

SCIKIT-LEARN PIPELINES AND CROSS VALIDATION

# ML Pipelines using scikit-learn and GridSearchCV

Managing ML workflows with Pipelines and using GridSearch Cross validation techniques for parameter tuning



Nikhil pentapalli

Jun 7, 2020 · 6 min read



Image from [Unsplash](#)

*ML calculations and algorithms generally process enormous information. A pipeline is an approach to chain those information handling ventures as required in an organized manner. Not just that, it likewise helps in making incredible work process and Reproducible code.*

## What is a ML Pipeline?

A pipeline is a progression of steps where information is changed. It originates from the "pipe and filter" plan design. In this way, you may have a class for each filter and afterward another class to join those means into the pipeline and make a complete final pipeline.

## Methods of a Scikit-Learn Pipeline

Pipelines **must** have those two methods:

- The word “fit” is to learn on the data and acquire its state
- The word “transform” (or “predict”) to actually process the data and generate a prediction.

It's also possible to call this method or to chain both of them:

- The word “fit\_transform” is to fit and then transform the data, but all in one go, which allows for great code optimizations when these two methods must be done simultaneously.

First, we shall define the model pipelines and then we do Grid search cross validation technique to find the optimal model for our problem statement.

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.feature_extraction.text import CountVectorizer
from sklearn import feature_extraction, linear_model,
model_selection, preprocessing
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
```

Here I am using the bbc news dataset from kaggle for building the pipelines.

```
news=pd.read_csv('bbc-text.csv')
news.head()
```

Out[4]:

category

text

	category	text
0	tech	tv future in the hands of viewers with home th...
1	business	worldcom boss left books alone former worldc...
2	sport	tigers wary of farrell gamble leicester say ...
3	sport	yeading face newcastle in fa cup premiership s...
4	entertainment	ocean s twelve raids box office ocean s twelve...

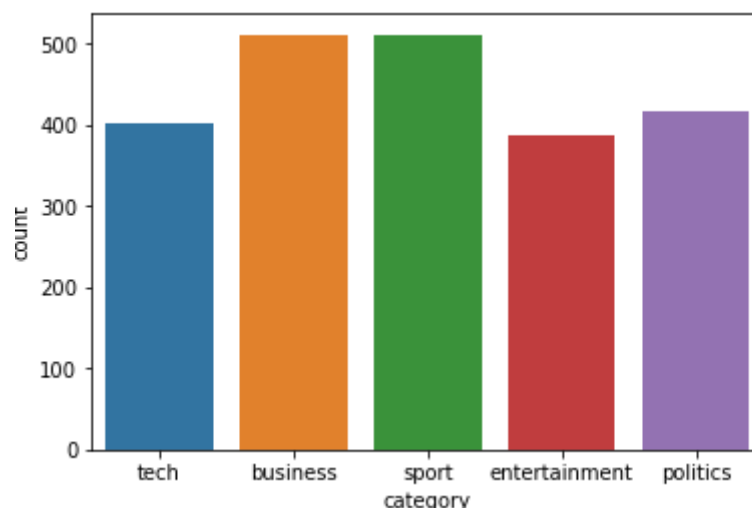
sample data

```
news['category'].value_counts()
```

```
Out[7]: sport      511
        business    510
        politics    417
        tech        401
        entertainment 386
        Name: category, dtype: int64
```

Value Counts of each category

```
sns.countplot(x="category", data=news)
```



Plot of each category

## Logistic Regression Classifier

```
x_train,x_test,y_train,y_test = train_test_split(news['text'],
news.category, test_size=0.2, random_state=2020)
```

