Identifying fine population structure using genomic scale data

Daniel Lawson Heilbronn Institute for Mathematical Research School of mathematics University of Bristol

dan.lawson@bristol.ac.uk

<u>Work with:</u> Garrett Hellenthal Daniel Falush Simon Myers

www.paintmychromosomes.com

## Talk outline

- 1 Ancestry Overview
- 2 Previous population model motivation (STRUCTURE)
- 3 Summary Statistic approaches
- 4 ChromoPainter ancestry summary model
- 5 FineSTRUCTURE population model
- 6 Applications and benefits









## Outline: The process



Step 1: SNPs are converted to similarity matrices Step 2: Analyse the population structure

#### Ancestry process



Daniel Lawson, University of Bristol dan.lawson@bristol.ac.uk Each generation randomly chooses parent
Red individuals ancestors to those sampled
Take limit N→∞

keeping N/T constant • Rate of *coalescence* between pairs  $\rightarrow 1$ • Lots known about this *Coalescent Tree* distribution.

#### Ancestry process with recombination



Daniel Lawson, University of Briston dan.lawson@bristol.ac.uk • Probability  $\rho$  of getting DNA from two parents

- Creates a graph structure
- Coalescence rates unchanged...
- But now a birth/death process
- Easily simulated but few analytical results

### Ancestral Recombination Graph



Hein, Schierup and Wiuf 'Gene Genealogies, Variation and Evolution', OUP 2005

## Ancestral Recombination Graph -Summary

Ancestral Recombination Graph (ARG) model
 backwards in time, ignore unobserved ancestors

is equivalent to the

- Forwards in time model
  - Random mating, within known size populations
  - No selection
- Inference under the ARG is impossible for reasonable datasets
- But when recombination is large, each SNP is independently drawn from a random tree

## Genetic drift



#### Present populations 'Independent drift' of SNP distribution

Kimura derived exact distribution for this case... (nasty)
Wright studied a related model, leads to Beta distribution of SNP frequency

$$p_k \sim \text{Beta}(p_0, \nu)$$

• Normal distribution approximation exists... simpler to handle covariance effects

$$p_k \sim \mathcal{N}\left(p_0, \sigma^2 p_0(1-p_0)\right)$$

# STRUCTURE(AMA) Model

 $\begin{aligned} x_{il} | q_i, \mathbf{p}_l \sim \operatorname{Binom}(p_{q_i l}) \\ \text{(SNPs depend on pop SNP frequency)} \\ p_{kl} \sim \operatorname{Beta}(\beta_l) \\ \text{(Pop SNP frequency drifts independently)} \\ \beta_l \sim P(\beta_l) \\ \text{(Some ancestral frequency model)} \end{aligned}$ 

 $\mathbf{q}, K \sim \mathrm{DP}(\alpha, G_0)$ (Some prior on assigning individuals to populations) p is a vector of length K
The Dirichlet Process prior is for K and the number of individuals in each population
STRUCTURE\* uses a different prior for q, K

13/43

- Pella & Masuda 2006
- STRUCTURAMA (Huelsenbeck & Andolfatto 2007)

Daniel Lawson, University of Bristol dan.lawson@bristol.ac.uk

\*Pritchard, Stephens and Donnelly, Genetics, 155:945-959, 2

#### Extension to large datasets

- Lots of latent parameters (SNP frequencies, etc) makes MCMC slow and sticky for large datasets
- Also, deal with correlations between SNPs
- Return to model?
- Start with summary statistic approaches

# Clustering

• K clusters (which should be inferred) • Find 'best' assignment of individuals into populations • 'Best' is some score, potentially using: Raw Data (Direct)  $Y_{il}$ Dimensionality Distances between  $X_{ii}$  $V_{id}$ reduction points (Similarity-(Spectral/PCA) based)  $L \times N$  $N \times N$  $D \times N$ Daniel Lawson, University of Bristol  $L \gg N \gg D$ 15/43dan.lawson@bristol.ac.uk

