014. Landscape genomics and a common garden trial reveal adaptive differentiation to temperature across Europe in the tree species *Alnus glutinosa*. *Molecular ecology*, pages 4709–4721.

de Villemereuil, P. and Gaggiotti, O. E. 2015. A new FST-based method to uncover local adaptation using environmental variables. *Methods in Ecology and Evolution*, 6(11): 1248–1258.

Duforet-Frebourg, N., Bazin, E., and Blum, M. G. B. 2014. Genome scans for detecting footprints of local adaptation using a Bayesian factor model. *Molecular biology and evolution*, 31(9): 1–13.

Eckert, A. J., Bower, A. D., González-Martínez, S. C., Wegrzyn, J. L., Coop, G., and Neale, D. B. 2010. Back to nature: Ecological genomics of loblolly pine (*Pinus taeda*, Pinaceae). *Molecular Ecology*, 19(17): 3789–3805.

- 446 Foll, M. and Gaggiotti, O. 2008. A genome-scan method to 484
447 identify selected loci appropriate for both dominant and 485
448 codominant markers: A Bayesian perspective. *Genetics*, 486
449 180(2): 977–993.
- 450 Frichot, E., Schoville, S. D., Bouchard, G., and François,
451 O. 2013. Testing for Associations between Loci and
452 Environmental Gradients Using Latent Factor Mixed
453 Models. *Molecular biology and evolution*, 30(7): 1687–
454 99.
- 455 Hecht, B. C., Matala, A. P., Hess, J. E., and Narum, S. R.
456 2015. Environmental adaptation in Chinook salmon
457 (*Oncorhynchus tshawytscha*) throughout their North
458 American range. *Molecular Ecology*, 24(22): 5573–5595.
- 459 Lachenbruch, P. A. 2011. Variable selection when missing
460 values are present: a case study. *Statistical Methods in*
461 *Medical Research*, 20(4): 429–444.
- 462 Lasky, J. R., Des Marais, D. L., McKay, J. K.,
463 Richards, J. H., Juenger, T. E., and Keitt, T. H.
464 2012. Characterizing genomic variation of *Arabidopsis*
465 *thaliana*: The roles of geography and climate. *Molecular*
466 *Ecology*, 21(22): 5512–5529.
- 467 Legendre, P. and Legendre, L. 2012. *Numerical ecology*.
468 Elsevier.
- 469 Luu, K., Bazin, E., Blum, M. G., Bazin, É., and Blum,
470 M. G. 2016. pcadapt: an R package to perform genome
471 scans for selection based on principal component
472 analysis. *bioRxiv*, 33: 056135.
- 473 Peng, B. and Kimmel, M. 2005. simuPOP: A forward-
474 time population genetics simulation environment.
475 *Bioinformatics*, 21(18): 3686–3687.
- 476 Rao, C. R. 1964. The Use and Interpretation of Principal
477 Component Analysis in Applied Research. *Sankhy: The*
478 *Indian Journal of Statistics, Series A*, 26: 329–358.
- 479 Steane, D. a., Potts, B. M., McLean, E., Prober, S. M.,
480 Stock, W. D., Vaillancourt, R. E., and Byrne, M.
481 2014. Genome-wide scans detect adaptation to aridity
482 in a widespread forest tree species. *Molecular ecology*,
483 23(10): 2500–13.
- 484 Vatsiou, A. I., Bazin, E., and Gaggiotti, O. 2015. A
485 comparison of recent methods for the detection of
486 selective sweeps. *Mol Ecol*, Accepted.