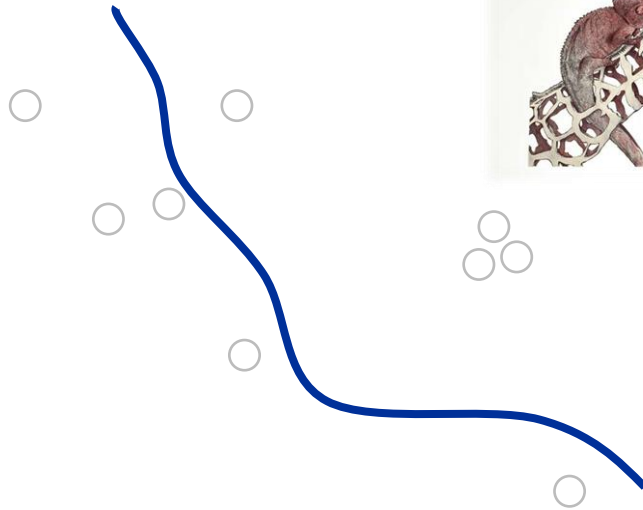


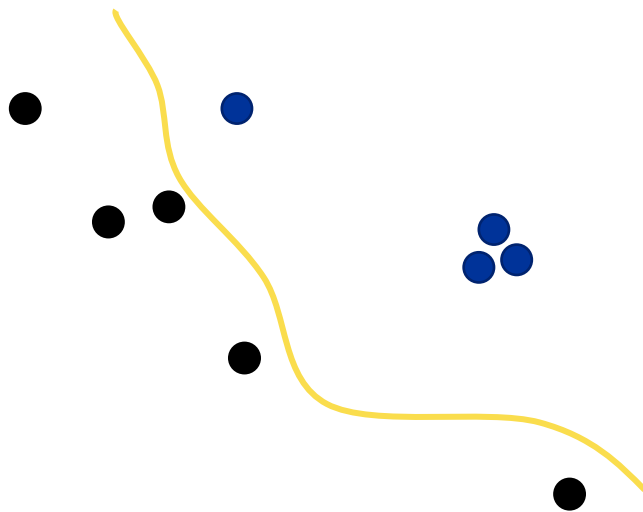
Assumption for population structure analysis:

- **neutral loci** = no effect of selection included
- **classical population genetics approach** = populations are (*thought to be*) known (e.g. we want to quantify level of genetic differentiation between two localities / ?populations)
- BUT populations are **not usually known** (e.g. due to no obvious spatial heterogeneity over the distribution range)
 - we want to **reveal any potential population differentiation/structure according to our genetic data**

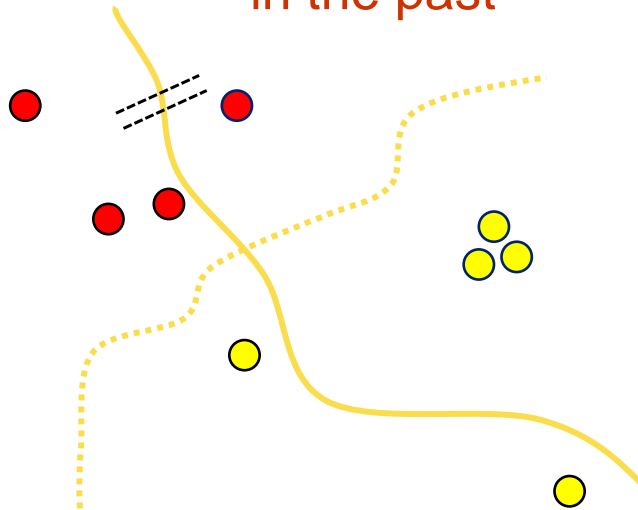
We have sampled animals in nature –
Is it one or several populations???



