

Pied PyPler

**Why packaging is important for both open and
close data science projects**

Motivation

How many times did you...

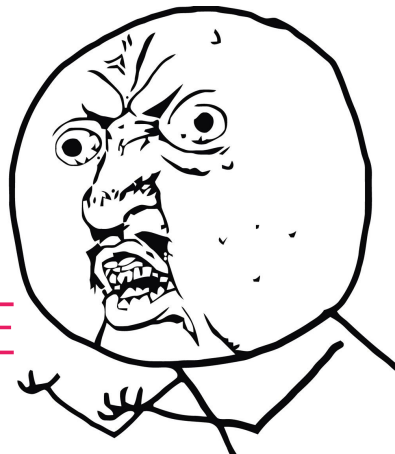
- Looked for an old piece of code
- Searched Dropbox & private repos
- Started Jupyter to search notebooks
- **Gave up and rewritten it**

How many times did you...

- Looked for an old piece of code
- Searched Dropbox & private repos
- Started Jupyter to search notebooks
- **Gave up and rewritten it**

CODE

Y U NO HERE

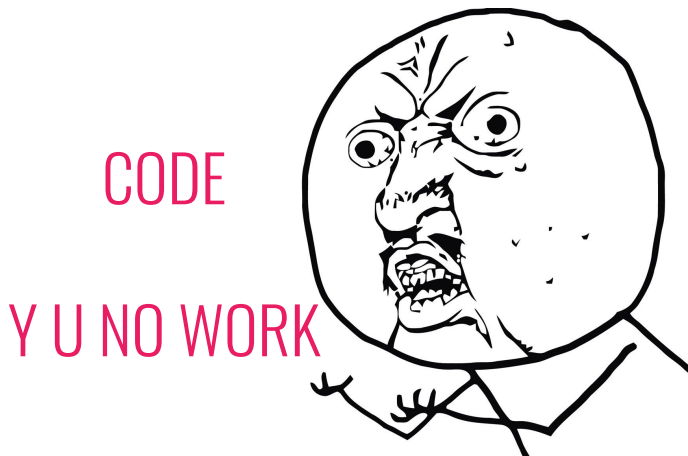


How many times did you...

- Looked for an old piece of code
- Found it!
- Copied it over
- What does this line do?!
- **Re-written half of it**

How many times did you...

- Looked for an old piece of code
- Found it!
- Copied it over
- What does this line do?!
- **Re-written half of it**



The solution?
Packaging!

What's in it for me?

- Write **less code** (over time)
- Code is **used more**
- All shall love you and **despair**

What's in it for me?

- Wr
- Co
- All



What's in it for me?

- Wr
- Co
- All



ALSO, NO MATH!

You will learn:

- Benefits to your code
- Benefits to your colleagues
- Benefits to you
- Which DS code fits packaging
- Designing a simple Python package

You will learn:

- Benefits to your code
- Benefits to your colleagues
- **Benefits to you**
- What DS codes fits packaging
- Designing a simple Python package



Benefits

Benefits to your code

- Built-in incentive to:
 - **Simplify & clarify**
 - **Document**
 - **Test**
- Gains **developer scrutiny**

Benefits to your code

- Built-in incentive to:
 - **Simplify & clarify**
 - **Document**
 - **Test**
- Gains **developer scrutiny**

AKA,

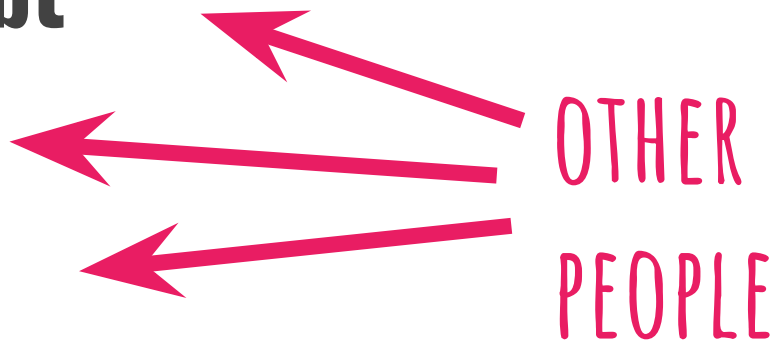
IT WORKS...

