# Issue #16710: Pipeline requires both fit and transform method to be available instead of only fit_transform

**Link to Issue**: https://github.com/scikit-learn/scikit-learn/issues/16710

**Summary of issue**: Calling a pipeline with a nonparametric function raises an error since the function transform() is missing. The pipeline itself calls the function fit_transform() if it's present. For nonparametric functions (the most prominent being t-SNE) a regular transform() method does not exist since there is no projection or mapping that is learned. But it could still be used for dimensionality reduction.

## Testing Environment

// Assuming user is running python 3.9+ and sklearn 1.0+ in a similar environment

```
>>> import sklearn; sklearn.show_versions()
System:
    python: 3.9.10 | packaged by conda-forge | (main, Feb  1 2022, 21:21:54) [MSC v.1929 64 bit
(AMD64)]
executable:/sklearn-env/python.exe
   machine: Windows-10-10.0.19044-SP0

Python dependencies:
        pip: 22.0.3
  setuptools: 60.9.3
    sklearn: 1.1.dev0
     numpy: 1.22.2
     scipy: 1.8.0
    Cython: 0.29.28
    pandas: None
  matplotlib: None
     joblib: 1.1.0
threadpoolctl: 3.1.0
     pytest: 7.0.1

Built with OpenMP: True
```

## Resources

```
from sklearn.decomposition import PCA
from sklearn.manifold import TSNE
from sklearn.pipeline import make_pipeline
import numpy as np
```

**Test Case**

// Input in python console

```
>>> X = np.array([[1], [2]])
>>> y = np.array([1, 2])
>>> pipe = make_pipeline(TSNE(), PCA())
>>> pipe.fit(X, y)
```

**Expected Output**

// A pipeline is created with warning message

```
…/course-project-keycap-guardians/scikit-learn/sklearn/pipeline.py:215: UserWarning: 'TSNE()' (type
<class 'sklearn.manifold._t_sne.TSNE'>) does not implement transform.
  Pipeline(steps=[('tsne', TSNE()), ('pca', PCA())])
```

## Implementation

**File modified:** scikit-learn/sklearn/pipeline.py

Modified logic of the code to allow pipelines with missing transformer implementation

```
204 -            if not (hasattr(t, "fit") or          205 +            if not (hasattr(t, "fit") or
        hasattr(t, "fit_transform")) or not hasattr(          (hasattr(t, "t")) and (hasattr(t, "fit") or
205 -                t, "transform"                      206 +                    (hasattr(t, "fit_transform")))
                                                          and (hasattr(t, "fit_transform") or
206 -            ):                                      207 +                    hasattr(t, "transform"))):
```

Added warnings for when transform is not implemented

```
                                                        214 +            if not hasattr(t, "transform"):
                                                        215 +                warnings.warn("'%s' (type %s) does
                                                                  not implement"
                                                        216 +                    " transform." % (t, type(t)))
```

Adjusted error messages to inform users that the creation of a pipeline is now allowed.

```
         raise TypeError(                               208         raise TypeError(
             "All intermediate steps should             209             "All intermediate steps should
    be "                                                        be "
-            "transformers and implement                210 +            "transformers and implement
    fit and transform "                                        fit and transform, fit_transform "
-            "or be the string                          211 +            "or be the string
    'passthrough' "                                            'passthrough'. "
             "'%s' (type %s) doesn't" % (t,             212             "'%s' (type %s) doesn't" % (t,
    type(t))                                                    type(t))
         )                                              213         )
```

```
176     msg = (                                         178     msg = (
177 -        "Last step of Pipeline should implement fit 179 +        "Last step of Pipeline should implement
    "                                                       fit, fit_transform "
178         "or be the string 'passthrough'"            180         "or be the string 'passthrough'"
179         ".*NoFit.*"                                 181         ".*NoFit.*"
180     )                                               182     )
```

```
1027              raise TypeError(                    1030              raise TypeError(
1028                  "All estimators should           1031                  "All estimators should
        implement fit and "                                   implement fit and "
1029 -               "transform. '%s' (type %s)        1032 +               "transform, fit_transform.
        doesn't" % (t, type(t))                               '%s' (type %s) doesn't" % (t, type(t))
1030                  )                                1033                  )
```

**File modified:** scikit-learn/sklearn/tests/test_pipeline.py

<u>Added test</u> to check that a warning is raised when a pipe is trying to fit when there is a missing transform implementation.

```
255  + def test_pipeline_tsne_pca():
256  +     X = np.array([[1], [2]])
257  +     y = np.array([1, 2])
258  +     pipe = make_pipeline(TSNE(), PCA())
259  +     with pytest.warns(UserWarning, match="does not
       implement transform"):
260  +         pipe.fit(X, y)
261  +
262  +
```