We are interested in genetic
structure of populations



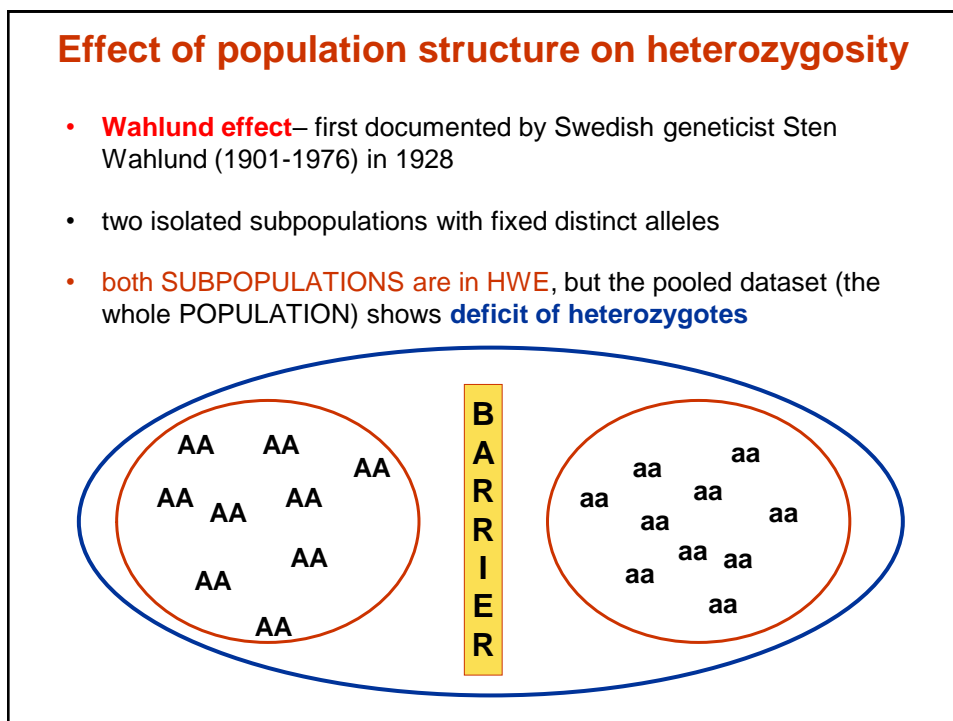
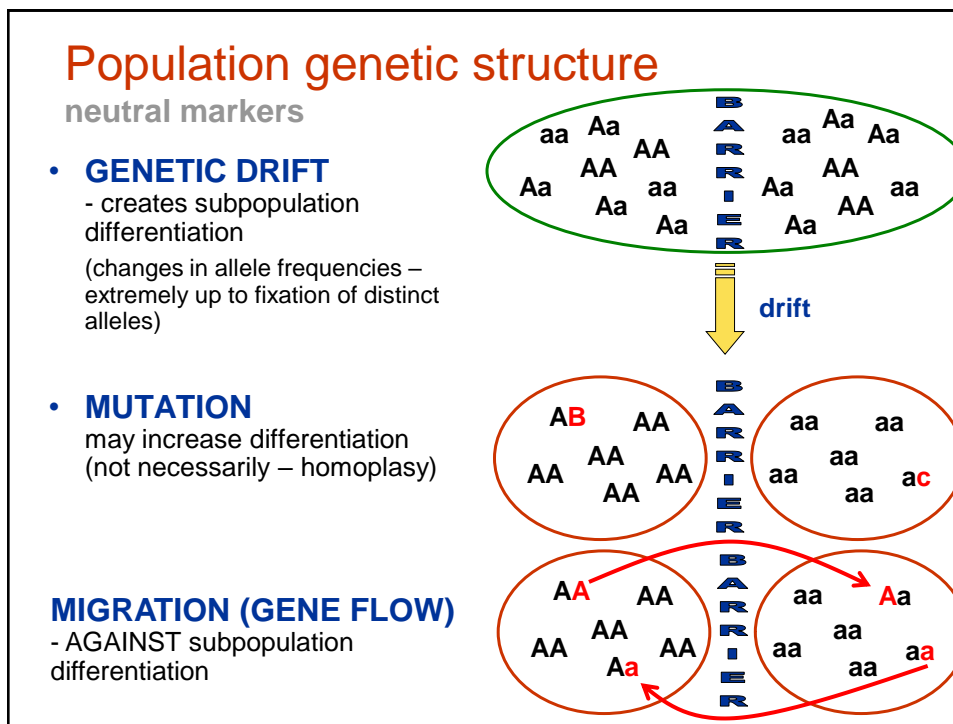
Recently observed genetic structure indicates what happened in the past



Genetic structure – any pattern in the genetic make-up of individuals within a population

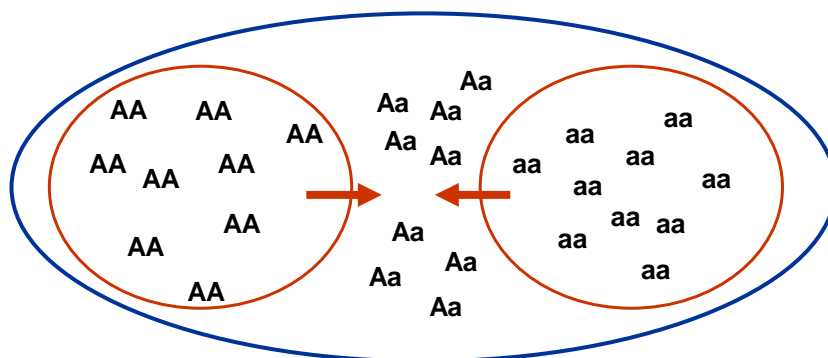
AIMS:

- Detection of **any** genetic structure (subdivision) in a population (in my dataset)
- Are there any **differences** between „different“ (in space and time) populations?
- Quantification of such differences = **description of genetic structure in population**
- What factors shape (have shaped) these differences? e.g. **population history**
- Is there any migration/connection between different populations? = detection and quantification of **gene flow**, what influences gene flow (e.g. **spatial heterogeneity**)
- What happens during migration/connection of populations? = **hybridisation**



