

LEA: Un paquetage R pour la génomique des populations

Olivier François
Université Grenoble-Alpes

Grenoble, Juin 2015

LEA : Landscape and Ecological genomewide Association studies

Methods in Ecology and Evolution



Methods in Ecology and Evolution 2015

doi: 10.1111/2041-210X.12382

APPLICATION

LEA: An R package for landscape and ecological association studies

Eric Fricot¹ and Olivier François^{1*}

LEA: Un programme R pour l'écologie moléculaire (à haut-débit)

- ▶ **Objectifs généraux:** Evaluer les effets des interactions de longue durée entre les organismes d'une population et leurs environnements (climat, régimes alimentaires, exposition aux pathogènes).

LEA: Un programme R pour l'écologie moléculaire (à haut-débit)

- ▶ Analyser la structure génétique des populations.
- ▶ Effectuer des cribles génomiques pour détecter les régions du génome répondant à la sélection naturelle.

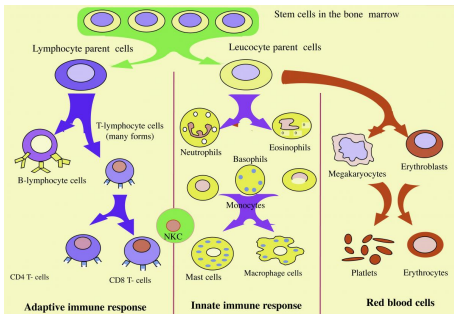
Exemples chez l'homme

- ▶ Des mutations des gènes *EGLN1* and *PPARA* confèrent une **tolérance à l'hypoxie** et permettent l'adaptation à la **haute altitude** chez les populations tibétaines (Simonsen et al. Science 2011).



Exemples chez l'homme

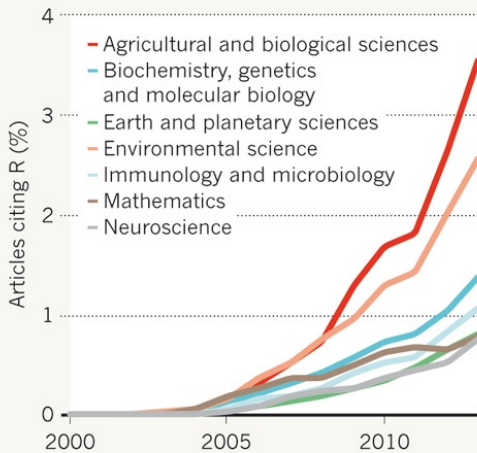
- ▶ Excès de gènes associés à des **maladies auto-immunes** en réponse à des pressions sélectives exercées par les pathogènes (Fumagalli et al. Plos Genetics 2012).



Pourquoi choisir R?

A RISING TIDE OF R

An increasing proportion of research articles explicitly reference R or an R package.



Bases de données

- ▶ Projets de génomique humaine : **HGDP** (Li et al. Science 2008), **The 1000 Genome project** (Nature 2012)
- ▶ Matrices de génotypes de SNPs¹ de **très grandes** tailles : $(Y_{i\ell}) \sim$ **1–5 Giga** entrées ($Y_{i\ell} = 0,1,2$).

¹SNP = Single Nucleotide Polymorphism (locus de l'ADN montrant une variation dans les populations)

Bases de données

- ▶ Bases de données environnementales: **Worldclim**, **WHO**
- ▶ Bases de données bioinformatiques: dbSNP, Genome browsers, NHI-EBI GWAS catalog, etc

Modèles statistiques implantés dans LEA

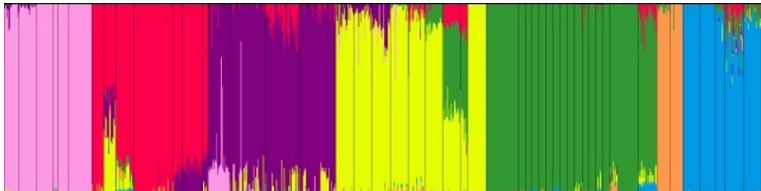
- ▶ **snmf** : Factorisation de matrice non-négative (avec contraintes de convexité)
 - ▶ Estimation de la structure génétique des populations et des coefficients individuels d'ascendance génétique.
 - ▶ Détection de régions génomiques répondant à la sélection naturelle (adaptation).

Modèles statistiques implantés dans LEA

- ▶ **lfmm** : Modèles mixtes à facteurs latents
 - ▶ Détection de régions génomiques associées à des variables écologiques
 - ▶ Correction des facteurs de confusion induits par la structure génétique des populations.

