**CULTURAL HERITAGE OF SMALL HOMELANDS**
**Food heritage – Seminar I.**
**4.3.2019 Nitra, Slovakia**

# Introduction to genetic identification of animals using low and high density data

Nina Moravčíková

Slovak University of Agriculture in Nitra
Faculty of Agrobiology and Food Resources
Department of Animal Genetics and Breeding Biology

---

# About this presentation

- **Part 1:** Analysis of diversity based on microsatellites
    - **1.1** Data manipulation
    - **1.2** Intro to diversity analysis
    - **1.3** Analysis of population structure and visualisation in R
    - **1.4** Practical exercise

- **Part 2:** Analysis of diversity on genome-wide level using SNPs
    - **2.1** Intro to PLINK
    - **2.2** Data processing in PLINK
    - **2.3** Analysis of population structure and visualisation in R
    - **2.4** Practical exercises

# Part 1: Analysis of diversity based on microsatellites

# 1.1 Data manipulation
## From biological samples to genetic analysis

1) isolation of genomic DNA from biological samples (blood, hair roots, semen, ….) in lab

2) lab identification of animals' genotypes by molecular-genetic methods based on use of genetic markers (microsatellites, SNPs, etc.) showing polymorphism
   - genetic polymorphism – occurrence of two or more genetically determined phenotypes in the same population

3) the quality control of genotyping data and development of input databases depending on the type of data

4) analysis of genetic diversity

# 1.2 Intro to diversity analysis

- frequency of alleles and genotypes
- Hardy – Weinberg equilibrium in population ($\chi^2$ test, ....)
- heterozygosity and homozygosity
- effective number of alleles
- polymorphic information content
- Wright's F – statistics: $F_{IS}$ – (f, molecular equivalent of pedigree inbreeding), $F_{ST}$ (genetic relationship between pops) a $F_{IT}$ (total inbreeding in metapopulation)
- genetic distance (Nei's $D_a$, ....)
- analysis of molecular variance (AMOVA)
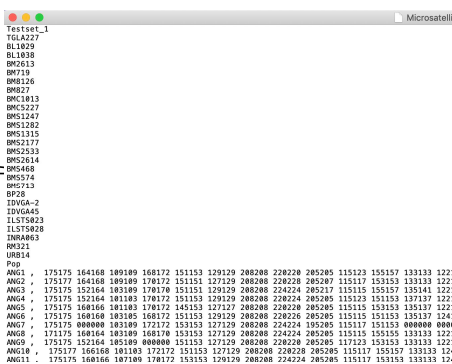- analysis of principal components (PCA)

# 1.2 Intro to diversity analysis
**Input data in genepop format**

**Genepop**
- one of the most commonly used format of data
- Genetix, Genepop, Arlequine, various R packages

Microsatellite_testset_1 ☞

## 1.2 Intro to diversity analysis
### Analysis based on the frequency of alleles

**Frequency by Pop:** outputs allele frequencies at each locus by population

**Frequency by Locus:** outputs allele frequencies in each population with loci in columns

## 1.2 Intro to diversity analysis
### Analysis based on the frequency of alleles

**Heterozygosity**
- Observed
- Expected

**Wright's F statistics ($F_{IS}$, $F_{ST}$ and $F_{IT}$)**

# 1.2 Intro to diversity analysis
**Analysis of diversity on intra- and inter-population level**

**Basic indices:**

- **MNA** – mean number of alleles per populations and loci
- **ENA** – effective allele number per populations and loci
- **AR** – total number of alleles per loci within each population
- **Heterozygosity** – average number within and across populations
- **Koeficient of inbreeding (F resp. $F_{IS}$)** – within and across populations
- **F – statistics**
- **AMOVA** (analysis of molecular variance)

# 1.2 Intro to diversity analysis
**Analysis of population genetic structure**

**Nei's genetic distances**

**Wright's $F_{ST}$ index**

**Analysis of principal components (PCA)**
- multivariate technique that allows one to find and plot the major patterns within a multivariate dataset e.g. multiple loci and multiple samples