Wahlund effect (isolate breaking)

Homozygosity reduction when subpopulations merge



Wahlund, S. (1928) Zusammensetzung von Population und Korrelationserscheinung vom Standpunkt der Vererbungslehre aus betrachtet. *Hereditas*, 11: 65–106

Wahlund effect – an example

- Bunnarsjöarna lake (northern Sweden) – „brown trout“
- one trait with 2 alleles

	170/170	170/172 (= Ho)	172/172	Total	p	2pq (=He)
Přítok	50	0 (0)	0	50	1.000	0.000
Odtok	1	13 (0.26)	36	50	0.150	0.255
Whole lake (expected)	51 (33.1)	13 (0.13) (48.9)	36 (18.1)	100	0.575	0.489

$$p^2 = 0.575^2$$

$$q^2 = 0.425^2$$



Ryman et al. 1979

Wright's F-statistics

$$F_{IS}, F_{ST}, F_{IT}$$



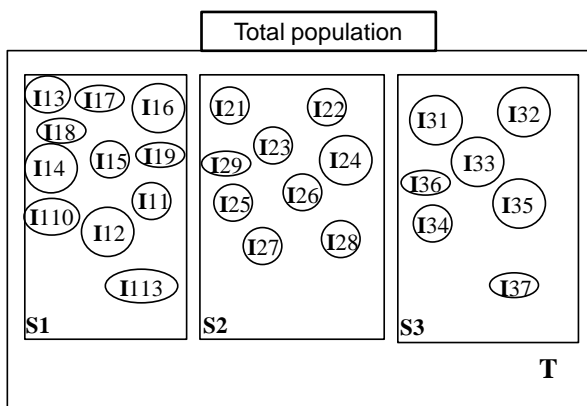
Masatoshi Nei
*1931



Sewall Wright
1889 - 1988

- Wright (1950), Nei (e.g. 1987)
- **detecting and describing population structure**
- describe heterozygosity (i.e. deviation from HWE) at different levels

Estimate of population structure effect on genetic diversity



- 3 levels (Total, Subpopulation, Individual)
- x subpopulations ($x = 1$ to k ; here $k = 3$)
- each subpopulation has N_x individuals
- AA, AB, BB – genotypes with different symbols
- e.g. I1-13 = 13st individual from the 1st subpopulation

F-statistics and heterozygosity

H_I – averaged observed heterozygosity of an individual in a subpopulation
 H_S – expected heterozygosity of an individual in a subpopulation **under HWE**
 H_T – expected heterozygosity of an individual over the total population under HWE

$$H_I = \sum_{x=1}^k H_x / k \quad H_x = \text{observed heterozygosity in subpopulation } x$$

$$H_S = 1 - \sum_{i=1}^j p_{i,x}^2 \quad p_{i,x}^2 = \text{frequency of } i\text{-th allele in subpopulation } x \quad \bar{H}_S = \sum_{x=1}^k H_S / k \quad \text{averaged expected heterozygosity in subpopulation}$$

$$H_T = 2p_0q_0 \quad p_o = \text{allele frequency in the total population}$$

- for two alleles at a single locus (Wright 1950)
- more complicated for more alleles (Nei 1987)

F-statistics

$$F_{IS} = \frac{\bar{H}_S - H_I}{\bar{H}_S} \quad \text{Heterozygosity decrease of an individual due to non-random mating in a subpopulation (vs. HWE)}$$

Heterozygosity over all populations →

$$F_{ST} = \frac{H_T - \bar{H}_S}{H_T} \quad \text{Influence of division of the total population in subpopulations (i.e. heterozygosity decrease due to Wahlund effect)}$$

$$F_{IT} = \frac{H_T - H_I}{H_T} \quad \text{Total coefficient of inbreeding } F_{IT} \text{ - measures heterozygosity decrease of an individual in relation to the total population}$$

$$(1 - F_{IT}) = (1 - F_{ST})(1 - F_{IS})$$

Weir & Cockerham (1984) f ($\sim F_{IS}$), θ ($\sim F_{ST}$), F ($\sim F_{IT}$)
 Correction for sample size and number of subpopulations

Computation of F-statistics

Locus	Subpopulation 1 ($N_1=40$)				Subpopulation 2 ($N_2=20$)				$p_{O(i)}$	Note
	AA	AB	BB	$p_{1(i)}$	AA	AB	BB	$p_{2(i)}$		
Loc I	10	20	10	0.5	5	10	5	0.5	0.5	HWE
Loc II	16	8	16	0.5	4	4	12	0.3	0.4	heterozygote deficit
Loc III	12	28	0	0.65	6	12	2	0.6	0.625	heterozygote excess
Loc IV	0	0	40	0.0	20	0	0	1.0	0.5	alternatively fixed alleles

