

# How to improve PCA based methods of genome scan using ecological data: detecting selection using RDA.

Eric Bazin,<sup>\*,1</sup> Keurcien Luu,<sup>2</sup> Michael G. B. Blum,<sup>2</sup>

<sup>1</sup>LECA, Université de Grenoble

<sup>2</sup>TIMC, Université de Grenoble

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

Associate Editor:

## Abstract

Ordination is a common tool in Ecology that aims at representing complex biological information on a reduced space. For instance, it is frequently used to study geographic distribution pattern of species diversity and to study the link between ecological variable such as temperature, drought, etc, on the species turnover. Recently, these methodologies are becoming quite popular in Landscape Genomic where one wants to study the link between environmental variable and the distribution pattern of genome wide diversity. However, it remains unclear what are the expected outcome of such approaches since genetic diversity has presumably a very different dynamic from species diversity. Simulations studies could help to shed light on this problem but they are still lacking whereas it tends to be broadly accepted as a pertinent approach. Furthermore, recent development have proposed to use ordination methods such as PCA to detect genes under selection. Simulations tend to support this idea as it seems to be quite robust to the underlying population structure and dynamic. Some authors have proposed to use other ordination approaches such as RDA, taking advantage of using environmental data. However no clear statistical framework have been developed to efficiently implement this idea in a robust and efficient test and once again, we don't know what is expected from the outcome of such approaches: which genes will be detected under which selective pressures? This paper aims at proposing a new test based on RDA approaches to search for genes under selection and to compare it to a classical PCA method. Thanks to individual based simulation, we compare both performance and robustness. Additionally, we test the efficiency of constrained ordination method such as RDA to detect relevant selective gradient since this was lacking in the Landscape Genomic literature. Finally, to illustrate the pertinence of such method in concrete example, we apply it to a real dataset.

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  
of a beneficial allele and its spread by the

© 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

mean of selective sweep. Others aim at picking up alleles with strong correlation with some environmental variable (e.g. Temperature, drought) with the idea that these alleles may confer a selective advantage to the individuals (??). Finally, other methods aim at detecting genomic region involved in local adaptation process. These region should have an increased differentiation between population because different alleles tend to be beneficial in each environment. Differentiation between population is excepted under the hypothesis of geographical isolation. Therefore, this region can be detected by quantifying the level of differentiation using some statistics and detecting the regions with unexpectedly high values. A common statistic and very easily comprehensible in population genetic is  $F_{st}$ . Many methods use this parameter as a basis in many different implementation of genome scans (???). These are model based method where parameters such as  $F_{st}$  are usually inferred using likelihood or Bayesian methods. This mean that users must have some a priori on their parameter value and the best model that fits their data in order to expect the best from their analysis. However, it is often difficult to get a satisfactory a priori picture of the demographic and population structure of the species one is interested in. Indeed many species are not clearly structured in different populations but more or less show a pattern of isolation by distance without clear geographical

barrier to gene flow. One solution would be to use more complex model that better reflect reality but these are difficult to implement in Bayesian framework. Additionally, these latter methods are very time consuming and the increase of both model complexity and the amount of data to analyze in terms of the number of individuals and loci makes them more and more difficult to use. A new path has opened recently with the use of multivariate methods. The idea is to capture the whole genome geographic structure using an ordination method such as ACP. Following this analysis, outliers loci are detected if they have extremely high correlation with one or more ordination axis (??). These are very efficient methods and simulations have shown that while they are very fast, they show similar efficiency than classical Bayesian method and sometimes perform better when the simulated demographic model drift from the model implemented in bayesian method, usually the island model. For instance ? have shown their method to be better when population are structured in hierarchical set or in isolation by distance pattern. Nevertheless, one conundrum of such approaches is the difficulty to interpret ordination axis in term of ecological meanings. These are usually tight to geographical axis (latitudinal or longitudinal) but they are not necessary linked to an environmental variable such as Temperature, drought, diet habit, etc. Therefore, when this information exists, it has to be a posteriori used as a mean

