

# *Genomalicious* tutorial 1: Basic ingredients

*J.A. Thia*

*29 July 2019*

## Preamble

Every good recipe is built around simple ingredients. Likewise, bioinformatic pipelines, while complex, require some fundamental methods and data types that are the bare necessities to any study implementing population genomic approaches.

In this tutorial, you will:

1. Become familiar with some basic features of *genomalicious*.
2. Develop familiarity with SNP data structures.
3. Learn to use *genomalicious* functions to import and manipulate SNP data structures.

## Getting to know *genomalicious*

Let's start by loading the *genomalicious* library into your R session.

```
library(genomalicious)
```

*Genomalicious* contains a number of demonstrative toy datasets that you can experiment on. These need to be loaded in with the `data()` function and follow the naming convention, `genomalicious_[data name]`, where `[data name]` is a unique identifier.

For example, let's take a look at a toy dataset of allele frequencies:

```
data("genomalicious_Freqs")
```

```
# Take a look at the structure of this data, the dim() function reports the dimensions.
dim(genomalicious_Freqs)
```

```
## [1] 4 8
```

```
# Print the data to screen.
genomalicious_Freqs
```

```
##      Chrom_1_8 Chrom_2_16 Chrom_3_24 Chrom_4_32 Chrom_5_40 Chrom_6_48
## Pop1 0.5882639 0.6985712 0.9281353 0.7154336 0.2783644 0.5004546
## Pop2 0.5143548 0.6962864 0.9250919 0.3973361 0.3866906 0.4006825
## Pop3 0.7224595 0.6567121 0.8721565 0.2434911 0.6373447 0.9227758
## Pop4 0.6771910 0.5523064 0.9278697 0.4081590 0.6456204 0.8862767
##      Chrom_7_56 Chrom_8_64
## Pop1 0.4777536 0.06605165
## Pop2 0.3632988 0.06384908
## Pop3 0.2049917 0.19425487
## Pop4 0.2442950 0.15228191
```

You will notice that `genomalicious_Freqs` is a matrix with 4 rows and 8 columns. The rows contain the populations, whereas the columns contain the loci, and the cells contain allele frequencies. This data is in **wide format**, but we will dig into what this means a little later.

You can learn more about a dataset by using `?`  and the data object name, e.g. `?genomalicious_Freqs`. We will look at other datasets in later tutorials, but you can peruse the various datasets by typing `genomalicious_` and hitting the TAB key to get a list of options.

## Loading VCFs into R

Variant call files (VCFs) are one of the typical end products from read assembly/variant calling pipelines. These files contain information about genotypes attributed to each sample and associated mapping and genotype calling statistics.

*Genomalicious* offers a very simple way to import VCFs into R as **long format** data tables, whereby loci and populations are both in rows. C.f. with the wide format data described above.

We will now import a demo VCF into R using the `vcf2DT()` function. First, we need to find where *genomalicious* is installed and make a path to the demo file.

```
# Create a link to raw external datasets in genomalicious
genomaliciousExtData <- paste0(find.package('genomalicious'), '/extdata')

# This command here shows you the VCF file that comes with genomalicious
list.files(genomaliciousExtData, pattern='_poolseq.vcf')
```

```
## [1] "genomalicious_poolseq.vcf"
```

```
# Use this to create a path to that file
vcfPath <- paste0(genomaliciousExtData, '/genomalicious_poolseq.vcf')

# The value of vcfPath with depend on your system
vcfPath
```

```
## [1] "C:/Users/Josh/Documents/R/win-library/3.4/genomalicious/extdata/genomalicious_poolseq.vcf"
```