Mean allele A frequency in the whole population

Computation of allele frequencies

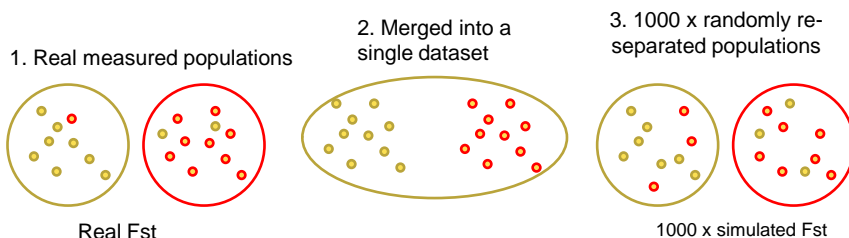
Locus	Observed heterozygosity		Expected heterozygosity			Wright's F-statistics		
	$H_{1(i)}$	$H_{2(i)}$	$H_{1(i)}$	$H_{S(i)}$	$H_{T(i)}$	$F_{IS(i)}$	$F_{ST(i)}$	$F_{IT(i)}$
Loc I	0.5	0.5	0.5	0.5	0.5	0.0	0.0	0.0
Loc II	0.2	0.2	0.2	0.46	0.48	0.565	0.042	0.583
Loc III	0.7	0.6	0.65	0.4675	0.46875	-0.39	0.0027	-0.387
Loc IV	0.0	0.0	0.0	0.0	0.5	---	1.0	1.0
Mean						0.058	0.261	0.300

Mean values of F-statistics may hide distinct evolution history of different loci

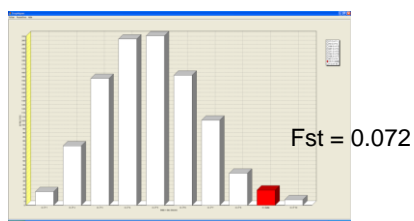
F-statistics

- F_{IS} decrease of heterozygosity in local subpopulation
high values – inbreeding
- F_{IT} summary measure – limited use
- F_{ST} = **subdivision measure** = limited gene flow between subpopulations (i.e. existence of a barrier – Wahlundeffect)
 - originally developed for estimation of the amount of allelic fixation due to genetic drift (**fixation index**)

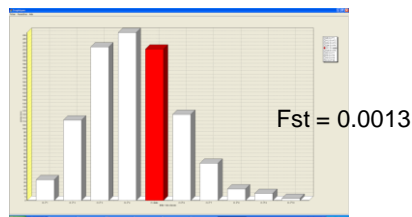
Permutation test of Fst significance



TWO DIFFERENT CASES:



0.80 % simulated values higher than real Fst
 $p = 0.008$ (i.e. significant difference)



35.40 % simulated values higher than real Fst
 $p = 0.354$ (e.g. non-significant difference)

F_{ST} computation – an example

	A/A	A/B (=H _o)	B/B	Total	p	2pq (=H _e)
Přítok	50	0 (0)	0	50	1.000	0.000
Odtok	1	13 (0.26)	36	50	0.150	0.255
Whole lake	51	13 (0.13)	36	100	0.575	0.489
(expected)	(33.1)	(48.9)	(18.1)			

$$F_{ST} = \frac{H_T - \bar{H}_s}{H_T} = \frac{0.489 - 0.128}{0.489} = 0.728$$

As a consequence of gene flow barrier:
Heterozygosity is about 72.8% lower than would be under HWE

Ryman et al. 1979



F_{ST} analysis – BE AWARE

Global vs. pairwise indices

Absolute values depends on heterozygosity level of used loci!!!

(i.e. microsatellite-based F_{ST} cannot be compared to allozyme-based F_{ST})

Demands standardization: $F_{ST}' = F_{ST}/F_{STmax}$ (Hedrick 2005)

– e.g. GenAlEx

In case of null alleles presence: needs to be corrected!
(increased F_{ST} – increase of homozygosity); FreeNA software



Giant Panda

- 192 feces samples → 136 genotypes → 53 unique genotypes
- separation by a river (ca 26 ky ago) and by roads (recently)
- even the roads are important barriers, even if less



(Zhu et al., 2011)

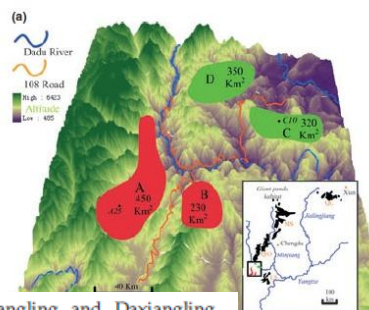


Table 3 Pairwise F_{ST} in the Xiaoxiangling and Daxiangling populations

Patch	A	B	C	D
A				
B	0.033*			
C	0.107*	0.062*		
D	0.107*	0.097*	0.037*	

*Significant level after Bonferroni correction ($P < 0.01$).