Benefits to people

- Code used **more often**, in more places
- **Other people** fixing & improving
- Manage **technical debt**
- Easier to **deploy**
- Save others writing it

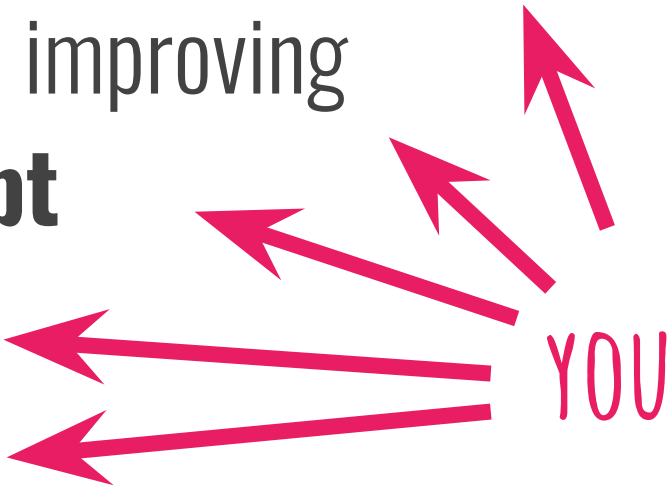
Benefits to people

- Code used **more often**, in more places
- **Other people** fixing & improving
- Manage **technical debt**
- Easier to **deploy**
- Save others writing it



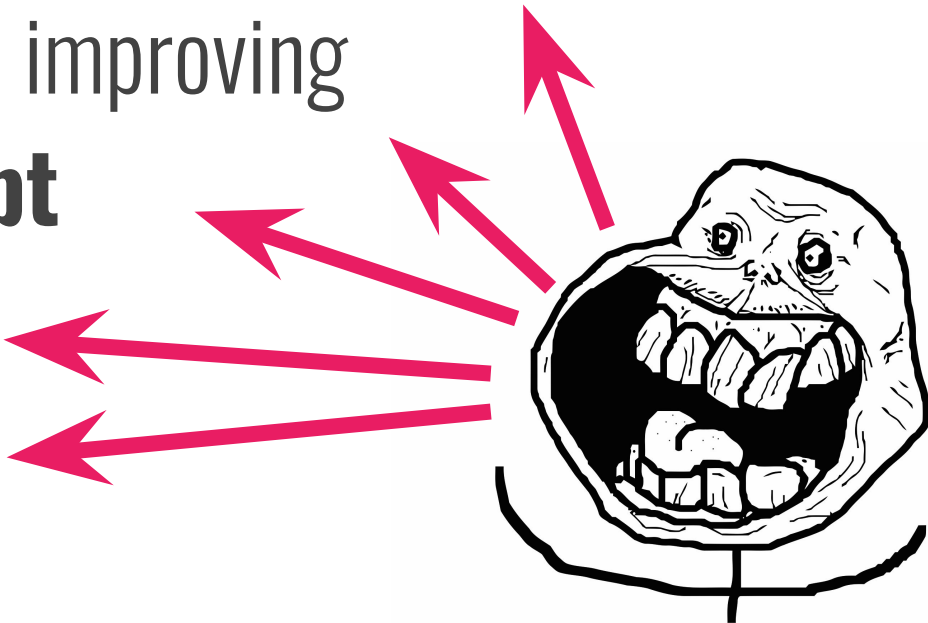
Benefits to people

- Code used **more often**, in more places
- **Other people** fixing & improving
- Manage **technical debt**
- Easier to **deploy**
- Save others writing it



Benefits to people

- Code used **more often**, in more places
- **Other people** fixing & improving
- Manage **technical debt**
- Easier to **deploy**
- Save others writing it



What to package?

