

d073a706b5 ▾

...

scikit-multiflow / src / skmultiflow / core / base.py / &lt;&gt; Jump to ▾



smastelini Fix PEP8 problems (#247) ... ✓

History

5 contributors



679 lines (547 sloc) 23.5 KB

...

```
1  """ Base class for all estimators in scikit-multiflow """
2
3  # Some code has been copied from sklearn for compatibility.
4  # Following is the commit information to track the last matching version with sklearn's base.py:
5  # https://github.com/scikit-learn/scikit-learn/
6  # blob/1fe00b58949a6b9bce45e9e15eb8b9c138bd6a2e/sklearn/base.py
7
8  from abc import ABCMeta, abstractmethod
9
10 import copy
11 import warnings
12 from collections import defaultdict
13 import inspect
14 import re
15
16 import numpy as np
17
18 from .._version import __version__
19
20 _DEFAULT_TAGS = {
21     'non_deterministic': False,
22     'requires_positive_data': False,
23     'X_types': ['2darray'],
24     'poor_score': False,
25     'no_validation': False,
26     'multioutput': False,
27     "allow_nan": False,
28     'stateless': False,
29     'multilabel': False,
30     '_skip_test': False,
31     'multioutput_only': False}
32
33
34 def clone(estimator, safe=True):
```

```

35 """Constructs a new estimator with the same parameters.
36
37 Clone does a deep copy of the model in an estimator
38 without actually copying attached data. It yields a new estimator
39 with the same parameters that has not been fit on any data.
40
41 Parameters
42 -----
43 estimator : estimator object, or list, tuple or set of objects
44             The estimator or group of estimators to be cloned
45
46 safe : boolean, optional
47       If safe is false, clone will fall back to a deep copy on objects
48       that are not estimators.
49
50 Notes
51 ----
52 Taken from sklearn for compatibility.
53 """
54 estimator_type = type(estimator)
55 # XXX: not handling dictionaries
56 if estimator_type in (list, tuple, set, frozenset):
57     return estimator_type([clone(e, safe=safe) for e in estimator])
58 elif not hasattr(estimator, 'get_params') or isinstance(estimator, type):
59     if not safe:
60         return copy.deepcopy(estimator)
61     else:
62         raise TypeError("Cannot clone object '%s' (type %s): "
63                         "it does not seem to be a valid estimator "
64                         "as it does not implement a 'get_params' methods."
65                         % (repr(estimator), type(estimator)))
66 klass = estimator.__class__
67 new_object_params = estimator.get_params(deep=False)
68 for name, param in new_object_params.items():
69     new_object_params[name] = clone(param, safe=False)
70 new_object = klass(**new_object_params)
71 params_set = new_object.get_params(deep=False)
72
73 # quick sanity check of the parameters of the clone
74 for name in new_object_params:
75     param1 = new_object_params[name]
76     param2 = params_set[name]
77     if param1 is not param2:
78         raise RuntimeError('Cannot clone object %s, as the constructor '
79                             "'either does not set or modifies parameter %s' %"
80                             (estimator, name))
81
82 return new_object
83
84 def _pprint(params, offset=0, printer=repr):
85     """Pretty print the dictionary 'params'
86
87 Parameters
88 -----

```

```

89     params : dict
90         The dictionary to pretty print
91
92     offset : int
93         The offset in characters to add at the begin of each line.
94
95     printer : callable
96         The function to convert entries to strings, typically
97         the builtin str or repr
98
99     Notes
100     -----
101     Taken from sklearn for compatibility.
102     """
103     # Do a multi-line justified repr:
104     options = np.get_printoptions()
105     np.set_printoptions(precision=5, threshold=64, edgeitems=2)
106     params_list = list()
107     this_line_length = offset
108     line_sep = ',\n' + (1 + offset // 2) * ' '
109     for i, (k, v) in enumerate(sorted(params.items())):
110         if type(v) is float:
111             # use str for representing floating point numbers
112             # this way we get consistent representation across
113             # architectures and versions.
114             this_repr = '%s=%s' % (k, str(v))
115         else:
116             # use repr of the rest
117             this_repr = '%s=%s' % (k, printer(v))
118         if len(this_repr) > 500:
119             this_repr = this_repr[:300] + '...' + this_repr[-100:]
120         if i > 0:
121             if (this_line_length + len(this_repr) >= 75 or '\n' in this_repr):
122                 params_list.append(line_sep)
123                 this_line_length = len(line_sep)
124             else:
125                 params_list.append(', ')
126                 this_line_length += 2
127         params_list.append(this_repr)
128         this_line_length += len(this_repr)
129
130     np.set_printoptions(**options)
131     lines = ''.join(params_list)
132     # Strip trailing space to avoid nightmare in doctests
133     lines = '\n'.join(line.rstrip(' ') for line in lines.split('\n'))
134     return lines
135
136
137 def _update_if_consistent(dict1, dict2):
138     common_keys = set(dict1.keys()).intersection(dict2.keys())
139     for key in common_keys:
140         if dict1[key] != dict2[key]:
141             raise TypeError("Inconsistent values for tag {}: {} != {}".format(
142                 key, dict1[key], dict2[key]