```

pipe_lr = Pipeline([('vect', CountVectorizer()),
                    ('tfidf', TfidfTransformer()),
                    ('model', LogisticRegression())])

model = pipe_lr.fit(x_train, y_train)
prediction = model.predict(x_test)
print("\naccuracy: {}".format(round(accuracy_score(y_test,
prediction)*100,2)))

print('\n',confusion_matrix(y_test, prediction))

print('\n',classification_report(y_test, prediction))

```

accuracy: 96.63%

```

[[ 97   0   0   1   3]
 [  1  80   1   0   0]
 [  3   0  70   1   1]
 [  0   0   0 105   0]
 [  4   0   0   0  78]]

```

	precision	recall	f1-score	support
business	0.92	0.96	0.94	101
entertainment	1.00	0.98	0.99	82
politics	0.99	0.93	0.96	75
sport	0.98	1.00	0.99	105
tech	0.95	0.95	0.95	82
avg / total	0.97	0.97	0.97	445

complete report of the model

We imported pipeline from scikit-learn library. In the pipeline we can define the function that has to be performed in a sequence. Here we first split our data into train, test using train\_test\_split.

Then we defined CountVectorizer, Tf-Idf, Logistic regression in an order in our pipeline. This way it reduces the amount of code and pipelining the model helps in comparing it with different models and getting an optimal model for our choice.

Similarly let's create pipelines for different models with the same data as input and select the best model out of everything.

### Creation of different Pipelines:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.pipeline import Pipeline

```

```
3 from sklearn.model_selection import GridSearchCV
4 from sklearn.metrics import accuracy_score
5 from sklearn.externals import joblib
6 from sklearn.tree import DecisionTreeClassifier
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn import svm
10
11 x_train,x_test,y_train,y_test = train_test_split(news['text'], news.category, test_size=0.2, ra
12
13
14 #Now let us construct some pipelines of our choice.
15
16 pipe_lr = Pipeline([('vect', CountVectorizer()),
17                     ('tfidf', TfidfTransformer()),
18                     ('clf', LogisticRegression(random_state=42))])
19
20 pipe_dt = Pipeline([('vect', CountVectorizer()),
21                     ('tfidf', TfidfTransformer()),
22                     ('model', DecisionTreeClassifier(random_state=42))])
23
24 pipe_rf = Pipeline([('vect', CountVectorizer()),
25                     ('tfidf', TfidfTransformer()),
26                     ('clf', RandomForestClassifier(random_state=42))])
27
28
29 pipe_svm = Pipeline([('vect', CountVectorizer()),
30                     ('tfidf', TfidfTransformer()),
31                     ('clf', svm.SVC(random_state=42))])
32
33
34
35
```

pipeline\_creation.py hosted with ❤ by GitHub

[view raw](#)

In the above code i constructed four pipelines for Four different models.

1. Logistic Regression
2. Decision Tree
3. Random Forest
4. Support Vector Machine

You can absolutely change the architecture of the single pipeline. For example, in Logistic Regression if you feel the target values needs to be scaled you can use standard scaler or you can absolutely use some preprocessing techniques.

Let us assume a case where preprocessing is needed and how to handle such case.

```
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler, OneHotEncoder

num_transformer =
Pipeline([('imputer', SimpleImputer(strategy='mean')),
          ('scaler', StandardScaler())])

cat_transformer = Pipeline([
    ('imputer', SimpleImputer(strategy='constant',
                               fill_value='missing')),
    ('onehot', OneHotEncoder(handle_unknown='ignore'))])
```

So i have created two separate pipelines to handle different types of data that is num\_transformer handles the missing numbers using simple imputer while the standard scaler will transform your data such that its distribution will have a mean value 0 and standard deviation of 1.

In cat\_transformer i have used simple imputer and one hot encoder to encode the categories.

Now let us take the columns of the numerical features as well as categorical features and apply the pipeline of preprocessing steps described above. To do that i take the train\_data and select the specific columns as shown below.

```
num_features = train_data.select_dtypes(include=['int64']).columns
cat_features = train_data.select_dtypes(include=['object']).columns
```

Once we are ready we can now push everything into our Column Transformer which can be imported from sklearn.compose

```
from sklearn.compose import ColumnTransformer
```

```
preprocessor = ColumnTransformer([
    ('num', num_transformer, num_features),
    ('cat', cat_transformer, cat_features)])
```

So Now unlike above i just did all the preprocessing in one step which can be used for multiple models.

```
pipe_rf = Pipeline([('preprocess',preprocessor),
    ('clf', RandomForestClassifier(random_state=42))])
```

Let us now fit the models using GridSearchCV which helps us in model selection by passing many different params for each pipeline and getting the best model as well as best params with which the model was fit using. So let's get started by defining some params for grid search.

