# MLE for individual ancestries, population covariances, and selection
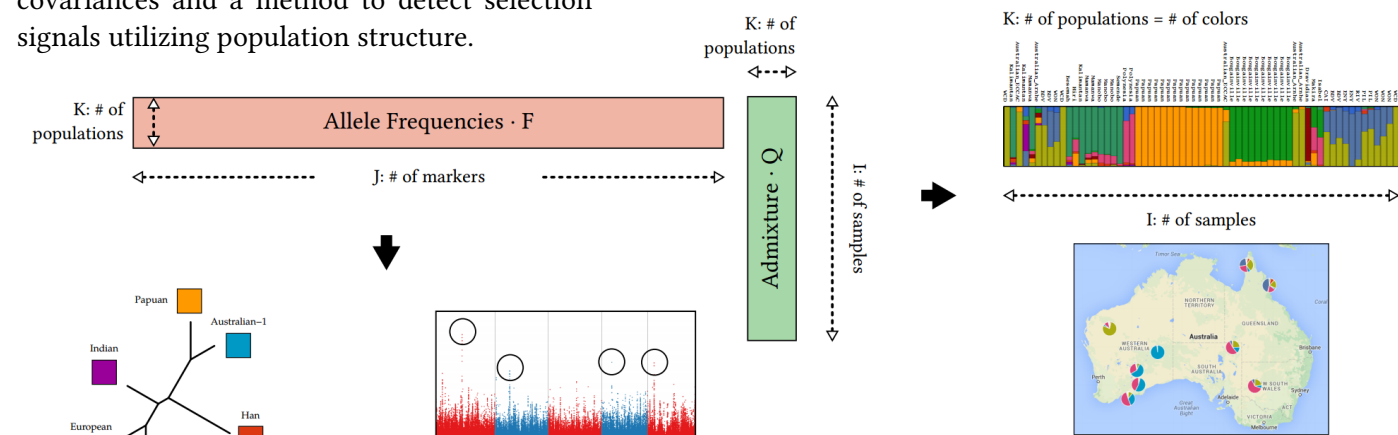
Jade Yu Cheng[1,2*], Rasmus Nielsen[1,3]

1. Departments of Integrative Biology and Statistics, University of California, Berkeley, CA, USA
2. Bioinformatics Research Centre, Aarhus University, Aarhus, Denmark
3. Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Denmark
* Email: ycheng@cs.au.dk

## Overview

Questions about population structure arise in many contexts such as human evolution and disease genetics association studies. The initial focus of our project has been to study population structure by developing a model-based inference strategy using quadratic programming with the active set algorithm.

Inspired by a number of recent methods that use Gaussian distributions to approximate the distribution of allele frequencies among populations, our project has also focused on developing a method to infer population covariances and a method to detect selection signals utilizing population structure.



## Population Covariances Inference

$$P(f_j \mid \Omega,\ \mu_j) \sim \mathcal{N}(\mu_j,\ \mu_j(1-\mu_j)\,\Omega).$$

**Likelihood Model**

$$\ln[P_2(F)]$$

$$= \ln\left\{\prod_j^J \left[\frac{1}{\sqrt{|2\pi c_j \Omega'|}}\right.\right.$$

$$\left.\left.\exp\left(-\frac{1}{2}\cdot f_j'^{T}\cdot(c_j\Omega')^{-1}\cdot f_j'\right)\right]\right\}$$

$$= -\frac{1}{2}\cdot\sum_j^J\left\{(K-1)\cdot\ln(2\pi c_j) + \ln[\det(\Omega')] + \frac{1}{c_j}\cdot f_j'^{T}\cdot\Omega'^{-1}\cdot f_j'\right\}$$

