

---

# **PyMethylProcess Documentation**

***Release 0.1***

**Joshua Levy**

**Mar 22, 2019**



## CONTENTS:

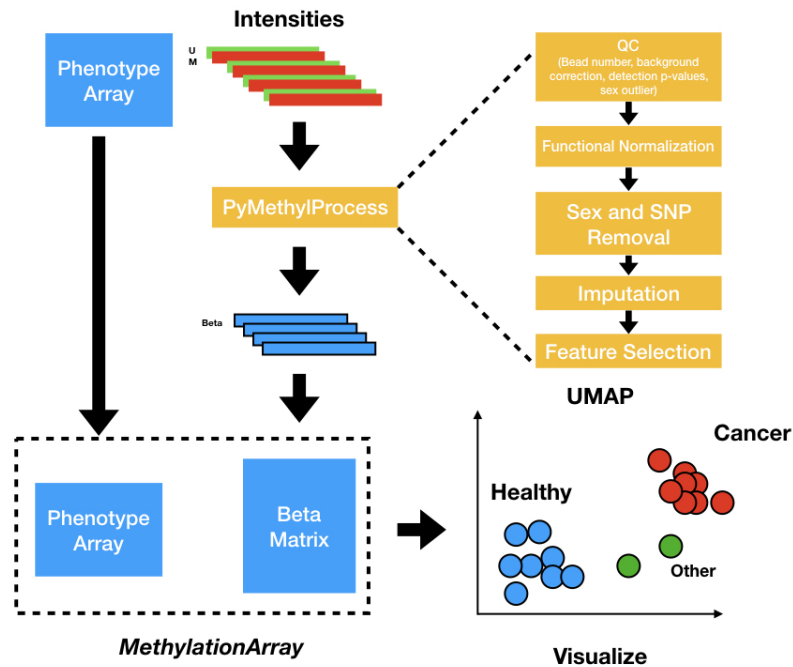
<b>1</b>	<b>PreProcessDataTypes.py</b>	<b>5</b>
<b>2</b>	<b>MethylationDataTypes.py</b>	<b>9</b>
<b>3</b>	<b>meffil_functions.py</b>	<b>13</b>
<b>4</b>	<b>general_machine_learning.py</b>	<b>15</b>
<b>5</b>	<b>pymethyl-install</b>	<b>17</b>
5.1	change_gcc_path . . . . .	17
5.2	install_bioconductor . . . . .	17
5.3	install_custom . . . . .	17
5.4	install_meffil . . . . .	18
5.5	install_minfi_others . . . . .	18
5.6	install_r_packages . . . . .	18
5.7	install_some_deps . . . . .	18
5.8	install_tcga_biolinks . . . . .	18
<b>6</b>	<b>pymethyl-visualize</b>	<b>19</b>
6.1	plot_cell_type_results . . . . .	19
6.2	plot_heatmap . . . . .	19
6.3	transform_plot . . . . .	20
<b>7</b>	<b>pymethyl-preprocess</b>	<b>23</b>
7.1	batch_deploy_preprocess . . . . .	23
7.2	combine_methylation_arrays . . . . .	24
7.3	concat_sample_sheets . . . . .	24
7.4	create_sample_sheet . . . . .	24
7.5	download_clinical . . . . .	25
7.6	download_geo . . . . .	25
7.7	download_tcga . . . . .	26
7.8	feature_select . . . . .	26
7.9	get_categorical_distribution . . . . .	26
7.10	imputation_pipeline . . . . .	27
7.11	meffil_encode . . . . .	28
7.12	merge_sample_sheets . . . . .	28
7.13	na_report . . . . .	29
7.14	preprocess_pipeline . . . . .	29
7.15	remove_diseases . . . . .	30
7.16	split_preprocess_input_by_subtype . . . . .	30

<b>8</b>	<b>pymethyl-utils</b>	<b>33</b>
8.1	backup_pkl . . . . .	33
8.2	bin_column . . . . .	33
8.3	concat_csv . . . . .	34
8.4	counts . . . . .	34
8.5	create_external_validation_set . . . . .	34
8.6	feature_select_train_val_test . . . . .	35
8.7	fix_key . . . . .	35
8.8	modify_pheno_data . . . . .	36
8.9	move_jpg . . . . .	36
8.10	overwrite_pheno_data . . . . .	36
8.11	pkl_to_csv . . . . .	37
8.12	print_number_sex_cpgs . . . . .	37
8.13	print_shape . . . . .	37
8.14	ref_estimate_cell_counts . . . . .	38
8.15	remove_sex . . . . .	38
8.16	remove_snps . . . . .	39
8.17	set_part_array_background . . . . .	39
8.18	stratify . . . . .	39
8.19	subset_array . . . . .	40
8.20	train_test_val_split . . . . .	40
8.21	write_cpgs . . . . .	40
<b>9</b>	<b>Indices and tables</b>	<b>43</b>
	<b>Python Module Index</b>	<b>45</b>
	<b>Index</b>	<b>47</b>

<https://github.com/Christensen-Lab-Dartmouth/PyMethylProcess>

To get started, download pymethylprocess using Docker (joshualevy44/pymethylprocess) or PIP (pymethylprocess) and run pymethyl-install\_r\_dependencies.

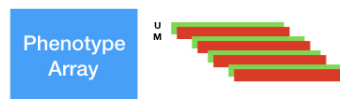
There is both an API and CLI available for use. Examples for CLI usage can be found in ./example\_scripts.



# Pipeline

`pymethyl-preprocess download_geo -g GSE87571`

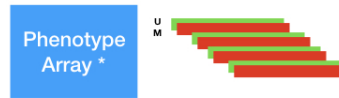
- Download
- Format
- Preprocess
- Visualize



# Pipeline

```
pymethyl-preprocess create_sample_sheet -is ./geo_idats/
GSE87571_clinical_info.csv -s geo -i geo_idats/ -os
geo_idats/samplesheet.csv -d "disease state:ch1" -c
include_col.txt
```

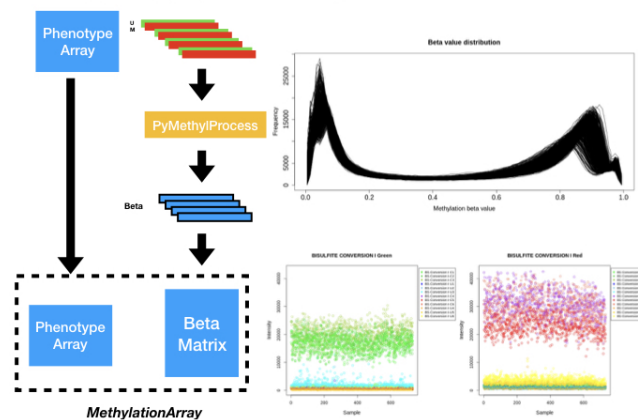
- Download
- **Format**
- Preprocess
- Visualize



# Pipeline

```
pymethyl-preprocess preprocess_pipeline -i geo_idats/ -p minfi -noob
pymethyl-utils remove_sex -i preprocess_outputs/methyl_array.pkl
pymethyl-preprocess imputation_pipeline -i ./autosomal/methyl_array.pkl -s fancyimpute -m KNN -k 15
pymethyl-preprocess feature_select -n 300000
```

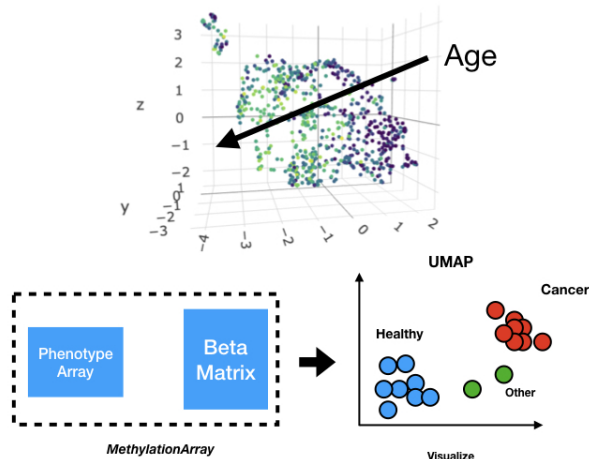
- Download
- Format
- **Preprocess**
- Visualize



# Pipeline

```
pymethyl-visualize transform_plot -o visualizations/pre_vae_umap.html -c Age -nn 8
```

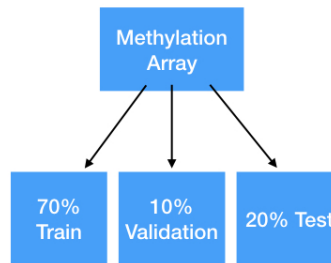
- Download
- Format
- Preprocess
- Visualize



# Pipeline

```
pymethyl-utils train_test_val_split -tp .8 -vp .125
```

- Download
- Format
- Preprocess
- Visualize







## PREPROCESSDATATYPES.PY

Contains datatypes core to downloading IDATs, preprocessing IDATs and samplesheets.

```
class pymethylprocess.PreProcessDataTypes.PreProcessIDAT (idat_dir, minfi=None, en-  
mix=None, base=None,  
meffil=None)
```

Class that will preprocess IDATs using R pipelines.

**idat\_dir** Location of idats or samplesheet csv.

**minfi** Rpy2 importr minfi library, default to None will load through rpy2

**enmix** Rpy2 importr enmix library, default to None will load through rpy2

**base** Rpy2 importr base library, default to None will load through rpy2

**meffil** Rpy2 importr meffil library, default to None will load through rpy2

**export\_csv** (*output\_dir*)

Export pheno and beta dataframes to CSVs

**output\_dir** Where to store csvs.

**export\_pickle** (*output\_pickle, disease=""*)

Export pheno and beta dataframes to pickle, stored in python dict that can be loaded into MethylationArray

**output\_pickle** Where to store MethylationArray.

**disease** Custom naming scheme for data.

**export\_sql** (*output\_db, disease=""*)

Export pheno and beta dataframes to SQL

**output\_db** Where to store data, sqlite db.

**disease** Custom naming scheme for data.

**extract\_manifest** ()

Get manifest from RGSet.

**extract\_pheno\_data** (*methylset=False*)

Extract pheno data from MSet or RGSet, minfi.

**methylset** If MSet has been created, set to True, else extract from original RGSet.

**filter\_beta** ()

After creating beta, filter out outliers.

**get\_beta** ()

Get beta value matrix from minfi after finding RSet.

**get\_meth()**  
Get methylation intensity matrix from MSet

**get\_unmeth()**  
Get unmethylated intensity matrix from MSet

**load\_idats()**  
For minfi pipeline, load IDATs from specified idat\_dir.

**move\_jpg()**  
Move jpeg files from current working directory to the idat directory.

**output\_pheno\_beta** (*meffil=False*)  
Get pheno and beta dataframe objects stored as attributes for input to MethylationArray object.  
**meffil** True if ran meffil pipeline.

**plot\_original\_qc** (*output\_dir*)  
Plot QC results from ENmix pipeline and possible minfi. Still experimental.  
**output\_dir** Where to store plots.

**plot\_qc\_metrics** (*output\_dir*)  
Plot QC results from ENmix pipeline and possible minfi. Still experimental.  
**output\_dir** Where to store plots.

**preprocessENmix** (*n\_cores=6*)  
Run ENmix preprocessing pipeline.  
**n\_cores** Number of CPUs to use.

**preprocessMeffil** (*n\_cores=6, n\_pcs=4, qc\_report\_fname='qc/report.html', normalization\_report\_fname='norm/report.html', pc\_plot\_fname='qc/pc\_plot.pdf', useCache=True, qc\_only=True, qc\_parameters={'p.beadnum.cpgs': 0.1, 'p.beadnum.samples': 0.1, 'p.detection.cpgs': 0.1, 'p.detection.samples': 0.1}, rm\_sex=False*)  
Run meffil preprocessing pipeline with functional normalization.  
**n\_cores** Number of CPUs to use.  
**n\_pcs** Number of principal components to use for functional normalization, set to -1 to autoselect via kneedle algorithm.  
**qc\_report\_fname** HTML filename to store QC report.  
**normalization\_report\_fname** HTML filename to store normalization report  
**pc\_plot\_fname** PDF file to store principal components plot.  
**useCache** Use saved QC objects instead of running through QC again.  
**qc\_only** Perform QC, then save and quit before normalization.  
**qc\_parameters** Python dictionary with parameters for qc.  
**rm\_sex** Remove non-autosomal cpgs?