of interpretation but are not involved in the genome scan but in order to quantify multilocus inference process. It should be recalled that adaptation to an environmental gradient (????). natural selection is the result of a complex set of These studies whereby relationships between environmental pressures and that it most often environmental data and large multilocus data is acts on several characters simultaneously and explored are becoming more and more popular that these characters are encoded by several and are often coined as Ecological Genomics genes which generally have weak effects. In or Landscape Genomics studies. However the order to extract the maximum of all available concept of using constrained ordination methods information, it seems therefore necessary to use to analyse genomic data has never been tested approaches that are able to compile all kind of on simulated datasets. This paper aims at filling variable (e.g. alleles, phenotypic measurement, this gap. First, we show how one can make biotic and abiotic variables). One natural way use of a constrained ordination method namely to overcome this limitation would be to use Redundancy Analysis (RDA) as an efficient and more sophisticated ordination method than ACP robust genome scan method. We discarded the like methods. Constrained ordination methods other constraint ordination methods such as CCA (i.e. Redundancy Analysis, RDA, Canonical since they are very similar in their principles. RDA Correspondence Analysis, CCA) are well-known has already been used for instance by ? to perform set of approaches in Ecology for instance to genome scan in order to detect loci involved in explain the species distribution pattern by the adaptation to climate in *Arabidopsis thaliana*. the mean of environmental data. They have Outliers were identified as SNPs with the greatest specifically been designed in order to deal with squared scores along the first RDA axis (i.e. those biological complexity. In the population genomic in the 0.5 % tail). We build on this idea to era, it seems that data amount, complexity develop a comprehensive and robust statistical and heterogeneity is often a limitation to the test that allows to search for outliers on an use of inference methods based on classical arbitrary number of RDA axis simultaneously and population genetic models. Although they are allows to control precisely for the false discovery more difficult to interpret, such approaches rate. Using simulations, we show that it has better would be complementary to the model based results than PCA-based method. Second, thanks method because of their long-term use in ecology to these simulations, we show that RDA can and their efficiency on complex and large indeed help to identify important environmental datasets. These method have sometimes been gradient that better explain the adaptive variation used in population genomic studies, not as a 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  $\beta_{jk}$  is the regression coefficient corresponding to the  $j$ -th SNP regressed by the  $k$ -th ordination axis, and  $\epsilon_j$  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  $\text{Khi2}$  distributed after a correction with inflation factor (Luu et al., 2016). Pvalues are computed using  $K$  degree of freedom. We use the FDR approach to control for false positives. Qvalue are computed with `qvalueR` package and a loci is considered as an outlier if its qvalue 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

201 using scatter plots. To control for false positive, 233 must be very smooth and genetic differentiation  
 202 we used the same qvalue threshold (i.e.  $q = 10\%$ ). 234 must show an isolation by distance pattern over  
 203 Environmental genomic 235 the 64 populations. This is where pcadapt is  
 204 Once outliers have been identified, we isolate 236 best designed for. Loci are biallelic (0 or 1) like  
 205 them in a separate matrix A defining an 237 SNPs. Allele frequency of the whole population is  
 206 ”adaptively enriched genetic space” as coined by 238 initialized at 0.5. 1000 loci are defined. They are  
 207 ?. Following their methodology, we perform a 239 separated in 200 chunks of 5 SNPs in physical  
 208 second constrained ordination (RDA) on matrix 240 linkage with recombination rate between adjacent  
 209 A against environmental data. The rational of 241 loci fixed at 0.1. 3 different Traits are coded by  
 210 this analysis is to remove neutral variation before 242 a group of 10 different loci. The first trait is  
 211 performing ordination in order to have a better 243 coded by loci 1, 11, 21, ..., 91. The trait value is  
 212 picture of which environmental gradients have the 244 simply the sum of genotype value and therefore  
 213 strongest association with the adaptive genetic 245 can take value between 0 and 20. For the sake  
 214 space. On the simulated dataset, we report the 246 of realism, we add to each trait a random noise  
 215  $R^2$  statistics between env1, env2 and env3 and 247 (non heritable variation) drawn from a normal  
 216 the first three ordination axis to have an idea 248 distribution  $N(0,2)$ . The second trait is coded by  
 217 of which they are better associated with and if 249 loci 101, 111, ..., 191 and the third is coded by loci  
 218 the ordination space succeed in separating the 250 201, 211, ..., 291. Each trait is therefore coded by  
 219 environmental effect on different axis. 251 free recombining SNP loci. In other words, there  
 220 Simulations 252 are 30 coding SNPs among 1000. Selection can  
 221 To test for the efficiency of RDA in population 253 have an effect on linked loci, for instance, loci 2,  
 222 genomic, we performed simulations using simuPop 254 3, 4 and 5 can be impacted by selection on locus  
 223 python library (?). We compared our approach to 255 1. However, recombination is high enough (0.1) to  
 224 PCAdapt method to perform genome scans. Both 256 expect a limited linkage effect. We have defined  
 225 approach are equivalent except their ordination 257 10 different environmental variables. The first one  
 226 method. Finally we use these simulations to 258 determines the selective pressure on trait 1, the  
 227 evaluate RDA approach as a mean to detect 259 second one on trait 2 and the third one on trait  
 228 selective environmental gradient. A lattice of 260 3. The first environment variable is a quadratic  
 229 8x8 populations is simulated (i.e. 64 populations 261 gradient coded by function  $env1 = -(\cos(\theta) * (i - 3.5))^2 - (\sin(\theta) * (j - 3.5))^2 + 18, \theta = \pi/2$ , i and  
 230 in total). Each population is initialized with 262  $j$  being the population indicator on the 8x8  
 231 200 diploid individual with random genotypes. 263 lattice. The second one is a linear plan gradient  
 232 Migration is set to 0.5 so that population structure 264