where $c_j = \mu_j(1-\mu_j)$

$$f_j' = f_j - f_{j0}.$$

**NM**

**Repeat** until a stopping criteria is reached
  Evaluate each point in the simplex using the objective function
  Determine the point $p_{min}$ with the lowest score
  Reflect $p_{min}$ through the centroid of the remaining vertices to $p_r$
  **If** the score at $p_r$ is the highest score in the simplex **Then**
    Expand $p_r$ away from the centroid to $p_e$
    Use $p_e$ in place of $p_{min}$
  **Else If** the score at $p_r$ is still the lowest score **Then**
    Contract $p_r$ toward the centroid to point $p_c$
    **If** the score at $p_c$ is no longer the lowest score **Then**
      Use $p_c$ to replace $p_{min}$
    **Else**
      Determine the point $p_{max}$ with the highest score
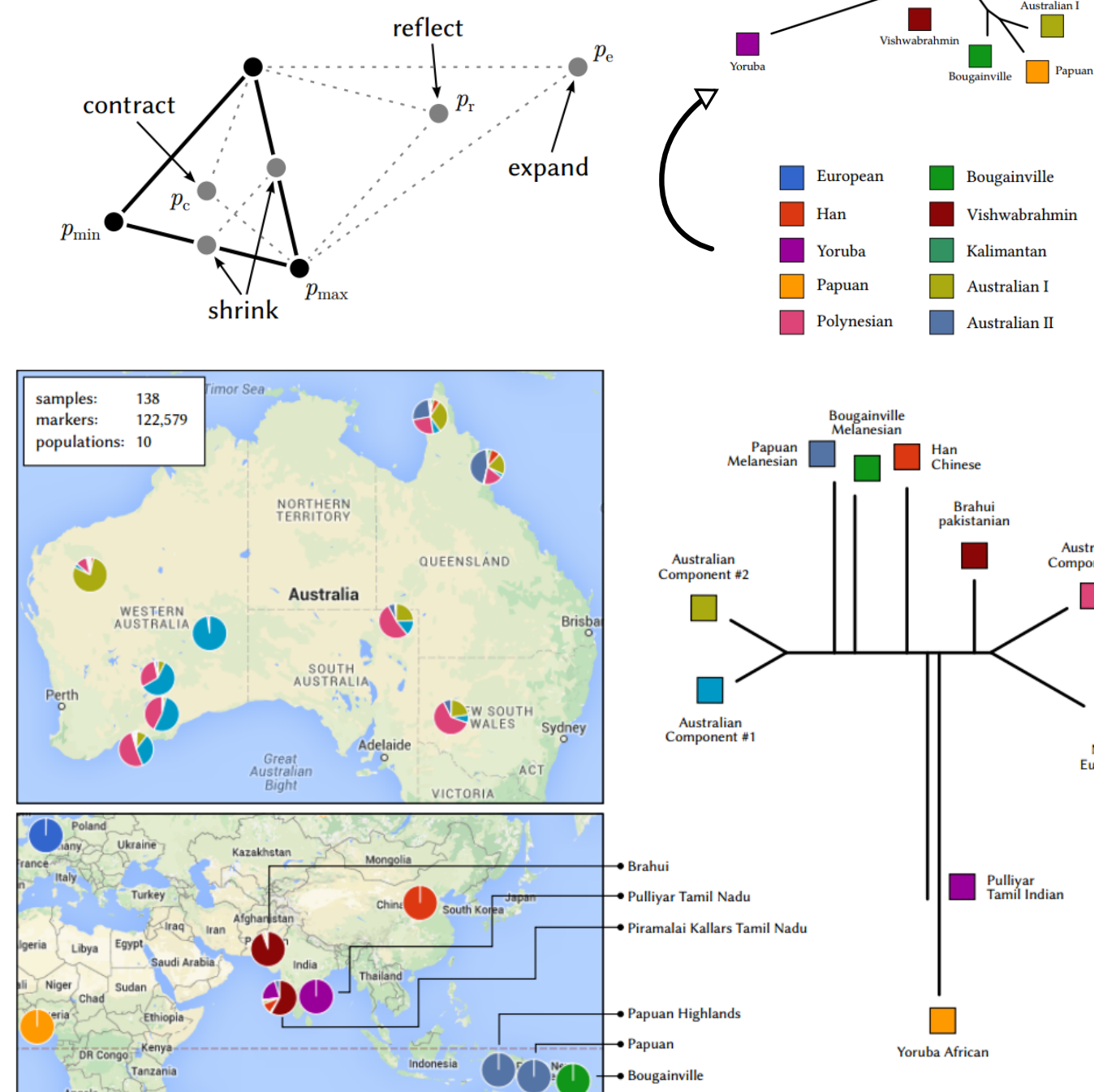      Shrink all points in the simplex around $p_{max}$
    **End If**
  **Else**
    Use $p_r$ in place of $p_{min}$
  **End If**
**End Repeat**