You can navigate yourself to the path stored in `vcfPath` and open the VCF using a text editor. However, we can simply read the lines of the file and print them to screen:

```
# Read in and print the first 10 lines of the demo VCF
head(readLines(vcfPath), 10)
```

```
## [1] "##This is a toy dataset for the R package genomalicious - it emulates a VCF file"
## [2] "##INFO=<ID=DP,Number=1,Type=Integer,Description='The total depth across samples'>"
## [3] "##FORMAT=<ID=DP,Number=1,Type=Integer,Description='The total depth in a sample'>"
## [4] "##FORMAT=<ID=RO,Number=1,Type=Integer,Description='The reference allele counts in a sample'>"
## [5] "##FORMAT=<ID=AO,Number=1,Type=Integer,Description='The alternate allele counts in a sample'>"
## [6] "#CHROM\tPOS\tID\tREF\tALT\tQUAL\tFILTER\tINFO\tFORMAT\tPop1_Rep1\tPop1_Rep2\tPop1_Rep3\tPop2_Rep1\tPop2_Rep2\tPop2_Rep3\tPop2_Rep4\tPop2_Rep5\tPop2_Rep6\tPop2_Rep7\tPop2_Rep8\tPop2_Rep9\tPop2_Rep10\tPop2_Rep11\tPop2_Rep12\tPop2_Rep13\tPop2_Rep14\tPop2_Rep15\tPop2_Rep16\tPop2_Rep17\tPop2_Rep18\tPop2_Rep19\tPop2_Rep20\tPop2_Rep21\tPop2_Rep22\tPop2_Rep23\tPop2_Rep24\tPop2_Rep25\tPop2_Rep26\tPop2_Rep27\tPop2_Rep28\tPop2_Rep29\tPop2_Rep30\tPop2_Rep31\tPop2_Rep32\tPop2_Rep33\tPop2_Rep34\tPop2_Rep35\tPop2_Rep36\tPop2_Rep37\tPop2_Rep38\tPop2_Rep39\tPop2_Rep40\tPop2_Rep41\tPop2_Rep42\tPop2_Rep43\tPop2_Rep44\tPop2_Rep45\tPop2_Rep46\tPop2_Rep47\tPop2_Rep48\tPop2_Rep49\tPop2_Rep50\tPop2_Rep51\tPop2_Rep52\tPop2_Rep53\tPop2_Rep54\tPop2_Rep55\tPop2_Rep56\tPop2_Rep57\tPop2_Rep58\tPop2_Rep59\tPop2_Rep60\tPop2_Rep61\tPop2_Rep62\tPop2_Rep63\tPop2_Rep64\tPop2_Rep65\tPop2_Rep66\tPop2_Rep67\tPop2_Rep68\tPop2_Rep69\tPop2_Rep70\tPop2_Rep71\tPop2_Rep72\tPop2_Rep73\tPop2_Rep74\tPop2_Rep75\tPop2_Rep76\tPop2_Rep77\tPop2_Rep78\tPop2_Rep79\tPop2_Rep80\tPop2_Rep81\tPop2_Rep82\tPop2_Rep83\tPop2_Rep84\tPop2_Rep85\tPop2_Rep86\tPop2_Rep87\tPop2_Rep88\tPop2_Rep89\tPop2_Rep90\tPop2_Rep91\tPop2_Rep92\tPop2_Rep93\tPop2_Rep94\tPop2_Rep95\tPop2_Rep96\tPop2_Rep97\tPop2_Rep98\tPop2_Rep99\tPop2_Rep100\tPop2_Rep101\tPop2_Rep102\tPop2_Rep103\tPop2_Rep104\tPop2_Rep105\tPop2_Rep106\tPop2_Rep107\tPop2_Rep108\tPop2_Rep109\tPop2_Rep110\tPop2_Rep111\tPop2_Rep112\tPop2_Rep113\tPop2_Rep114\tPop2_Rep115\tPop2_Rep116\tPop2_Rep117\tPop2_Rep118\tPop2_Rep119\tPop2_Rep120\tPop2_Rep121\tPop2_Rep122\tPop2_Rep123\tPop2_Rep124\tPop2_Rep125\tPop2_Rep126\tPop2_Rep127\tPop2_Rep128\tPop2_Rep129\tPop2_Rep130\tPop2_Rep131\tPop2_Rep132\tPop2_Rep133\tPop2_Rep134\tPop2_Rep135\tPop2_Rep136\tPop2_Rep137\tPop2_Rep138\tPop2_Rep139\tPop2_Rep140\tPop2_Rep141\tPop2_Rep142\tPop2_Rep143\tPop2_Rep144\tPop2_Rep145\tPop2_Rep146\tPop2_Rep147\tPop2_Rep148\tPop2_Rep149\tPop2_Rep150\tPop2_Rep151\tPop2_Rep152\tPop2_Rep153\tPop2_Rep154\tPop2_Rep155\tPop2_Rep156\tPop2_Rep157\tPop2_Rep158\tPop2_Rep159\tPop2_Rep160\tPop2_Rep161\tPop2_Rep162\tPop2_Rep163\tPop2_Rep164\tPop2_Rep165\tPop2_Rep166\tPop2_Rep167\tPop2_Rep168\tPop2_Rep169\tPop2_Rep170\tPop2_Rep171\tPop2_Rep172\tPop2_Rep173\tPop2_Rep174\tPop2_Rep175\tPop2_Rep176\tPop2_Rep177\tPop2_Rep178\tPop2_Rep179\tPop2_Rep180\tPop2_Rep181\tPop2_Rep182\tPop2_Rep183\tPop2_Rep184\tPop2_Rep185\tPop2_Rep186\tPop2_Rep187\tPop2_Rep188\tPop2_Rep189\tPop2_Rep190\tPop2_Rep191\tPop2_Rep192\tPop2_Rep193\tPop2_Rep194\tPop2_Rep195\tPop2_Rep196\tPop2_Rep197\tPop2_Rep198\tPop2_Rep199\tPop2_Rep200\tPop2_Rep201\tPop2_Rep202\tPop2_Rep203\tPop2_Rep204\tPop2_Rep205\tPop2_Rep206\tPop2_Rep207\tPop2_Rep208\tPop2_Rep209\tPop2_Rep210\tPop2_Rep211\tPop2_Rep212\tPop2_Rep213\tPop2_Rep214\tPop2_Rep215\tPop2_Rep216\tPop2_Rep217\tPop2_Rep218\tPop2_Rep219\tPop2_Rep220\tPop2_Rep221\tPop2_Rep222\tPop2_Rep223\tPop2_Rep224\tPop2_Rep225\tPop2_Rep226\tPop2_Rep227\tPop2_Rep228\tPop2_Rep229\tPop2_Rep230\tPop2_Rep231\tPop2_Rep232\tPop2_Rep233\tPop2_Rep234\tPop2_Rep235\tPop2_Rep236\tPop2_Rep237\tPop2_Rep238\tPop2_Rep239\tPop2_Rep240\tPop2_Rep241\tPop2_Rep242\tPop2_Rep243\tPop2_Rep244\tPop2_Rep245\tPop2_Rep246\tPop2_Rep247\tPop2_Rep248\tPop2_Rep249\tPop2_Rep250\tPop2_Rep251\tPop2_Rep252\tPop2_Rep253\tPop2_Rep254\tPop2_Rep255\tPop2_Rep256\tPop2_Rep257\tPop2_Rep258\tPop2_Rep259\tPop2_Rep260\tPop2_Rep261\tPop2_Rep262\tPop2_Rep263\tPop2_Rep264\tPop2_Rep265\tPop2_Rep266\tPop2_Rep267\tPop2_Rep268\tPop2_Rep269\tPop2_Rep270\tPop2_Rep271\tPop2_Rep272\tPop2_Rep273\tPop2_Rep274\tPop2_Rep275\tPop2_Rep276\tPop2_Rep277\tPop2_Rep278\tPop2_Rep279\tPop2_Rep280\tPop2_Rep281\tPop2_Rep282\tPop2_Rep283\tPop2_Rep284\tPop2_Rep285\tPop2_Rep286\tPop2_Rep287\tPop2_Rep288\tPop2_Rep289\tPop2_Rep290\tPop2_Rep291\tPop2_Rep292\tPop2_Rep293\tPop2_Rep294\tPop2_Rep295\tPop2_Rep296\tPop2_Rep297\tPop2_Rep298\tPop2_Rep299\tPop2_Rep300\tPop2_Rep301\tPop2_Rep302\tPop2_Rep303\tPop2_Rep304\tPop2_Rep305\tPop2_Rep306\tPop2_Rep307\tPop2