

Creation and simulation of personalized microbiota models through metagenomic data integration

Author: Federico Baldini, Molecular Systems Physiology Group, University of Luxembourg.

INTRODUCTION

This tutorial shows the steps that *MyPipe* automatically performs to create and simulate personalized microbiota models through metagenomic data integration. Please note that this tutorial uses as an example a small dataset (4 columns and 30 rows) with the purpose of demonstrating the functionalities of the pipeline. We recommend using high-performance computing clusters when assembling and simulating from bigger datasets.

The pipeline is divided into 3 parts:

1. **[PART 1]** Analysis of individuals' specific microbes abundances is computed. Individuals' metabolic diversity in relation to microbiota size and disease presence, as well as, classical multidimensional scaling (PCoA) on individuals' reaction repertoire are examples.
2. **[PART 2]** 1. Constructing a global metabolic model (petri) containing all the microbes listed in the study. 2. Building individuals' specific models integrating abundance data retrieved from metagenomics. For each organism, reactions are coupled to their objective function.
3. **[PART 3]** A specific range of growth is imposed for each microbiota model and simulations under specific diet regimes are carried. Set of standard analysis to apply to the personalized models. PCA of computed MHPCs of individuals as for example.

USAGE

Normally, once provided all the input variables in the driver (*StartMyPipe*), the only action required is to run the driver itself. However, for this tutorial, we will disable the automon functionality and compute each section manually.

DRIVER

This file has to be modified by the user to launch the pipeline and to define inputs and outputs files and locations.

```
% We first set the paths to input and output files
inputDataPath=[0];

% path to where the CBRA Toolbox is located
global CBRA

% path to microbe models (download ROME models from https://www.vib.uzh.ch/rome/romeindex)
romePath = [getenv('HOME') filesep 'rome' filesep 'rot'];

% path where to save results
resPath = [CBRA filesep '.log'];

% path to and name of the file with dietary information. Here,
% we will use an "Average European" diet that is located in the
% DietaryRepresentation folder.
dietaryFilePath=[CBRA filesep 'papers' filesep '2018_microbiomeModelingToolbox' filesep 'resources' filesep 'AverageEuropeanDiet'];
```

Then we set the path and the name of the file from which to load the abundances. For this tutorial, to reduce the time of computations, we will use a reduced version of the example file (*normCoverageReduced.csv*) provided in the folder Resources: only 4 individuals and 30 strains will be considered. Please, note that abundances are normalized to a total sum of one.

```
abundancePath=[CBRA filesep 'tutorial' filesep 'analysis' filesep 'microbiomeModelingToolbox' filesep 'normCoverageReduced.csv'];
```

Next inputs will define:

1. name of the objective function of organisms
2. format to use to save images
3. number of cores to use for the pipeline execution
4. if to enable automatic detection and correction of possible bugs
5. if to enable compatibility mode
6. if stratification criteria are available
7. if to simulate also a rich diet
8. if to use an external solver and save models with diet
9. the type of PVA function to use to solve

The following setting should work for almost any system, but please check carefully to be sure these options are valid for you. A more detailed description of these variables is available in the documentation.

The same inputs needs to be set in the driver file *StartMyPipe* when running *myPipe* outside of this tutorial or directly in the *startMyPipe* function.

```
% name of the objective function of organisms
objFnc=[ 'ss_biomass(e)' ];

% the output is a vectorized picture, change to '-png' for .png
figFnc = '-degplot';

% number of cores dedicated for parallelization
numWorkers = 2;

% autoFix for names mismatch
autoFix = true;

% if outputs in open formats should be produced for each section
debug = false;

% if documentation (.css) on stratification criteria is available
loadOfFilePath='none';
```

```
% To enable also rich diet simulations
rtol = false;

% If to use an external solver and case models with diet
external = false;

% The type of PBA function to use to solve
funcType = true;

% To turn off the autotask to be able to manually execute each part of the pipeline
autotask = false;

[subject, basePath,~, resPath, dietFolderPath, subFolderPath, objFile, TugPath, numWorkers, autoFile, compMod, rtol, external, funcType, autoFile]=
```

PIPELINE: [PART 1]

The number of organisms, their names, the number of samples and their identifiers are automatically detected from the input file.

```
[pathObj, sampleName, strIndex] = getIndividualsSettings(subFolderPath)
```

Now we detect from the content of the results folder if PART1 was already computed: if the associated file is already present in the results folder its execution is skipped else its execution starts

```
[magP] = detectOutput(resPath, 'magInfo.mat');

if isempty(magP)
    % Mapping file found: loading from resPath and skipping [PART1] analysis!
    disp('')
    load(strcat(resPath, 'magInfo.mat'))
end
```

In case PART 1 was not computed we will compute it now. We will first load the models and create a cell array containing them. This cell array will be used as input by many functions in the pipeline. Any possible constraint from each model reactions will be removed. Moreover we will run and subsequently plot the results of some analysis that are computed. The main outputs are:

- Metabolic diversity** The number of mapped organisms for each individual compared to the total number of unique reactions (extrapolated by the number of reactions of each organism). Please, note that bigger circles with a number inside represent overlapping individuals for metabolic diversity.
- Classical multidimensional scaling of each individual reactions repertoire**