## Pairwise summary statistics

Allele Sharing distance/Identity by state  
Edit distance/Norm' distance e.g. L1, L2Count the number of shared  
SNPs
$$\sum_{l} |Y_{il} - Y_{jl}|$$
Covariance $\sum_{l} (Y_{il} - f_l) (Y_{jl} - f_l)$ PCA uses the covariance $\sum_{l} \frac{(Y_{il} - f_l) (Y_{jl} - f_l)}{f_l(1 - f_l)}$ EIGENSTRAT use this normalisation  
to do PCA

Daniel Lawson, University of Bristol dan.lawson@bristol.ac.uk

16/43

 $\bullet \bullet \bullet \bullet \bullet$ 

#### **Dimensionality** reduction



# How many dimensions?

#### • MAP (Minimum Average Partial <u>Correlation), Velicer 1976</u>

- Remove largest eigenvalue
- Compute (partial) correlation between remaining eigenvectors and the data (accounting for previous eigenvectors)
- Repeat
- Tracy-Widom Distribution (TW 1994, Patterson et al 2006)
  - Theoretical distribution of the EVs of an *L* by *N* matrix (*L* is the (effective) number of SNPs)
  - Remove biggest EV if bigger than some quantile
  - Repeat

#### • Parallel Analysis (Horn 1965)

• Simulate many matrices the same shape, with the same mean and variance as the data

N

• Keep all components bigger than some quantile of the simulated values



### Clustering: MVN

Multivariate Normal ("Soft K-Means", implemented by MCLUST in R)
Infer mean and variance for each cluster
BIC model selection for *K*(*Alternatives: K-means, UPGMA, etc*)



#### The coancestry matrix

• (unlinked) coancestry matrix:

$$Y_{ij} = \sum_{l=1}^{L} \frac{x_{il} x_{jl}}{\sum_{j \neq i} x_{jl}} + \frac{(1 - x_{il})(1 - x_{jl})}{\sum_{j \neq i} (1 - x_{jl})}$$

- To O(N), Coancestry matrix is a rotation of the (EIGENSTRAT PCA) eigenvector matrix
  - If SNPs are uncorrelated
  - and the number of individuals is large
- To O(N), Coancestry matrix is a *sufficient statistic* for the STRUCTURE likelihood
  - If additionally, drift is small

#### Local genealogies



ChromoPainter 'Coancestry' similarity matrix
Unlinked limit: normalised allele sharing



Time to MRCA with haplotype 1

See: Li and Stephens, Genetics 165:2213-2233, 2003

1.1

1.24

0.52

0.52

0.06

0.01

0.06

0.08

0

0.09

Haplotype 1

### Li and Stephens Hidden Markov Model

#### Closest haplotype ('painting')



Switches occur at constant rate ρgl (Scaled Genetic distance)
Mutations occur at constant rate θ
Can efficiently infer X<sub>1</sub>...<sub>L</sub> (painting)
And X<sub>l</sub>|X<sub>l-1</sub> (switches, the number of 'chunks')

#### Similarity measures of simulated data





#### FineSTRUCTURE Population structure model

Individuals exchangeable within populations

$$x_{ab} = \sum_{i \in a, j \in b} x_{ij}$$

• Populations donate chunks independently at a characteristic rate  $P_{ab}$ 

$$p(X|P) = \prod_{a,b=1}^{K} \left(\frac{P_{ab}}{\hat{n}_{b}}\right)^{x_{ab}}$$
Population assignment
Number of individuals to donate from
population

Coancestry matrix

#### Probability of a partition

• Dirichlet Process prior for partition  $\eta$  :

 $\eta \sim \alpha^K \prod_{k=1}^K \Gamma(\hat{n}_k) \qquad \{P_1, \cdots, P_K\} | \eta = \prod_{b=1}^K G_0$ 

• Rows of  $P_{ab}$  (i.e.  $G_0$ ) are Dirichlet (containing hidden biological parameters)...

• ... so conjugate, and we integrate out  $P_{ab}$ 

• No individual/population level parameters

#### Additional results

- To O(N), FineSTRUCTURE likelihood is equivalent to the STRUCTURE\* likelihood
  - *if SNPs are uncorrelated,*
  - -drift is weak,
  - -genotyped SNPs are not very rare
- Empirically, with linkage model we do better.