_Rep308\tPop2_Rep309\tPop2_Rep310\tPop2_Rep311\tPop2_Rep312\tPop2_Rep313\tPop2_Rep314\tPop2_Rep315\tPop2_Rep316\tPop2_Rep317\tPop2_Rep318\tPop2_Rep319\tPop2_Rep320\tPop2_Rep321\tPop2_Rep322\tPop2_Rep323\tPop2_Rep324\tPop2_Rep325\tPop2_Rep326\tPop2_Rep327\tPop2_Rep328\tPop2_Rep329\tPop2_Rep330\tPop2_Rep331\tPop2_Rep332\tPop2_Rep333\tPop2_Rep334\tPop2_Rep335\tPop2_Rep336\tPop2_Rep337\tPop2_Rep338\tPop2_Rep339\tPop2_Rep340\tPop2_Rep341\tPop2_Rep342\tPop2_Rep343\tPop2_Rep344\tPop2_Rep345\tPop2_Rep346\tPop2_Rep347\tPop2_Rep348\tPop2_Rep349\tPop2_Rep350\tPop2_Rep351\tPop2_Rep352\tPop2_Rep353\tPop2_Rep354\tPop2_Rep355\tPop2_Rep356\tPop2_Rep357\tPop2_Rep358\tPop2_Rep359\tPop2_Rep360\tPop2_Rep361\tPop2_Rep362\tPop2_Rep363\tPop2_Rep364\tPop2_Rep365\tPop2_Rep366\tPop2_Rep367\tPop2_Rep368\tPop2_Rep369\tPop2_Rep370\tPop2_Rep371\tPop2_Rep372\tPop2_Rep373\tPop2_Rep374\tPop2_Rep375\tPop2_Rep376\tPop2_Rep377\tPop2_Rep378\tPop2_Rep379\tPop2_Rep380\tPop2_Rep381\tPop2_Rep382\tPop2_Rep383\tPop2_Rep384\tPop2_Rep385\tPop2_Rep386\tPop2_Rep387\tPop2_Rep388\tPop2_Rep389\tPop2_Rep390\tPop2_Rep391\tPop2_Rep392\tPop2_Rep393\tPop2_Rep394\tPop2_Rep395\tPop2_Rep396\tPop2_Rep397\tPop2_Rep398\tPop2_Rep399\tPop2_Rep400\tPop2_Rep401\tPop2_Rep402\tPop2_Rep403\tPop2_Rep404\tPop2_Rep405\tPop2_Rep406\tPop2_Rep407\tPop2_Rep408\tPop2_Rep409\tPop2_Rep410\tPop2_Rep411\tPop2_Rep412\tPop2_Rep413\tPop2_Rep414\tPop2_Rep415\tPop2_Rep416\tPop2_Rep417\tPop2_Rep418\tPop2_Rep419\tPop2_Rep420\tPop2_Rep421\tPop2_Rep422\tPop2_Rep423\tPop2_Rep424\tPop2_Rep425\tPop2_Rep426\tPop2_Rep427\tPop2_Rep428\tPop2_Rep429\tPop2_Rep430\tPop2_Rep431\tPop2_Rep432\tPop2_Rep433\tPop2_Rep434\tPop2_Rep435\tPop2_Rep436\tPop2_Rep437\tPop2_Rep438\tPop2_Rep439\tPop2_Rep440\tPop2_Rep441\tPop2_Rep442\tPop2_Rep443\tPop2_Rep444\tPop2_Rep445\tPop2_Rep446\tPop2_Rep447\tPop2_Rep448\tPop2_Rep449\tPop2_Rep450\tPop2_Rep451\tPop2_Rep452\tPop2_Rep453\tPop2_Rep454\tPop2_Rep455\tPop2_Rep456\tPop2_Rep457\tPop2_Rep458\tPop2_Rep459\tPop2_Rep460\tPop2_Rep461\tPop2_Rep462\tPop2_Rep463\tPop2_Rep464\tPop2_Rep465\tPop2_Rep466\tPop2_Rep467\tPop2_Rep468\tPop2_Rep469\tPop2_Rep470\tPop2_Rep471\tPop2_Rep472\tPop2_Rep473\tPop2_Rep474\tPop2_Rep475\tPop2_Rep476\tPop2_Rep477\tPop2_Rep478\tPop2_Rep479\tPop2_Rep480\tPop2_Rep481\tPop2_Rep482\tPop2_Rep483\tPop2_Rep484\tPop2_Rep485\tPop2_Rep486\tPop2_Rep487\tPop2_Rep488\tPop2_Rep489\tPop2_Rep490\tPop2_Rep491\tPop2_Rep492\tPop2_Rep493\tPop2_Rep494\tPop2_Rep495\tPop2_Rep496\tPop2_Rep497\tPop2_Rep498\tPop2_Rep499\tPop2_Rep500\tPop2_Rep501\tPop2_Rep502\tPop2_Rep503\tPop2_Rep504\tPop2_Rep505\tPop2_Rep506\tPop2_Rep507\tPop2_Rep508\tPop2_Rep509\tPop2_Rep510\tPop2_Rep511\tPop2_Rep512\tPop2_Rep513\tPop2_Rep514\tPop2_Rep515\tPop2_Rep516\tPop2_Rep517\tPop2_Rep518\tPop2_Rep519\tPop2_Rep520\tPop2_Rep521\tPop2_Rep522\tPop2_Rep523\tPop2_Rep524\tPop2_Rep525\tPop2_Rep526\tPop2_Rep527\tPop2_Rep528\tPop2_Rep529\tPop2_Rep530\tPop2_Rep531\tPop2_Rep532\tPop2_Rep533\tPop2_Rep534\tPop2_Rep535\tPop2_Rep536\tPop2_Rep537\tPop2_Rep538\tPop2_Rep539\tPop2_Rep540\tPop2_Rep541\tPop2_Rep542\tPop2_Rep543\tPop2_Rep544\tPop2_Rep545\tPop2_Rep546\tPop2_Rep547\tPop2_Rep548\tPop2_Rep549\tPop2_Rep550\tPop2_Rep551\tPop2_Rep552\tPop2_Rep553\tPop2_Rep554\tPop2_Rep555\tPop2_Rep556\tPop2_Rep557\tPop2_Rep558\tPop2_Rep559\tPop2_Rep560\tPop2_Rep561\tPop2_Rep562\tPop2_Rep563\tPop2_Rep564\tPop2_Rep565\tPop2_Rep566\tPop2_Rep567\tPop2_Rep568\tPop2_Rep569\tPop2_Rep570\tPop2_Rep571\tPop2_Rep572\tPop2_Rep573\tPop2_Rep574\tPop2_Rep575\tPop2_Rep576\tPop2_Rep577\tPop2_Rep578\tPop2_Rep579\tPop2_Rep580\tPop2_Rep581\tPop2_Rep582\tPop2_Rep583\tPop2_Rep584\tPop2_Rep585\tPop2_Rep586\tPop2_Rep587\tPop2_Rep588\tPop2_Rep589\tPop2_Rep590\tPop2_Rep591\tPop2_Rep592\tPop2_Rep593\tPop2_Rep594\tPop2_Rep595\tPop2_Rep596\tPop2_Rep597\tPop2_Rep598\tPop2_Rep599\tPop2_Rep600\tPop2_Rep601\tPop2_Rep602\tPop2_Rep603\tPop2_Rep604\tPop2_Rep605\tPop2_Rep606\tPop2_Rep607\tPop2_Rep608\tPop2_Rep609\tPop2_Rep610\tPop2_Rep611\tPop2_Rep612\tPop2_Rep613\tPop2_Rep614\tPop2_Rep615\tPop2_Rep616\tPop2_Rep617\tPop2_Rep618\tPop2_Rep619\tPop2_Rep620\tPop2_Rep621\tPop2_Rep622\tPop2_Rep623\tPop2_Rep624\tPop2_Rep625\tPop2_Rep626\tPop2_Rep627\tPop2_Rep628\tPop2_Rep629\tPop2_Rep630\tPop2_Rep631\tPop2_Rep632\tPop2_Rep633\tPop2_Rep634\tPop2_Rep635\tPop2_Rep636\tPop2_Rep637\tPop2_Rep638\tPop2_Rep639\tPop2_Rep640\tPop2_Rep641\tPop2_Rep642\tPop2_Rep643\tPop2_Rep644\tPop2_Rep645\tPop2_Rep646\tPop2_Rep647\tPop2_Rep648\tPop2_Rep649\tPop2_Rep650\tPop2_Rep651\tPop2_Rep652\tPop2_Rep653\tPop2_Rep654\tPop2_Rep655\tPop2_Rep656\tPop2_Rep657\tPop2_Rep658\tPop2_Rep659\tPop2_Rep660\tPop2_Rep661\tPop2_Rep662\tPop2_Rep663\tPop2_Rep664\tPop2_Rep665\tPop2_Rep666\tPop2_Rep667\tPop2_Rep668\tPop2_Rep669\tPop2_Rep670\tPop2_Rep671\