Coefficients d'ascendance génétique

- ▶ Les modèles de métissage génétique supposent que les **individus ont hérité leurs gènes de K populations ancestrales** (K n'est pas connu).
- ▶ Coefficient d'ascendance génétique : la fraction de genome de l'individual i provenant de la population ancestrale k , Q_{ik} .



Méthode **snmf**

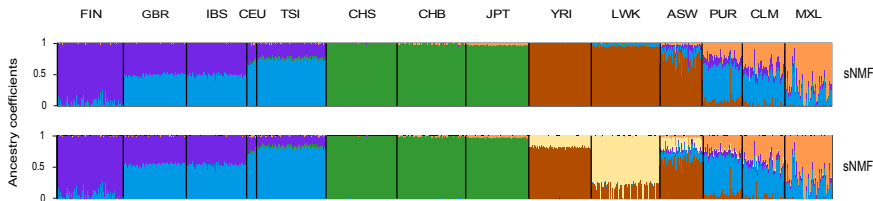
- ▶ La fonction **snmf** calcule des estimateurs de moindres carrés de la matrice **Q** (de rang K)

$$\text{LS}(\mathbf{Q}, \mathbf{F}) = \|\mathbf{Y} - \mathbf{Q}\mathbf{F}\|_{\text{F}}^2 + \sqrt{\alpha} \sum_{i=1}^n \|\mathbf{Q}_i\|_1^2, \quad \mathbf{Q}, \mathbf{F} \geq 0,$$

- ▶ sous les contraintes

$$\sum_{k=1}^K \mathbf{Q}_{ik} = 1, \quad \sum_{j=0}^2 \mathbf{F}_{kj}(j) = 1, \quad j = 0, 1, 2.$$

Résultat pour les données 1000 Genomes



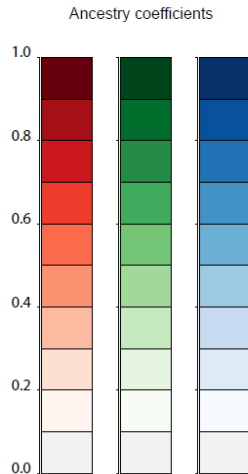
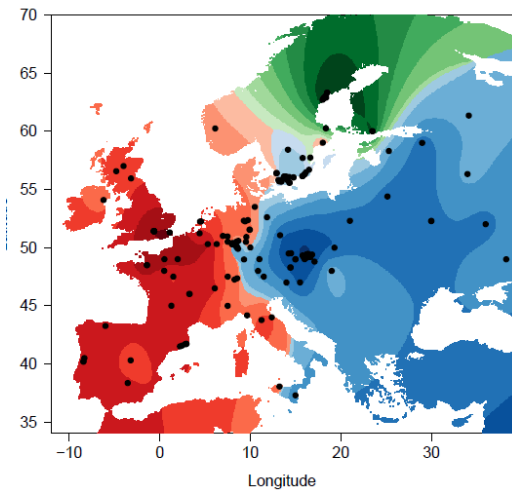
Dimension de la matrice : $n = 2k$, $p = 3M$,
Temps de calcul CPU ≈ 30 min
(Frichot et al. Genetics 2014)

Résultats pour la plante *A. thaliana*

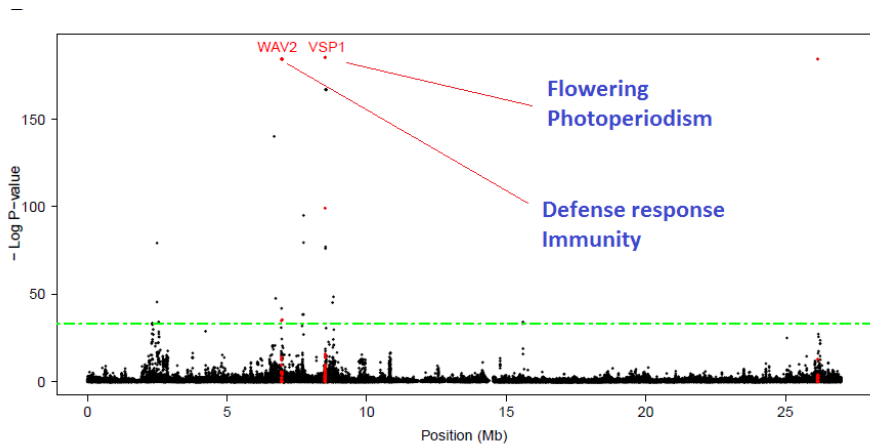
- ▶ Lignées européennes de la plante modèle *A. thaliana* (216K SNPs, 170 individus).
- ▶ La valeur du nombre de populations, K , est choisie selon un critère d'entropie croisée $\rightarrow K = 3$.