```

```

143         ))
144     dict1.update(dict2)
145     return dict1
146
147
148 class BaseEstimator:
149     """Base Estimator class for compatibility with scikit-learn.
150
151     Notes
152     -----
153     * All estimators should specify all the parameters that can be set
154       at the class level in their ``__init__`` as explicit keyword
155       arguments (no ``*args`` or ``**kwargs``).
156     * Taken from sklearn for compatibility.
157     """
158
159     @classmethod
160     def _get_param_names(cls):
161         """Get parameter names for the estimator"""
162         # fetch the constructor or the original constructor before
163         # deprecation wrapping if any
164         init = getattr(cls.__init__, 'deprecated_original', cls.__init__)
165         if init is object.__init__:
166             # No explicit constructor to introspect
167             return []
168
169         # introspect the constructor arguments to find the model parameters
170         # to represent
171         init_signature = inspect.signature(init)
172         # Consider the constructor parameters excluding 'self'
173         parameters = [p for p in init_signature.parameters.values()
174                       if p.name != 'self' and p.kind != p.VAR_KEYWORD]
175         for p in parameters:
176             if p.kind == p.VAR_POSITIONAL:
177                 raise RuntimeError("scikit-multiflow estimators should always "
178                                "specify their parameters in the signature"
179                                " of their __init__ (no varargs).")
180             "%s with constructor %s doesn't "
181             " follow this convention."
182             % (cls, init_signature))
183         # Extract and sort argument names excluding 'self'
184         return sorted([p.name for p in parameters])
185
186     def get_params(self, deep=True):
187         """Get parameters for this estimator.
188
189         Parameters
190         -----
191         deep : boolean, optional
192             If True, will return the parameters for this estimator and
193             contained subobjects that are estimators.
194
195         Returns
196         -----

```

```

197     params : mapping of string to any
198           Parameter names mapped to their values.
199     """
200     out = dict()
201     for key in self._get_param_names():
202         value = getattr(self, key, None)
203         if deep and hasattr(value, 'get_params'):
204             deep_items = value.get_params().items()
205             out.update((key + '__' + k, val) for k, val in deep_items)
206         out[key] = value
207     return out
208
209 def set_params(self, **params):
210     """Set the parameters of this estimator.
211
212     The method works on simple estimators as well as on nested objects
213     (such as pipelines). The latter have parameters of the form
214     ``<component>__<parameter>`` so that it's possible to update each
215     component of a nested object.
216
217     Returns
218     -----
219     self
220     """
221     if not params:
222         # Simple optimization to gain speed (inspect is slow)
223         return self
224     valid_params = self.get_params(deep=True)
225
226     nested_params = defaultdict(dict) # grouped by prefix
227     for key, value in params.items():
228         key, delim, sub_key = key.partition('__')
229         if key not in valid_params:
230             raise ValueError('Invalid parameter %s for estimator %s. '
231                               'Check the list of available parameters '
232                               'with `estimator.get_params().keys()`.`' %
233                               (key, self))
234
235         if delim:
236             nested_params[key][sub_key] = value
237         else:
238             setattr(self, key, value)
239             valid_params[key] = value
240
241     for key, sub_params in nested_params.items():
242         valid_params[key].set_params(**sub_params)
243
244     return self
245
246 def __repr__(self, N_CHAR_MAX=700):
247     # N_CHAR_MAX is the (approximate) maximum number of non-blank
248     # characters to render. We pass it as an optional parameter to ease
249     # the tests.
250