\*Pritchard, Stephens and Donnelly, Genetics, 155:945-959, 2000 Calculations due to Simon Myers

### More details

- MCMC based Bayesian inference
- Post-processing step for a tree relating populations
   Identify & merge close relatives first
- Infer *K*: Structure bar plots don't make sense
- New ways to visualise results needed...

#### HGDP dataset

#### 650K SNPs on 938 individuals from 53 pops (5-45 inds/pop)



website: http://hgdp.uchicago.edu/ (picture from Cvalli-Sforza (2005) Nat Rev Genet)

Slide courtesy of GH

### MAP tree: whole world HGDP data







### Other uses

- Other quantities:
  - Matrix  $L_{ij}$  of the length of haplotype chunks
  - Matrix  $M_{ij}$  of mutations on haplotype chunks
  - $-P(X_{l+d} = b | X_l = a)$  Correlation between 'origin' of SNPs as a function of genetic distance
- Applications
  - Population inference
  - Identifying groupable individuals (e.g for GWAS, demographic inference, other modelling)
  - ... Haplotype-based association studies?
  - Admixture dating

### Admixture dating

- Expected length of haplotypes between populations halves each generation
- Exponential tract length distribution





#### Brahui-Yoruba 95/5% admixture (30 gens)









#### Former Mongolian Empire



Daniel Lawson, University of Bristol dan.lawson@bristol.ac.uk

#### Slide courtesy of GH

### Personal genomics (e.g. 23andMe)

(Sharing approx 250K SNPs)

- http://fennoscandia.blogspot.co.uk (~200 Scandanavian individuals)
- http://www.harappadna.org (~700 South Asian individuals)
- http://dodecad.blogspot.co.uk (~400 Balkans/West Asia individuals)
- http://eurogenes.blogspot.co.uk (~500 European individuals)
- http://magnusducatus.blogspot.co.uk (~100 Lithuanian individuals)



#### Summary statistics (requiem)

• FastIBD (Browning & Browning 2011) is another useful linked summary statistic • Different features to ChromoPainter • If you know the correct number of PCA components to retain, you can do very well with PCA and simple models • But it goes very wrong in some (real) circumstances • New linked approaches will be developed • ChromoPainter/FineSTRUCTURE pipeline is a lot more robust than model-free alternatives

### Further info

 ChromoPainter/FineSTRUCTURE software/code/info:

www.paintmychromosomes.com

- Includes GUI: Windows/Linux/Mac (cmd line only)
- ChromoPainter/FineSTRUCTURE publication: Lawson et al 2012 PLoS Genetics
- Summary statistics review: Lawson & Falush 2012 An. Rev. Hum. Genet. (to appear)
- Admixture dating paper in preparation
- POBI paper in preparation
- 'Admixture' model in the works

#### Acknowledgements:

#### fineSTRUCTURE



Garrett Hellenthal (Oxford) (CP algorithm, admixture dating)



Simon Myers (Oxford) (theory, admixture dating)



Daniel Falush (Max Planck Institute) (CP/FS concept)

- Peter Green (Bristol) Grant, support
- Bluecrystal HPC facilities @ Bristol
- <u>www.paintmychromosomes.com</u>





(Individual labels not shown)

#### FastIBD (Browning and Browning 2011)

- Alternative linked model: Identify *r* closest segments of DNA for each pair of individuals
- Genetic lengths of each are related to time since common ancestor
- Similarity measure: sum of the genetic lengths found for each pair
- Somewhat heuristic, has some tuning parameters, but empirically works well



#### Comparison of clusterings







## Weak Biological Model for prior

'Correct' Ancestral Recombination Graph for the limit of large populations at large time with simple population structure



#### Posterior evaluation

- MCMC update of hyperparameters and partitions
- Partition moves:
  - Move an individual
  - -Merge
  - Split
  - -Merge and resplit
- Merge/split 'nearly Gibbs' move:

$$p(q_{m};a,b) = p(q_{1}) p(q_{2}|q_{1}) \cdots p(q_{m}|q_{1:m-1})$$
  
$$p(q_{m}=a) \approx \hat{n}_{a} \int F(x_{m}|P_{m}) dH_{$$

(Not exact as the 'unsplit' population interacts with the remaining dataset)