Other outputs computed during this phase are saved together with the previous ones into the .mat file called **magInfo.mat** if the **compMod** option is enabled (disabled here and by default in the **magPipe** pipeline) these results are outputted as different .csv files. For simplicity reasons we will not discuss these additional outputs in this tutorial for a description of them, please refer to the documentation.

```
[magP] = detectOutput(resPath, 'magInfo.mat')
```

```
if isempty(magP)
    % Loading models
    model = loadModel(resPath, strIndex, objFile);

    % Computing genetic information
    [reac, genes, kIndex, gIndex, resPath, reactObj, reactTab, reactTab, reactObj, reactObj] = getMagInfo(model, subFolderPath, pathObj);
    writeTable(cell2table(reactObj, 'VariableName', {'Reactions'; sampleName}), [strcat(resPath, 'Reactions.csv')]);

    % Plotting genetic information
    [PCoA] = plotMappingToFa(resPath, gIndex, reactObj, reactTab, reactObj, subFolderPath, figure, sampleName, strIndex);

    if compMod == 1
        andir(strcat(resPath, 'compfile'))
        writetable(arraytable(reactObj, arraytable(reactTab, 'VariableName', sampleName))), [resPath 'compfile' filesep 'reactTab.csv'];
        writetable(arraytable(reactObj, 'VariableName', sampleName), [resPath 'compfile' filesep 'reactObj.csv']);
        writetable(arraytable(strIndex, arraytable(reactObj, 'VariableName', sampleName))), [resPath 'compfile' filesep 'reactObj.csv'];
        coverIndex(strcat(resPath, 'compfile/PCoA_tab.csv'), PCoA);
    end

    % Make all the created variables
    % Create tables and save all the created variables
    reactTab(arraytable(reactObj, arraytable(reactTab, 'VariableName', sampleName))), [resPath 'compfile' filesep 'reactTab.csv'];
    reactObj(arraytable(reactObj, 'VariableName', sampleName));
    reactObj(arraytable(strIndex, arraytable(reactObj, 'VariableName', sampleName)));

    save(strcat(resPath, 'magInfo.mat'))
end
end % trigger for autotask
```

PIPELINE: [PART 2.1]

Checking consistency of inputs: if autoFile == 0 halts execution with error mag if inconsistencies are detected, otherwise it really tries hard to fix the problem and continues execution when possible.

```
[autoStat, fixVec, strIndex] = checkConsistency(strIndex, autoFile);
```

Now we detect from the content of the results folder if PART2 was already computed: if the associated file is already present in the results folder its execution is skipped else its execution starts

```
[magP] = detectOutput(resPath, 'setup_allData.mat');

if isempty(magP)
```

```

        modbuild = 1)
    else
        modbuild = 0)
    in "global setup file found: loading from repaths and skipping [PART.1] analysis"
    disp(s)
end
end of Trigger For Autoload

```

A model joining all the reconstructions contained in the study will be created in this section. This model will be later used, integrating abundances coming from the metagenomic sequencing, to derive the different microbiota models. The result of this section will be automatically saved in the results folder.

```

17 modbuild = 1
   setupFactSetupReactor(modile, strains, {}),obj(re)
   setupNames "Global reconstruction with lumen / fecal compartments as host"
   setupReactor
   save(ct:cat(repaths,"setup_allfact.mat"), "setup")
end

17 modbuild=0
   load(ct:cat(repaths,"setup_allfact.mat"))
end

```

PIPELINE: [PART 2.2]

Now we will create the different microbiota models integrating the given abundances. Coupling constraints and personalized "cumulative-biomass" objective functions are also added. Models that are already existent will not be recreated, and new microbiota models will be saved in the results folder.

```
[createdModels]=createMicrobiotaModels(absolutePaths, repaths, setup, sampleName, strains, pathWeb)
```

PIPELINE: [PART 3]

In this phase, for each microbiota model, a diet, in the form of set constraints to the exchanges reactions of the diet compartment, is integrated. Flux Variability Analysis for all the exchange reactions of the diet and fecal compartment is also computed and saved in a file called "simRes". Specifically what computed and saved are:

1. **ID** a vector containing the names of metabolites for which PVA of exchange reactions was computed
2. **fact** a cell array containing min flux trough uptake and max trough secretion exchanges (later used for computing NMPCs)
3. **fact** a cell array containing min flux trough uptake and min trough secretion exchanges
4. **goal** an array containing the value of objectives for each microbiota model with rich and selected diet
5. **infeasible** cell array containing the names of the microbiota models that reported an infeasible status when solved for their objective

```
[ID, fact, fact, pRes1, infeasible]=microbiotaModelsLauncher(repaths, setup, sampleName, dietPath, rDiet, ID, extraInfo, pathWeb, factType)
```

Finally, NMPCs (net maximal production capability) are computed in a metabolite resolved manner and saved in a comma delimited file in the results folder. NMPCs indicate the maximal production of each metabolite and are computed as the absolute value of the sum of the maximal secretion flux with the maximal uptake flux. The similarity of metabolic profiles (using the different NMPCs as features) between individuals is also evaluated with classical multidimensional scaling.

```
[Fap,T]= getMetabolicFact(repaths, ID, sampleName, rDiet, ID, pathWeb, loadToPath, fact, f, figure)
```

Additionally, it is possible to retrieve and export, comprehensively, all the results (fluxes) computed during the simulations for a specified diet. Since PVA is computed on diet and fecal exchanges, every metabolite will have four different values for each individual, values corresponding min and max of uptake and secretion.

```
[fIndex] = extractFullFact(repaths, ID, "id diet", sampleName, fact, fact)
```