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automating your flow
- Data science infrastructure
- Specific technique implementation

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automating your flow
- Data science infrastructure
- Specific technique implementation

imbutil

pypi v0.0.6

python 3.5, 3.6

build passing

codecov 100%

license MIT

Additions to the `imbalanced-learn` package.

```
from imbutil.combine import MinMaxRandomSampler; from imblearn import pipeline;
# oversampling minority classes to 100 and undersampling majority classes to 800
sampler = MinMaxRandomSampler(min_freq=100, max_freq=800)
sampling_clf = pipeline.make_pipeline(sampler, inner_clf)
```

Contents

- 1 Installation
- 2 Basic Use
 - 2.1 combine
- 3 Contributing
 - 3.1 Installing for development

96 lines (83 sloc) | 3.44 KB

```
1  """Randomly samples data to bring all class frequencies into a range."""
2
3  import numpy as np
4  from sklearn.utils import check_X_y
5
6  from imblearn.base import SamplerMixin
7  from imblearn.utils import check_target_type, hash_X_y
8  from imblearn.under_sampling import RandomUnderSampler
9  from imblearn.over_sampling import RandomOverSampler
10
11
12  class MinMaxRandomSampler(SamplerMixin):
13      """Random samples data to bring all class frequencies into a range.
14
15      Parameters
```



```
class MinMaxRandomSampler(SamplerMixin):
```

```
    """Random samples data to bring all class frequencies into a range.
```

```
    Parameters
```

```
    -----
```

```
    min_freq : int
```

```
        The minimum frequency for a class after sampling. All classes with  
        fewer samples are over-sampled to have this number of samples.
```

```
    max_freq : int
```

```
        The maximum frequency for a class after sampling. All classes with  
        more samples are under-sampled to have this number of samples.
```

```
    random_state : int, RandomState instance or None, optional (default=None)
```

```
        If int, random_state is the seed used by the random number generator;
```

```
        If RandomState instance, random_state is the random number generator;
```

```
        If None, the random number generator is the RandomState instance used  
        by np.random.
```

```
    """
```

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automating your flow
- Data science infrastructure
- Specific technique implementation



pypi v0.0.11

python 3.5, 3.6

build passing

codecov 100%

license MIT

scikit-learn wrappers for Python fastText .

```
>>> from skift import FirstColFtClassifier
>>> df = pandas.DataFrame([['woof', 0], ['meow', 1]], columns=['txt', 'lbl'])
>>> sk_clf = FirstColFtClassifier(lr=0.3, epoch=10)
>>> sk_clf.fit(df[['txt']], df[['lbl']])
>>> sk_clf.predict([['woof']])
[0]
```

Contents

- [1 Installation](#)
- [2 Features](#)
- [3 Wrappers](#)
 - [3.1 Standard wrappers](#)
 - [3.2 pandas-dependent wrappers](#)

```
2  """scikit-learn classifier wrapper for fasttext."""
3
4  import os
5  import abc
6
7  import numpy as np
8  from fastText import train_supervised
9  from sklearn.base import BaseEstimator, ClassifierMixin
10 from sklearn.utils.multiclass import unique_labels
11 from sklearn.exceptions import NotFittedError
12
13 from .util import (
14     temp_dataset_fpath,
15     dump_xy_to_fasttext_format,
16     python_fasttext_model_to_bytes,
17     bytes_to_python_fasttext_model,
18 )
19
20
21 class FtClassifierABC(BaseEstimator, ClassifierMixin, metaclass=abc.ABCMeta):
```

<> Code

! Issues 0

Pull requests 0

Projects 0

Wiki

Insights

Settings

Return ndarrays instead of lists while predicting #2

Closed

uniaz opened this issue on Mar 13 · 4 comments



uniaz commented on Mar 13

Contributor



The functions `predict`, `predict_proba` return `lists` instead of `numpy` arrays which makes them unusable with classifiers like `sklearn.multiclass.OneVsRestClassifier`. `GridSearch` and other similar functionality also don't work.