```

```

251 from ..utils._pprint import _EstimatorPrettyPrinter
252
253 N_MAX_ELEMENTS_TO_SHOW = 30 # number of elements to show in sequences
254
255 # use ellipsis for sequences with a lot of elements
256 pp = _EstimatorPrettyPrinter(
257     compact=True, indent=1, indent_at_name=True,
258     n_max_elements_to_show=N_MAX_ELEMENTS_TO_SHOW)
259
260 repr_ = pp.pformat(self)
261
262 # Use brute force ellipsis when there are a lot of non-blank characters
263 n_nonblank = len(''.join(repr_.split()))
264 if n_nonblank > N_CHAR_MAX:
265     lim = N_CHAR_MAX // 2 # approx number of chars to keep on both ends
266     regex = r'^(\s*\S){%d}' % lim
267     # The regex '^(\s*\S){%d}' % n
268     # matches from the start of the string until the nth non-blank
269     # character:
270     # - ^ matches the start of string
271     # - (pattern){n} matches n repetitions of pattern
272     # - \s*\S matches a non-blank char following zero or more blanks
273     left_lim = re.match(regex, repr_).end()
274     right_lim = re.match(regex, repr[::-1]).end()
275
276     if '\n' in repr_[left_lim:-right_lim]:
277         # The left side and right side aren't on the same line.
278         # To avoid weird cuts, e.g.:
279         # categoric...ore',
280         # we need to start the right side with an appropriate newline
281         # character so that it renders properly as:
282         # categoric...
283         # handle_unknown='ignore',
284         # so we add [^\n]*\n which matches until the next \n
285         regex += r'[^\n]*\n'
286         right_lim = re.match(regex, repr[::-1]).end()
287
288     ellipsis = '...'
289     if left_lim + len(ellipsis) < len(repr_) - right_lim:
290         # Only add ellipsis if it results in a shorter repr
291         repr_ = repr_[:left_lim] + '...' + repr_[-right_lim:]
292
293     return repr_
294
295 def __getstate__(self):
296     try:
297         state = super().__getstate__()
298     except AttributeError:
299         state = self.__dict__.copy()
300
301     if type(self).__module__.startswith('skmultiflow.'):
302         return dict(state.items(), _skmultiflow_version=__version__)
303     else:
304         return state

```

```

305
306 def __setstate__(self, state):
307     if type(self).__module__.startswith('skmultiflow.'):
308         pickle_version = state.pop("_skmultiflow_version", "pre-0.18")
309         if pickle_version != __version__:
310             warnings.warn(
311                 "Trying to unpickle estimator {0} from version {1} when "
312                 "using version {2}. This might lead to breaking code or "
313                 "invalid results. Use at your own risk.".format(
314                     self.__class__.__name__, pickle_version, __version__),
315                 UserWarning)
316     try:
317         super().__setstate__(state)
318     except AttributeError:
319         self.__dict__.update(state)
320
321 def _get_tags(self):
322     collected_tags = {}
323     for base_class in inspect.getmro(self.__class__):
324         if (hasattr(base_class, '_more_tags') and base_class != self.__class__):
325             more_tags = base_class._more_tags(self)
326             collected_tags = _update_if_consistent(collected_tags,
327                                                     more_tags)
328     if hasattr(self, '_more_tags'):
329         more_tags = self._more_tags()
330         collected_tags = _update_if_consistent(collected_tags, more_tags)
331     tags = _DEFAULT_TAGS.copy()
332     tags.update(collected_tags)
333     return tags
334

```

```

335
336 class BaseSKMObject(BaseEstimator):

```

```

337     """Base class for most objects in scikit-multiflow
338
339     Notes
340     -----
341     This class provides additional functionality not available in the base estimator
342     from scikit-learn
343     """
344     def reset(self):
345         """ Resets the estimator to its initial state.
346
347     Returns
348     -----
349     self
350
351     """
352     # non-optimized default implementation; override if a better
353     # method is possible for a given object
354     command = ''.join([line.strip() for line in self.__repr__().split()])
355     command = command.replace(str(self.__class__.__name__), 'self.__init__')
356     exec(command)
357
358     def get_info(self):

```

```

359         """ Collects and returns the information about the configuration of the estimator
360
361     Returns
362     -----
363     string
364         Configuration of the estimator.
365     """
366     return self.__repr__()
367
368
369 class ClassifierMixin(metaclass=ABCMeta):
370     """Mixin class for all classifiers in scikit-multiflow."""
371     _estimator_type = "classifier"
372
373     def fit(self, X, y, classes=None, sample_weight=None):
374         """ Fit the model.
375
376         Parameters
377         -----
378         X : numpy.ndarray of shape (n_samples, n_features)
379             The features to train the model.
380
381         y: numpy.ndarray of shape (n_samples, n_targets)
382             An array-like with the class labels of all samples in X.
383
384         classes: numpy.ndarray, optional (default=None)
385             Contains all possible/known class labels. Usage varies depending
386             on the learning method.
387
388         sample_weight: numpy.ndarray, optional (default=None)
389             Samples weight. If not provided, uniform weights are assumed.
390             Usage varies depending on the learning method.
391
392     Returns
393     -----
394     self
395
396     """
397     # non-optimized default implementation; override if a better
398     # method is possible for a given classifier
399     self.partial_fit(X, y, classes=classes, sample_weight=sample_weight)
400
401     return self
402
403     @abstractmethod
404     def partial_fit(self, X, y, classes=None, sample_weight=None):
405         """ Partially (incrementally) fit the model.
406
407         Parameters
408         -----
409         X : numpy.ndarray of shape (n_samples, n_features)
410             The features to train the model.
411
412         y: numpy.ndarray of shape (n_samples)