tPop2_Rep672\tPop2_Rep673\tPop2_Rep674\tPop2_Rep675\tPop2_Rep676\tPop2_Rep677\tPop2_Rep678\tPop2_Rep679\tPop2_Rep680\tPop2_Rep681\tPop2_Rep682\tPop2_Rep683\tPop2_Rep684\tPop2_Rep685\tPop2_Rep686\tPop2_Rep687\tPop2_Rep688\tPop2_Rep689\tPop2_Rep690\tPop2_Rep691\tPop2_Rep692\tPop2_Rep693\tPop2_Rep694\tPop2_Rep695\tPop2_Rep696\tPop2_Rep697\tPop2_Rep698\tPop2_Rep699\tPop2_Rep700\tPop2_Rep701\tPop2_Rep702\tPop2_Rep703\tPop2_Rep704\tPop2_Rep705\tPop2_Rep706\tPop2_Rep707\tPop2_Rep708\tPop2_Rep709\tPop2_Rep710\tPop2_Rep711\tPop2_Rep712\tPop2_Rep713\tPop2_Rep714\tPop2_Rep715\tPop2_Rep716\tPop2_Rep717\tPop2_Rep718\tPop2_Rep719\tPop2_Rep720\tPop2_Rep721\tPop2_Rep722\tPop2_Rep723\tPop2_Rep724\tPop2_Rep725\tPop2_Rep726\tPop2_Rep727\tPop2_Rep728\tPop2_Rep729\tPop2_Rep730\tPop2_Rep731\tPop2_Rep732\tPop2_Rep733\tPop2_Rep734\tPop2_Rep735\tPop2_Rep736\tPop2_Rep737\tPop2_Rep738\tPop2_Rep739\tPop2_Rep740\tPop2_Rep741\tPop2_Rep742\tPop2_Rep743\tPop2_Rep744\tPop2_Rep745\tPop2_Rep746\tPop2_Rep747\tPop2_Rep748\tPop2_Rep749\tPop2_Rep750\tPop2_Rep751\tPop2_Rep752\tPop2_Rep753\tPop2_Rep754\tPop2_Rep755\tPop2_Rep756\tPop2_Rep757\tPop2_Rep758\tPop2_Rep759\tPop2_Rep760\tPop2_Rep761\tPop2_Rep762\tPop2_Rep763\tPop2_Rep764\tPop2_Rep765\tPop2_Rep766\tPop2_Rep767\tPop2_Rep768\tPop2_Rep769\tPop2_Rep770\tPop2_Rep771\tPop2_Rep772\tPop2_Rep773\tPop2_Rep774\tPop2_Rep775\tPop2_Rep776\tPop2_Rep777\tPop2_Rep778\tPop2_Rep779\tPop2_Rep780\tPop2_Rep781\tPop2_Rep782\tPop2_Rep783\tPop2_Rep784\tPop2_Rep785\tPop2_Rep786\tPop2_Rep787\tPop2_Rep788\tPop2_Rep789\tPop2_Rep790\tPop2_Rep791\tPop2_Rep792\tPop2_Rep793\tPop2_Rep794\tPop2_Rep795\tPop2_Rep796\tPop2_Rep797\tPop2_Rep798\tPop2_Rep799\tPop2_Rep800\tPop2_Rep801\tPop2_Rep802\tPop2_Rep803\tPop2_Rep804\tPop2_Rep805\tPop2_Rep806\tPop2_Rep807\tPop2_Rep808\tPop2_Rep809\tPop2_Rep810\tPop2_Rep811\tPop2_Rep812\tPop2_Rep813\tPop2_Rep814\tPop2_Rep815\tPop2_Rep816\tPop2_Rep817\tPop2_Rep818\tPop2_Rep819\tPop2_Rep820\tPop2_Rep821\tPop2_Rep822\tPop2_Rep823\tPop2_Rep824\tPop2_Rep825\tPop2_Rep826\tPop2_Rep827\tPop2_Rep828\tPop2_Rep829\tPop2_Rep830\tPop2_Rep831\tPop2_Rep832\tPop2_Rep833\tPop2_Rep834\tPop2_Rep835\tPop2_Rep836\tPop2_Rep837\tPop2_Rep838\tPop2_Rep839\tPop2_Rep840\tPop2_Rep841\tPop2_Rep842\tPop2_Rep843\tPop2_Rep844\tPop2_Rep845\tPop2_Rep846\tPop2_Rep847\tPop2_Rep848\tPop2_Rep849\tPop2_Rep850\tPop2_Rep851\tPop2_Rep852\tPop2_Rep853\tPop2_Rep854\tPop2_Rep855\tPop2_Rep856\tPop2_Rep857\tPop2_Rep858\tPop2_Rep859\tPop2_Rep860\tPop2_Rep861\tPop2_Rep862\tPop2_Rep863\tPop2_Rep864\tPop2_Rep865\tPop2_Rep866\tPop2_Rep867\tPop2_Rep868\tPop2_Rep869\tPop2_Rep870\tPop2_Rep871\tPop2_Rep872\tPop2_Rep873\tPop2_Rep874\tPop2_Rep875\tPop2_Rep876\tPop2_Rep877\tPop2_Rep878\tPop2_Rep879\tPop2_Rep880\tPop2_Rep881\tPop2_Rep882\tPop2_Rep883\tPop2_Rep884\tPop2_Rep885\tPop2_Rep886\tPop2_Rep887\tPop2_Rep888\tPop2_Rep889\tPop2_Rep890\tPop2_Rep891\tPop2_Rep892\tPop2_Rep893\tPop2_Rep894\tPop2_Rep895\tPop2_Rep896\tPop2_Rep897\tPop2_Rep898\tPop2_Rep899\tPop2_Rep900\tPop2_Rep901\tPop2_Rep902\tPop2_Rep903\tPop2_Rep904\tPop2_Rep905\tPop2_Rep906\tPop2_Rep907\tPop2_Rep908\tPop2_Rep909\tPop2_Rep910\tPop2_Rep911\tPop2_Rep912\tPop2_Rep913\tPop2_Rep914\tPop2_Rep915\tPop2_Rep916\tPop2_Rep917\tPop2_Rep918\tPop2_Rep919\tPop2_Rep920\tPop2_Rep921\tPop2_Rep922\tPop2_Rep923\tPop2_Rep924\tPop2_Rep925\tPop2_Rep926\tPop2_Rep927\tPop2_Rep928\tPop2_Rep929\tPop2_Rep930\tPop2_Rep931\tPop2_Rep932\tPop2_Rep933\tPop2_Rep934\tPop2_Rep935\tPop2_Rep936\tPop2_Rep937\tPop2_Rep938\tPop2_Rep939\tPop2_Rep940\tPop2_Rep941\tPop2_Rep942\tPop2_Rep943\tPop2_Rep944\tPop2_Rep945\tPop2_Rep946\tPop2_Rep947\tPop2_Rep948\tPop2_Rep949\tPop2_Rep950\tPop2_Rep951\tPop2_Rep952\tPop2_Rep953\tPop2_Rep954\tPop2_Rep955\tPop2_Rep956\tPop2_Rep957\tPop2_Rep958\tPop2_Rep959\tPop2_Rep960\tPop2_Rep961\tPop2_Rep962\tPop2_Rep963\tPop2_Rep964\tPop2_Rep965\tPop2_Rep966\tPop2_Rep967\tPop2_Rep968\tPop2_Rep969\tPop2_Rep970\tPop2_Rep971\tPop2_Rep972\tPop2_Rep973\tPop2_Rep974\tPop2_Rep975\tPop2_Rep976\tPop2_Rep977\tPop2_Rep978\tPop2_Rep979\tPop2_Rep980\tPop2_Rep981\tPop2_Rep982\tPop2_Rep983\tPop2_Rep984\tPop2_Rep985\tPop2_Rep986\tPop2_Rep987\tPop2_Rep988\tPop2_Rep989\tPop2_Rep990\tPop2_Rep991\tPop2_Rep992\tPop2_Rep993\tPop2_Rep994\tPop2_Rep995\tPop2_Rep996\tPop2_Rep997\tPop2_Rep998\tPop2_Rep999\tPop2_Rep1000\tPop2_Rep1001\tPop2_Rep1002\tPop2_Rep1003\tPop2_Rep1004\tPop2_Rep1005\tPop2_Rep1006\tPop2_Rep1007\tPop2_Rep1008\tPop2_Rep1009\tPop2_Rep1010\tPop2_Rep1011\tPop2_Rep1012\tPop2_Rep1013\tPop2_Rep1014\tPop2_Rep1015\tPop2_Rep1016\tPop2_Rep1017\tPop2_Rep1018\tPop2_Rep1019\tPop2_Rep1020\tPop2_Rep1021\tPop2_Rep1022\tPop2_Rep1023\tPop2_Rep1024\tPop2_Rep1025\tPop2_Rep1026\tPop2_Rep1027\tPop2_Rep1028\tPop2_Rep1029\tPop2_Rep1030\tPop2_Rep1031\tPop2_Rep1032\tPop2_Rep1033\tPop2_Rep1034\tPop2_Rep1035\tPop2_Rep1036\tPop2_Rep1037\tPop2_Rep1038\tPop2_Rep1039\tPop2_Rep1040\tPop2_Rep1041\tPop2_Rep1042\tPop2_Rep1043\tPop2_Rep1044\tPop2_Rep1045\tPop2_Rep1046\tPop2_Rep1047\tPop2_Rep1048\tPop2_Rep1049\tPop2_Rep1050\tPop2_Rep1051\tPop2_Rep1052\tPop2_Rep1053\tPop2_Rep1054\tPop2_Rep1055\tPop2_Rep1056\tPop2_Rep1057\tPop2_Rep1058\tPop2_Rep1059\tPop2_Rep1060\tPop2_Rep1061\tPop2_Rep1062\tPop2_Rep1063\tPop2_Rep1064\tPop2_Rep1065\tPop2_Rep1066\tPop2_Rep1067\tPop2_Rep1068\tPop2_Rep1069\tPop2_Rep1070\tPop2_Rep1071\tPop2_Rep1072\tPop2_Rep1073\tPop2_Rep1074\tPop2_Rep1075\tPop2_Rep1076\tPop2_Rep1077\tPop2_Rep1078\tPop2_Rep1079\tPop2_Rep1080\tPop2_Rep1081\tPop2_Rep1082\tPop2_Rep1083\tPop2_Rep1084\tPop2_Rep1085\tPop2_Rep1086\tPop2_Rep1087\tPop2_Rep1088\tPop2_Rep1089\tPop2_Rep1090\tPop2_Rep1091\tPop2_Rep1092\tPop2_Rep1093\tPop2_Rep1094\tPop2_Rep1095\tPop2_Rep1096\tPop2_Rep1097\tPop2_Rep1098\tPop2_Rep1099\tPop2_Rep1100\tPop2_Rep1101\tPop2_Rep1102\tPop2_Rep1103\tPop2_Rep1104\tPop2_Rep1105\tPop2_Rep1106\tPop2_Rep1107\tPop2_Rep1108\tPop2_Rep1109\tPop2_Rep1110\tPop2_Rep1111\tPop2_Rep1112\tPop2_Rep1113\tPop2_Rep1114\tPop2_Rep1115\tPop2_Rep1116\tPop2_Rep1117\tPop2_Rep1118\tPop2_Rep1119\tPop2_Rep1120\tPop2_Rep1121\tPop2_Rep1122\tPop2_Rep1123\tPop2_Rep1124\tPop2_Rep1125\tPop2_Rep1126\tPop2_Rep1127\tPop2_Rep1128\tPop2_Rep1129\tPop2_Rep1130\tPop2_Rep1131\tPop2_Rep1132\tPop2_Rep1133\tPop2_Rep1134\tPop2_Rep1135\tPop2_Rep1136\tPop2_Rep1137\tPop2_Rep1138\tPop2_Rep1139\tPop2_Rep1140\tPop2_Rep1141\tPop2_Rep1142\tPop2_Rep1143\tPop2_Rep1144\tPop2_Rep1145\tPop2_Rep1146\tPop2_Rep1147\tPop2_Rep1148\tPop2_Rep1149\tPop2_Rep1150\tPop2_Rep1151\tPop2_Rep1152\tPop2_Rep1153\tPop2_Rep1154\tPop2_Rep1155\tPop2_Rep1156\tPop2_Rep1157\tPop2_Rep1158\tPop2_Rep1159\tPop2_Rep1160\tPop2_Rep1161\tPop2_Rep1162\tPop2_Rep1163\tPop2_Rep1164\tPop2_Rep1165\tPop2_Rep1166\tPop2_Rep1167\tPop2_Rep1168\tPop2_Rep1169\tPop2_Rep1170\tPop2_Rep1171\tPop2_Rep1172\tPop2_Rep1173\tPop2_Rep1174\tPop2_Rep1175\tPop2_Rep1176\tPop2_Rep1177\tPop2_Rep1178\tPop2_Rep1179\tPop2_Rep1180\tPop2_Rep1181\tPop2_Rep1182\tPop2_Rep1183\tPop2_Rep1184\tPop2_Rep1185\tPop2_Rep1186\tPop2_Rep1187\tPop2_Rep1188\tPop2_Rep1189\tPop2_Rep1190\tPop2_Rep1191\tPop2_Rep1192\tPop2_Rep1193\tPop2_Rep1194\tPop2_Rep1195\tPop2_Rep1196\tPop2_Rep1197\tPop2_Rep1198\tPop2_Rep1199\tPop2_Rep1200\tPop2_Rep1201\tPop2_Rep1202\tPop2_Rep1203\tPop2_Rep1204\tPop2_Rep1205\tPop2_Rep1206\tPop2_Rep1207\tPop2_Rep1208\tPop2_Rep1209\tPop2_Rep1210\tPop2_Rep1211\tPop2_Rep1212\tPop2_Rep1213\tPop2_Rep1214\tPop2_Rep1215\tPop2_Rep1216\tPop2_Rep1217\tPop2_Rep1218\tPop2_Rep1219\tPop2_Rep1220\tPop2_Rep1221\tPop2_Rep1222\tPop2_Rep1223\tPop2_Rep1224\tPop2_Rep1225\tPop2_Rep1226\tPop2_Rep1227\tPop2_Rep1228\tPop2_Rep1229\tPop2_Rep1230\tPop2_Rep1231\tPop2_Rep1232\tPop2_Rep1233\tPop2_Rep1234\tPop2_Rep1235\tPop2_Rep1236\tPop2_Rep1237\tPop2_Rep1238\tPop2_Rep1239\tPop2_Rep1240\tPop2_Rep1241\tPop2_Rep1242\tPop2_Rep1243\tPop2_Rep1244\tPop2_Rep1245\tPop2_Rep1246\tPop2_Rep1247\tPop2_Rep1248\tPop2_Rep1249\tPop2_Rep1250\tPop2_Rep1251\tPop2_Rep1252\tPop2_Rep1253\tPop2_Rep1254\tPop2_Rep1255\tPop2_Rep1256\tPop2_Rep1257\tPop2_Rep1258\tPop2_Rep1259\tPop2_Rep1260\tPop2_Rep1261\tPop2_Rep1262\tPop2_Rep1263\tPop2_Rep1264\tPop2_Rep1265\tPop2_Rep1266\tPop2_Rep1267\tPop2_Rep1268\tPop2_Rep1269\tPop2_Rep1270\tPop2_Rep1271\tPop2_Rep1272\tPop2_Rep1273\tPop2_Rep1274\tPop2_Rep1275\tPop2_Rep1276\tPop2_Rep1277\tPop2_Rep1278\tPop2_Rep1279\tPop2_Rep1280\tPop2_Rep1281\tPop2_Rep1282\tPop2_Rep1283\tPop2_Rep1284\tPop2_Rep1285\tPop2_Rep1286\tPop2_Rep1287\tPop2_Rep1288\tPop2_Rep1289\tPop2_Rep1290\tPop2_Rep1291\tPop2_Rep1292\tPop2_Rep1293\tPop2_Rep1294\tPop2_Rep1295\tPop2_Rep1296\tPop2_Rep1297\tPop2_Rep1298\tPop2_Rep1299\tPop2_Rep1300\tPop2_Rep1301\tPop2_Rep1302\tPop2_Rep1303\tPop2_Rep1304\tPop2_Rep1305\tPop2_Rep1306\tPop2_Rep1307\tPop2_Rep1308\tPop2_Rep1309\tPop2_Rep1310\tPop2_Rep1311\tPop2_Rep1312\tPop2_Rep1313\tPop2_Rep1314\tPop2_Rep1315\tPop2_Rep1316\tPop2_Rep1317\tPop2_Rep1318\tPop2_Rep1319\tPop2_Rep1320\tPop2_Rep1321\tPop2_Rep1322\tPop2_Rep1323\tPop2_Rep1324\tPop2_Rep1325\tPop2_Rep1326\tPop2_Rep1327\tPop2_Rep1328\tPop2_Rep1329\tPop2_Rep1330\tPop2_Rep1331\tPop2_Rep1332\tPop2_Rep1333\tPop2_Rep1334\tPop2_Rep1335\tPop2_Rep1336\tPop2_Rep1337\tPop2_Rep1338\tPop2_Rep1339\tPop2_Rep1340\tPop2_Rep1341\tPop2_Rep1342\tPop2_Rep1343\tPop2_Rep1344\tPop2_Rep1345\tPop2_Rep1346\tPop2_Rep1347\tPop2_Rep1348\tPop2_Rep1349\tPop2_Rep1350\tPop2_Rep1351\tPop2_Rep1352\tPop2_Rep1353\tPop2_Rep1354\tPop2_Rep1355\tPop2_Rep1356\tPop2_Rep1357\tPop2_Rep1358\tPop2_Rep1359\tPop2_Rep1360\tPop2_Rep1361\tPop2_Rep1362\tPop2_Rep1363\tPop2_Rep1364\tPop2_Rep1365\tPop2_Rep1366\tPop2_Rep1367\tPop2_Rep1368\tPop2_Rep1369\tPop2
```

