# Package "JePhCort"

May . 2012

**Title :** JePhCort - A genotype-phenotype correlation tool based on phylogenetic analysis.

**Version :** 1.0

**Author :** Dr. Farhat Habib, Amol Kolte

**Maintainer :** Amol Kolte <amolkolte1989@gmail.com>

**Description :** A purely Unix-based bioinformatics tool to locate phenotype associated genotypic markers (SNPs) by taking into account the phylogenetic relationship among the species. JePhCort has been built using R and python.

**Dependencies:**

| | |
|---|---|
| python 2.7 (not 3.0) | R (>=2.14.0) |
| ete (2.1 alpha) | ape (<=2.8) |
| python-tk | igraph |
| numpy (1.6.1) | matrix (1.0) |
| scipy (0.9) | quadprog (1.5) |
| rpy (2.25) | phangorn (1.5) |

**Package Compilation:** *JePhCort* do not require package compilation. It mainly uses 2 standalone scripts (resurrect.R and reanimate.py) :

- ressurect.R (Performs ancestral sequence reconstruction)

- reanimate.py (Performs phenotype reconstruction and permutation test)

Keep all the scripts under single location, make sure all the dependancies are satisfied

**Running *JePhCort***

**Step 1.** Perform ancestral sequence reconstruction

resurrect.R ⟨alignemnt_sequence_file⟩ ⟨newick_tree_file⟩ ⟨fasta/phylip⟩

**Step 2.** Perform ancestral phenotype reconstruction and correlation

reanimate.py -s ⟨alignemnt_sequence_file⟩ -t ⟨newick_tree_file⟩ -f ⟨fasta/phylip⟩ -i ⟨No_iterations_for_permutation_test⟩ -p ⟨phenotype_file⟩ -o ⟨output_file⟩

**Files and formats :** *JePhCort* obligatorily requires three user input text files. The proper formats are discussed below :

i) **Sequence file** – Nucleotide sequences (SNPs) can be submitted in the standard *fasta* or *phylip* format.

ii) **Phylogeny/Tree file** – In the standard *newick* format (with branch lengths) as shown in the adjacent figure.
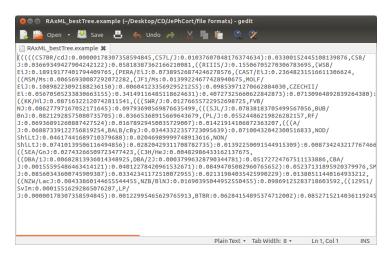


Figure 1: *Newick* file format

iii) **Phenotype file** – This is a simple tab-separated text file in a customized format. Depending upon the availability of the phenotypic data, user can choose between two-column or the three-column file format (as shown in Figure 3).

      The first column in Two-column file format represents 'name of the species' and other represents the 'mean phenotypic value'. On other hand, Three-column file format possesses an extra column of 'standard-deviation in the phenotypic values' (as shown in Figure 4).

**Result Interpretation :** *JePhCort* generates a tab-separated text file once the operation is successfully completed. The file consists of five entities, which are explained below.

1. **SNP serial ID** – Every row in the result file represents a genotypic marker *i.e.* an SNP. They are alloted serial numbers (starting from *Zero*) as per their position in the original sequence file. Thus, the $n^{th}$ SNP serial ID in the result represents the $(n+1)^{th}$ SNP in the original file.

2. **p-value** – Lower the p-value, higher is the significance of association between a given SNP and the phenotype.
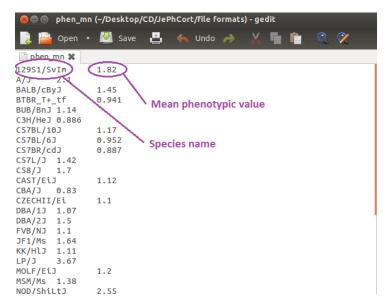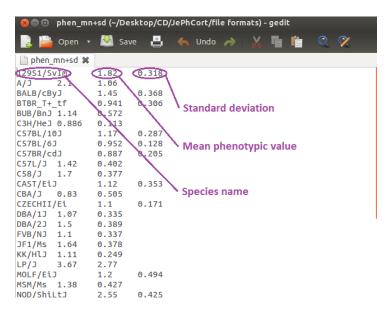
Figure 2: Two-column phenotype file format



Figure 3: Three-column phenotype file format