This is a quick fix.



```
@@ -156,10 +156,10 @@ def predict(self, X):
```

```
156 156         y : array of int of shape = [n_samples]
157 157         Predicted labels for the given input samples.
158 158         """"
```

```
159 -         return [
```

```
159 +         return np.array([
```

```
160 160             self._clean_label(res[0][0])
```

```
161 161             for res in self._predict(X)
```

```
162 -         ]
```

```
162 +         ], dtype=np.float_)
```

```
163 163
```

```
164 164     def _format_probas(self, result):
```

```
165 165         lbl_prob_pairs = zip(result[0], result[1])
```



```
@@ -181,10 +181,10 @@ def predict_proba(self, X):
```

```
181 181         The class probabilities of the input samples. The order of the
182 182         classes corresponds to that in the attribute classes_.
183 183         """"
```

```
184 -         return [
```

```
184 +         return np.array([
```

```
185 185             self._format_probas(res)
```

```
186 186             for res in self._predict(X, self.num_classes_)
```

```
187 -         ]
```

```
187 +         ], dtype=np.float_)
```

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automatizing your flow
- Data science infrastructure
- Specific technique implementation

pdpipe

pypi v0.0.27

python 3.5, 3.6

build passing

codecov 100%

License MIT

Easy pipelines for pandas DataFrames.

```
>>> df = pd.DataFrame(  
    data=[[4, 165, 'USA'], [2, 180, 'UK'], [2, 170, 'Greece']],  
    index=['Dana', 'Jane', 'Nick'],  
    columns=['Medals', 'Height', 'Born']  
)  
>>> pipeline = pdp.ColDrop('Medals').Binarize('Born')  
>>> pipeline(df)
```

	Height	Born_UK	Born_USA
Dana	165	0	1
Jane	180	1	0
Nick	170	0	0

3.2.2 Printing Pipelines

A pipeline structure can be clearly displayed by printing the object:

```
>>> drop_name = pdp.ColDrop("Name")
>>> binar_label = pdp.Binarize("Label")
>>> map_job = pdp.MapColVals("Job", {"Part": True, "Full": True, "No": False})
>>> pipeline = pdp.PdPipeline([drop_name, binar_label, map_job])
>>> print(pipeline)
A pdpipe pipeline:
[ 0] Drop column Name
[ 1] Binarize Label
[ 2] Map values of column Job with {'Part': True, 'Full': True, 'No': False}.
```

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automating your flow
- Data science infrastructure
- Specific technique implementation

barn

Simple local/remote dataset store for Python.

```
from barn import Dataset
twitter_usa = Dataset(name='twitter_usa', task='NER')
# download from an azure block blob storage and load into a dataframe
twitter_usa.download(tags=['preprocessed'], version='20180305')
df = twitter_usa.df(tags=['preprocessed'], version='20180305')
```

Contents

- [1 Installation](#)
- [2 Features](#)

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automating your flow
- Data science infrastructure
- Specific technique implementation

Build text classification models

See tests/ folder for usage.

NOT MY PACKAGE!

Word based models

When dataset represented as (docs, words) word based models can be created using Token

```
from keras_text.models import TokenModelFactory
from keras_text.models import YoonKimCNN, AttentionRNN, StackedRNN

# RNN models can use `max_tokens=None` to indicate variable length words per mini
factory = TokenModelFactory(1, tokenizer.token_index, max_tokens=100, embedding_t
word_encoder_model = YoonKimCNN()
model = factory.build_model(token_encoder_model=word_encoder_model)
model.compile(optimizer='adam', loss='categorical_crossentropy')
model.summary()
```

Types of code worth packaging

- Extensions to existing libraries
- Adapting existing code to a common API
- Automating your flow
- Data science infrastructure
- Specific technique implementation

How to package?