**preprocessNoob()**  
Run minfi preprocessing with Noob normalization

**preprocessRAW()**  
Run minfi preprocessing with RAW normalization

**preprocess\_enmix\_pipeline** (*n\_cores=6, pipeline='enmix', noob=False, qc\_only=False, use\_cache=False*)

Run complete ENmix or minfi preprocessing pipeline.

**n\_cores** Number CPUs.

**pipeline** Run enmix or minfi

**noob** Noob norm or RAW if minfi running.

**qc\_only** Save and quit after only running QC?

**use\_cache** Load preexisting RGSet instead of running QC again.

**return\_beta** ()

Return minfi RSet after having created MSet.

**to\_methyl\_array** (*disease=""*)

Convert results from preprocessing into MethylationArray, and directly return MethylationArray object.

**disease** Custom naming scheme for data.

**class** pymethylprocess.PreProcessDataTypes.PreProcessPhenoData (*pheno\_sheet, idat\_dir, header\_line=0*)

Class that will manipulate phenotype samplesheet before preprocessing of IDATs.

**pheno\_sheet** Location of clinical info csv.

**idat\_dir** Location of idats

**header\_line** Where to start reading clinical csv

**concat** (*other\_formatted\_sheet*)

Concat multiple PreProcessPhenoData objects, concat their dataframes to accept more than one samplesheet/dataset.

**other\_formatted\_sheet** Other PreProcessPhenoData to concat.

**export** (*output\_sheet\_name*)

Export pheno data to csv after done with manipulation.

**output\_sheet\_name** Output csv name.

**format\_custom** (*basename\_col, disease\_class\_column, include\_columns={}*)

Custom format clinical sheet if user supplied idats.

**basename\_col** Column name of sample names.

**disease\_class\_column** Disease column of clinical info csv.

**include\_columns** Dictionary specifying other columns to include, and new names to assign them to.

**format\_geo** (*disease\_class\_column='methylation class:ch1', include\_columns={}*)

Format clinical sheets if downloaded geo idats.

**disease\_class\_column** Disease column of clinical info csv.

**include\_columns** Dictionary specifying other columns to include, and new names to assign them to.

**format\_tcga** (*mapping\_file='idat\_filename\_case.txt'*)

Format clinical sheets if downloaded tcga idats.

**mapping\_file** Maps uuids to proper tcga sample names, should be downloaded with tcga clinical information.

**get\_categorical\_distribution** (*key, disease\_only=False, subtype\_delimiter=', '*)

Print categorical distribution, counts for each unique value in phenotype column.

**key** Phenotype Column.

**disease\_only** Whether to split phenotype column entries by delimiter.

**subtype\_delimiter** Subtype delimiter to split on.

**merge** (*other\_formatted\_sheet, use\_second\_sheet\_disease=True, no\_disease\_merge=False*)

Merge multiple PreProcessPhenoData objects, merge their dataframes to accept more than one samplesheet/dataset or add more pheno info.

**other\_formatted\_sheet** Other PreProcessPhenoData to merge.

**use\_second\_sheet\_disease** Change disease column to that of second sheet instead of first.

**no\_disease\_merge** Keep both disease columns from both sheets.

**remove\_diseases** (*exclude\_disease\_list, low\_count, disease\_only, subtype\_delimiter*)

Remove samples with certain diseases from disease column.

**exclude\_disease\_list** List containing diseases to remove.

**low\_count** Remove samples that have less than x disease occurrences in column.

**disease\_only** Whether to split phenotype column entries by delimiter.

**subtype\_delimiter** Subtype delimiter to split on.

**split\_key** (*key, subtype\_delimiter*)

Split pheno column by key, with subtype delimiter, eg. entry S1,s2 -> S1 with delimiter “,”.

**key** Pheno column name.

**subtype\_delimiter** Subtype delimiter to split on.

**class** pymethylprocess.PreProcessDataTypes.TCGADownloader

Downloads TCGA and GEO IDAT and clinical data

**download\_clinical** (*output\_dir*)

Download TCGA Clinical Data.

**output\_dir** Where to output clinical data csv.

**download\_geo** (*query, output\_dir*)

Download GEO IDATs.

**query** GEO accession number to query, must be 450k/850k.

**output\_dir** Output directory to store idats and clinical information csv

**download\_tcg** (*output\_dir*)

Download TCGA IDATs.

**output\_dir** Where to output idat files.

## METHYLATIONDATATYPES.PY

Contains datatypes core to storing beta and phenotype methylation data, and imputation.

```
class pymethylprocess.MethylationDataTypes.ImputerObject (solver, method, opts={})
    Class that stores and accesses different types of imputers. Construct sklearn-like imputer given certain input
    arguments.

    solver Library for imputation, eg. sklearn, fancyimpute.

    method Imputation method in library, named.

    opts Additional options to assign to imputer.

    return_imputer ()
        Return initialized sklearn-like imputer.

class pymethylprocess.MethylationDataTypes.MethylationArray (pheno_df, beta_df,
                                                             name="")
    Stores beta and phenotype information and performs various operations. Initialize MethylationArray object by
    inputting dataframe of phenotypes and dataframe of beta values with samples as index.

    pheno_df Phenotype dataframe (samples x covariates)

    beta_df Beta Values Dataframe (samples x cpGs)

    bin_column (col, n_bins)
        Turn continuous variable/covariate into categorical bins. Returns name of new column and updates phe-
        notype matrix to reflect this change.

    col Continuous column of phenotype array to bin.

    n_bins Number of bins to create.

    categorical_breakdown (key)
        Print categorical distribution, counts for each unique value in phenotype column.

    key Phenotype Column.

    feature_select (n_top_cpGs, feature_selection_method='mad', metric='correlation', nn=10)
        Perform unsupervised feature selection on MethylationArray.

    n_top_cpGs Number of CpGs to retain.

    feature_selection_method Method to perform selection.

    metric If considering structural feature selection like SPEC, use this distance metric.

    nn Number of nearest neighbors.

    classmethod from_pickle (input_pickle)
        Load MethylationArray stored in pickle.
```

Usage: `MethylationArray.from_pickle([input_pickle])`

**input\_pickle** Stored MethylationArray pickle.

**groupby** (*key*)  
Groupby for Methylation Array. Returns generator of methylation arrays grouped by key.

**preprocess\_sample\_df** New phenotype dataframe.

**impute** (*imputer*)  
Perform imputation on NaN beta values. Input imputer returned from `ImputerObject`.

**imputer** Type of imputer object, in sklearn type interface.

**merge\_preprocess\_sheet** (*preprocess\_sample\_df*)  
Feed in another phenotype dataframe that will be merged with existing phenotype array.

**preprocess\_sample\_df** New phenotype dataframe.

**overwrite\_pheno\_data** (*preprocess\_sample\_df*)  
Feed in another phenotype dataframe that will overwrite overlapping keys of existing phenotype array.

**preprocess\_sample\_df** New phenotype dataframe.

**remove\_missingness** (*cpg\_threshold=None, sample\_threshold=None*)  
Remove samples and CpGs with certain level of missingness..

**cpg\_threshold** If more than fraction of Samples for this CpG are missing, remove cpg.

**sample\_threshold** If more than fraction of CpGs for this sample are missing, remove sample.

**remove\_na\_samples** (*outcome\_cols*)  
Remove samples of MethylationArray who have missing values in phenotype column.

**outcome\_cols** Phenotype columns, if any rows contain missing values, samples are removed.

**remove\_whitespace** (*key*)  
Remove whitespaces from phenotype column.

**key** Phenotype column.

**return\_cpgs** ()  
Return list of cpGs of MethylationArray

**return\_idx** ()  
Return sample names of MethylationArray.

**return\_raw\_beta\_array** ()  
Return numpy array of methylation beta values.

**return\_shape** ()  
Return dimensionality and number of samples of beta matrix.

**split\_by\_subtype** (*disease\_only, subtype\_delimiter*)  
Split MethylationArray into generator of MethylationArrays by phenotype column. Much akin to groupby. Only splits from disease column.

**disease\_only** Consider disease superclass.

**subtype\_delimiter** How to break up disease column if using `disease_only`.

**split\_key** (*key, subtype\_delimiter*)  
Manipulate an entire phenotype column, splitting each element up by some delimiter.

**key** Phenotype column.

**subtype\_delimiter** How to break up strings in columns. S1,s2 -> S1 for instance.

**split\_train\_test** (*train\_p=0.8, stratified=True, disease\_only=False, key='disease', subtype\_delimiter=',', val\_p=0.0*)

Split MethylationArray into training and test sets, with option to stratify by categorical covariate.

**train\_p** Fraction of methylation array to use as training set.

**stratified** Whether to stratify by categorical variable.

**disease\_only** Consider disease superclass by some delimiter. For instance if disease is S1,s2, superclass would be S1.

**key** Column to stratify on.

**subtype\_delimiter** How to split disease column into super/subclass.

**val\_p** If set greater than 0, will create additional validation set, fraction of which is broken off from training set.

**subsample** (*key='disease', n\_samples=None, frac=None, categorical=False*)

Subsample MethylationArray, make the set randomly smaller.

**key** If stratifying, use this column of pheno array.

**n\_samples** Number of samples to consider overall, or per stratum.

**frac** Alternative to n\_samples, where x frac of array or stratum is considered.

**categorical** Whether to stratify by column.

**subset\_cpgs** (*cpgs*)

Subset beta matrix by list of CpGs. Parameters ——— cpgs

CpGs to subset by.

**subset\_index** (*index*)

Subset MethylationArray by samples.

**index** Sample names to subset by.

**write\_csvs** (*output\_dir*)

Write phenotype data and beta values to csvs.

**output\_dir** Directory to output csv files.

**write\_db** (*conn, disease=""*)

Store phenotype data and beta values in SQL database.

**conn** SQLite connection.

**disease** Create new tables in db that are related to disease state by this name.

**write\_pickle** (*output\_pickle, disease=""*)

Store phenotype data and beta values in pickle file. Is default file format for storing MethylationArray objects.

**output\_pickle** Pickle file to store MethylationArray data.

**class** pymethylprocess.MethylationDataTypes.**MethylationArrays** (*list\_methylation\_arrays*)

Literally a list of methylation arrays, with methods operate on these arrays that is memory efficient. Initialize with list of methylation arrays. Can optionally leave list empty or with one element.

**list\_methylation\_arrays** List of methylation arrays.

**combine** (*array\_generator=None*)

Combine the list of methylation arrays into one array via concatenation of beta matrices and phenotype arrays.

**array\_generator** Generator of additional methylation arrays for computational memory minimization.

**impute** (*imputer*)

Impute all methylation arrays.

**imputer** Type of imputation, sklearn-like.

**write\_dbs** (*conn*)

Write list of methylation arrays to SQL database. Recommend naming MethylationArray.

**conn** SQL connection.

**write\_pkls** (*pkl*)

Write list of methylation arrays to single pickle. Recommend naming each MethylationArray.

**pkl** Pickle file to write to.

`pymethylprocess.MethylationDataTypes.extract_pheno_beta_df_from_folder` (*folder*)

Return phenotype and beta dataframes from specified folder with csv.

**folder** Input folder.

`pymethylprocess.MethylationDataTypes.extract_pheno_beta_df_from_pickle_dict` (*input\_dict*,  
*dis-*  
*ease=""*)

Return phenotype and beta dataframes from specified dictionary storing MethylationArray python dictionary.

**input\_dict** Python dictionary storing pheno/beta information.

`pymethylprocess.MethylationDataTypes.extract_pheno_beta_df_from_sql` (*conn*,  
*dis-*  
*ease=""*)

Return phenotype and beta dataframes from SQL tables storing MethylationArray info.