265 coded by function  $env2 = h * \cos(\theta) * (i - 1) + h * \sin(\theta) * (j - 1) + k$  with  $h = 2$ ,  $\theta = \pi/4$  and  $k =$  297 are relative and selection arises on parents and  
266 3. The third environment variable simulates 298 determine their number of offsprings. Simulations  
267 a coarse environment with value  $env3 = 2$  for 299 are made across 500 generations. At the end  
268 all populations except population  $(i, j) = (2, 2)$ , 300 of simulation, we sample 10 individuals per  
269  $(2, 3)$ ,  $(3, 2)$ ,  $(3, 3)$ ,  $(6, 2)$ ,  $(6, 3)$ ,  $(7, 2)$ ,  $(7, 3)$ ,  $(2, 6)$ , 301 population. Therefore, we have a sample of 640  
270  $(2, 7)$ ,  $(3, 6)$ ,  $(3, 7)$ ,  $(6, 6)$ ,  $(6, 7)$ ,  $(7, 6)$ ,  $(7, 7)$  for 302 individuals with 1000 SNP-like loci.  
271 which  $env3 = 18$ . Env4, env5 and env6 have 303 Real dataset  
272 exactly the same equation than env1, env2 and 304 The Loblolly pine dataset is a sample of 682  
273 env3 respectively. The remaining 4 environment 305 individuals genotyped on 1,730 SNPs selected in  
274 variable are similar to env2 but with different 306 ESTs (?). 60 climatic variables were available and  
275 value of  $h$  and  $\theta$ . Env7 has  $h = 2$ ,  $\theta = 0$  and  $k =$  307 summarized by the authors in the five first axis  
276 3. Env8 has  $h = 2$ ,  $\theta = \pi/4$  and  $k = 0$ . Env9 has 308 of a PCA. The first axis, PC1 is mainly linked  
277  $h = 1$ ,  $\theta = \pi/4$  and  $k = 4$ . Env10 has  $h = 0.5$ ,  $\theta =$  309 to latitude, longitude, temperature, and winter  
278  $\pi/4$  and  $k = 8$ . Graphical representation of mean 310 aridity. PC2 is linked to longitude, spring-fall  
279 environmental value for environment 1, 2 and 311 aridity, and precipitation. We inputed the missing  
280 3 is given in Fig. ???. Environment 4, 5 and 6 312 data using a very simple algorithm implement in  
281 have respectively the same mean value spatial 313 function `sing.im` of the R package `linkim` (?). It  
282 distribution. For a graphical representation of 314 imputes the missing value based on the observed  
283 environment 7 to 10, see supplementary material. 315 data proportions. We used  $K = 4$  axis to compute  
284 Environmental equation gives a mean value of 316 Malahanobis distances.  
285 the environment variable. To avoid colinearity 317 The Chinook salmon consists of 19 703 SNP  
286 between environments variable, we added noise by 318 loci genotyped on 1956 total individuals pooled  
287 drawing an environment value within a normal 319 in 46 collections. ? have estimated that between  
288 distribution  $N(\mu = env, \sigma = 1)$ . Fitness for each 320 5.8 and 21.8% of genomic variation can be  
289 trait is set to be  $-e^{((x - env)^2 / (2 * \omega^2))}$ ,  $x$  being the 321 accounted for by environmental features, and  
290 quantitative trait value,  $env$  the environmental 322 566 putatively adaptive loci were identified as  
291 value and  $\omega$  is defining selection strength and 323 targets of environmental adaptation. Therefore  
292 has been set to 10 which in our experience seems 324 this dataset is a good candidate to test ACP  
293 sufficient for loci to be often detected. To get 325 and RDA approaches to detect outliers and  
294 the overall fitness for a given individual, fitness 326 selective gradients. Five variables (`MigDistKM`,  
295 associated to each trait are multiplied. Fitness 327 `StreamOrder`, `bio03`, `bio17` and `bio18`) have  
296 been used among 24 different climate and 328