G_{ST} (Nei 1973)

- Analogy of F_{ST} for **haploid (haplodiploid) organisms, mtDNA sequences**
- Takes into account **haplotype (gene) diversity** instead of heterozygosity
- *Haplotype diversity* = probability that any two randomly chosen sequences in a population will be different
- Pracuje tedy jen s frekvencemi alel, ne s procentem heterozygotů

R_{ST}

- Analogy of F_{ST}
- Takes into account **the size of alleles** (number of repeats in microsatellite loci)
- Assumption of a known mutation model
assumption of SMM (stepwise mutation model)
- Indicates traces of mutations
 - $R_{ST} > F_{ST}$ **higher effect of mutations**
 - $R_{ST} = F_{ST}$ **higher effect of genetic drift**
- Randomisation tests for R_{ST} significance (Hardy et al. 2003, program SPAGeDi 1.1)


AMOVA

Excoffier et al. 1992

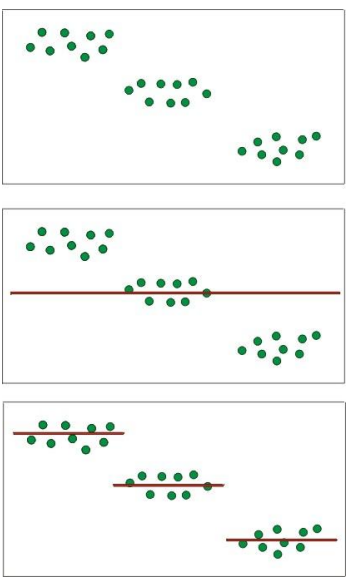
Arlequin ver. 2.000
A software for population genetics data analysis

Authors:
Stefan Schneider
David Roessli
Laurent Excoffier

Contact Arlequin:
Url: <http://anthropologie.unige.ch/arlequin/>
Mail: arlequin@sc2a.unige.ch



- **A**nalysis of **M**olecular **V**ariance
- Analysis of allele frequencies variance (before in *Cockerham & Weir 1987, 1993*)
- **Quantifies population differentiation**
- Takes into account difference between alleles – allelic state (mutations)
- Program ARLEQUIN
- Data:
sequences
microsatellites (assuming SMM *stepwise mutation model*)



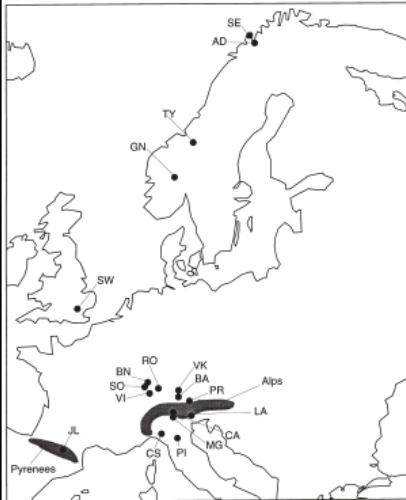
Hierarchical AMOVA

How much variation may be explained by:

- differentiation in big **groups of populations**
- differentiation in **populations** within the groups
- differentiation between **individuals** within the populations

Bombus pascuorum

Widmer & Schmid-Hempel 1999



	F / Φ	d.f.	SSD†	Variance component	% Total variance*
Among populations	F	17	77.71	0.07	4.51*
	Φ	17	5198.20	5.02	8.74*
Among regions	F	4	56.15	0.08	5.16*
	Φ	4	3464.94	4.58	7.49*
Among populations within regions	F	11	24.35	0.02	1.11*
	Φ	11	1773.71	2.16	3.53*
Between north and south of Alps	F	1	38.57	0.11	7.12*
	Φ	1	2622.89	7.25	11.74*
Among populations north and south of the Alps, respectively	F	16	39.14	0.02	1.46*
	Φ	16	2575.31	2.18	3.53*

†Sum of squared deviations.

* $P < 0.001$.

Microsatellites, AMOVA
Most explained by the Alps

AMOVA and F-statistics

description of results, not causes → possible alternative explanations
(use of population history analyses – based on coalescence and allele phylogenetics)

Recent separation,
no gene flow

a ← d → b

Old separation, but
continuous (low)
gene flow

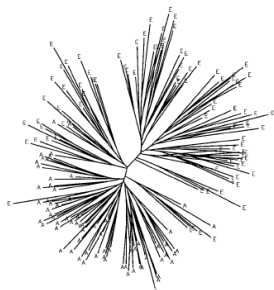
a ← d → b

Time

Clustering methods

DISTANCE-BASED methods

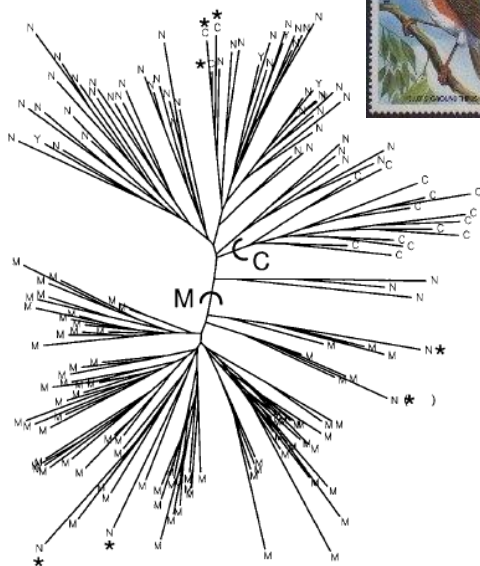
- a tree or a plot is constructed according to a **pairwise distance matrix**
- clusters then may be defined **visually**