**conn** SQL connection.



## MEFFIL\_FUNCTIONS.PY

Contains a few R functions that interact with meffil and minfi.

`pymethylprocess.meffil_functions.est_cell_counts_IDOL (rgset, library)`

Given RGSet object, estimate cell counts for 450k/850k using reference approach via IDOL library.

**rgset** RGSet object stored in python via rpy2

**library** What type of CpG library to use.

`pymethylprocess.meffil_functions.est_cell_counts_meffil (qc_list,  
cell_type_reference)`

Given QCObject list R object, estimate cell counts using reference approach via meffil.

**qc\_list** R list containing qc objects.

**cell\_type\_reference** Reference blood/tissue set.

`pymethylprocess.meffil_functions.est_cell_counts_minfi (rgset)`

Given RGSet object, estimate cell counts using reference approach via minfi.

**rgset** RGSet object stored in python via rpy2

`pymethylprocess.meffil_functions.load_detection_p_values_beadnum (qc_list,  
n_cores)`

Return list of detection p-value matrix and bead number matrix.

**qc\_list** R list containing qc objects.

**n\_cores** Number of cores to use in computation.

`pymethylprocess.meffil_functions.r_autosomal_cpgs (array_type='450k')`

Return list of autosomal cpg probes per platform.

**array\_type** 450k/850k array?

`pymethylprocess.meffil_functions.r_snp_cpgs (array_type='450k')`

Return list of SNP cpg probes per platform.

**array\_type** 450k/850k array?

`pymethylprocess.meffil_functions.remove_sex (beta, array_type='450k')`

Remove non-autosomal cpGs from beta matrix.

**array\_type** 450k/850k array?

`pymethylprocess.meffil_functions.set_missing (beta, pval_beadnum, detection_val=1e-06)`

Set missing beta values to NA, taking into account detection values and bead number thresholds.

**pval\_beadnum** Detection pvalues and number of beads per cpg/samples

**detection\_val** If threshold to set site to missingness based on p-value detection.



## GENERAL\_MACHINE\_LEARNING.PY

Contains a machine learning class to perform scikit-learn like operations, along with held-out hyperparameter grid search.

```
class pymethylprocess.general_machine_learning.MachineLearning(model, options,  
                                                                grid={}, labe-  
                                                                lencode=False,  
                                                                n_eval=0)
```

Machine learning class to run sklearn-like pipeline on MethylationArray data. Initialize object with scikit-learn model, and optionally supply a hyperparameter search grid.

**model** Scikit-learn-like model, classification, regression, dimensionality reduction, clustering etc.

**options** Options to supply model in form of dictionary.

**grid** Alternatively, supply search grid to search for best hyperparameters.

**labelencode** T/F encode string labels.

**n\_eval** Number of evaluations for randomized grid search, if set to 0, perform exhaustive grid search

**assign\_results\_to\_pheno\_col** (*methyl\_array, new\_col, output\_pkl*)  
Assign results to new phenotype column.

**methyl\_array** MethylationArray.

**new\_col** New column name.

**output\_pkl** Output pickle to dump MethylationArray to.

**fit** (*train\_methyl\_array, val\_methyl\_array=None, outcome\_cols=None*)  
Fit data to model.

**train\_methyl\_array** Training MethylationArray.

**val\_methyl\_array** Validation MethylationArray. Can set to None.

**outcome\_cols** Set to none if not needed, but phenotype column to train on, can be multiple.

**fit\_predict** (*train\_methyl\_array, outcome\_cols=None*)  
Fit and predict training data.

**train\_methyl\_array** Training MethylationArray.

**outcome\_cols** Set to none if not needed, but phenotype column to train on, can be multiple.

**fit\_transform** (*train\_methyl\_array, outcome\_cols=None*)  
Fit and transform to training data.

**train\_methyl\_array** Training MethylationArray.

**outcome\_cols** Set to none if not needed, but phenotype column to train on, can be multiple.

**predict** (*test\_methyl\_array*)  
Make new predictions on test methylation array.

**test\_methyl\_array** Testing MethylationArray.

**return\_outcome\_metric** (*methyl\_array, outcome\_cols, metric, run\_bootstrap=False*)  
Supply metric to evaluate results.

**methyl\_array** MethylationArray to evaluate.

**outcome\_cols** Outcome phenotype columns.

**metric** Sklearn evaluation metric.

**run\_bootstrap** Make 95% CI from 1k bootstraps.

**store\_results** (*output\_pkl, results\_dict={}*)  
Store results in pickle file.

**output\_pkl** Output pickle to dump results to.

**results\_dict** Supply own results dict to be dumped.

**transform** (*test\_methyl\_array*)  
Transform test methylation array.

**test\_methyl\_array** Testing MethylationArray.

**transform\_results\_to\_beta** (*methyl\_array, output\_pkl*)  
Transform beta matrix into reduced beta matrix and store.

**methyl\_array** MethylationArray.

**output\_pkl** Output pickle to dump MethylationArray to.

## PYMETHYL-INSTALL

```
pymethyl-install [OPTIONS] COMMAND [ARGS]...
```

### Options

#### **--version**

Show the version and exit.

### 5.1 change\_gcc\_path

Change GCC and G++ paths if don't have version 7.2.0. [Experimental]

```
pymethyl-install change_gcc_path [OPTIONS]
```

### 5.2 install\_bioconductor

Installs bioconductor.

```
pymethyl-install install_bioconductor [OPTIONS]
```

### 5.3 install\_custom

Installs bioconductor packages.

```
pymethyl-install install_custom [OPTIONS]
```

### Options

**-p, --package** <package>  
Custom packages. [default: ENmix]

**-m, --manager**  
Use BiocManager (recommended).

## 5.4 install\_meffil

Installs meffil (update!).

```
pymethyl-install install_meffil [OPTIONS]
```

## 5.5 install\_minfi\_others

Installs minfi and other dependencies.

```
pymethyl-install install_minfi_others [OPTIONS]
```

## 5.6 install\_r\_packages

Installs r packages.

```
pymethyl-install install_r_packages [OPTIONS]
```

### Options

**-p, --package** <package>  
Custom packages. [default: ]

## 5.7 install\_some\_deps

Installs bioconductor, minfi, enmix, tcga biolinks, and meffil.

```
pymethyl-install install_some_deps [OPTIONS]
```

## 5.8 install\_tcga\_biolinks

Installs tcga biolinks.

```
pymethyl-install install_tcga_biolinks [OPTIONS]
```

## PYMETHYL-VISUALIZE

```
pymethyl-visualize [OPTIONS] COMMAND [ARGS]...
```

### Options

**--version**  
Show the version and exit.

### 6.1 plot\_cell\_type\_results

Plot csv containing cell type results into side by side boxplots.

```
pymethyl-visualize plot_cell_type_results [OPTIONS]
```

### Options

**-i, --input\_csv** <input\_csv>  
Input csv. [default: cell\_type\_estimates.csv]

**-o, --outfilename** <outfilename>  
Output png. [default: visualizations/cell\_type\_results.png]

**-cols, --plot\_cols** <plot\_cols>  
Plot columns. [default: Gran, CD4T, CD8T, Bcell, Mono, NK, gMDSC]

**-fs, --font\_scale** <font\_scale>  
Font scaling [default: 1.0]

### 6.2 plot\_heatmap

Plot heatmap from CSV file.

```
pymethyl-visualize plot_heatmap [OPTIONS]
```

## Options

**-i, --input\_csv** <input\_csv>  
Input csv. [default: ]

**-o, --outfilename** <outfilename>  
Output png. [default: output.png]

**-idx, --index\_col** <index\_col>  
Index load dataframe [default: 0]

**-fs, --font\_scale** <font\_scale>  
Font scaling [default: 1.0]

**-min, --min\_val** <min\_val>  
Min heat val [default: 0.0]

**-max, --max\_val** <max\_val>  
Max heat val, if -1, defaults to None [default: 1.0]

**-a, --annot**  
Annotate heatmap [default: False]

**-n, --norm**  
Normalize matrix data [default: False]

**-c, --cluster**  
Cluster matrix data [default: False]

**-m, --matrix\_type** <matrix\_type>  
Type of matrix supplied [default: none]

**-x, --xticks**  
Show x ticks [default: False]

**-y, --yticks**  
Show y ticks [default: False]

**-t, --transpose**  
Transpose matrix data [default: False]

**-col, --color\_column** <color\_column>  
Color column. [default: color]

## 6.3 transform\_plot

Dimensionality reduce VAE or original beta values using UMAP and plot using plotly.

```
pymethyl-visualize transform_plot [OPTIONS]
```

## Options

**-i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]

**-c, --column\_of\_interest** <column\_of\_interest>  
Column extract from phenotype data. [default: disease]



**-o, --output\_file** <output\_file>  
Output visualization. [default: ./visualization.html]

**-nn, --n\_neighbors** <n\_neighbors>  
Number of neighbors UMAP. [default: 5]

**-a, --axes\_off**  
Whether to turn axes on or off.

**-s, --supervised**  
Supervise umap embedding.

**-d, --min\_dist** <min\_dist>  
UMAP min distance. [default: 0.1]

**-m, --metric** <metric>  
Reduction metric. [default: euclidean]

**-cc, --case\_control\_override**  
Add controls from case\_control column and override current disease for classification tasks. [default: False]



## PYMETHYL-PREPROCESS

```
pymethyl-preprocess [OPTIONS] COMMAND [ARGS]...
```

### Options

**--version**  
Show the version and exit.

### 7.1 batch\_deploy\_preprocess

Deploy multiple preprocessing jobs in series or parallel.

```
pymethyl-preprocess batch_deploy_preprocess [OPTIONS]
```

### Options

**-n, --n\_cores** <n\_cores>  
Number cores to use for preprocessing. [default: 6]

**-i, --subtype\_output\_dir** <subtype\_output\_dir>  
Output subtypes pheno csv. [default: ./preprocess\_outputs/]

**-m, --meffil**  
Preprocess using meffil.

**-t, --torque**  
Job submission torque.

**-r, --run**  
Actually run local job or just print out command.

**-s, --series**  
Run commands in series.

**-p, --pc\_qc\_parameters\_csv** <pc\_qc\_parameters\_csv>  
For meffil, qc parameters and pcs for final qc and functional normalization. [default: ./preprocess\_outputs/pc\_qc\_parameters.csv]

**-u, --use\_cache**  
If this is selected, loads qc results rather than running qc again. Only works for meffil selection.

**-qc, --qc\_only**

Only perform QC for meffil pipeline, caches results into rds file for loading again, only works if use\_cache is false.

**-c, --chunk\_size <chunk\_size>**

If not series, chunk up and run these number of commands at once.. -1 means all commands at once.