Simple case: Pella and Masuda Canadian J. Fish. Aquatic Science 63:576-596, 2003 Daniel Lawson, University of Bristol dan.lawson@bristol.ac.uk

51/43

#### Clustering into k clusters

- Find "similar" individuals l∈ [1, L]
  Three main approaches:

  Cluster on raw data Y<sub>il</sub>
  Cluster on similarity matrix X<sub>ij</sub>
  Cluster on dimensionality reduced version of data, e.g. MDS/PCA/SVD V<sub>id</sub>
  d ∈ [1, D]
  - Recall:  $L \gg N \gg D$  O(k) = O(d)?
- Lowest dimension description usually best...
  Raw data approach terrible here (without good model)

#### The future – Admixture model

- Pure population structure is not correct recent mixing leads to admixture
  - Seek conjugate mixture model for individuals
  - Hierarchical Dirichlet Process!
  - Interpretation: Pure populations created by drift, we see mixtures
- Better model:
  - Allow drift and admixture to both occur in real time
  - Requires more sophisticated model, can we keep conjugacy?

53/43

-(Matrix Coalescent\* results available)

Daniel Lawson, Initichlet diffusion treas and Conceptetics, 161:1641-1650, 2002 dan.lawson@bristol.ac.uk \*\*Neal, in J. M. Bernardo, et al. (ed.), Bayesian Statistics 7, pp. 619-629, 2003

#### Comparison of linked methods



#### Posterior evaluation: building block

• Sample from posterior  $p(q_m; a, b) = p(q_1) p(q_2|q_1) \cdots p(q_m|q_{1:m-1})$ 

Metropolis-Hastings proposal for a split:
 – Random individuals creates population *a* and *b* from *c*

#### – Move rest from *c* with probability

 $p(m;a) \propto \hat{n}_{a} \int F(x_{m}|p_{m}) dH_{<m, S(p_{m})} \\ \approx n_{a} \frac{P(S_{a}, \{i=1, \cdots, m\}) P(S_{c}, \{i=1, \cdots, m\})}{P(S_{a}, \{i=1, \cdots, m-1\}) P(S_{c}, \{i=1, \cdots, m-1\})}$ 

(Not exact as the 'unsplit' population interacts with the remaining dataset)

Daniel Lawson, University of Bristol

dan.lawson@bristol.aExdact case: Pella and Masuda Canadian J. Fish. Aquatic Science 63:576-596, 200

55/43

#### Probability of a partition

Rows of  $P_{ab}$  are Dirichlet – Conjugate to multinomial, sum to 1 – Weak prior

Compute posterior incrementally due to conjugacy

$$p(x_a|q) = \prod_{m \in a} \int F(x_m|P_a, q) dH_{$$

$$dH_{\langle m, S_a}(P_a) = Dirichlet(P_a; \{\beta_{ab} + x_{\langle m, b}\}_{b=1, \cdots, K})$$

(Idea: add each individual, update Dirichlet posterior, use as prior for the next individual)

#### Final model

• Posterior  

$$p(\eta|X) \propto \alpha^{K} \prod_{a=1}^{K} \Gamma(\hat{n}_{a}) \frac{\Gamma(\beta_{a})}{\Gamma(x_{a}+\beta_{a})} \prod_{b=1}^{K} \frac{\Gamma(x_{ab}/c+\beta_{ab})}{\Gamma(\beta_{ab})\hat{n}_{b}^{x_{ab}}}$$

• Prior for hyperparameters

$$\beta_{ab} = \begin{cases} \gamma V_b & \text{if } a \neq b \\ \gamma (1 + \delta) V_b & \text{if } a = b \\ Prift \text{ due to mutation} & \text{Ancestral dot} \end{cases}$$

Ancestral donation frequency

 $\gamma = (1 - F)/F$   $\blacksquare$  Drift in allele frequency

#### Posterior visualisation

- Too many populations!
- Pairwise coincidence matrix
- Create MAP (maximum a posteriori) tree from MAP partition
- Show partition split posterior support
- (Population summary of data matrix *X*)

#### Comparison of linked methods