Résultat pour *A. thaliana*



Crible génomique pour *A. thaliana*



lfmm: Méthode d'association génomique

- ▶ **lfmm: latent factor mixed model**
- ▶ Estimer la **corrélation entre la fréquence d'allèle** pour un SNP donné et une **variable d'intérêt écologique**.
- ▶ Les cribles génomiques cherchent à identifier les SNPs montrant des différences significatives en comparaison avec le fond génomique.
- ▶ **Caveat**: La structure de population et d'autres facteurs de confusion, créent de nombreux faux positifs.

lfmm: Méthode d'association génomique

- ▶ **Latent factor mixed models**

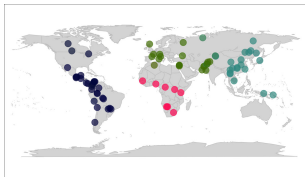
$$Y_{il} = \mu_l + \beta_l^T X_i + U_i^T V_l + \epsilon_{il},$$

où β_l est un vecteur de coefficients de régression, U_i est un facteur latent, et la matrice V_l contient les poids correspondants (*loadings*) (Frichot et al. Mol. Biol. Evol. 2013).

- ▶ K facteurs latents.

Human Genome Diversity Project (660k SNP arrays)

- ▶ Echantillons d'ADN de 1,043 individus répartis mondialement dans 52 populations
- ▶ Données climatiques pour chacune des 52 populations (WorldClim database, résolution de 1km^2)



- ▶ 11 variables bioclimatiques interpolées, collectées sur une période de 50 ans (1950-2000)
- ▶ Nombre de facteurs latents choisi par un test de Tracy-Widom.

GWAS-SNPs associated with environmental predictors.

Gene	Trait association	$-\log_{10} P$ -value
<i>OCA2/HERC2</i>	Eye and hair color, pigmentation	9.15
<i>DHCR7</i>	Vitamin D levels	7.78
<i>SLC45A2</i>	Hair color	6.90
Intergenic <i>MUC7</i>	Alcoholism	8.91
<i>ZMIZ1</i>	Crohn's disease	8.77
<i>KLK3</i>	Prostate Cancer	8.61
<i>ICOSLG</i>	Celiac disease	7.02
<i>HLA-DRA</i>	Systemic sclerosis	6.97
<i>NCAPG-LCORL</i>	Height	9.43
<i>BOK</i>	Brain structure and development	9.43

Genic SNPs associated with environmental predictors.

Gene	Annotation (dbSNPs)	$-\log_{10} P\text{-value}$
<i>EPHB4</i>	Heart morphogenesis and angiogenesis	16.54
<i>NRG1</i>	Nervous system development, cell proliferation	16.21
<i>RBM19</i>	Regulation of embryonic development	15.98
<i>EYA2</i>	Eye development and DNA repair	15.9
<i>POLA1</i>	Mitotic cell cycle and cell proliferation	15.87

Résumé

- ▶ Des algorithmes rapides de factorisation matricielle (modèles à facteurs latents)
- ▶ Permettant d'étudier la structure génétique des populations
- ▶ De lancer des cribles génomiques tout en contrôlant les fausses découvertes.

[Home](#) » [Bioconductor 3.1](#) » [Software Packages](#) » [LEA](#)

LEA

platforms all

downloads available

posts 0

in Bioc < 6 months

build ok

commits 0.50

LEA: an R package for Landscape and Ecological Association Studies

Bioconductor version: Release (3.1)

LEA is an R package dedicated to landscape genomics and ecological association tests. LEA can run analyses of population structure and genome scans for local adaptation. It includes statistical methods for estimating ancestry coefficients from large genotypic matrices and evaluating the number of ancestral populations (`snmf`, `pca`); and identifying genetic polymorphisms that exhibit high correlation with some environmental gradient or with the variables used as proxies for ecological pressures (`lfmm`), and controlling the false discovery rate. LEA is mainly based on optimized C programs that can scale with the dimension of very large data sets.

Author: Eric Frichot <eric.frichot at gmail.com>, Olivier Francois <olivier.francois at imag.fr>