## 7.2 combine\_methylation\_arrays

If split MethylationArrays by subtype for either preprocessing or imputation, can use to recombine data for downstream step.

```
pymethyl-preprocess combine_methylation_arrays [OPTIONS]
```

### Options

**-i, --input\_pkls <input\_pkls>**

Input pickles for beta and phenotype data. [default: ./preprocess\_outputs/methyl\_array.pkl]

**-d, --optional\_input\_pkl\_dir <optional\_input\_pkl\_dir>**

Auto grab input pkls. [default: ]

**-o, --output\_pkl <output\_pkl>**

Output database for beta and phenotype data. [default: ./combined\_outputs/methyl\_array.pkl]

**-e, --exclude <exclude>**

If -d selected, these diseases will be excluded from study. [default: ]

## 7.3 concat\_sample\_sheets

Concat two sample files for more fields for minfi+ input, adds more samples.

```
pymethyl-preprocess concat_sample_sheets [OPTIONS]
```

### Options

**-s1, --sample\_sheet1 <sample\_sheet1>**

Clinical information downloaded from tcga/geo/custom, formatted using create\_sample\_sheet. [default: ./tcga\_idats/clinical\_info1.csv]

**-s2, --sample\_sheet2 <sample\_sheet2>**

Clinical information downloaded from tcga/geo/custom, formatted using create\_sample\_sheet. [default: ./tcga\_idats/clinical\_info2.csv]

**-os, --output\_sample\_sheet <output\_sample\_sheet>**

CSV for minfi input. [default: ./tcga\_idats/minfiSheet.csv]

## 7.4 create\_sample\_sheet

Create sample sheet for input to minfi, meffil, or enmix.

```
pymethyl-preprocess create_sample_sheet [OPTIONS]
```

### Options

- is, --input\_sample\_sheet** <input\_sample\_sheet>  
Clinical information downloaded from tcga/geo/custom. [default: ./tcga\_idats/clinical\_info.csv]
- s, --source\_type** <source\_type>  
Source type of data. [default: tcga]
- i, --idat\_dir** <idat\_dir>  
Idat directory. [default: ./tcga\_idats/]
- os, --output\_sample\_sheet** <output\_sample\_sheet>  
CSV for minfi input. [default: ./tcga\_idats/minfiSheet.csv]
- m, --mapping\_file** <mapping\_file>  
Mapping file from uuid to TCGA barcode. Downloaded using download\_tcga. [default: ./idat\_filename\_case.txt]
- l, --header\_line** <header\_line>  
Line to begin reading csv/xlsx. [default: 0]
- d, --disease\_class\_column** <disease\_class\_column>  
Disease classification column, for custom and geo datasets. [default: methylation class:ch1]
- b, --basename\_col** <basename\_col>  
Basename classification column, for custom datasets. [default: Sentrix ID (.idat)]
- c, --include\_columns\_file** <include\_columns\_file>  
Custom columns file containing columns to keep, separated by n. Add a tab for each line if you wish to rename columns: original\_name t new\_column\_name [default: ]

## 7.5 download\_clinical

Download all TCGA 450k clinical info.

```
pymethyl-preprocess download_clinical [OPTIONS]
```

### Options

- o, --output\_dir** <output\_dir>  
Output directory for exported idats. [default: ./tcga\_idats/]

## 7.6 download\_geo

Download geo methylation study idats and clinical info.

```
pymethyl-preprocess download_geo [OPTIONS]
```

## Options

- g, --geo\_query** <geo\_query>  
GEO study to query. [default: ]
- o, --output\_dir** <output\_dir>  
Output directory for exported idats. [default: ./geo\_idats/]

## 7.7 download\_tcga

Download all tcga 450k data.

```
pymethyl-preprocess download_tcga [OPTIONS]
```

## Options

- o, --output\_dir** <output\_dir>  
Output directory for exported idats. [default: ./tcga\_idats/]

## 7.8 feature\_select

Filter CpGs by taking x top CpGs with highest mean absolute deviation scores or via spectral feature selection.

```
pymethyl-preprocess feature_select [OPTIONS]
```

## Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./imputed\_outputs/methyl\_array.pkl]
- o, --output\_pkl** <output\_pkl>  
Output database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]
- n, --n\_top\_cpgs** <n\_top\_cpgs>  
Number cpgs to include with highest variance across population. [default: 300000]
- f, --feature\_selection\_method** <feature\_selection\_method>
- mm, --metric** <metric>
- nn, --n\_neighbors** <n\_neighbors>  
Number neighbors for feature selection, default enacts rbf kernel. [default: 0]
- m, --mad\_top\_cpgs** <mad\_top\_cpgs>  
Number cpgs to apply mad filtering first before more sophisticated feature selection. If 0 or primary feature selection is mad, no mad pre-filtering. [default: 0]

## 7.9 get\_categorical\_distribution

Get categorical distribution of columns of sample sheet.

```
pymethyl-preprocess get_categorical_distribution [OPTIONS]
```

### Options

- is, --formatted\_sample\_sheet** <formatted\_sample\_sheet>  
Clinical information downloaded from tcga/geo/custom, formatted using create\_sample\_sheet. [default: ./tcga\_idats/minfiSheet.csv]
- k, --key** <key>  
Column of csv to print statistics for. [default: disease]
- d, --disease\_only**  
Only look at disease, or text before subtype\_delimiter.
- sd, --subtype\_delimiter** <subtype\_delimiter>  
Delimiter for disease extraction. [default: ,]

## 7.10 imputation\_pipeline

Imputation of subtype or no subtype using various imputation methods.

```
pymethyl-preprocess imputation_pipeline [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./combined\_outputs/methyl\_array.pkl]
- ss, --split\_by\_subtype**  
Imputes CpGs by subtype before combining again.
- m, --method** <method>  
Method of imputation. [default: KNN]
- s, --solver** <solver>  
Imputation library. [default: fancyimpute]
- k, --n\_neighbors** <n\_neighbors>  
Number neighbors for imputation if using KNN. [default: 5]
- r, --orientation** <orientation>  
Impute CpGs or samples. [default: Samples]
- o, --output\_pkl** <output\_pkl>  
Output database for beta and phenotype data. [default: ./imputed\_outputs/methyl\_array.pkl]
- n, --n\_top\_cpgs** <n\_top\_cpgs>  
Number cpGs to include with highest variance across population. Greater than 0 allows for mad filtering during imputation to skip mad step. [default: 0]
- f, --feature\_selection\_method** <feature\_selection\_method>
- mm, --metric** <metric>
- nfs, --n\_neighbors\_fs** <n\_neighbors\_fs>  
Number neighbors for feature selection, default enacts rbf kernel. [default: 0]

- d, --disease\_only**  
Only look at disease, or text before subtype\_delimiter.
- sd, --subtype\_delimiter** <subtype\_delimiter>  
Delimiter for disease extraction. [default: ,]
- st, --sample\_threshold** <sample\_threshold>  
Value between 0 and 1 for NaN removal. If samples has sample\_threshold proportion of cpgs missing, then remove sample. Set to -1 to not remove samples. [default: -1.0]
- ct, --cpg\_threshold** <cpg\_threshold>  
Value between 0 and 1 for NaN removal. If cpgs has cpg\_threshold proportion of samples missing, then remove cpg. Set to -1 to not remove samples. [default: -1.0]

## 7.11 meffil\_encode

Reformat file for meffil input.

```
pymethyl-preprocess meffil_encode [OPTIONS]
```

### Options

- is, --input\_sample\_sheet** <input\_sample\_sheet>  
CSV for minfi input. [default: ./tcga\_idats/minfiSheet.csv]
- os, --output\_sample\_sheet** <output\_sample\_sheet>  
CSV for minfi input. [default: ./tcga\_idats/minfiSheet.csv]

## 7.12 merge\_sample\_sheets

Merge two sample files for more fields for minfi+ input.

```
pymethyl-preprocess merge_sample_sheets [OPTIONS]
```

### Options

- s1, --sample\_sheet1** <sample\_sheet1>  
Clinical information downloaded from tcga/geo/custom, formatted using create\_sample\_sheet. [default: ./tcga\_idats/clinical\_info1.csv]
- s2, --sample\_sheet2** <sample\_sheet2>  
Clinical information downloaded from tcga/geo/custom, formatted using create\_sample\_sheet. [default: ./tcga\_idats/clinical\_info2.csv]
- os, --output\_sample\_sheet** <output\_sample\_sheet>  
CSV for minfi input. [default: ./tcga\_idats/minfiSheet.csv]
- d, --second\_sheet\_disease**  
Use second sheet's disease column.
- nd, --no\_disease\_merge**  
Don't merge disease columns.



## 7.13 na\_report

Print proportion of missing values throughout dataset.

```
pymethyl-preprocess na_report [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./preprocess\_outputs/methyl\_array.pkl]
- o, --output\_dir** <output\_dir>  
Output database for na report. [default: ./na\_report/]
- r, --head\_directory**  
-i option becomes directory, and searches there for multiple input pickles.

## 7.14 preprocess\_pipeline

Perform preprocessing of idats using enmix or meffil.

```
pymethyl-preprocess preprocess_pipeline [OPTIONS]
```

### Options

- i, --idat\_dir** <idat\_dir>  
Idat dir for one sample sheet, alternatively can be your phenotype sample sheet. [default: ./tcga\_idats/]
- n, --n\_cores** <n\_cores>  
Number cores to use for preprocessing. [default: 6]
- o, --output\_pkl** <output\_pkl>  
Output database for beta and phenotype data. [default: ./preprocess\_outputs/methyl\_array.pkl]
- m, --meffil**  
Preprocess using meffil.
- pc, --n\_pcs** <n\_pcs>  
For meffil, number of principal components for functional normalization. If set to -1, then PCs are selected using elbow method. [default: -1]
- p, --pipeline** <pipeline>  
If not meffil, preprocess using minfi or enmix. [default: enmix]
- noob, --noob\_norm**  
Run noob normalization of minfi selected.
- u, --use\_cache**  
If this is selected, loads qc results rather than running qc again and update with new qc parameters. Only works for meffil selection. Minfi and enmix just loads RG Set.
- qc, --qc\_only**  
Only perform QC for meffil pipeline, caches results into rds file for loading again, only works if use\_cache is false. Minfi and enmix just saves the RGSet before preprocessing.

- bns, --p\_beadnum\_samples** <p\_beadnum\_samples>  
From meffil documentation, “fraction of probes that failed the threshold of 3 beads”. [default: 0.05]
- pds, --p\_detection\_samples** <p\_detection\_samples>  
From meffil documentation, “fraction of probes that failed a detection.pvalue threshold of 0.01”. [default: 0.05]
- bnc, --p\_beadnum\_cpgs** <p\_beadnum\_cpgs>  
From meffil documentation, “fraction of samples that failed the threshold of 3 beads”. [default: 0.05]
- pdc, --p\_detection\_cpgs** <p\_detection\_cpgs>  
From meffil documentation, “fraction of samples that failed a detection.pvalue threshold of 0.01”. [default: 0.05]
- sc, --sex\_cutoff** <sex\_cutoff>  
From meffil documentation, “difference of total median intensity for Y chromosome probes and X chromosome probes”. [default: -2]
- sd, --sex\_sd** <sex\_sd>  
From meffil documentation, “sex detection outliers if outside this range”. [default: 5]

## 7.15 remove\_diseases

Exclude diseases from study by count number or exclusion list.

```
pymethyl-preprocess remove_diseases [OPTIONS]
```

### Options

- is, --formatted\_sample\_sheet** <formatted\_sample\_sheet>  
Clinical information downloaded from tcga/geo/custom, formatted using create\_sample\_sheet. [default: ./tcga\_idats/clinical\_info.csv]
- e, --exclude\_disease\_list** <exclude\_disease\_list>  
List of conditions to exclude, from disease column, comma delimited. [default: ]
- os, --output\_sheet\_name** <output\_sheet\_name>  
CSV for minfi input. [default: ./tcga\_idats/minfiSheet.csv]
- l, --low\_count** <low\_count>  
Remove diseases if they are below a certain count, default this is not used. [default: 0]
- d, --disease\_only**  
Only look at disease, or text before subtype\_delimiter.
- sd, --subtype\_delimiter** <subtype\_delimiter>  
Delimiter for disease extraction. [default: ,]

## 7.16 split\_preprocess\_input\_by\_subtype

Split preprocess input samplesheet by disease subtype.

```
pymethyl-preprocess split_preprocess_input_by_subtype [OPTIONS]
```

## Options

- i, --idat\_csv** <idat\_csv>  
Idat csv for one sample sheet, alternatively can be your phenotype sample sheet. [default: ./tcga\_idats/minfiSheet.csv]
- d, --disease\_only**  
Only look at disease, or text before subtype\_delimiter.
- sd, --subtype\_delimiter** <subtype\_delimiter>  
Delimiter for disease extraction. [default: ,]
- o, --subtype\_output\_dir** <subtype\_output\_dir>  
Output subtypes pheno csv. [default: ./preprocess\_outputs/]



## PYMETHYL-UTILS

```
pymethyl-utils [OPTIONS] COMMAND [ARGS]...
```

### Options

**--version**  
Show the version and exit.

## 8.1 backup\_pkl

Copy methylarray pickle to new location to backup.

```
pymethyl-utils backup_pkl [OPTIONS]
```

