	Importation of example data and expected data formats  In the following, we import a Holstein-Friesian cattle dataset with 265 cows from five half-sib families from Musa and Reinsch [1], a subset of data from previous studies [2, 3]. The dataset contains 39780 markers for 29 (autosomal) chromosomes and marker effects for three example milk traits.  The genetic map, marker effects, phased genotypes, and phenotypic information can be imported using the function msq.example_data():
In [2]:	The main information required by the functions in the msq module and how they should be provided are explained below:  1. Phenotypic information: This information should be provided as a pandas data frame with two
In [3]:	group ID 0 F 10001
	1 F 10002 2 F 10003 3 F 10004 4 F 10005 260 F 10261 261 F 10262 262 F 10263 263 F 10264
in [4]:	264 F 10265  [265 rows x 2 columns]  All individuals in this data are females (F). In order to demonstrate the use of some functions for zygotes, let us assume the first 130 individuals are males (M) like:
	group.iloc[0:130, 0] = "M"  print(group)  group ID  M 10001  M 10002  M 10003  M 10004  M 10005
	4 M 10005
	1. Genetic map: The genetic map should be provided as a pandas data frame. The first three columns must be the chromosome number, marker name, and marker position in base pairs. The remaining column(s) should be maker distance/position (cM) or recombination rates in the order listed in the group column of the phenotypic information (1). Let us assume the group classification is sex. If males are listed first, then the fourth column should be the marker distance/position (or recombination rates)
In [5]:	for males. The fifth column should be the marker distance/position (or recombination rates) for females. If an average map for all groups is provided, the data frame will contain four columns only, with the average marker distance/position as the fourth column (or recombination rates). Below is an example of how the genetic map should be provided:
[0] •	Print(gmap)  CHR SNPName Position group1  0 1 SNP1 113641 0.113641  1 1 SNP2 244698 0.244698  2 1 SNP3 369418 0.369418  3 1 SNP4 447277 0.447277  4 1 SNP5 487653 0.487653
	39775 29 SNP39776 51899151 51.899151 39776 29 SNP39777 51920849 51.920849 39777 29 SNP39778 51986600 51.986600 39778 29 SNP39779 52030414 52.030414 39779 29 SNP39780 52112161 52.112161  [39780 rows x 4 columns]  Let us add another column to the data frame to simulate a data frame with group-specific maps.
in [6]:	gmap.insert(4, "group2", gmap.iloc[:, 3], True) print(gmap)  CHR SNPName Position group1 group2 0 1 SNP1 113641 0.113641 0.113641 1 1 SNP2 244698 0.244698 0.244698
	2 1 SNP3 369418 0.369418 0.369418 3 1 SNP4 447277 0.447277 0.447277 4 1 SNP5 487653 0.487653 0.487653 39775 29 SNP39776 51899151 51.899151 51.899151 39776 29 SNP39777 51920849 51.920849 39777 29 SNP39778 51986600 51.986600 39778 29 SNP39779 52030414 52.030414 39779 29 SNP39780 52112161 52.112161
	[39780 rows x 5 columns]  PyMSQ will use group1 as the marker distance/position (or recombination rates) for males and group2 for females because males are first in the data frame for phenotypic information.  Note: When using recombination rates, the first marker on a chromosome must be zero, followed by
	recombination rate between markers m and m+1. However, the first element on each chromosome does not have to be zero when using genetic position/distance (specified in cM) because the function utilizes the differences between given values. This allows any subset of markers from a genetic map to be conveniently used by copying their physical (base-pairs) and genetic map (cM) positions.  1. Allele substitution effects: The allele substitution or marker effects should be provided as a
	pandas data frame with named columns. This data frame provides the trait names used in PyMSQ functions. As a result, users should name these columns according to how they want them to appear in the results. The number of rows of this data frame should be equal to the number of markers, and the number of columns should be equal to the number of traits plus one. The first column should be the marker name (same as in the genetic map), and the remaining column(s) should be allele substitution or marker effects for each trait. Here is an example:
in [7]:	print(meff)  SNPName fat pH protein  SNP1 0.000059 -0.000163 -0.000211  SNP2 -0.000051 -0.000773 -0.000006  SNP3 0.000034 -0.000047 -0.000075  SNP4 0.000026 0.000101 -0.000118
	3 SNP4 0.000026 0.000118 4 SNP5 0.000021 0.000066 0.000075 39775 SNP39776 -0.000104 -0.000078 -0.000017 39776 SNP39777 0.000205 -0.000145 -0.000053 39777 SNP39778 0.000034 0.000025 0.000043 39778 SNP39779 -0.000148 -0.000033 -0.000003 39779 SNP39780 -0.000059 -0.000004 0.000013
	<ol> <li>Phased genotypes: Phased genotypes: Individuals' paternal and maternal haplotypes should be arranged in rows with their IDs. The phased genotypes should also be provided as a pandas data frame in one of the following ways:</li> <li>it as strings array – The data frame should have two columns (ID and haplotypes as strings without</li> </ol>
in [8]:	<ul> <li>i: as strings array – The data frame should have two columns (ID and haplotypes as strings without spaces), and the number of rows should be 2 * the number of individuals, as follows:</li> <li>print (gmat)</li> <li>1</li> <li>10001</li> <li>0222222222222222222222222222222222222</li></ul>
	<pre>1    10001    2222020222222222222222222</pre>
	<ul> <li>529 10265 22002222222222222222222222222222222</li></ul>
	Note: Since genotypes are already phased into haplotypes, PyMSQ does not expect missing values and, therefore, does not accept data with missing values. It does, however, accept any coding for phased haplotypes as long as they are biallelic (i.e., two integers must be used for coding) and taken from integers ranging from 0 to 9. As a result, outputs from several genotype phasing software can be used with ease.  Checking the input information for errors
	PyMSQ requires that the main information, including index weights, be checked for errors using the function msq.Datacheck to enhance the smooth running of other functions. Briefly, it performs checks to ensure the following:  • The number of traits matches the length of index weights.  • The ID of individuals in phenotypic information and phased genotypes are ordered, and the same.
	<ul> <li>There are no missing allele substitution or marker effects, and all entries are numeric.</li> <li>Marker names in the genetic map and marker effects are ordered, and the same.</li> <li>The number of maps is the same as the number of groups if the map number is more than one.</li> <li>The number of markers in phased genotypes, genetic map, and marker effects is the same.</li> <li>The marker distance/position or recombination rates on each chromosome are ordered.</li> </ul>
In [9]:	index_wt = [1, 1, 1]
	<pre>start = time.time() data = msq.Datacheck(gmap=gmap, meff=meff, gmat=gmat, group=group,</pre>
	Number of groups: 2: ['M' 'F']  Number of specific maps: 2  Number of chromosomes: 29  Total no. markers: 39780  Number of trait(s): 3  Trait name(s) and Index weight(s)  fat: 1  pH: 1  protein: 1
n [10]:	Time taken: 2.9 secs  The main information may be deleted because they are already stored as a class object.
	The population covariance matrix is derived using function msq.popcovmat according to Bonk et al. [4] if parameter method=1 or Santos et al. [5] if method=2. The map type must be specified using mposunit= "cM" for genetic position/distance (centimorgan), or mposunit= "reco" for recombination rates. If mposunit='reco', an additional check is run to ensure that the first marker on each chromosome has a value of zero. An error message is displayed on the screen if the value is not zero. When group-specific
n [11]:	maps are available, this function returns a named list of lists containing population covariance matrices (for each chromosome) for each group.
	Estimation of Mendelian sampling (co)variance for gametes produced by an individual  The function msq.msvarcov_g can be used to estimate the Mendelian sampling variance (MSV) for each
n [12]:	<pre>msvmsc_g = msq.msvarcov_g(info=data, covmat=covmatrices,</pre>
	sub_id=None, progress=True)  print('Time taken: ', round(time.time() - start, 2), 'secs')  print(msvmsc_g)  Progress:
	2 10003 M 0.021665 0.013760 0.010475 0.015662 0.009849 3 10004 M 0.000949 0.000317 0.004836 -0.000044 -0.000158 4 10005 M 0.021892 0.014447 0.012676 0.015576 0.010120 260 10261 F 0.022112 0.012953 0.010050 0.017090 0.009524 261 10262 F 0.023853 0.016183 0.012896 0.017691 0.012189 262 10263 F 0.000604 0.000096 0.002694 0.000025 0.000126 263 10264 F 0.022809 0.015301 0.012748 0.017123 0.011174
	264 10265 F 0.022491 0.016226 0.013892 0.016460 0.011973  protein AG_fat AG_pH AG_protein AG 0 0.002321 0.000756 0.001802 0.002300 0.004858 1 0.015098 0.055475 0.039215 0.043486 0.138176 2 0.013710 0.051087 0.034084 0.039221 0.124392 3 0.001352 0.001223 0.004995 0.001149 0.007367 4 0.013318 0.051915 0.037243 0.039013 0.128172
	260 0.015024 0.052155 0.032527 0.041637 0.126319 261 0.015089 0.057727 0.041268 0.044968 0.143963 262 0.001841 0.000725 0.002916 0.001992 0.005633 263 0.014931 0.055232 0.039223 0.043227 0.137683 264 0.013582 0.055177 0.042091 0.042014 0.139282  [265 rows x 12 columns]
ı [13]:	As seen above, Mendelian sampling (co-)variance will be estimated for all individuals if parameter sub_id is None. There may be cases where a user might be interested in a subset of individuals. In this case, the user can provide individuals of interest as a pandas data frame with one column. Assuming the candidates of interest are  aaa = np.arange(0, 265, 40) df gam = data.group.iloc[aaa, 1]
	print (df_gam)  0    10001 40    10041 80    10081 120    10121 160    10161 200    10201 240    10241
n [14]:	Name: ID, dtype: int64  Then the Mendelian sampling (co-)variances of traits for these individuals are given below:
	<pre>print(msvmsc_gsub)  Time taken: 0.66 secs</pre>
	0 10101 11 0.022001 0.0111/0 0.012002 0.01012
	4 10161 F 0.000469 -0.000055 0.002331 0.000368 -0.000489 5 10201 F 0.000466 -0.000022 0.002190 0.000254 -0.000319 6 10241 F 0.000566 0.000213 0.001609 0.000415 0.000152  protein AG_fat AG_pH AG_protein AG 0 0.002321 0.000756 0.001802 0.002300 0.004858 1 0.001163 0.000676 0.003282 0.001089 0.005048 2 0.003116 0.001162 0.003546 0.003382 0.008090
	4 10161 F 0.000469 -0.000055 0.002331 0.000368 -0.000489 5 10201 F 0.000466 -0.000022 0.002190 0.000254 -0.000319 6 10241 F 0.000566 0.000213 0.001609 0.000415 0.000152  protein AG_fat AG_pH AG_protein AG 0 0.002321 0.000756 0.001802 0.002300 0.004858 1 0.001163 0.000676 0.003282 0.001089 0.005048
n [15]:	### 10161
n [15]:	<pre>4 10161  F  0.000469 -0.000055  0.002331  0.000368  -0.000489 5 10201  F  0.000466 -0.000022  0.002190  0.000254  -0.000319 6 10241  F  0.000566  0.000213  0.001609  0.000415  0.000152  protein</pre>
ı [15]:	1 10161
n [15]:	10161   F   0.000466 -0.000025   0.002331   0.000368   -0.000489     5   10201   F   0.000466 -0.000022   0.000250   0.000254   -0.000519     6   10241   F   0.000566   0.000213   0.001609   0.000415   0.000152
n [15]:	10161   F   0.000469 -0.000055   0.002331   0.000368   -0.000489     5   10201   F   0.000466 -0.00022   0.002190   0.000254   -0.000319     6   10241   F   0.000566   0.000213   0.001609   0.000415   0.000152
n [15]:	10161
	10161
	1016  F
	1.0016  F 0.000469 -0.000025 0.002331 0.000366 -0.000489
n [16]:	1
n [16]:	1   101
n [16]:	1
n [16]:	1
n [16]:	10.00
n [16]:	1316
n [16]:	### 1
n [16]:	10.00
n [16]:	1016
n [16]:	# 16.56   17.0 0.000   10.0000   10.0000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.000000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.000000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.000000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.000000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.00000   10.000000   10.000000   10.000000   10.000000   10.000000   10.000000   10.000000   10.000000   10.000000   10.00000000   10.000000000   10.0000000000
n [16]:	# 10.51
n [16]:	March   Process   Proces
n [16]:	## Section of Selection criteria    Section of Selection selection criteria
n [16]:	## Section of Selection criteria  ## Sec
n [16]:	The content of the
n [16]: n [17]:	Section   Control of
n [16]:	## Committee   Com
n [16]: n [17]:	### SECRETARY PROPRIES THE CONTROL OF THE CONTROL O
n [16]: n [27]:	### STATES   10 months of States   10 months
n [16]: n [27]:	## Secret Secretary Control of Co
n [16]: n [27]:	## STATE OF THE PROPERTY OF TH
n [16]: n [27]:	### Security of Control of Contro
1 [16]:  1 [17]:  1 [21]:	# Section 1 - Company of the Company