## Admixture Inference

$$\ln[P_1(Q,\ F)] = \sum_i^I \sum_j^J \left\{ g_{ij}\cdot\ln\left[\sum_k^K q_{ik}\cdot f_{kj}\right] + (2-g_{ij})\cdot\ln\left[\sum_k^K q_{ik}\cdot(1-f_{kj})\right]\right\}.$$

**Likelihood Model**

$$\max_{\Delta_{Q_i}}\left\{\frac{1}{2}\Delta_{Q_i}^T H_{Q_i}\Delta_{Q_i} + D_{Q_i}^T\Delta_{Q_i}\right\}$$

s.t. $\quad A\Delta_{Q_i} \leq a$

$\qquad B\Delta_{Q_i} = b$

$$\max_{\Delta_{F_j}}\left\{\frac{1}{2}\Delta_{F_j}^T H_{F_j}\Delta_{F_j} + D_{F_j}^T\Delta_{F_j}\right\}$$

s.t. $\quad A\Delta_{F_j} \leq a$

**QPAS**

Find a feasible starting point
Initialize the corresponding active set
**Repeat**
  Solve the equality problem defined by the active set
  Compute the Lagrange multipliers of the active set
  **If** the solved approximation is within the feasible region
    **If** all Lagrange multipliers are negative
      Return the solved approximation
    **Else**
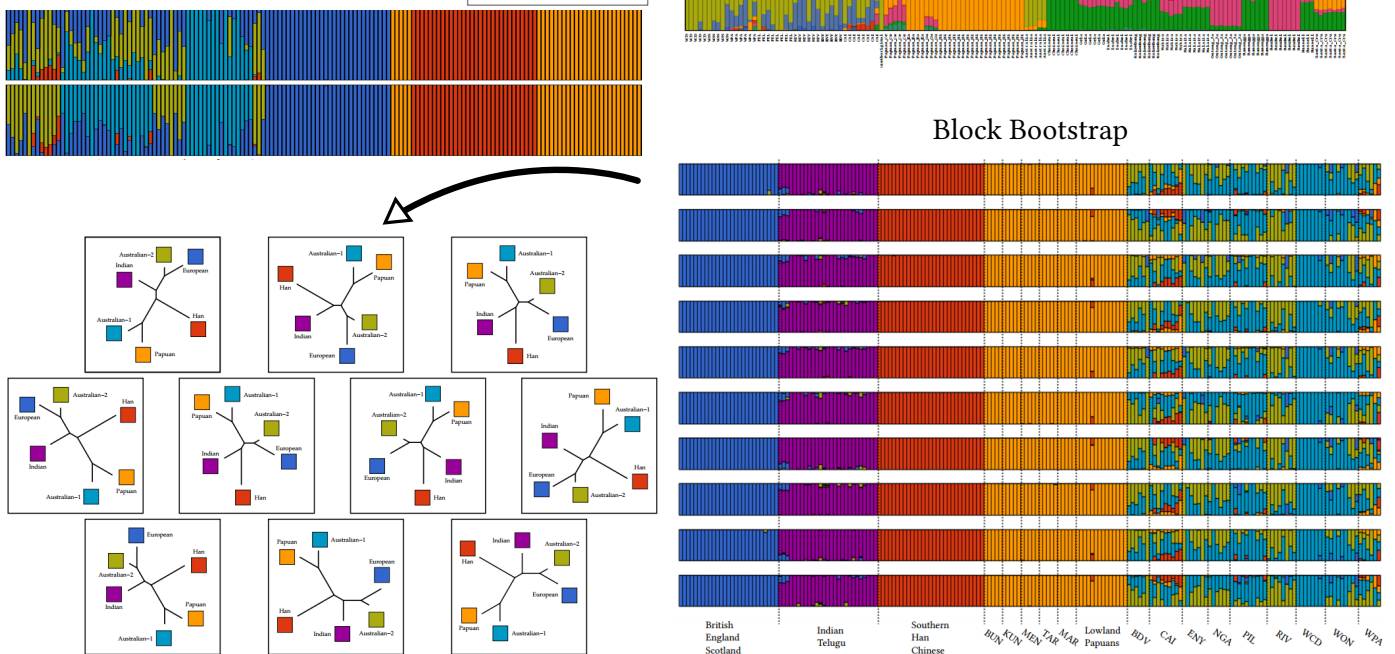      Remove the constraint with the largest Lagrange multiplier
    **End If**
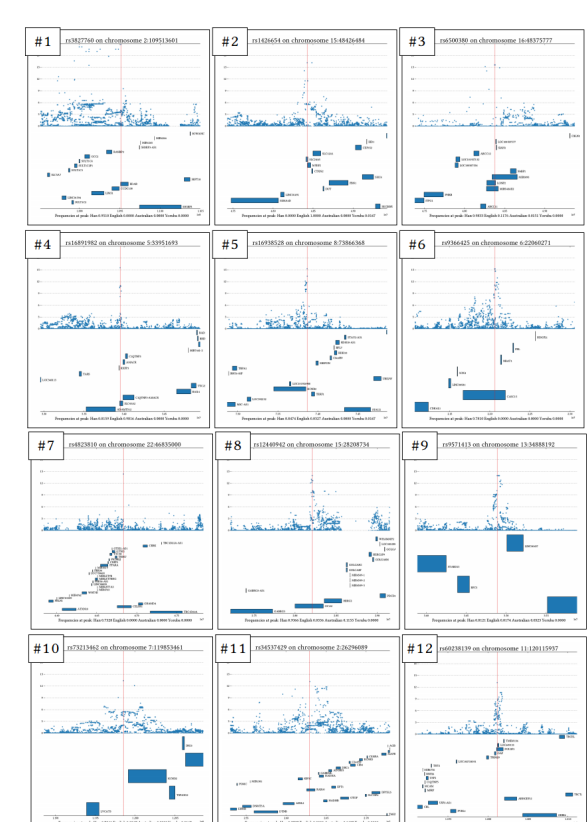  **Else**
    Take the shortest step back into the feasible region
    Insert the corresponding constraint into the active set
  **End If**
**End Repeat**