### Options

**-i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]

**-o, --output\_pkl** <output\_pkl>  
Output database for beta and phenotype data. [default: ./backup/methyl\_array.pkl]

## 8.2 bin\_column

Convert continuous phenotype column into categorical by binning.

```
pymethyl-utils bin_column [OPTIONS]
```

### Options

**-t, --test\_pkl** <test\_pkl>  
Pickle containing testing set. [default: ./train\_val\_test\_sets/test\_methyl\_array.pkl]

**-c, --col** <col>  
Column to turn into bins. [default: age]

**-n, --n\_bins** <n\_bins>  
Number of bins. [default: 10]

**-ot, --output\_test\_pkl** <output\_test\_pkl>  
Binned shap pickle for further testing. [default: ./train\_val\_test\_sets/test\_methyl\_array\_shap\_binned.pkl]

## 8.3 concat\_csv

Concatenate two csv files together.

```
pymethyl-utils concat_csv [OPTIONS]
```

### Options

**-i1, --input\_csv** <input\_csv>  
Beta csv. [default: ./beta1.csv]

**-i2, --input\_csv2** <input\_csv2>  
Beta/other csv 2. [default: ./cell\_estimates.csv]

**-o, --output\_csv** <output\_csv>  
Output csv. [default: ./beta.concat.csv]

**-a, --axis** <axis>  
Axis to merge on. Columns are 0, rows are 1. [default: 1]

## 8.4 counts

Return categorical breakdown of phenotype column.

```
pymethyl-utils counts [OPTIONS]
```

### Options

**-i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]

**-k, --key** <key>  
Key to split on. [default: disease]

## 8.5 create\_external\_validation\_set

Create external validation set containing same CpGs as training set.

```
pymethyl-utils create_external_validation_set [OPTIONS]
```

## Options

- t, --train\_pkl** <train\_pkl>  
Input methyl array. [default: ./train\_val\_test\_sets/train\_methyl\_array.pkl]
- q, --query\_pkl** <query\_pkl>  
Input methylation array to add/subtract cpGs to. [default: ./final\_preprocessed/methyl\_array.pkl]
- o, --output\_pkl** <output\_pkl>  
Output methyl array external validation. [default: ./external\_validation/methyl\_array.pkl]
- c, --cpg\_replace\_method** <cpg\_replace\_method>  
What to do for missing CpGs. [default: mid]

## 8.6 feature\_select\_train\_val\_test

Filter CpGs by taking x top CpGs with highest mean absolute deviation scores or via spectral feature selection.

```
pymethyl-utils feature_select_train_val_test [OPTIONS]
```

## Options

- i, --input\_pkl\_dir** <input\_pkl\_dir>  
Input database for beta and phenotype data. [default: ./train\_val\_test\_sets/]
- o, --output\_dir** <output\_dir>  
Output database for beta and phenotype data. [default: ./train\_val\_test\_sets\_fs/]
- n, --n\_top\_cpGs** <n\_top\_cpGs>  
Number cpGs to include with highest variance across population. [default: 300000]
- f, --feature\_selection\_method** <feature\_selection\_method>
- mm, --metric** <metric>
- nn, --n\_neighbors** <n\_neighbors>  
Number neighbors for feature selection, default enacts rbf kernel. [default: 0]
- m, --mad\_top\_cpGs** <mad\_top\_cpGs>  
Number cpGs to apply mad filtering first before more sophisticated feature selection. If 0 or primary feature selection is mad, no mad pre-filtering. [default: 0]

## 8.7 fix\_key

Format certain column of phenotype array in MethylationArray.

```
pymethyl-utils fix_key [OPTIONS]
```

## Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]

**-k, --key <key>**  
Key to split on. [default: disease]

**-d, --disease\_only**  
Only look at disease, or text before subtype\_delimiter.

**-sd, --subtype\_delimiter <subtype\_delimiter>**  
Delimiter for disease extraction. [default: ,]

**-o, --output\_pkl <output\_pkl>**  
Input database for beta and phenotype data. [default: ./fixed\_preprocessed/methyl\_array.pkl]

## 8.8 modify\_pheno\_data

Use another spreadsheet to add more descriptive data to methylarray.

```
pymethyl-utils modify_pheno_data [OPTIONS]
```

### Options

**-i, --input\_pkl <input\_pkl>**  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]

**-is, --input\_formatted\_sample\_sheet <input\_formatted\_sample\_sheet>**  
Information passed through function create\_sample\_sheet, has Basename and disease fields. [default: ./tcga\_idats/minfi\_sheet.csv]

**-o, --output\_pkl <output\_pkl>**  
Output database for beta and phenotype data. [default: ./modified\_processed/methyl\_array.pkl]

## 8.9 move\_jpg

Move preprocessing jpegs to preprocessing output directory.

```
pymethyl-utils move_jpg [OPTIONS]
```

### Options

**-i, --input\_dir <input\_dir>**  
Directory containing jpg. [default: ./]

**-o, --output\_dir <output\_dir>**  
Output directory for images. [default: ./preprocess\_output\_images/]

## 8.10 overwrite\_pheno\_data

Use another spreadsheet to add more descriptive data to methylarray.

```
pymethyl-utils overwrite_pheno_data [OPTIONS]
```



## Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]
- is, --input\_formatted\_sample\_sheet** <input\_formatted\_sample\_sheet>  
Information passed through function create\_sample\_sheet, has Basename and disease fields. [default: ./tcga\_idats/minfi\_sheet.csv]
- o, --output\_pkl** <output\_pkl>  
Output database for beta and phenotype data. [default: ./modified\_processed/methyl\_array.pkl]
- c, --index\_col** <index\_col>  
Index col when reading csv. [default: 0]

## 8.11 pkl\_to\_csv

Output methylarray pickle to csv.

```
pymethyl-utils pkl_to_csv [OPTIONS]
```

## Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]
- o, --output\_dir** <output\_dir>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/]
- c, --col** <col>  
Column to color. [default: ]

## 8.12 print\_number\_sex\_cpgs

Print number of non-autosomal CpGs.

```
pymethyl-utils print_number_sex_cpgs [OPTIONS]
```

## Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]
- a, --array\_type** <array\_type>  
Array Type. [default: 450k]

## 8.13 print\_shape

Print dimensions of beta matrix.

```
pymethyl-utils print_shape [OPTIONS]
```

### Options

**-i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]

## 8.14 ref\_estimate\_cell\_counts

Reference based cell type estimates.

```
pymethyl-utils ref_estimate_cell_counts [OPTIONS]
```

### Options

**-ro, --input\_r\_object\_dir** <input\_r\_object\_dir>  
Input directory containing qc data. [default: ./preprocess\_outputs/]

**-a, --algorithm** <algorithm>  
Algorithm to run cell type. [default: meffil]

**-ref, --reference** <reference>  
Cell Type Reference. [default: cord blood gse68456]

**-l, --library** <library>  
IDOL Library. [default: IDOLOptimizedCpGs450klegacy]

**-o, --output\_csv** <output\_csv>  
Output cell type estimates. [default: ./added\_cell\_counts/cell\_type\_estimates.csv]

## 8.15 remove\_sex

Remove non-autosomal CpGs.

```
pymethyl-utils remove_sex [OPTIONS]
```

### Options

**-i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./preprocess\_outputs/methyl\_array.pkl]

**-o, --output\_pkl** <output\_pkl>  
Output methyl array autosomal. [default: ./autosomal/methyl\_array.pkl]

**-a, --array\_type** <array\_type>  
Array Type. [default: 450k]

## 8.16 remove\_snps

Remove SNPs from methylation array.

```
pymethyl-utils remove_snps [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./autosomal/methyl\_array.pkl]
- o, --output\_pkl** <output\_pkl>  
Output methyl array autosomal. [default: ./no\_snp/methyl\_array.pkl]
- a, --array\_type** <array\_type>  
Array Type. [default: 450k]

## 8.17 set\_part\_array\_background

Set subset of CpGs from beta matrix to background values.

```
pymethyl-utils set_part_array_background [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input methyl array. [default: ./final\_preprocessed/methyl\_array.pkl]
- c, --cpg\_pkl** <cpg\_pkl>  
Pickled numpy array for subsetting. [default: ./subset\_cpgs.pkl]
- o, --output\_pkl** <output\_pkl>  
Output methyl array external validation. [default: ./removal/methyl\_array.pkl]

## 8.18 stratify

Split methylation array by key and store.

```
pymethyl-utils stratify [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]
- k, --key** <key>  
Key to split on. [default: disease]
- o, --output\_dir** <output\_dir>  
Output directory for stratified. [default: ./stratified/]

## 8.19 subset\_array

Only retain certain number of CpGs from methylation array.

```
pymethyl-utils subset_array [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input methyl array. [default: ./final\_preprocessed/methyl\_array.pkl]
- c, --cpg\_pkl** <cpg\_pkl>  
Pickled numpy array for subsetting. [default: ./subset\_cpgs.pkl]
- o, --output\_pkl** <output\_pkl>  
Output methyl array external validation. [default: ./subset/methyl\_array.pkl]

## 8.20 train\_test\_val\_split

Split methylation array into train, test, val.

```
pymethyl-utils train_test_val_split [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input database for beta and phenotype data. [default: ./final\_preprocessed/methyl\_array.pkl]
- o, --output\_dir** <output\_dir>  
Output directory for training, testing, and validation sets. [default: ./train\_val\_test\_sets/]
- tp, --train\_percent** <train\_percent>  
Percent data training on. [default: 0.8]
- vp, --val\_percent** <val\_percent>  
Percent of training data that comprises validation set. [default: 0.1]
- cat, --categorical**  
Multi-class prediction. [default: False]
- do, --disease\_only**  
Only look at disease, or text before subtype\_delimiter.
- k, --key** <key>  
Key to split on. [default: disease]
- sd, --subtype\_delimiter** <subtype\_delimiter>  
Delimiter for disease extraction. [default: .]

## 8.21 write\_cpgs

Write CpGs in methylation array to file.

```
pymethyl-utils write_cpgs [OPTIONS]
```

### Options

- i, --input\_pkl** <input\_pkl>  
Input methyl array. [default: ./final\_preprocessed/methyl\_array.pkl]
- c, --cpg\_pkl** <cpg\_pkl>  
Pickled numpy array for subsetting. [default: ./subset\_cpgs.pkl]



## INDICES AND TABLES

- `genindex`
- `modindex`
- `search`





## PYTHON MODULE INDEX

### p

`pymethylprocess.general_machine_learning`,  
13  
`pymethylprocess.meffil_functions`, 12  
`pymethylprocess.MethylationDataTypes`, 8  
`pymethylprocess.PreProcessDataTypes`, 3



## INDEX

### Symbols

- version
  - pymethyl-install command line option, [17](#)
  - pymethyl-preprocess command line option, [23](#)
  - pymethyl-utils command line option, [33](#)
  - pymethyl-visualize command line option, [19](#)
- a, -algorithm <algorithm>
  - pymethyl-utils-ref\_estimate\_cell\_counts command line option, [38](#)
- a, -annot
  - pymethyl-visualize-plot\_heatmap command line option, [20](#)
- a, -array\_type <array\_type>
  - pymethyl-utils-print\_number\_sex\_cpgs command line option, [37](#)
  - pymethyl-utils-remove\_sex command line option, [38](#)
  - pymethyl-utils-remove\_snps command line option, [39](#)
- a, -axes\_off
  - pymethyl-visualize-transform\_plot command line option, [21](#)
- a, -axis <axis>
  - pymethyl-utils-concat\_csv command line option, [34](#)
- b, -basename\_col <basename\_col>
  - pymethyl-preprocess-create\_sample\_sheet command line option, [25](#)
- bnc, -p\_beadnum\_cpgs <p\_beadnum\_cpgs>
  - pymethyl-preprocess-preprocess\_pipeline command line option, [30](#)
- bns, -p\_beadnum\_samples <p\_beadnum\_samples>
  - pymethyl-preprocess-preprocess\_pipeline command line option, [29](#)
- c, -chunk\_size <chunk\_size>
  - pymethyl-preprocess-batch\_deploy\_preprocess command line option, [24](#)
- c, -cluster
  - pymethyl-visualize-plot\_heatmap command line option, [20](#)
- c, -col <col>
  - pymethyl-utils-bin\_column command line option, [33](#)
  - pymethyl-utils-pkl\_to\_csv command line option, [37](#)
- c, -column\_of\_interest <column\_of\_interest>
  - pymethyl-visualize-transform\_plot command line option, [20](#)
- c, -cpg\_pkl <cpg\_pkl>
  - pymethyl-utils-set\_part\_array\_background command line option, [39](#)
  - pymethyl-utils-subset\_array command line option, [40](#)
  - pymethyl-utils-write\_cpgs command line option, [41](#)
- c, -cpg\_replace\_method <cpg\_replace\_method>
  - pymethyl-utils-create\_external\_validation\_set command line option, [35](#)
- c, -include\_columns\_file <include\_columns\_file>
  - pymethyl-preprocess-create\_sample\_sheet command line option, [25](#)
- c, -index\_col <index\_col>
  - pymethyl-utils-overwrite\_pheno\_data command line option, [37](#)
- cat, -categorical
  - pymethyl-utils-train\_test\_val\_split command line option, [40](#)
- cc, -case\_control\_override
  - pymethyl-visualize-transform\_plot command line option, [21](#)
- col, -color\_column <color\_column>
  - pymethyl-visualize-plot\_heatmap command line option, [20](#)
- cols, -plot\_cols <plot\_cols>
  - pymethyl-visualize-plot\_cell\_type\_results command line option, [19](#)