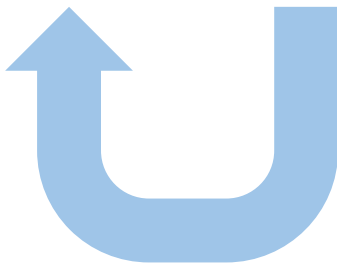
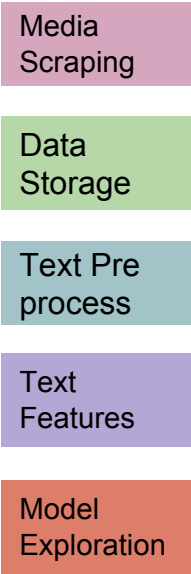
Good practices

- Do one thing
- Generalize the use case
- Keep it simple
- Advertise professionalism

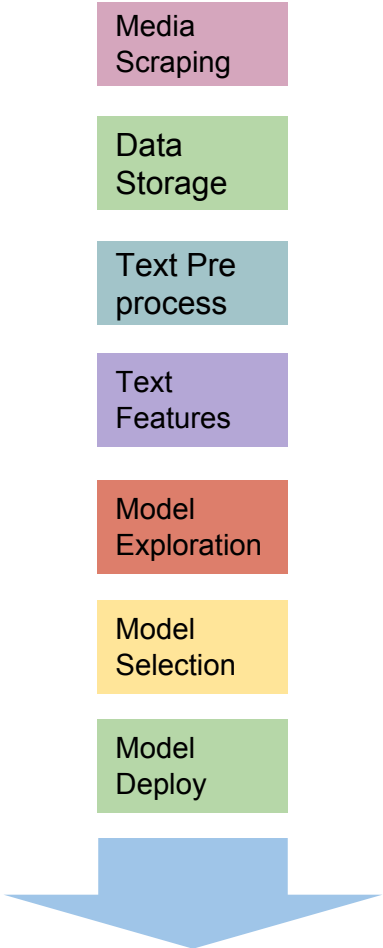
Do one thing

- NO: Data processing + model deployment + feature selection
- Don't recreate a flow; rather, a slice.
(horizontal, not vertical)

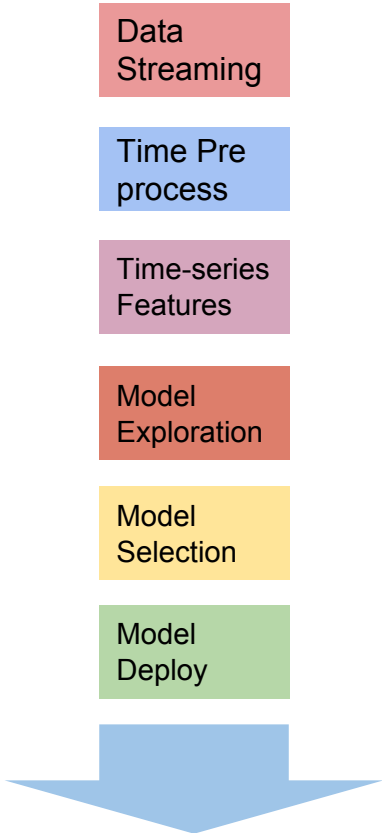
Sentiment Analysis (Research)



Sentiment Analysis (Prod)



Trend Pred. (Prod)



Sentiment Analysis (Research)

Sentiment Analysis (Prod)

Trend Pred. (Prod)

Media
Scraping

Media
Scraping

Data
Storage

Data
Storage

Data
Streaming

Text Pre
process

Text Pre
process

Time Pre
process

Text
Features

Text
Features

Time-series
Features

Model
Exploration

Model
Exploration

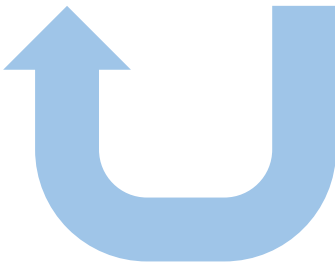
Model
Exploration

Model
Selection

Model
Selection

Model
Deploy

Model
Deploy



Generalize the use case

- Parameterize
- Split when X ways forward
 - If it's not too much work...