MODEL-BASED methods

- observations from each cluster are random draws from some parametric **model**
- **inference for the parameters** corresponding to each cluster is done jointly with **inference for the cluster membership** of each individual
- standard statistical methods are used (e.g. maximum-likelihood in Bayesian methods)

Turdus helleri



- Fragments of humid tropical forest
- Localities Chawia, Ngangao, Mbololo, Yale (Kenya)
- 7 microsatellite loci
- Neighbour-joining
- * wrongly clustered individuals

Clustering method based on microsatellite distances

Factorial correspondence analysis

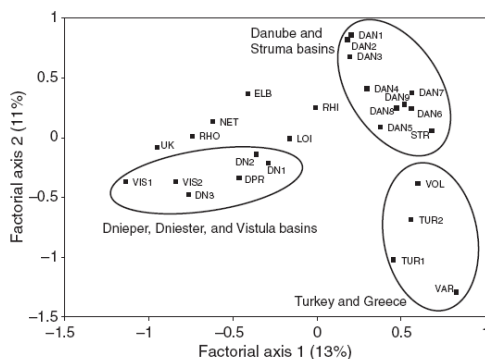


Fig. 2 A two-dimensional plot of the factorial correspondence analysis performed using GENETIX based on 12 microsatellite loci. Three geographical groups are bounded by grey lines.

- each locus as one variable, reduction of number of variables
- **Genetix** – orientační zjištění strukturovanosti populace
- individuals vs. populations

STRUCTURE program

Pritchard, Stephens and Donnelly 2000, Genetics

- a model-based Bayesian clustering method
- uses multilocus genotype data (e.g. microsatellites, RFLPs, SNPs; various levels of ploidy)
- MCMC algorithm
- INFERS POPULATION STRUCTURE:
 - presence of population structure
 - assignment of individuals to populations
 - identification of migrants or admixed individuals (parameter Q – individual membership coefficient)

Model implemented in STRUCTURE assumes:

- **K populations/clusters (K may be unknown)**
- each of K populations is characterized by **a set of allele frequencies** at each locus
- **within each of K populations** marker loci are at LINKAGE EQUILIBRIUM with each other and in HARDY-WEINBERG EQUILIBRIUM

under these assumptions each allele at each locus in each genotype is an independent draw from the appropriate frequency distribution, and this is completely specified by the **probability distribution** $P(X|Z,P)$

X – genotypes of the sampled individuals

Z – unknown populations of origin of the individuals

P – unknown allele frequencies in all populations

MODELS in STRUCTURE



ANCESTRY MODELS

- no admixture model
- admixture model
- linkage model
- models with informative priors



ALLELE FREQUENCY MODELS

- independent frequencies model
- correlated frequencies model

Ancestry models:

NO ADMIXTURE MODEL

- each individual is discretely from one of the K populations
- the output reports the posterior probability that individual i is from population K
- the prior probability for each population is $1/K$

This model is appropriate for studying fully discrete populations and is often more powerful than the admixture model at **detecting subtle structure**.

Ancestry models:

ADMIXTURE MODEL

- individuals may have mixed ancestry
- each individual has inherited **some proportion** of its genome from each of the K populations = Q
- the output records **the posterior mean estimates** of these proportions

Recommended as a starting point for most populations.

“It is a reasonably flexible model for dealing with many of the complexities of real populations. Admixture is a common feature of real data, and you probably won’t find it if you use the no-admixture model.”

Allele frequency models:

INDEPENDENT FREQUENCIES MODEL

- the allele frequencies in each population are independent draws from a distribution that is specified by a **parameter λ**
- this prior says that we expect allele frequencies in different populations to be **reasonably different** from each other

Allele frequency models:

CORRELATED FREQUENCIES MODEL

- frequencies in the different populations are likely **to be similar** (probably due to migration or shared ancestry)
- this prior says that the allele frequencies in different populations may be **quite similar** between the populations
- better clustering for **closely related populations**
- but may increase the risk of over-estimating K
- *If one population is quite divergent from the others, the correlated model can sometimes achieve better inference if that population is removed.*