```
-ct, -cpg_threshold <cpg_threshold>      -i, -idat_csv <idat_csv>
    pymethyl-preprocess-imputation_pipeline pymethyl-preprocess-split_preprocess_input_by_s
        command line option,28                command line option,31
-d, -disease_class_column                 -i, -idat_dir <idat_dir>
    <disease_class_column>                 pymethyl-preprocess-create_sample_sheet
    pymethyl-preprocess-create_sample_sheet command line option,25
        command line option,25                pymethyl-preprocess-preprocess_pipeline
-d, -disease_only                         command line option,29
    pymethyl-preprocess-get_categorical_distribution_csv <input_csv>
        command line option,27                pymethyl-visualize-plot_cell_type_results
    pymethyl-preprocess-imputation_pipeline command line option,19
        command line option,27                pymethyl-visualize-plot_heatmap
    pymethyl-preprocess-remove_diseases    command line option,20
        command line option,30                -i, -input_dir <input_dir>
    pymethyl-preprocess-split_preprocess_input_by_cell_type move_jpg command
        command line option,31                line option,36
    pymethyl-utils-fix_key command          -i, -input_pkl <input_pkl>
        line option,36                        pymethyl-preprocess-feature_select
-d, -min_dist <min_dist>                  command line option,26
    pymethyl-visualize-transform_plot       pymethyl-preprocess-imputation_pipeline
        command line option,21                command line option,27
-d, -optional_input_pkl_dir               pymethyl-preprocess-na_report
    <optional_input_pkl_dir>                command line option,29
    pymethyl-preprocess-combine_methylation_pymethyl-utils-backup_pkl command
        command line option,24                line option,33
-d, -second_sheet_disease                 pymethyl-utils-counts command line
    pymethyl-preprocess-merge_sample_sheets option,34
        command line option,28                pymethyl-utils-fix_key command
-d, -do, -disease_only                    line option,35
    pymethyl-utils-train_test_val_split     pymethyl-utils-modify_pheno_data
        command line option,40                command line option,36
-e, -exclude <exclude>                    pymethyl-utils-overwrite_pheno_data
    pymethyl-preprocess-combine_methylation_array command line option,37
        command line option,24                pymethyl-utils-pkl_to_csv command
-e, -exclude_disease_list                  line option,37
    <exclude_disease_list>                 pymethyl-utils-print_number_sex_cpgs
    pymethyl-preprocess-remove_diseases    command line option,37
        command line option,30                pymethyl-utils-print_shape command
-f, -feature_selection_method              line option,38
    <feature_selection_method>             pymethyl-utils-remove_sex command
    pymethyl-preprocess-feature_select     line option,38
        command line option,26                pymethyl-utils-remove_snps command
    pymethyl-preprocess-imputation_pipeline line option,39
        command line option,27                pymethyl-utils-set_part_array_background
    pymethyl-utils-feature_select_train_val_test command line option,39
        command line option,35                pymethyl-utils-stratify command
-fs, -font_scale <font_scale>              line option,39
    pymethyl-visualize-plot_cell_type_results pymethyl-utils-subset_array
        command line option,19                command line option,40
    pymethyl-visualize-plot_heatmap        pymethyl-utils-train_test_val_split
        command line option,20                command line option,40
-g, -geo_query <geo_query>                pymethyl-utils-write_cpgs command
    pymethyl-preprocess-download_geo        line option,41
        command line option,26                pymethyl-visualize-transform_plot
```

command line option, <a href="#">20</a>	command line option, <a href="#">25</a>
-i, -input_pkl_dir <input_pkl_dir>	-l, -library <library>
pymethyl-utils-feature_select_train_val_test	pymethyl-utils-ref_estimate_cell_counts
command line option, <a href="#">35</a>	command line option, <a href="#">38</a>
-i, -input_pkls <input_pkls>	-l, -low_count <low_count>
pymethyl-preprocess-combine_methylation	pymethyl-preprocess-remove_diseases
command line option, <a href="#">24</a>	command line option, <a href="#">30</a>
-i, -subtype_output_dir	-m, -mad_top_cpgs <mad_top_cpgs>
<subtype_output_dir>	pymethyl-preprocess-feature_select
pymethyl-preprocess-batch_deploy_preprocess	command line option, <a href="#">26</a>
command line option, <a href="#">23</a>	pymethyl-utils-feature_select_train_val_test
-il, -input_csv <input_csv>	command line option, <a href="#">35</a>
pymethyl-utils-concat_csv	-m, -manager
command line option, <a href="#">34</a>	pymethyl-install-install_custom
-i2, -input_csv2 <input_csv2>	command line option, <a href="#">17</a>
pymethyl-utils-concat_csv	-m, -mapping_file <mapping_file>
command line option, <a href="#">34</a>	pymethyl-preprocess-create_sample_sheet
-idx, -index_col <index_col>	command line option, <a href="#">25</a>
pymethyl-visualize-plot_heatmap	-m, -matrix_type <matrix_type>
command line option, <a href="#">20</a>	pymethyl-visualize-plot_heatmap
-is, -formatted_sample_sheet	command line option, <a href="#">20</a>
<formatted_sample_sheet>	-m, -meffil
pymethyl-preprocess-get_categorical_distribution	pymethyl-preprocess-batch_deploy_preprocess
command line option, <a href="#">27</a>	command line option, <a href="#">23</a>
pymethyl-preprocess-remove_diseases	pymethyl-preprocess-preprocess_pipeline
command line option, <a href="#">30</a>	command line option, <a href="#">29</a>
-is, -input_formatted_sample_sheet	-m, -method <method>
<input_formatted_sample_sheet>	pymethyl-preprocess-imputation_pipeline
pymethyl-utils-modify_pheno_data	command line option, <a href="#">27</a>
command line option, <a href="#">36</a>	-m, -metric <metric>
pymethyl-utils-overwrite_pheno_data	pymethyl-visualize-transform_plot
command line option, <a href="#">37</a>	command line option, <a href="#">21</a>
-is, -input_sample_sheet	-max, -max_val <max_val>
<input_sample_sheet>	pymethyl-visualize-plot_heatmap
pymethyl-preprocess-create_sample_sheet	command line option, <a href="#">20</a>
command line option, <a href="#">25</a>	-min, -min_val <min_val>
pymethyl-preprocess-meffil_encode	pymethyl-visualize-plot_heatmap
command line option, <a href="#">28</a>	command line option, <a href="#">20</a>
-k, -key <key>	-mm, -metric <metric>
pymethyl-preprocess-get_categorical_distribution	pymethyl-preprocess-feature_select
command line option, <a href="#">27</a>	command line option, <a href="#">26</a>
pymethyl-utils-counts	pymethyl-preprocess-imputation_pipeline
command line option, <a href="#">34</a>	command line option, <a href="#">27</a>
pymethyl-utils-fix_key	pymethyl-utils-feature_select_train_val_test
command line option, <a href="#">35</a>	command line option, <a href="#">35</a>
pymethyl-utils-stratify	-n, -n_bins <n_bins>
command line option, <a href="#">39</a>	pymethyl-utils-bin_column
pymethyl-utils-train_test_val_split	command line option, <a href="#">33</a>
command line option, <a href="#">40</a>	-n, -n_cores <n_cores>
-k, -n_neighbors <n_neighbors>	pymethyl-preprocess-batch_deploy_preprocess
pymethyl-preprocess-imputation_pipeline	command line option, <a href="#">23</a>
command line option, <a href="#">27</a>	pymethyl-preprocess-preprocess_pipeline
-l, -header_line <header_line>	command line option, <a href="#">29</a>
pymethyl-preprocess-create_sample_sheet	-n_top_cpgs <n_top_cpgs>

```
pymethyl-preprocess-feature_select -o, -output_file <output_file>
    command line option,26 pymethyl-visualize-transform_plot
pymethyl-preprocess-imputation_pipeline command line option,20
    command line option,27 -o, -output_pkl <output_pkl>
pymethyl-utils-feature_select_train_val_test pymethyl-preprocess-combine_methylation_arrays
    command line option,35 command line option,24
-n, -norm pymethyl-preprocess-feature_select
    pymethyl-visualize-plot_heatmap command line option,26
    command line option,20 pymethyl-preprocess-imputation_pipeline
    command line option,27
-nd, -no_disease_merge pymethyl-preprocess-merge_sample_sheets pymethyl-preprocess-preprocess_pipeline
    command line option,28 command line option,29
-nfs, -n_neighbors_fs <n_neighbors_fs> pymethyl-utils-backup_pkl command
    pymethyl-preprocess-imputation_pipeline line option,33
    command line option,27 pymethyl-utils-create_external_validation_set
    command line option,35
-nn, -n_neighbors <n_neighbors> pymethyl-utils-fix_key command
    pymethyl-preprocess-feature_select line option,36
    command line option,26 pymethyl-utils-modify_pheno_data
    pymethyl-utils-feature_select_train_val_test command line option,36
    command line option,35 pymethyl-utils-overwrite_pheno_data
    pymethyl-visualize-transform_plot command line option,37
    command line option,21 pymethyl-utils-remove_sex command
    line option,38
-noob, -noob_norm pymethyl-preprocess-preprocess_pipeline
    command line option,29 pymethyl-utils-remove_snps command
    line option,39
-o, -outfilename <outfilename> pymethyl-utils-set_part_array_background
    pymethyl-visualize-plot_cell_type_result command line option,39
    command line option,19 pymethyl-utils-subset_array
    command line option,20 command line option,40
-o, -output_csv <output_csv> -o, -subtype_output_dir
    pymethyl-utils-concat_csv command <subtype_output_dir>
    line option,34 pymethyl-preprocess-split_preprocess_input_by_s
    command line option,31
pymethyl-utils-ref_estimate_cell_counts -os, -output_sample_sheet
    command line option,38 <output_sample_sheet>
-o, -output_dir <output_dir> pymethyl-preprocess-concat_sample_sheets
    command line option,25 command line option,24
pymethyl-preprocess-download_clinical pymethyl-preprocess-create_sample_sheet
    command line option,26 command line option,25
pymethyl-preprocess-download_geo pymethyl-preprocess-meffil_encode
    command line option,26 command line option,28
pymethyl-preprocess-na_report pymethyl-preprocess-merge_sample_sheets
    command line option,29 command line option,28
pymethyl-utils-feature_select_train_val_test -o, -output_sheet_name
    command line option,35 <output_sheet_name>
pymethyl-utils-move_jpg command pymethyl-preprocess-remove_diseases
    line option,36 command line option,30
pymethyl-utils-pkl_to_csv command -ot, -output_test_pkl
    line option,37 <output_test_pkl>
pymethyl-utils-stratify command pymethyl-utils-bin_column command
    line option,39 line option,34
pymethyl-utils-train_test_val_split -p, -package <package>
    command line option,40 pymethyl-install-install_custom
```

---