```



```

413         An array-like with the class labels of all samples in X.
414
415     classes: numpy.ndarray, optional (default=None)
416         Array with all possible/known class labels. Usage varies depending
417         on the learning method.
418
419     sample_weight: numpy.ndarray of shape (n_samples), optional (default=None)
420         Samples weight. If not provided, uniform weights are assumed.
421         Usage varies depending on the learning method.
422
423     Returns
424     -----
425     self
426
427     """
428     raise NotImplementedError
429
430 @abstractmethod
431 def predict(self, X):
432     """ Predict classes for the passed data.
433
434     Parameters
435     -----
436     X : numpy.ndarray of shape (n_samples, n_features)
437         The set of data samples to predict the class labels for.
438
439     Returns
440     -----
441     A numpy.ndarray with all the predictions for the samples in X.
442
443     """
444     raise NotImplementedError
445
446 @abstractmethod
447 def predict_proba(self, X):
448     """ Estimates the probability of each sample in X belonging to each of the class-labels.
449
450     Parameters
451     -----
452     X : numpy.ndarray of shape (n_samples, n_features)
453         The matrix of samples one wants to predict the class probabilities for.
454
455     Returns
456     -----
457     A numpy.ndarray of shape (n_samples, n_labels), in which each outer entry is associated
458     with the X entry of the same index. And where the list in index [i] contains
459     len(self.target_values) elements, each of which represents the probability that
460     the i-th sample of X belongs to a certain class-label.
461
462     """
463     raise NotImplementedError
464
465 def score(self, X, y, sample_weight=None):
466     """Returns the mean accuracy on the given test data and labels.

```

```

467
468     In multi-label classification, this is the subset accuracy
469     which is a harsh metric since you require for each sample that
470     each label set be correctly predicted.
471
472     Parameters
473     -----
474     X : array-like, shape = (n_samples, n_features)
475         Test samples.
476
477     y : array-like, shape = (n_samples) or (n_samples, n_outputs)
478         True labels for X.
479
480     sample_weight : array-like, shape = [n_samples], optional
481         Sample weights.
482
483     Returns
484     -----
485     score : float
486         Mean accuracy of self.predict(X) wrt. y.
487
488     """
489     from sklearn.metrics import accuracy_score
490     return accuracy_score(y, self.predict(X), sample_weight=sample_weight)
491
492
493 class RegressorMixin(metaclass=ABCMeta):
494     """Mixin class for all regression estimators in scikit-multiflow."""
495     _estimator_type = "regressor"
496
497     def fit(self, X, y, sample_weight=None):
498         """ Fit the model.
499
500         Parameters
501         -----
502         X : numpy.ndarray of shape (n_samples, n_features)
503             The features to train the model.
504
505         y: numpy.ndarray of shape (n_samples, n_targets)
506             An array-like with the target values of all samples in X.
507
508         sample_weight: numpy.ndarray, optional (default=None)
509             Samples weight. If not provided, uniform weights are assumed. Usage varies
510             depending on the learning method.
511
512         Returns
513         -----
514         self
515
516         """
517         # non-optimized default implementation; override if a better
518         # method is possible for a given regressor
519         self.partial_fit(X, y, sample_weight=sample_weight)
520