Falush, Stephens and Pritchard 2003, Genetics

MODELS in STRUCTURE



ANCESTRY MODELS

ALLELE FREQUENCY MODELS

- no admixture model

- admixture model

- linkage model
- models with informative priors

- independent frequencies model

- correlated frequencies model

How long to run it

it is not possible to determine suitable run-lengths theoretically
this requires some experimentation on the part of the user

burnin length: how long to run the simulation before collecting data to minimize the effect of the starting configuration

- typically a burnin of 10,000—100,000 is more than adequate

run length: how long to run the simulation after the burnin to get accurate parameter estimates

- several runs at each K , possibly of different lengths, and see whether you get consistent answers
- you can get good estimates of the parameter values (P and Q) with runs of 10,000–100,000 steps, but accurate estimation of $\Pr(X|K)$ may require longer runs
- at least 500,000

In practice your run length may be determined by your computer speed and patience as much as anything else.

STRUCTURE program

Pritchard, Stephens et Donnelly 2000, Genetics

The screenshot shows the STRUCTURE software interface. The 'Project Data' table is visible, with columns for Locus 1 through Locus 8. The rows represent different parameter sets, such as 'paramset_run_10 (K=2)', 'paramset_run_11 (K=3)', etc. The data in the table consists of numerical values representing genotypes at each locus for each parameter set.

Data format: genotypes of an individual in TWO rows

		loc_a	loc_b	loc_c	loc_d	loc_e
George	1	-9	145	66	0	92
George	1	-9	-9	64	0	94
Paula	1	106	142	68	1	92
Paula	1	106	148	64	0	94
Matthew	2	110	145	-9	0	92
Matthew	2	110	148	66	1	-9
Bob	2	108	142	64	1	94
Bob	2	-9	142	-9	0	94
Anja	1	112	142	-9	1	-9
Anja	1	114	142	66	1	94
Peter	1	-9	145	66	0	-9
Peter	1	110	145	-9	1	-9
Carsten	2	108	145	62	0	-9
Carsten	2	110	145	64	1	92

Needs to be specified:

number of individuals, ploidy of the data, number of loci, missing value symbol (integer)

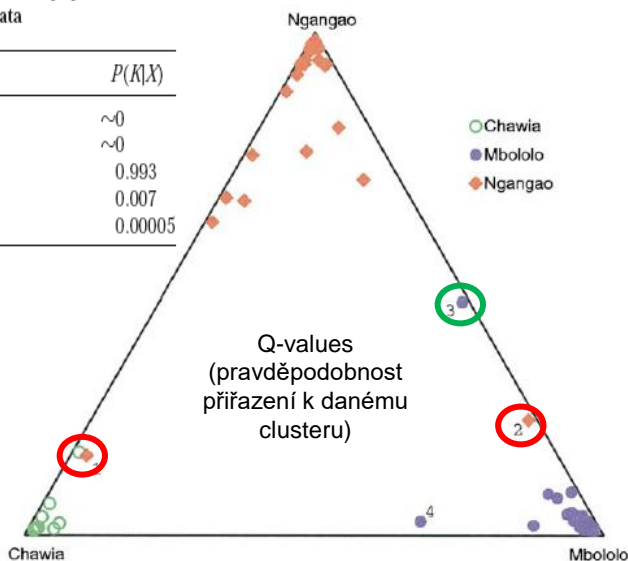
Program STRUCTURE – graphical output

Inferring the value of K , the number of populations,
for the *T. helleri* data

K	$\log P(\lambda K)$	$P(K \lambda)$
1	-3144	~ 0
2	-2769	~ 0
3	-2678	0.993
4	-2683	0.007
5	-2688	0.00005

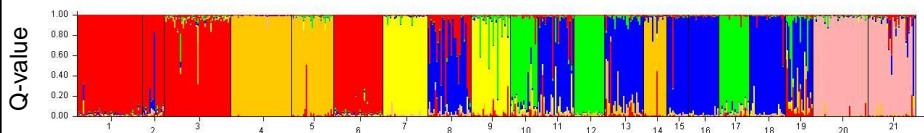
recent migrants

a hybrid?