```

    command line option, 17
pymethyl-install-install_r_packages
    command line option, 18
-p, -pc_qc_parameters_csv
    <pc_qc_parameters_csv>
pymethyl-preprocess-batch_deploy_preprocess
    command line option, 23
-p, -pipeline <pipeline>
pymethyl-preprocess-preprocess_pipeline2, -sample_sheet2 <sample_sheet2>
    command line option, 29
-pc, -n_pcs <n_pcs>
pymethyl-preprocess-preprocess_pipeline
    command line option, 29
-pdc, -p_detection_cpgs
    <p_detection_cpgs>
pymethyl-preprocess-preprocess_pipeline
    command line option, 30
-pds, -p_detection_samples
    <p_detection_samples>
pymethyl-preprocess-preprocess_pipeline3, -sd, -sex_sd <sex_sd>
    command line option, 30
-q, -query_pkl <query_pkl>
pymethyl-utils-create_external_validation_set
    command line option, 35
-qc, -qc_only
pymethyl-preprocess-batch_deploy_preprocess4
    command line option, 23
pymethyl-preprocess-preprocess_pipeline
    command line option, 29
-r, -head_directory
pymethyl-preprocess-na_report
    command line option, 29
-r, -orientation <orientation>
pymethyl-preprocess-imputation_pipeline5, -split_by_subtype
    command line option, 27
-r, -run
pymethyl-preprocess-batch_deploy_preprocess5
    command line option, 23
-ref, -reference <reference>
pymethyl-utils-ref_estimate_cell_counts
    command line option, 38
-ro, -input_r_object_dir
    <input_r_object_dir>
pymethyl-utils-ref_estimate_cell_counts6, -torque
    command line option, 38
-s, -series
pymethyl-preprocess-batch_deploy_preprocess6
    command line option, 23
-s, -solver <solver>
pymethyl-preprocess-imputation_pipeline6, -transpose
    command line option, 27
-s, -source_type <source_type>
pymethyl-preprocess-create_sample_sheet7, -train_percent <train_percent>
    command line option, 25
-s, -supervised
pymethyl-visualize-transform_plot
    command line option, 21
-sl, -sample_sheet1 <sample_sheet1>
pymethyl-preprocess-concat_sample_sheets
pymethyl-preprocess-merge_sample_sheets
    command line option, 28
-s2, -sample_sheet2 <sample_sheet2>
pymethyl-preprocess-concat_sample_sheets
    command line option, 24
pymethyl-preprocess-merge_sample_sheets
    command line option, 28
-sc, -sex_cutoff <sex_cutoff>
pymethyl-preprocess-preprocess_pipeline
    command line option, 30
-sd, -sex_sd <sex_sd>
pymethyl-preprocess-preprocess_pipeline
    command line option, 30
-sd, -subtype_delimiter
    <subtype_delimiter>
pymethyl-preprocess-get_categorical_distribution
    command line option, 27
pymethyl-preprocess-imputation_pipeline
    command line option, 28
pymethyl-preprocess-remove_diseases
    command line option, 30
pymethyl-preprocess-split_preprocess_input_by_subtype
    command line option, 31
pymethyl-utils-fix_key
    command line option, 36
pymethyl-utils-train_test_val_split
    command line option, 40
-s, -split_by_subtype
pymethyl-preprocess-imputation_pipeline
    command line option, 27
-sample_threshold
    <sample_threshold>
pymethyl-preprocess-imputation_pipeline
    command line option, 28
-t, -test_pkl <test_pkl>
pymethyl-utils-bin_column
    command line option, 33
-torque
pymethyl-preprocess-batch_deploy_preprocess7
    command line option, 23
-train_pkl <train_pkl>
pymethyl-utils-create_external_validation_set
    command line option, 35
-te, -transpose
pymethyl-visualize-plot_heatmap
    command line option, 20
-train_percent <train_percent>
pymethyl-utils-train_test_val_split

```

command line option, 40  
-u, -use\_cache  
pymethyl-preprocess-batch\_deploy\_preprocess()  
command line option, 23  
pymethyl-preprocess-preprocess\_pipeline  
command line option, 29  
-vp, -val\_percent <val\_percent>  
pymethyl-utils-train\_test\_val\_split  
command line option, 40  
-x, -xticks  
pymethyl-visualize-plot\_heatmap  
command line option, 20  
-y, -yticks  
pymethyl-visualize-plot\_heatmap  
command line option, 20

## A

assign\_results\_to\_pheno\_col() (pymethylprocess.general\_machine\_learning.MachineLearning method), 15

## B

bin\_column() (pymethylprocess.MethylationDataTypes.MethylationArray method), 9

## C

categorical\_breakdown() (pymethylprocess.MethylationDataTypes.MethylationArray method), 9

combine() (pymethylprocess.MethylationDataTypes.MethylationArrays method), 11

concat() (pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method), 7

## D

download\_clinical() (pymethylprocess.PreProcessDataTypes.TCGADownloader method), 8

download\_geo() (pymethylprocess.PreProcessDataTypes.TCGADownloader method), 8

download\_tcga() (pymethylprocess.PreProcessDataTypes.TCGADownloader method), 8

## E

est\_cell\_counts\_IDOL() (in module pymethylprocess.meffil\_functions), 13

est\_cell\_counts\_meffil() (in module pymethylprocess.meffil\_functions), 13

est\_cell\_counts\_minfi() (in module pymethylprocess.meffil\_functions), 13

export\_manifest() (pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method), 7

export\_csv() (pymethylprocess.PreProcessDataTypes.PreProcessIDAT method), 5

export\_pickle() (pymethylprocess.PreProcessDataTypes.PreProcessIDAT method), 5

export\_sql() (pymethylprocess.PreProcessDataTypes.PreProcessIDAT method), 5

extract\_manifest() (pymethylprocess.PreProcessDataTypes.PreProcessIDAT method), 5

extract\_pheno\_beta\_df\_from\_folder() (in module pymethylprocess.MethylationDataTypes), 12

extract\_pheno\_beta\_df\_from\_pickle\_dict() (in module pymethylprocess.MethylationDataTypes), 12

extract\_pheno\_beta\_df\_from\_sql() (in module pymethylprocess.MethylationDataTypes), 12

extract\_pheno\_data() (pymethylprocess.PreProcessDataTypes.PreProcessIDAT method), 5

## F

feature\_select() (pymethylprocess.MethylationDataTypes.MethylationArray method), 9

filter\_beta() (pymethylprocess.PreProcessDataTypes.PreProcessIDAT method), 5

fit() (pymethylprocess.general\_machine\_learning.MachineLearning method), 15

fit\_predict() (pymethylprocess.general\_machine\_learning.MachineLearning method), 15

fit\_transform() (pymethylprocess.general\_machine\_learning.MachineLearning method), 15

format\_custom() (pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method), 7

format\_geo() (pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method), 7

format\_tcga() (pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method), 7



`from_pickle()` (*pymethylprocess.MethylationDataTypes.MethylationArray class method*), 9

## G

`get_beta()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 5

`get_categorical_distribution()` (*pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method*), 7

`get_meth()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 5

`get_unmeth()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`groupby()` (*pymethylprocess.MethylationDataTypes.MethylationArray method*), 10

## I

`impute()` (*pymethylprocess.MethylationDataTypes.MethylationArray method*), 10

`impute()` (*pymethylprocess.MethylationDataTypes.MethylationArrays method*), 12

`ImputerObject` (class in *pymethylprocess.MethylationDataTypes*), 9

## L

`load_detection_p_values_beadnum()` (in module *pymethylprocess.meffil\_functions*), 13

`load_idats()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

## M

`MachineLearning` (class in *pymethylprocess.general\_machine\_learning*), 15

`merge()` (*pymethylprocess.PreProcessDataTypes.PreProcessPhenoData method*), 8

`merge_preprocess_sheet()` (*pymethylprocess.MethylationDataTypes.MethylationArray method*), 10

`MethylationArray` (class in *pymethylprocess.MethylationDataTypes*), 9

`MethylationArrays` (class in *pymethylprocess.MethylationDataTypes*), 11

`move_jpg()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

## O

`output_pheno_beta()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`overwrite_pheno_data()` (*pymethylprocess.MethylationDataTypes.MethylationArray method*), 10

## P

`plot_original_qc()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`plot_qc_metrics()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`predict()` (*pymethylprocess.general\_machine\_learning.MachineLearning method*), 15

`preprocess_enmix_pipeline()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`preprocessENmix()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`PreProcessIDAT` (class in *pymethylprocess.PreProcessDataTypes*), 5

`preprocessMeffil()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`preprocessNoob()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`PreProcessPhenoData` (class in *pymethylprocess.PreProcessDataTypes*), 7

`preprocessRAW()` (*pymethylprocess.PreProcessDataTypes.PreProcessIDAT method*), 6

`pymethyl-install` command line option `-version`, 17

`pymethyl-install-install_custom` command line option `-m, -manager`, 17

`-p, -package <package>`, 17

`pymethyl-install-install_r_packages` command line option

`-p, -package <package>`, 18

`pymethyl-preprocess` command line option `-version`, 23

`pymethyl-preprocess-batch_deploy_preprocess` command line option

`-c, -chunk_size <chunk_size>`, 24

`-i, -subtype_output_dir <subtype_output_dir>`, 23