# 1.3 Analysis of population structure and visualisation in R

➢ R Studio

➢Free …

…case sensitive!
…opposite slashes than Windows defaults

➢Package Adegenet

# 1.4 Practical excercise

- **Your turn!**

- Task description:

    **1. How many populations and markers are stored in testset 2?**

    **2. Compute basic diversity indices for each of analysed population**

    **3. Make AMOVA analysis**

    **4. Make PCA analysis.**

    **5. Make DAPC analysis and visualise the group assignment probability of individuals**

    **6. Plot the relationships between individuals using NJ unrooted tree**

# Part 2: Analysis of diversity on genome-wide level using SNPs

# 2.1 Intro to PLINK

- **PLINK** is a free, open-source whole genome association analysis toolset

- to perform a range of basic, large-scale analyses in a computationally efficient manner

## 2.1 Intro to PLINK

- Runs in the command line and/or using R

- In Windows, Linux and Mac

- Options preceded by double dash "--"

- All options:
  http://pngu.mgh.harvard.edu/~purcell/plink/reference.shtml#options

## 2.2 Data processing in PLINK using R
**Data manipulation**

- **Task:** Quality control of genotyping data

**QC criteria**

1. missing genotypes per sample max. 10 %
2. min. SNPs call rate 90 %
3. min. minor allele frequency (MAF) 0.01

# 2.3 Analysis of population structure and visualisation in R

- **Task:** load data in to R and run analysis

**load library for adegenet**

**load input data**
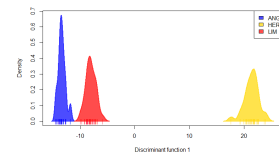
**run DAPC analysis**

---

# 2.3 Analysis of population structure and visualisation in R

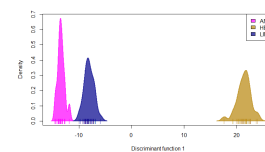- **Task:** load data in to R and run analysis

**visualise differentiation between pops based on first two DFs**

**visualise differentiation between pops based on first DF**
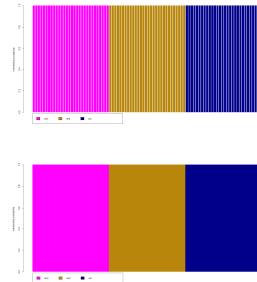
**change colours in picture**

## 2.3 Analysis of population structure and visualisation in R

- **Task:** load data in to R and run analysis

**make barplot to represent the group assignment probability of individuals to several groups**

**without spaces between columns in barplot**

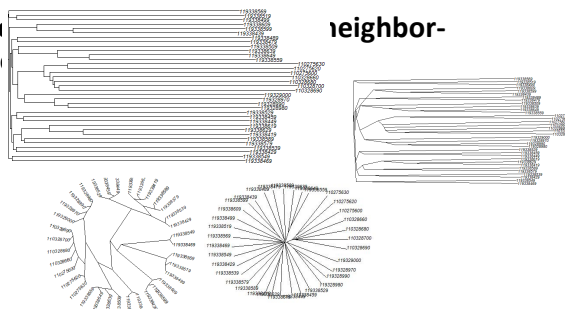## 2.3 Analysis of population structure and visualisation in R

- **Task:** plot simple neighbour-joining (NJ) tree

**load library for adegenet**

**perform the neighbor-j**          **neighbor-joining tree estimation**

**plot NJ tree**

**change type of tree**

# 2.3 Analysis of population structure and visualisation in R

- **Your turn!**

- Task description:

  **1. Make QC of data and create new dataQC .ped and .map file**

  - filter out all of animals and SNPs with more than 10% of missing data and MAF<0.01

  **2. Create input file for adegenet**

  **3. Make DAPC analysis and visualise the group assignment probability of individuals**

  **4. Plot the relationships between individuals using NJ unrooted tree**

Thank you for your attention!