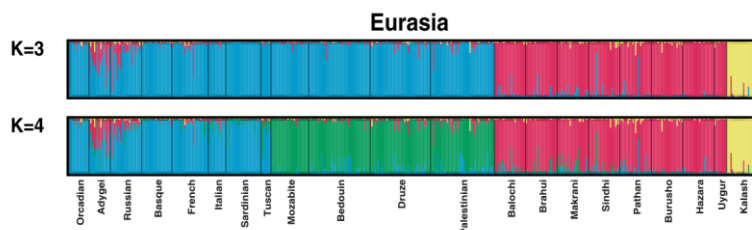


Admixture model – allows assignment of an individual to several clusters

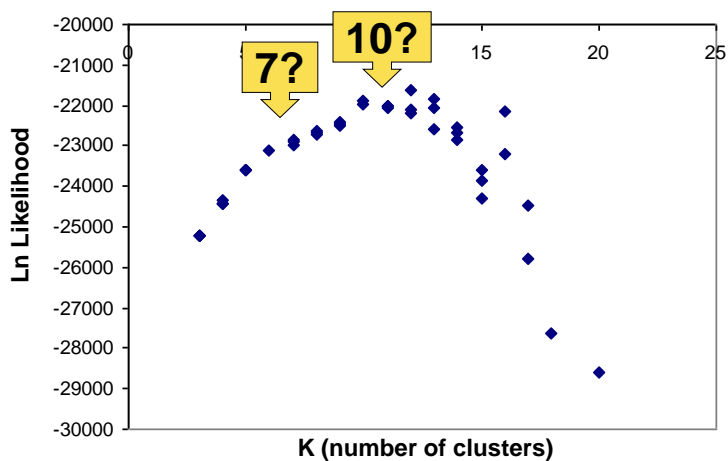
Barplot for $K = 7$



Genome proportion of each individual assigned to each of K clusters



What K is the best???



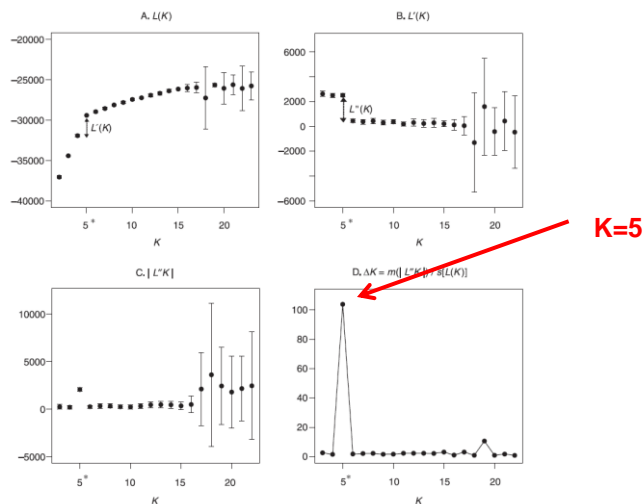
Molecular Ecology (2005) 14, 2611–2620

doi: 10.1111/j.1365-294X.2005.02553.x

Detecting the number of clusters of individuals using the software STRUCTURE: a simulation study

G. EVANNO, S. REGNAUT and J. GOUDET

Department of Ecology and Evolution, Biology building, University of Lausanne, CH 1015 Lausanne, Switzerland



Post-processing of the STRUCTURE outputs

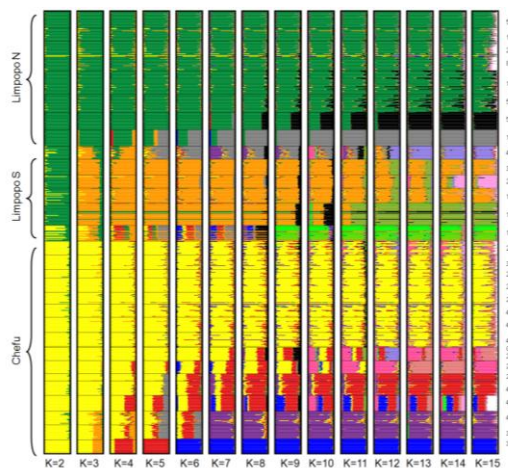
Main Pipeline Distruct for many K's Compare Best K Download Help Contact & Citing Issues

CLUMPAK - CLUSTER MARKOV PACKAGER ACROSS K

CLUMPAK was designed to aid users in four main objectives:

- Separate distinct solutions obtained from STRUCTURE-like programs.
- Compare and align solutions obtained for different K values.
- Compare results obtained using different models/data subsets/programs.
- Indicate the preferred value of K according to Evanno et al.

Graphical
output from
STRUCTURE –
a serie of
barplots with
increasing K

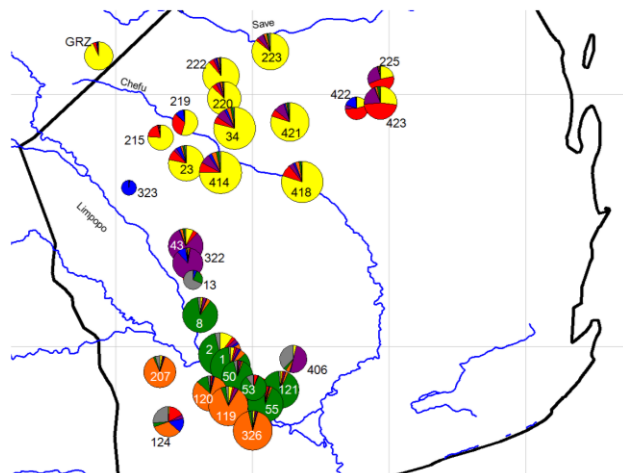


„forced clustering“

Picture of **hierarchical structure between clusters**

Bartáková et al. 2013

- Q-values for whole locality samples (not individuals)



Bartáková et al. 2013