environmental variables because they have been tested as significantly associated with the SNP variation rangewide (?). 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 analysis and measured the mean  $R^2$  between env1, env2 and env3 and each of the first three ordination axis. This is summarized in Fig. ??.

### Loblolly Pine

### Chinook Salmon

Our analysis of Chinook Salmon gave a list of 27 SNPs (Tab. ??). From the material of ?, we extracted their matching with coding sequences and the associated annotation.

## 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.

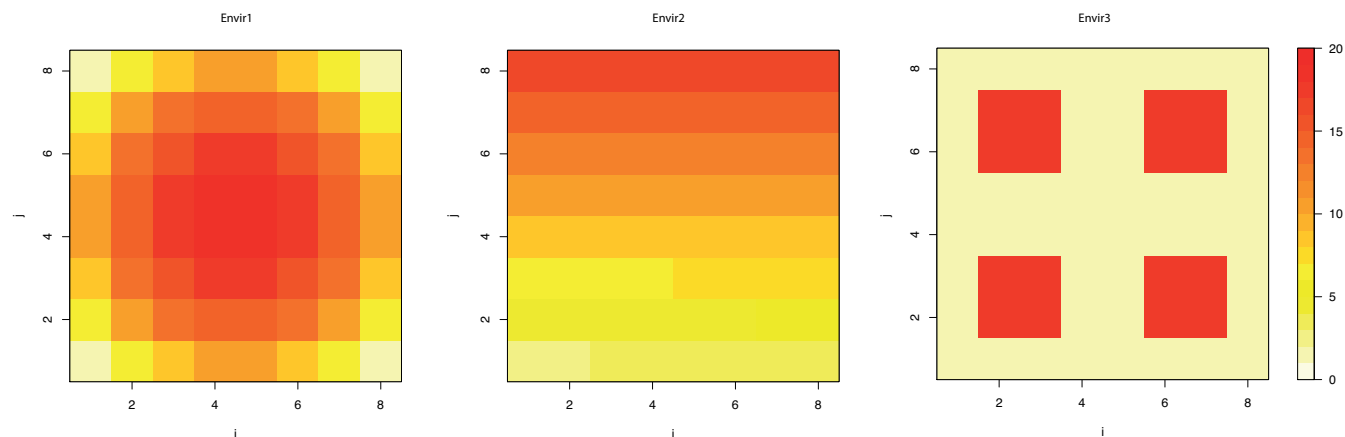


FIG. 1. Graphical representation of mean environmental value for environment 1, 2 and 3

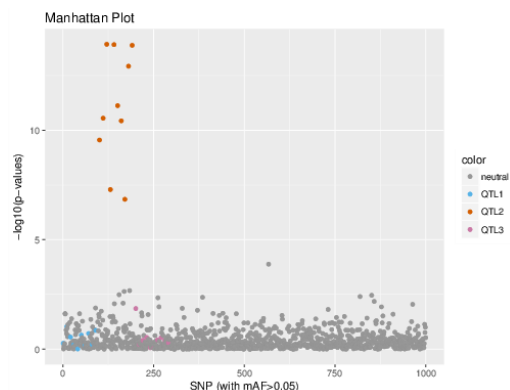


FIG. 2. Manhattan plot of the result of pcadapt on a simulated data set.