```

```

521         return self
522
523     @abstractmethod
524     def partial_fit(self, X, y, sample_weight=None):
525         """ Partially (incrementally) fit the model.
526
527         Parameters
528         -----
529         X : numpy.ndarray of shape (n_samples, n_features)
530             The features to train the model.
531
532         y: numpy.ndarray of shape (n_samples)
533             An array-like with the target values of all samples in X.
534
535         sample_weight: numpy.ndarray of shape (n_samples), optional (default=None)
536             Samples weight. If not provided, uniform weights are assumed. Usage varies
537             depending on the learning method.
538
539         Returns
540         -----
541         self
542
543         """
544         raise NotImplementedError
545
546     @abstractmethod
547     def predict(self, X):
548         """ Predict target values for the passed data.
549
550         Parameters
551         -----
552         X : numpy.ndarray of shape (n_samples, n_features)
553             The set of data samples to predict the target values for.
554
555         Returns
556         -----
557         A numpy.ndarray with all the predictions for the samples in X.
558
559         """
560         raise NotImplementedError
561
562     @abstractmethod
563     def predict_proba(self, X):
564         """ Estimates the probability for probabilistic/bayesian regressors
565
566         Parameters
567         -----
568         X : numpy.ndarray of shape (n_samples, n_features)
569             The matrix of samples one wants to predict the probabilities for.
570
571         Returns
572         -----
573         numpy.ndarray
574

```

```

575     """
576     raise NotImplementedError
577
578 def score(self, X, y, sample_weight=None):
579     """Returns the coefficient of determination R^2 of the prediction.
580
581     The coefficient R^2 is defined as (1 - u/v), where u is the residual
582     sum of squares ((y_true - y_pred) ** 2).sum() and v is the total
583     sum of squares ((y_true - y_true.mean()) ** 2).sum().
584     The best possible score is 1.0 and it can be negative (because the
585     model can be arbitrarily worse). A constant model that always
586     predicts the expected value of y, disregarding the input features,
587     would get a R^2 score of 0.0.
588
589     Parameters
590     -----
591     X : array-like, shape = (n_samples, n_features)
592         Test samples. For some estimators this may be a
593         precomputed kernel matrix instead, shape = (n_samples,
594         n_samples_fitted], where n_samples_fitted is the number of
595         samples used in the fitting for the estimator.
596
597     y : array-like, shape = (n_samples) or (n_samples, n_outputs)
598         True values for X.
599
600     sample_weight : array-like, shape = [n_samples], optional
601         Sample weights.
602
603     Returns
604     -----
605     score : float
606         R^2 of self.predict(X) wrt. y.
607
608     Notes
609     ----
610     The R2 score used when calling ``score`` on a regressor will use
611     ``multioutput='uniform_average'`` from version 0.23 to keep consistent
612     with ``metrics.r2_score``. This will influence the ``score`` method of
613     all the multioutput regressors (except for
614     ``multioutput.MultiOutputRegressor``). To specify the default value
615     manually and avoid the warning, please either call ``metrics.r2_score``
616     directly or make a custom scorer with ``metrics.make_scorer`` (the
617     built-in scorer ``'r2'`` uses ``multioutput='uniform_average'``).
618     """
619
620     from sklearn.metrics import r2_score
621     from sklearn.metrics.regression import _check_reg_targets
622     y_pred = self.predict(X)
623     # XXX: Remove the check in 0.23
624     y_type, _, _, _ = _check_reg_targets(y, y_pred, None)
625     if y_type == 'continuous-multioutput':
626         warnings.warn("The default value of multioutput (not exposed in "
627                       "score method) will change from 'variance_weighted' "
628                       "to 'uniform_average' in 0.23 to keep consistent ")

```

```

629         "with 'metrics.r2_score'. To specify the default "
630         "value manually and avoid the warning, please "
631         "either call 'metrics.r2_score' directly or make a "
632         "custom scorer with 'metrics.make_scorer' (the "
633         "built-in scorer 'r2' uses "
634         "multioutput='uniform_average').", FutureWarning)
635     return r2_score(y, y_pred, sample_weight=sample_weight,
636                    multioutput='variance_weighted')
637
638
639 class MetaEstimatorMixin(object):
640     """Mixin class for all meta estimators in scikit-multiflow."""
641     _required_parameters = ["estimator"]
642
643
644 class MultiOutputMixin(object):
645     """Mixin to mark estimators that support multioutput."""
646     def _more_tags(self):
647         return {'multioutput': True}
648
649
650 def is_classifier(estimator):
651     """Returns True if the given estimator is (probably) a classifier.
652
653     Parameters
654     -----
655     estimator : object
656         Estimator object to test.
657
658     Returns
659     -----
660     out : bool
661         True if estimator is a classifier and False otherwise.
662     """
663     return getattr(estimator, "_estimator_type", None) == "classifier"
664
665
666 def is_regressor(estimator):
667     """Returns True if the given estimator is (probably) a regressor.
668
669     Parameters
670     -----
671     estimator : object
672         Estimator object to test.
673
674     Returns
675     -----
676     out : bool
677         True if estimator is a regressor and False otherwise.
678     """
679     return getattr(estimator, "_estimator_type", None) == "regressor"

```