3.1 Standard wrappers

These wrappers do not make additional assumptions on input besides those commonly made by `scikit-learn` classifiers; i.e. that input is a 2d `ndarray` object and such.

- `FirstColFtClassifier` - An sklearn classifier adapter for fasttext that takes the first column of input `ndarray` objects as input.
- `IdxBasedFtClassifier` - An sklearn classifier adapter for fasttext that takes input by column index. This is set on object construction by providing the `input_ix` parameter to the constructor.

3.2 pandas-dependent wrappers

These wrappers assume the `X` parameter given to `fit`, `predict`, and `predict_proba` methods is a `pandas.DataFrame` object:

- `FirstObjFtClassifier` - An sklearn adapter for fasttext using the first column of `dtype == object` as input.
- `ColLblBasedFtClassifier` - An sklearn adapter for fasttext taking input by column label. This is set on object construction by providing the `input_col_lbl` parameter to the constructor.

```
@abc.abstractmethod
def _input_col(self, X):
    pass # pragma: no cover
```

FirstColFtClassifier

```
class FirstColFtClassifier(FtClassifierABC):
    """An sklearn classifier adapter for fast

    Parameters
    -----
    **kwargs
        Additional keyword arguments will be
        fasttext.train_supervised.
    """

    def _input_col(self, X):
        return np.array(X)[: , 0]
```

FirstObjFtClassifier

```
def _input_col(self, X):
    input_col_name = None
    for col_name, dtype in X.dtypes.items():
        if dtype == object:
            input_col_name = col_name
            break
    if input_col_name is not None:
        return X[input_col_name]
    raise ValueError("No object dtype column in input param X.")
```

Keep it simple

- Informative class/function names
- Minimal API
- A clear example with one/two imports

```
>>> from skift import FirstColFtClassifier
>>> df = pandas.DataFrame([['woof', 0], ['meow', 1]], columns=['txt', 'lbl'])
>>> sk_clf = FirstColFtClassifier(lr=0.3, epoch=10)
>>> sk_clf.fit(df[['txt']], df[['lbl']])
>>> sk_clf.predict([['woof']])
[0]
```

```
from skutil.estimators import ColumnIgnoringClassifier
# use a classifier that can't handle string data as
# an inner classifier in some stacked model, for example
```

```
>>> df = pd.DataFrame(
    data=[[4, 165, 'USA'], [2, 180, 'UK'], [2, 170, 'Greece']],
    index=['Dana', 'Jane', 'Nick'],
    columns=['Medals', 'Height', 'Born']
)
>>> pipeline = pdp.ColDrop('Medals').Binarize('Born')
>>> pipeline(df)
```

	Height	Born_UK	Born_USA
Dana	165	0	1
Jane	180	1	0
Nick	170	0	0

Advertise professionalism

- Continuous testing
- Test coverage
- Permissive License
- Documentation
- Encourage contributions

pdpipe

pypi v0.0.27

python 3.5, 3.6

build passing

codecov 100%

License MIT

Easy pipelines for pandas DataFrames.

pdpipeline

pypi v0.0.27

python 3.5, 3.6

build passing

codecov 100%

License MIT

Easy pipelines for pandas DataFrames.

pdpipeline 0.0.27

✓ Latest version

`pip install pdpipe`

Last released: May 28, 2018

Easy pipelines for pandas.