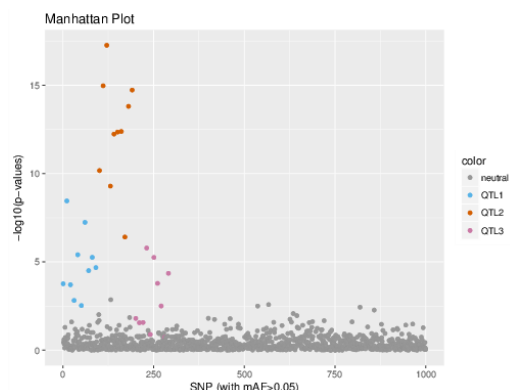


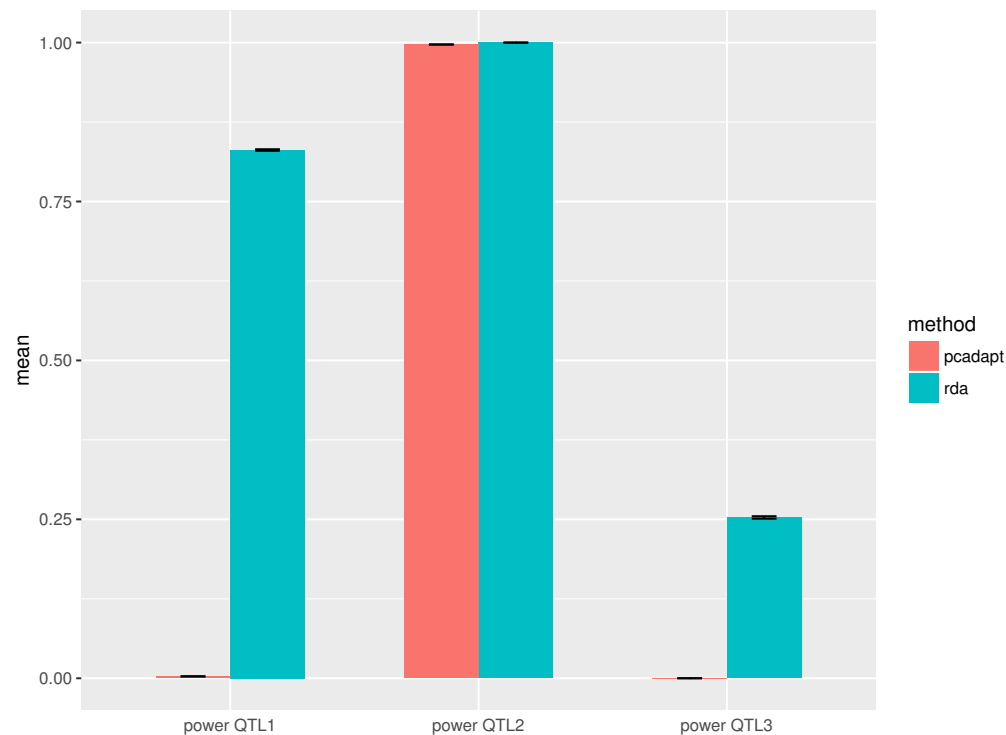
FIG. 3. Manhattan plot of the result of genome scan using RDA on a simulated data set.

Both methods have a good control of false discovery rate ( $8.36 \times 10^{-2}$  for pcadapt and  $8.51 \times 10^{-2}$  for RDA). Results summarized on Fig ?? is confirming 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 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 a metric that provides a holistic measure of





**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.

genomic adaptation. Indeed, in RDA1 is strongly associated with env1, RDA2 with env2 and RDA3 with env3 whereas poorly associated with the other axis. As expected, the correlated environment are also strongly associated with this respective axis. This is reflecting the fact that in reality it is difficult on an environmental gradient to distinguish among the covariable which one has a causal effect on the individual fitness. However, it is often sufficient for biologists when performing an exploratory analysis to identify combination of environment variable having a strong association with adaptive variation without knowing precisely the underlying mechanical process.

From the analysis of Chinook Salmon, we picked up some genes that can be interpreted regarding to the environmental variable. For instance, a