This demo VCF was constructed more simply for the purpose of this tutorial, but follows the same basic structure of that produced from variant calling software. You will see that the text is interspersed with `\t`: these indicate tabs that separate the various columns in the file.

A VCF has the following components:

1. A **comments section** marked with double hashtags, `##`. These typically contain details about the contents of the VCF or how the reads were called.
2. A **heading section** marked with a single hashtag, `#`. This is effectively the column names for the wide format SNP data.
3. All lines proceeding contain SNP data.

VCFs can be a bit tricky to interpret when you first see them. All columns **before** **FORMAT** contain information about the SNP (i.e. the chromosome, its position, the alleles). **FORMAT itself** contains a text string detailing the contents of the sample columns. Note, the contents of **FORMAT** are typically described in the comments section. We therefore know that `DP:RO:AO` represent the total read depth (DP), and the counts of the reference (RO) and alternate (AO) allele. Also note the difference in DP stored in the **INFO** column (read depth across all samples) and that in **FORMAT** (read depth within each sample).

All sample columns occur **after** **FORMAT** (e.g. `Pop1_Rep1`, `Pop1_Rep2`) and continue to the end of the line. Within each sample column, the values described in **FORMAT** are stored, with each value separated by a `:`. This demo VCF contains read data from a replicated pool-seq experiment. Individuals from four populations were pooled into a single library, and three replicate libraries were sequenced. The reads associated with each population replicate, for each SNP, are detailed in the *i*-by-*j* row-column coordinates.

For more in depth descriptions of VCFs, you should check out the resources provided by the developers of *samtools*, here.

Now that we understand the data we are working with, let's import it into R.

```
# Import VCF into R using the path name
poolSnps <- vcf2DT(vcfPath)

# First 8 rows of the imported SNP data
head(poolSnps, 8)
```

```
##      CHROM POS ID REF ALT QUAL FILTER  INFO      SAMPLE      LOCUS  DP RO
## 1: Chrom_1   8  .   A   T   50   PASS DP=854 Pop1_Rep1 Chrom_1_8 113 63
## 2: Chrom_2  16  .   C   G   40   PASS DP=785 Pop1_Rep1 Chrom_2_16  82 60
## 3: Chrom_3  24  .   G   A   45   PASS DP=864 Pop1_Rep1 Chrom_3_24 101 90
## 4: Chrom_4  32  .   T   C   43   PASS DP=792 Pop1_Rep1 Chrom_4_32 101 73
## 5: Chrom_5  40  .   A   C   30   PASS DP=825 Pop1_Rep1 Chrom_5_40 100 22
## 6: Chrom_6  48  .   T   G   37   PASS DP=854 Pop1_Rep1 Chrom_6_48 101 51
## 7: Chrom_7  56  .   C   G   38   PASS DP=820 Pop1_Rep1 Chrom_7_56 109 50
## 8: Chrom_8  64  .   T   A   52   PASS DP=838 Pop1_Rep1 Chrom_8_64 101  3
##      AO
## 1: 50
## 2: 22
## 3: 11
## 4: 28
## 5: 78
## 6: 50
## 7: 59
## 8: 98
```

```
# The imported data is stored as a data.table object
class(poolSnps)
```

```
## [1] "data.table" "data.frame"
```

As you can see, the function `vcf2DT()` converts the wide format VCF data into a long format data table. Instead of a single column for each sample (and loci in rows), all possible sample-by-locus combinations are stored in the rows of `poolSnps`, with a single column each for samples and loci. Using the `$` notation, you can access these columns as vectors.

```
# Column vector for sample and loci
head(poolSnps$SAMPLE)
```

```
## [1] "Pop1_Rep1" "Pop1_Rep1" "Pop1_Rep1" "Pop1_Rep1" "Pop1_Rep1" "Pop1_Rep1"
```

```
head(poolSnps$LOCUS)
```

```
## [1] "Chrom_1_8" "Chrom_2_16" "Chrom_3_24" "Chrom_4_32" "Chrom_5_40"
## [6] "Chrom_6_48"
```

Typically in these tutorials, and throughout the documentation for *genomalicious*, I will use `$` to indicate nested vectors of R objects (e.g. columns, list items).