```
-m, -meffil, 23
-n, -n_cores <n_cores>, 23
-p, -pc_qc_parameters_csv
    <pc_qc_parameters_csv>, 23
-qc, -qc_only, 23
-r, -run, 23
-s, -series, 23
-t, -torque, 23
-u, -use_cache, 23
pymethyl-preprocess-combine_methylation_arrays <formatted_sample_sheet>, 27
    command line option
-d, -optional_input_pkl_dir
    <optional_input_pkl_dir>, 24
-e, -exclude <exclude>, 24
-i, -input_pkls <input_pkls>, 24
-o, -output_pkl <output_pkl>, 24
pymethyl-preprocess-concat_sample_sheets
    command line option
-os, -output_sample_sheet
    <output_sample_sheet>, 24
-s1, -sample_sheet1
    <sample_sheet1>, 24
-s2, -sample_sheet2
    <sample_sheet2>, 24
pymethyl-preprocess-create_sample_sheet
    command line option
-b, -basename_col <basename_col>, 25
-c, -include_columns_file
    <include_columns_file>, 25
-d, -disease_class_column
    <disease_class_column>, 25
-i, -idat_dir <idat_dir>, 25
-is, -input_sample_sheet
    <input_sample_sheet>, 25
-l, -header_line <header_line>, 25
-m, -mapping_file <mapping_file>, 25
-os, -output_sample_sheet
    <output_sample_sheet>, 25
-s, -source_type <source_type>, 25
pymethyl-preprocess-download_clinical
    command line option
-o, -output_dir <output_dir>, 25
pymethyl-preprocess-download_geo
    command line option
-g, -geo_query <geo_query>, 26
-o, -output_dir <output_dir>, 26
pymethyl-preprocess-download_tcga
    command line option
-o, -output_dir <output_dir>, 26
pymethyl-preprocess-feature_select
    command line option
-f, -feature_selection_method
    <feature_selection_method>, 26
-i, -input_pkl <input_pkl>, 26
-m, -mad_top_cpgs <mad_top_cpgs>, 26
-mm, -metric <metric>, 26
-n, -n_top_cpgs <n_top_cpgs>, 26
-nn, -n_neighbors <n_neighbors>, 26
-o, -output_pkl <output_pkl>, 26
pymethyl-preprocess-get_categorical_distribution
    command line option
-d, -disease_only, 27
-is, -formatted_sample_sheet
-k, -key <key>, 27
-sd, -subtype_delimiter
    <subtype_delimiter>, 27
pymethyl-preprocess-imputation_pipeline
    command line option
-ct, -cpg_threshold
    <cpg_threshold>, 28
-d, -disease_only, 27
-f, -feature_selection_method
    <feature_selection_method>, 27
-i, -input_pkl <input_pkl>, 27
-k, -n_neighbors <n_neighbors>, 27
-m, -method <method>, 27
-mm, -metric <metric>, 27
-n, -n_top_cpgs <n_top_cpgs>, 27
-nfs, -n_neighbors_fs
    <n_neighbors_fs>, 27
-o, -output_pkl <output_pkl>, 27
-r, -orientation <orientation>, 27
-s, -solver <solver>, 27
-sd, -subtype_delimiter
    <subtype_delimiter>, 28
-ss, -split_by_subtype, 27
-st, -sample_threshold
    <sample_threshold>, 28
pymethyl-preprocess-meffil_encode
    command line option
-is, -input_sample_sheet
    <input_sample_sheet>, 28
-os, -output_sample_sheet
    <output_sample_sheet>, 28
pymethyl-preprocess-merge_sample_sheets
    command line option
-d, -second_sheet_disease, 28
-nd, -no_disease_merge, 28
-os, -output_sample_sheet
    <output_sample_sheet>, 28
-s1, -sample_sheet1
    <sample_sheet1>, 28
-s2, -sample_sheet2
    <sample_sheet2>, 28
pymethyl-preprocess-na_report
    command line option
-i, -input_pkl <input_pkl>, 29
```

-o, -output\_dir <output\_dir>, 29  
 -r, -head\_directory, 29  
 pymethyl-preprocess-preprocess\_pipeline command line option  
 -bnc, -p\_beadnum\_cpgs <p\_beadnum\_cpgs>, 30  
 -bns, -p\_beadnum\_samples <p\_beadnum\_samples>, 29  
 -i, -idat\_dir <idat\_dir>, 29  
 -m, -meffil, 29  
 -n, -n\_cores <n\_cores>, 29  
 -noob, -noob\_norm, 29  
 -o, -output\_pkl <output\_pkl>, 29  
 -p, -pipeline <pipeline>, 29  
 -pc, -n\_pcs <n\_pcs>, 29  
 -pdc, -p\_detection\_cpgs <p\_detection\_cpgs>, 30  
 -pds, -p\_detection\_samples <p\_detection\_samples>, 30  
 -qc, -qc\_only, 29  
 -sc, -sex\_cutoff <sex\_cutoff>, 30  
 -sd, -sex\_sd <sex\_sd>, 30  
 -u, -use\_cache, 29  
 pymethyl-preprocess-remove\_diseases command line option  
 -d, -disease\_only, 30  
 -e, -exclude\_disease\_list <exclude\_disease\_list>, 30  
 -is, -formatted\_sample\_sheet <formatted\_sample\_sheet>, 30  
 -l, -low\_count <low\_count>, 30  
 -os, -output\_sheet\_name <output\_sheet\_name>, 30  
 -sd, -subtype\_delimiter <subtype\_delimiter>, 30  
 pymethyl-preprocess-split\_preprocess\_input\_by\_subtype command line option  
 -d, -disease\_only, 31  
 -i, -idat\_csv <idat\_csv>, 31  
 -o, -subtype\_output\_dir <subtype\_output\_dir>, 31  
 -sd, -subtype\_delimiter <subtype\_delimiter>, 31  
 pymethyl-utils command line option  
 -version, 33  
 pymethyl-utils-backup\_pkl command line option  
 -i, -input\_pkl <input\_pkl>, 33  
 -o, -output\_pkl <output\_pkl>, 33  
 pymethyl-utils-bin\_column command line option  
 -c, -col <col>, 33  
 -n, -n\_bins <n\_bins>, 33  
 -ot, -output\_test\_pkl <output\_test\_pkl>, 34  
 -t, -test\_pkl <test\_pkl>, 33  
 pymethyl-utils-concat\_csv command line option  
 -a, -axis <axis>, 34  
 -i1, -input\_csv <input\_csv>, 34  
 -i2, -input\_csv2 <input\_csv2>, 34  
 -o, -output\_csv <output\_csv>, 34  
 pymethyl-utils-counts command line option  
 -i, -input\_pkl <input\_pkl>, 34  
 -k, -key <key>, 34  
 pymethyl-utils-create\_external\_validation\_set command line option  
 -c, -cpg\_replace\_method <cpg\_replace\_method>, 35  
 -o, -output\_pkl <output\_pkl>, 35  
 -q, -query\_pkl <query\_pkl>, 35  
 -t, -train\_pkl <train\_pkl>, 35  
 pymethyl-utils-feature\_select\_train\_val\_test command line option  
 -f, -feature\_selection\_method <feature\_selection\_method>, 35  
 -i, -input\_pkl\_dir <input\_pkl\_dir>, 35  
 -m, -mad\_top\_cpgs <mad\_top\_cpgs>, 35  
 -mm, -metric <metric>, 35  
 -n, -n\_top\_cpgs <n\_top\_cpgs>, 35  
 -nn, -n\_neighbors <n\_neighbors>, 35  
 -o, -output\_dir <output\_dir>, 35  
 pymethyl-utils-fix\_key command line option  
 -d, -disease\_only, 36  
 -i, -input\_pkl <input\_pkl>, 35  
 -k, -key <key>, 35  
 -o, -output\_pkl <output\_pkl>, 36  
 -sd, -subtype\_delimiter <subtype\_delimiter>, 36  
 pymethyl-utils-modify\_pheno\_data command line option  
 -i, -input\_pkl <input\_pkl>, 36  
 -is, -input\_formatted\_sample\_sheet <input\_formatted\_sample\_sheet>, 36  
 -o, -output\_pkl <output\_pkl>, 36  
 pymethyl-utils-move\_jpg command line option  
 -i, -input\_dir <input\_dir>, 36  
 -o, -output\_dir <output\_dir>, 36  
 pymethyl-utils-overwrite\_pheno\_data command line option  
 -c, -index\_col <index\_col>, 37  
 -i, -input\_pkl <input\_pkl>, 37  
 -is, -input\_formatted\_sample\_sheet

```
<input_formatted_sample_sheet>,
37
-o, -output_pkl <output_pkl>, 37
pymethyl-utils-pkl_to_csv command line
option
-c, -col <col>, 37
-i, -input_pkl <input_pkl>, 37
-o, -output_dir <output_dir>, 37
pymethyl-utils-print_number_sex_cpgs
command line option
-a, -array_type <array_type>, 37
-i, -input_pkl <input_pkl>, 37
pymethyl-utils-print_shape command
line option
-i, -input_pkl <input_pkl>, 38
pymethyl-utils-ref_estimate_cell_counts
command line option
-a, -algorithm <algorithm>, 38
-l, -library <library>, 38
-o, -output_csv <output_csv>, 38
-ref, -reference <reference>, 38
-ro, -input_r_object_dir
<input_r_object_dir>, 38
pymethyl-utils-remove_sex command line
option
-a, -array_type <array_type>, 38
-i, -input_pkl <input_pkl>, 38
-o, -output_pkl <output_pkl>, 38
pymethyl-utils-remove_snps command
line option
-a, -array_type <array_type>, 39
-i, -input_pkl <input_pkl>, 39
-o, -output_pkl <output_pkl>, 39
pymethyl-utils-set_part_array_background
command line option
-c, -cpg_pkl <cpg_pkl>, 39
-i, -input_pkl <input_pkl>, 39
-o, -output_pkl <output_pkl>, 39
pymethyl-utils-stratify command line
option
-i, -input_pkl <input_pkl>, 39
-k, -key <key>, 39
-o, -output_dir <output_dir>, 39
pymethyl-utils-subset_array command
line option
-c, -cpg_pkl <cpg_pkl>, 40
-i, -input_pkl <input_pkl>, 40
-o, -output_pkl <output_pkl>, 40
pymethyl-utils-train_test_val_split
command line option
-cat, -categorical, 40
-do, -disease_only, 40
-i, -input_pkl <input_pkl>, 40
-k, -key <key>, 40
-o, -output_dir <output_dir>, 40
-sd, -subtype_delimiter
<subtype_delimiter>, 40
-tp, -train_percent
<train_percent>, 40
-vp, -val_percent <val_percent>, 40
pymethyl-utils-write_cpgs command line
option
-c, -cpg_pkl <cpg_pkl>, 41
-i, -input_pkl <input_pkl>, 41
pymethyl-visualize command line option
-version, 19
pymethyl-visualize-plot_cell_type_results
command line option
-cols, -plot_cols <plot_cols>, 19
-fs, -font_scale <font_scale>, 19
-i, -input_csv <input_csv>, 19
-o, -outfilename <outfilename>, 19
pymethyl-visualize-plot_heatmap
command line option
-a, -annot, 20
-c, -cluster, 20
-col, -color_column <color_column>,
20
-fs, -font_scale <font_scale>, 20
-i, -input_csv <input_csv>, 20
-idx, -index_col <index_col>, 20
-m, -matrix_type <matrix_type>, 20
-max, -max_val <max_val>, 20
-min, -min_val <min_val>, 20
-n, -norm, 20
-o, -outfilename <outfilename>, 20
-t, -transpose, 20
-x, -xticks, 20
-y, -yticks, 20
pymethyl-visualize-transform_plot
command line option
-a, -axes_off, 21
-c, -column_of_interest
<column_of_interest>, 20
-cc, -case_control_override, 21
-d, -min_dist <min_dist>, 21
-i, -input_pkl <input_pkl>, 20
-m, -metric <metric>, 21
-nn, -n_neighbors <n_neighbors>, 21
-o, -output_file <output_file>, 20
-s, -supervised, 21
pymethylprocess.general_machine_learning
(module), 13
pymethylprocess.meffil_functions (mod-
ule), 12
pymethylprocess.MethylationDataTypes
(module), 8
```

`pymethylprocess.PreProcessDataTypes`  
(module), 3

## R

`r_autosomal_cpgs()` (in module `pymethylprocess.meffil_functions`), 13

`r_snp_cpgs()` (in module `pymethylprocess.meffil_functions`), 13

`remove_diseases()` (`pymethylprocess.PreProcessDataTypes.PreProcessPhenoData` method), 8

`remove_missingness()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`remove_na_samples()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`remove_sex()` (in module `pymethylprocess.meffil_functions`), 13

`remove_whitespace()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`return_beta()` (`pymethylprocess.PreProcessDataTypes.PreProcessIDAT` method), 7

`return_cpgs()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`return_idx()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`return_imputer()` (`pymethylprocess.MethylationDataTypes.ImputerObject` method), 9

`return_outcome_metric()` (`pymethylprocess.general_machine_learning.MachineLearning` method), 16

`return_raw_beta_array()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`return_shape()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

## S

`set_missing()` (in module `pymethylprocess.meffil_functions`), 13

`split_by_subtype()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`split_key()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`split_key()` (`pymethylprocess.PreProcessDataTypes.PreProcessPhenoData` method), 8

`split_train_test()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 10

`store_results()` (`pymethylprocess.general_machine_learning.MachineLearning` method), 16

`subsample()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 11

`subset_cpgs()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 11

`subset_index()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 11

## T

`TCGADownloader` (class in `pymethylprocess.PreProcessDataTypes`), 8

`to_methyl_array()` (`pymethylprocess.PreProcessDataTypes.PreProcessIDAT` method), 7

`transform()` (`pymethylprocess.general_machine_learning.MachineLearning` method), 16

`transform_results_to_beta()` (`pymethylprocess.general_machine_learning.MachineLearning` method), 16

## W

`write_csvs()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 11

`write_db()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 11

`write_dbs()` (`pymethylprocess.MethylationDataTypes.MethylationArrays` method), 12

`write_pickle()` (`pymethylprocess.MethylationDataTypes.MethylationArray` method), 11

`write_pkls()` (`pymethylprocess.MethylationDataTypes.MethylationArrays` method), 12