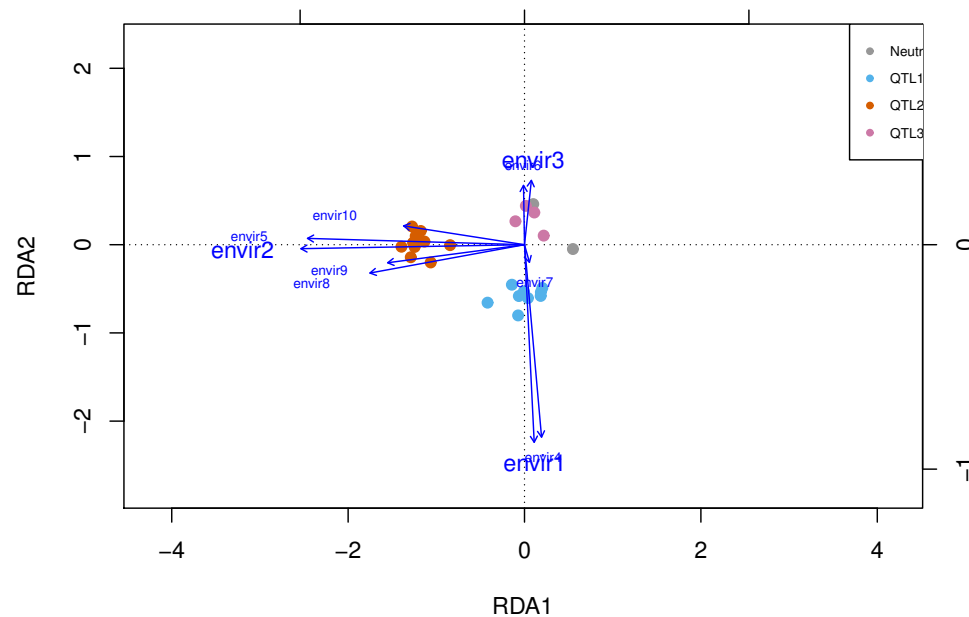
heat shock protein which are known to be involved in adaptation to temperature or lipolysis-stimulated lipoprotein receptor which are involved in regulation of lipid metabolic process. This latter process can reasonably thought to be involved in adaptation to food abundance and the need for salmon to migrate on a short or long distance.

## Supplementary Material

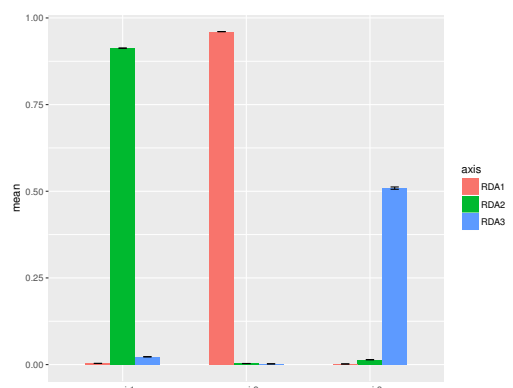
## Acknowledgments

## 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.



**FIG. 5.** RDA on the adaptively enriched genetic space. We discarded the individual points for readability. Dots represents outliers SNPs.  $R^2$  of envir1 with the first, second and third axis is (0.02%, 77.5%, 14.5%), envir2 is (99.3%, 0.003%, 0.001%) and envir3 is (0.009%, 0.82%, 64.7%)



**FIG. 6.**  $R^2$  between envir1, envir2 and envir3 and each of the first three ordination axis. Values are averaged across the 100 simulated datasets.

- 459 *Evolution*, 6(11): 1248–1258.
- 460 Duforet-Frebourg, N., Bazin, E., and Blum, M. G. B.
- 461 2014. Genome scans for detecting footprints of local
- 462 adaptation using a Bayesian factor model. *Molecular*
- 463 *biology and evolution*, 31(9): 1–13.
- 464 Eckert, A. J., Bower, A. D., González-Martínez, S. C.,
- 465 Wegrzyn, J. L., Coop, G., and Neale, D. B. 2010. Back
- 466 to nature: Ecological genomics of loblolly pine (*Pinus*
- 467 *taeda*, Pinaceae). *Molecular Ecology*, 19(17): 3789–3805.
- 468 Foll, M. and Gaggiotti, O. 2008. A genome-scan method to
- 469 identify selected loci appropriate for both dominant and
- 470 codominant markers: A Bayesian perspective. *Genetics*,
- 471 180(2): 977–993.
- 472 Frichot, E., Schoville, S. D., Bouchard, G., and François,
- 473 O. 2013. Testing for Associations between Loci and
- 474 Environmental Gradients Using Latent Factor Mixed
- 475 Models. *Molecular biology and evolution*, 30(7): 1687–
- 476 99.
- 477 Hecht, B. C., Matala, A. P., Hess, J. E., and Narum, S. R.
- 478 2015. Environmental adaptation in Chinook salmon
- 450 De Kort, H., Vandepitte, K., Bruun, H. H., Closset-Kopp,
- 451 D., Honnay, O., and Mergeay, J. 2014. Landscape
- 452 genomics and a common garden trial reveal adaptive
- 453 differentiation to temperature across Europe in the
- 454 tree species *Alnus glutinosa*. *Molecular ecology*, pages
- 455 4709–4721.
- 456 de Villemereuil, P. and Gaggiotti, O. E. 2015. A new
- 457 FST-based method to uncover local adaptation using
- 458 environmental variables. *Methods in Ecology and*