### Block Bootstrap



## Selection Detection



**SELECTION**

Obtain the full genotype dataset G with N markers and M samples
Sample N' markers with respect to LD (N' > 100,000) to form G'
QPAS over $\ln(P_1)$ using G'
  Produce admixture proportions Q' of size M by K
  Produce allele frequencies F' of size K by N'
Nelder-Mead over $\ln(P_2)$ using F'
  Produce variance covariance matrix $\Omega'$
QPAS over $\ln(P_1)$ using G while fixing Q'
  Produce allele frequencies F of size K by N
**Repeat** for each marker in F
  Set $l_{ratio}$ to zero
  **Repeat** for each $\alpha$ in a range of an even interval starting from 1.0
    Set $l_{new}$ to $\ln(P_2)$ calculated for this marker using $\alpha \times \Omega'$
    Set $l_{old}$ to $\ln(P_2)$ calculated for this marker using $\Omega'$
    **If** $2\times(l_{new}-l_{old})$ is greater than $l_{ratio}$
      Set $l_{ratio}$ to $2\times(l_{new}-l_{old})$
  **End Repeat**
  Emit $l_{ratio}$
**End Repeat**

**Top Twelve Peaks**

1. Hair thickness and curliness
2. Earwax moisture and underarm odor
3. Skin pigmentation
4. Skin pigmentation
5. 'maxDrinks' alcohol related
6. A melanoma tumor repressor
7. Neural tube defect and hair follicles
8. Skin pigmnetation
9. Intergenic
10. Intergenic
11. Insulin dependent regulation of glucose
12. Taste cells in the mouth