## Data tables for storing SNP data

*Genomalicious* is largely built around `data.table` classed objects. Data tables feel and more-or-less function as per standard R `data.frame` objects, but they differ in two major ways. Firstly, they are much more efficient at storing large volumes of data. Secondly, they have some nifty features that facilitate easy data manipulations. Though the purpose of this tutorial is not to provide a detailed demonstration of `data.table` object features (you can find that, [here](#)), let's just take a quick look at how you can harness the power of data tables:

```
# Subset rows by depth.
poolSnps[DP > 110,]
```

```
##      CHROM POS ID REF ALT QUAL FILTER  INFO  SAMPLE  LOCUS  DP RO
## 1: Chrom_1   8 .   A   T   50   PASS DP=854 Pop1_Rep1 Chrom_1_8 113 63
## 2: Chrom_1   8 .   A   T   50   PASS DP=854 Pop2_Rep1 Chrom_1_8 119 63
## 3: Chrom_6  48 .   T   G   37   PASS DP=854 Pop2_Rep1 Chrom_6_48 116 43
## 4: Chrom_1   8 .   A   T   50   PASS DP=854 Pop4_Rep1 Chrom_1_8 120 82
## 5: Chrom_5  40 .   A   C   30   PASS DP=825 Pop4_Rep1 Chrom_5_40 113 71
##      AO
## 1: 50
## 2: 56
## 3: 73
## 4: 38
## 5: 42
```

```
# You can apply functions to columns, for example,
# take the mean depth.
poolSnps[, mean(DP)]
```

```
## [1] 69.08333
```

```
# You can even specify subsets of the data that you want to
# apply the function to, for example, take the mean depth
# for each locus.
poolSnps[, mean(DP), by=LOCUS]
```

```
##      LOCUS      V1
## 1: Chrom_1_8 71.16667
## 2: Chrom_2_16 65.41667
## 3: Chrom_3_24 72.00000
## 4: Chrom_4_32 66.00000
## 5: Chrom_5_40 68.75000
## 6: Chrom_6_48 71.16667
## 7: Chrom_7_56 68.33333
## 8: Chrom_8_64 69.83333
```

```
# You can also use this feature of column manipulation to
# apply a function to the data and create a new column using
# the ':' notation. For example, add a new column
# containing the frequency of reads that contain the
# reference allele.
poolSnps[, R.FREQ:=R0/DP]
head(poolSnps, 4)
```

```
##      CHROM POS ID REF ALT QUAL FILTER INFO SAMPLE LOCUS DP R0
## 1: Chrom_1 8 . A T 50 PASS DP=854 Pop1_Rep1 Chrom_1_8 113 63
## 2: Chrom_2 16 . C G 40 PASS DP=785 Pop1_Rep1 Chrom_2_16 82 60
## 3: Chrom_3 24 . G A 45 PASS DP=864 Pop1_Rep1 Chrom_3_24 101 90
## 4: Chrom_4 32 . T C 43 PASS DP=792 Pop1_Rep1 Chrom_4_32 101 73
##      AO R.FREQ
## 1: 50 0.5575221
## 2: 22 0.7317073
## 3: 11 0.8910891
## 4: 28 0.7227723
```

Many functions in *genomalicious* take long format `data.table` objects as their direct input. From personal experience, I find data tables from 100s of individuals, and 1,000s to 10,000s of SNP loci are very manageable, though they may take a moment to load into R (especially when importing from a VCF). However, I believe the ease, simplicity, and versatility of working with data tables makes them an excellent object class for reduced representation SNPs.

Data tables are also easily transformed into matrices and pure data frames:

```
# Converting to a matrix
matDat <- as.matrix(poolSnps)
class(matDat)
```

```
## [1] "matrix"
```

```
head(matDat, 4)
```

```
##      CHROM      POS  ID  REF ALT QUAL FILTER INFO      SAMPLE
## [1,] "Chrom_1"    " 8" ". " "A" "T" "50" "PASS" "DP=854" "Pop1_Rep1"
## [2,] "Chrom_2"   "16" ". " "C" "G" "40" "PASS" "DP=785" "Pop1_Rep1"
## [3,] "Chrom_3"   "24" ". " "G" "A" "45" "PASS" "DP=864" "Pop1_Rep1"
## [4,] "Chrom_4"   "32" ". " "T" "C" "43" "PASS" "DP=792" "Pop1_Rep1"
##      LOCUS      DP    RO    AO  R.FREQ
## [1,] "Chrom_1_8"  "113" " 63" "50" "0.55752212"
## [2,] "Chrom_2_16" " 82" " 60" "22" "0.73170732"
## [3,] "Chrom_3_24" "101" " 90" "11" "0.89108911"
## [4,] "Chrom_4_32" "101" " 73" "28" "0.72277228"
```

```
# Converting to a data frame
dfDat <- as.data.frame(poolSnps)
class(dfDat)
```

```
## [1] "data.frame"
```

```
head(dfDat, 4)
```

```
##      CHROM POS  ID  REF ALT QUAL FILTER  INFO      SAMPLE      LOCUS DP RO AO
## 1 Chrom_1   8   .   A   T   50   PASS DP=854 Pop1_Rep1 Chrom_1_8 113 63 50
## 2 Chrom_2  16   .   C   G   40   PASS DP=785 Pop1_Rep1 Chrom_2_16  82 60 22
## 3 Chrom_3  24   .   G   A   45   PASS DP=864 Pop1_Rep1 Chrom_3_24 101 90 11
## 4 Chrom_4  32   .   T   C   43   PASS DP=792 Pop1_Rep1 Chrom_4_32 101 73 28
##      R.FREQ
## 1 0.5575221
## 2 0.7317073
## 3 0.8910891
## 4 0.7227723
```

## Data structures: Long-to-wide format and genotype values

Though many functions in R expect data structured in long format, there are lots of others that require data to be wide format, especially many population genetics/genomics packages. Remember, in long form, loci-by-sample combinations are all in rows, whereas in wide format, samples are in rows and loci are in columns.

There are two functions in *genomalicious* for long-to-wide conversions: `DT2Mat_freqs()` for **allele frequencies**, and `DT2Mat_genos()` for **individual genotypes**. Both return an R **matrix** object.

Let's first try with the pooled SNP data we imported at the start of this lesson. Earlier, you made a column in `poolSnps`, `$R.FREQ`, that contained the frequency of reads per sample of the reference allele. We will treat column `$R.FREQ` as an estimate of the sample allele frequency.

`DT2Mat_freqs()` requires specification of the sample, locus, and allele frequency columns in the long format data table as the arguments, `popCol`, `locusCol`, and `freqCol`, respectively:

```
# Convert long format data table of allele frequencies to a matrix
freqMat <- DT2Mat_freqs(poolSnps, popCol='SAMPLE', locusCol='LOCUS', freqCol='R.FREQ')
freqMat
```

```
##           Chrom_1_8 Chrom_2_16 Chrom_3_24 Chrom_4_32 Chrom_5_40 Chrom_6_48
## Pop1_Rep1 0.5575221 0.7317073 0.8910891 0.7227723 0.2200000 0.5049505
## Pop1_Rep2 0.6666667 0.8085106 0.9615385 0.7454545 0.3703704 0.5098039
## Pop1_Rep3 0.5652174 0.5416667 0.9791667 0.6727273 0.2558140 0.4800000
## Pop2_Rep1 0.5294118 0.7176471 0.8823529 0.4183673 0.4204545 0.3706897
## Pop2_Rep2 0.5090909 0.7068966 0.9615385 0.3846154 0.4375000 0.3650794
## Pop2_Rep3 0.5000000 0.6530612 0.9558824 0.3658537 0.2641509 0.5128205
## Pop3_Rep1 0.7674419 0.6769231 0.9313725 0.1666667 0.6739130 0.9176471
## Pop3_Rep2 0.7187500 0.7228916 0.8250000 0.3544304 0.6627907 0.9310345
## Pop3_Rep3 0.6739130 0.5517241 0.8611111 0.1875000 0.5588235 0.9393939
## Pop4_Rep1 0.6833333 0.5757576 0.9259259 0.4019608 0.6283186 0.9090909
## Pop4_Rep2 0.6666667 0.5142857 0.9583333 0.5000000 0.6290323 0.8620690
## Pop4_Rep3 0.6774194 0.5428571 0.9166667 0.3333333 0.7115385 0.8888889
##           Chrom_7_56 Chrom_8_64
## Pop1_Rep1 0.4587156 0.02970297
## Pop1_Rep2 0.4375000 0.08196721
## Pop1_Rep3 0.5476190 0.09090909
## Pop2_Rep1 0.2755102 0.06315789
## Pop2_Rep2 0.3770492 0.04838710
## Pop2_Rep3 0.4883721 0.08695652
## Pop3_Rep1 0.2048193 0.17500000
## Pop3_Rep2 0.2000000 0.21794872
## Pop3_Rep3 0.2142857 0.18918919
## Pop4_Rep1 0.2688172 0.20000000
## Pop4_Rep2 0.2166667 0.12698413
## Pop4_Rep3 0.2424242 0.09859155
```

```
# Check the class
class(freqMat)
```

```
## [1] "matrix"
```

```
# Sample names are stored in the rows of the matrix
rownames(freqMat)
```

```
## [1] "Pop1_Rep1" "Pop1_Rep2" "Pop1_Rep3" "Pop2_Rep1" "Pop2_Rep2"
## [6] "Pop2_Rep3" "Pop3_Rep1" "Pop3_Rep2" "Pop3_Rep3" "Pop4_Rep1"
## [11] "Pop4_Rep2" "Pop4_Rep3"
```

```
# Loci names are stored in the columns of the matrix
colnames(freqMat)
```

```
## [1] "Chrom_1_8" "Chrom_2_16" "Chrom_3_24" "Chrom_4_32" "Chrom_5_40"
## [6] "Chrom_6_48" "Chrom_7_56" "Chrom_8_64"
```