479 (Oncorhynchus tshawytscha) throughout their North  
480 American range. *Molecular Ecology*, 24(22): 5573–5595.

481 Lachenbruch, P. A. 2011. Variable selection when missing  
482 values are present: a case study. *Statistical Methods in*  
483 *Medical Research*, 20(4): 429–444.

484 Lasky, J. R., Des Marais, D. L., McKay, J. K.,  
485 Richards, J. H., Juenger, T. E., and Keitt, T. H.  
486 2012. Characterizing genomic variation of Arabidopsis  
487 thaliana: The roles of geography and climate. *Molecular*  
488 *Ecology*, 21(22): 5512–5529.

489 Legendre, P. and Legendre, L. 2012. *Numerical ecology*.  
490 Elsevier.

491 Luu, K., Bazin, E., Blum, M. G., Bazin, É., and Blum,  
492 M. G. 2016. pcadapt: an R package to perform genome  
493 scans for selection based on principal component  
494 analysis. *bioRxiv*, 33: 056135.

495 Peng, B. and Kimmel, M. 2005. simuPOP: A forward-  
496 time population genetics simulation environment.  
497 *Bioinformatics*, 21(18): 3686–3687.

498 Rao, C. R. 1964. The Use and Interpretation of Principal  
499 Component Analysis in Applied Research. *Sankhy: The*  
500 *Indian Journal of Statistics, Series A*, 26: 329–358.

501 Steane, D. a., Potts, B. M., McLean, E., Prober, S. M.,  
502 Stock, W. D., Vaillancourt, R. E., and Byrne, M.  
503 2014. Genome-wide scans detect adaptation to aridity  
504 in a widespread forest tree species. *Molecular ecology*,  
505 23(10): 2500–13.

506 Vatsiou, A. I., Bazin, E., and Gaggiotti, O. 2015. A  
507 comparison of recent methods for the detection of  
508 selective sweeps. *Mol Ecol*, Accepted.

**Table 1.** List of SNPs with *qvalue* < 0.1 and their matching with coding sequence when available.

	Locus	Sequence Description
1	8760.60	cell migration-inducing and hyaluronan-binding partial
2	11727.44	protein argonaute-1
3	15784.70	dna polymerase epsilon subunit 4
4	19372.14	protein fam122a-like isoform x1
5	19510.54	pantothenate kinase mitochondrial-like
6	19809.36	eukaryotic translation initiation factor 3 subunit j
7	22558.48	zinc finger protein gfi-1b-like
8	29912.62	g protein-activated inward rectifier potassium channel
9	30253.61	solute carrier family 1 (glial high affinity glutamate tra
10	30495.21	c-jun-amino-terminal kinase-interacting protein 4
11	33486.16	heat shock 70 kda protein 12a isoform x3
12	39480.19	ras association domain-containing protein 4
13	40284.30	rna-binding single-stranded-interacting protein 2-like i
14	41648.34	afadin- and alpha-actinin-binding protein
15	46982.22	unnamed protein product
16	50054.21	protein fam92a1-like isoform x1
17	54261.58	leukotriene b4 receptor 1-like
18	54497.54	ankyrin repeat domain-containing protein 50-like
19	56375.14	mms19 nucleotide excision repair protein homolog
20	60067.64	e3 ubiquitin-protein ligase trim37-like
21	66930.15	tubulin polyglutamylase ttl4-like isoform x2
22	69650.61	monocyte to macrophage differentiation factor 2
23	71287.48	lipolysis-stimulated lipoprotein receptor
24	74776.68	baculoviral iap repeat-containing protein 6 isoform x10
25	79151.39	guanine nucleotide-binding protein g g g subunit beta-
26	81519.68	nuclear receptor corepressor 1 isoform x3
27	89719.68	unnamed protein product, partial