Navigation

Project description

Release history

Download files

Project links

Homepage

Project description

pypi v0.0.27 python 3.5, 3.6 build passing codecov 100% License MIT

Easy pipelines for pandas DataFrames.

```
>>> df = pd.DataFrame(  
    data=[[4, 165, 'USA'], [2, 180, 'UK'], [2, 170, 'Greece']],  
    index=['Dana', 'Jane', 'Nick'],  
    columns=['Medals', 'Height', 'Born']  
)  
>>> pipeline = pdp.ColDrop('Medals').Binarize('Born')  
>>> pipeline(df)  
   Height  Born_UK  Born_USA  
Dana    165      0        1  
Jane    180      1        0  
Nick    170      0        0
```

pdpipe

pypi v0.0.27

python 3.5, 3.6

build passing

codecov 100%

License MIT

Easy pipelines for pandas DataFrames.

shaypal5 / pdpipe  build passing

Current Branches Build History Pull Requests

✓ master  Shay Palachy	CRON Update README.rst	👤 #143 passed - 45bc3d5 🔗	🕒 2 min 55 sec 📅 about 4 hours ago
✓ master  Shay Palachy	Update README.rst	👤 #142 passed - 45bc3d5 🔗	🕒 2 min 37 sec 📅 7 days ago
✓ v0.0.27  Shay Palachy	Merge branch 'master' of github.com:shaypal5/pdpipe	👤 #141 passed - 4b479e1 🔗	🕒 2 min 21 sec 📅 7 days ago
✓ master  Shay Palachy	Merge branch 'master' of github.com:shaypal5/pdpipe	👤 #140 passed - 4b479e1 🔗	🕒 2 min 27 sec 📅 7 days ago

pdpipeline

pypi v0.0.27 python 3.5, 3.6 build passing codecov 100% License MIT

Easy pipelines for pandas DataFrames.



- 1 Installation
- 2 Features
 - 2.1 Design Decisions
- 3 Basic Use
 - 3.1 Pipeline Stages
 - 3.1.1 Creating Pipeline Stages
 - 3.1.2 Applying Pipeline Stages
 - 3.1.3 Fittable Pipeline Stages
 - 3.2 Pipelines
 - 3.2.1 Creating Pipelines
 - 3.2.2 Printing Pipelines
 - 3.2.3 Pipeline Arithmetics
 - 3.2.4 Pipeline Chaining
 - 3.2.5 Pipeline Slicing
 - 3.2.6 Applying Pipelines
- 4 Types of Pipeline Stages
 - 4.1 Basic Stages
 - 4.2 Column Generation
 - 4.3 Scikit-learn-dependent Stages
 - 4.4 nltk-dependent Stages
- 5 Creating additional stages
 - 5.1 Extending PdPipelineStage
 - 5.2 Ad-Hoc Pipeline Stages
- 6 Contributing
 - 6.1 Installing for development
 - 6.2 Running the tests
 - 6.3 Adding documentation

6 Contributing

Package author and current maintainer is Shay Palachy (shay.palachy@gmail.com); You are more than welcome to approach him for help. Contributions are very welcomed, especially since this package is very much in its infancy and many other pipeline stages can be added.

6.1 Installing for development

Clone:

```
git clone git@github.com:shaypal5/pdpipe.git
```

Install in development mode with test dependencies:

```
cd pdpipe
pip install -e ".[test]"
```

6.2 Running the tests

To run the tests, use:

```
python -m pytest --cov=pdpipe
```

6.3 Adding documentation

This project is documented using the [numpy docstring conventions](#), which were chosen as they are perhaps the most widely-spread conventions that are both supported by common tools such as Sphinx and result in human-readable docstrings (in my personal opinion, of course). When documenting code you add to this project, please follow [these conventions](#).

Good practices

- Do one thing
- Generalize the use case
- Keep it simple
- Advertise professionalism

What's next?

Homework

Example repo:

<https://github.com/shaypal5/catlolzer>

- Go home
- Find code you can package
- Create & upload a first version
- **Email me** at shay.palachy@gmail.com
- I'll help (if needed)

Homework

- Go to your company
- Deploy a private PyPI server
- Create & upload a first version
- People catch up surprisingly fast
- **Email me** at shay.palachy@gmail.com

That's it!

Email me at shay.palachy@gmail.com