There is also a way to go in reverse — that is, convert the wide format matrix back into a long format data frame. This is done by specifying the argument `flip=TRUE` (default value is set to `FALSE`). Doing so

requires the input to be a wide format matrix, and specifying `popCol`, `locusCol`, and `freqCol` is required to determine the names of these columns in the new output long format data table.

```
freqDT <- DT2Mat_freqs(freqMat, popCol='SAMPLE', locusCol='LOCUS', freqCol='R.FREQ', flip=TRUE)
head(freqDT, 8)
```

```
##      SAMPLE      LOCUS      R.FREQ
## 1: Pop1_Rep1 Chrom_1_8 0.5575221
## 2: Pop1_Rep2 Chrom_1_8 0.6666667
## 3: Pop1_Rep3 Chrom_1_8 0.5652174
## 4: Pop2_Rep1 Chrom_1_8 0.5294118
## 5: Pop2_Rep2 Chrom_1_8 0.5090909
## 6: Pop2_Rep3 Chrom_1_8 0.5000000
## 7: Pop3_Rep1 Chrom_1_8 0.7674419
## 8: Pop3_Rep2 Chrom_1_8 0.7187500
```

Up to this point, we have worked with pooled sequence data and allele frequencies. Let's now take a look at individually genotyped data. We will import a *genomalicious* dataset of four populations:

```
# Import the dataset
data(genomalicious_4pops)
indSnps <- genomalicious_4pops

# Have a look at its structure
head(indSnps)
```

```
##      POP  SAMPLE      CHROM      LOCUS  GT
## 1: Pop1 Ind1.115 Chrom_3  Chrom_3_4 0/0
## 2: Pop1 Ind1.115 Chrom_5  Chrom_5_1 0/0
## 3: Pop1 Ind1.115 Chrom_7  Chrom_7_5 0/1
## 4: Pop1 Ind1.115 Chrom_8  Chrom_8_4 0/0
## 5: Pop1 Ind1.115 Chrom_12 Chrom_12_5 0/1
## 6: Pop1 Ind1.115 Chrom_15 Chrom_15_3 0/1
```

```
# The number of unique samples, per population
indSnps[, length(unique(SAMPLE)), by=POP]
```

```
##      POP V1
## 1: Pop1 30
## 2: Pop2 30
## 3: Pop3 30
## 4: Pop4 30
```

```
# The number of unique loci
length(unique(indSnps$LOCUS))
```

```
## [1] 1205
```

You would have noticed a column in `indSnps` called `$GT`: this column contains information on each sample's genotype at a particular locus. '0' is the reference allele and '1' is the alternate allele; each allele is separated



by a '/'. This is a **separated** format of genotype values. However, another way of representing genotypes is via **counts** of one of the alleles, e.g. '0', '1', or '2'.

Typically when genotypes are represented as allele counts, we assume biallelic data. If a diploid individual possess two alleles per locus, then it is only possible to keep track of a maximum of two unique alleles. In reality, there are other ways to work around this, but many population genomic analyses are constrained to be on biallelic SNPs.

The function `genoscore_converter()` can be used to convert between the different ways of scoring genotype values. The input into this function is simply a vector of genotypes. If inputting the separated format, the vector must be a **character** class object. If inputting allele count format, the vector is ideally **integer** class object. The function bases its counts of the Alternate allele, hence a genotype of '0/0', '0/1', and '1/1', are equivalent to 0, 1, and 2 (respectively).

```
# Practise run with simple vectors
genoscore_converter(c('0/0', '0/1', '1/1'))
```

```
## [1] 0 1 2
```

```
genoscore_converter(c(0L, 1, 2))
```

```
## [1] "0/0" "0/1" "1/1"
```

```
# Now manipulate the data table of genotypes
indSnps[, GT:=genoscore_converter(GT)]
```

```
head(indSnps, 8)
```

```
##      POP  SAMPLE    CHROM    LOCUS GT
## 1: Pop1 Ind1.115 Chrom_3  Chrom_3_4  0
## 2: Pop1 Ind1.115 Chrom_5  Chrom_5_1  0
## 3: Pop1 Ind1.115 Chrom_7  Chrom_7_5  1
## 4: Pop1 Ind1.115 Chrom_8  Chrom_8_4  0
## 5: Pop1 Ind1.115 Chrom_12 Chrom_12_5  1
## 6: Pop1 Ind1.115 Chrom_15 Chrom_15_3  1
## 7: Pop1 Ind1.115 Chrom_16 Chrom_16_2  0
## 8: Pop1 Ind1.115 Chrom_17 Chrom_17_3  2
```

Finally, let's practise converting a long data table of genotypes into a genotype matrix, using `DT2Mat_genos()`. The arguments are the same as `DT2Mat_freqs()`.

```
genoMat.sep <- DT2Mat_genos(indSnps, sampCol='SAMPLE', locusCol='LOCUS', genoCol='GT')
```

```
genoMat.sep[1:8, 1:4]
```

```
##      Chrom_1001_4 Chrom_1003_2 Chrom_1004_1 Chrom_1006_3
## Ind1.1040          2           0           0           2
## Ind1.1116          2           2           0           2
## Ind1.115           2           2           0           2
## Ind1.1153          2           1           0           2
## Ind1.1161          2           0           0           2
## Ind1.1184          2           1           0           1
## Ind1.1240          2           0           0           2
## Ind1.1315          0           2           0           2
```

## Postamble

This concludes the ‘Basic Ingredients’ tutorial. You should now be comfortable with the basic functionality of *genomalicious* for importing SNP data into R. You have also familiarised yourself with SNP data structures and learnt how *genomalicious* can be used to do some basic manipulations that are common in population genetic/genomic analyses.