Linear Regression takes l2 penalty by default.so i would like to experiment with l1 penalty.Similarly for Random forest in the selection criterion i could want to experiment on both 'gini' and 'entropy'. So i have passed both the values in to the clf\_criterion

You can also experiment on different kernels of your choice for the specific problem statement and use case.

```
1  # Set grid search params
2  param_range = [9, 10]
3  param_range_fl = [1.0, 0.5]
4
5  grid_params_lr = [{'clf__penalty': ['l1', 'l2'],
6                      'clf__C': param_range_fl,
7                      'clf__solver': ['liblinear']}]}
8
9
10 grid_params_rf = [{'clf__criterion': ['gini', 'entropy'],
11                     'clf__max_depth': param_range,
12                     'clf__min_samples_split': param_range[1:]}]
13
14 grid_params_svm = [{'clf__kernel': ['linear', 'rbf'],
15                     'clf__C': param_range}]
16
17 # Construct grid searches
18 jobs = -1
```

```
19
20 LR = GridSearchCV(estimator=pipe_lr,
21                   param_grid=grid_params_lr,
22                   scoring='accuracy',
23                   cv=10)
24
25
26
27 RF = GridSearchCV(estimator=pipe_rf,
28                   param_grid=grid_params_rf,
29                   scoring='accuracy',
30                   cv=10,
31                   n_jobs=n_jobs)
32
33
34 SVM = GridSearchCV(estimator=pipe_svm,
35                   param_grid=grid_params_svm,
36                   scoring='accuracy',
37                   cv=10,
38                   n_jobs=n_jobs)
39
40
41
42 # List of pipelines for iterating through each of them
43 grids = [LR, RF, SVM]
44
45 # Creating a dict for our reference
46 grid_dict = {0: 'Logistic Regression',
47              1: 'Random Forest',
48              2: 'Support Vector Machine'}
49
50 # Fit the grid search objects
51 print('Performing model optimizations...')
52 best_acc = 0.0
53 best_clf = 0
54 best_gs = ''
55 for idx, gs in enumerate(grids):
56     print('\nEstimator: %s' % grid_dict[idx])
57     gs.fit(x_train, y_train)
58     print('Best params are : %s' % gs.best_params_)
59     # Best training data accuracy
60     print('Best training accuracy: %.3f' % gs.best_score_)
61     # Predict on test data with best params
62     y_pred = gs.predict(x_test)
63     # Test data accuracy of model with best params
64     print('Test set accuracy score for best params: %.3f ' % accuracy_score(y_test, y_pred))
65     # Track best (highest test accuracy) model
66     if accuracy_score(y_test, y_pred) > best_acc:
```



```
67     best_acc = accuracy_score(y_test, y_pred)
68     best_gs = gs
69     best_clf = idx
70     print('\nClassifier with best test set accuracy: %s' % grid_dict[best_clf])
71
72     # Save best grid search pipeline to file
73     dump_file = 'best_grid_search_pipeline.pkl'
74     joblib.dump(best_gs, dump_file, compress=1)
75     print('\nSaved %s grid search pipeline to file: %s' % (grid_dict[best_clf], dump_file))
```

## OUTPUT:

Performing model optimizations...

Estimator: Logistic Regression

Best params: {'clf\_\_C': 1.0, 'clf\_\_penalty': 'l2', 'clf\_\_solver': 'liblinear'}

Best training accuracy: 0.969

Test set accuracy score for best params: 0.966

Estimator: Random Forest

Best params: {'clf\_\_criterion': 'gini', 'clf\_\_max\_depth': 10, 'clf\_\_min\_samples\_split': 10}

Best training accuracy: 0.844

Test set accuracy score for best params: 0.836

Estimator: Support Vector Machine

Best params: {'clf\_\_C': 9, 'clf\_\_kernel': 'linear'}

Best training accuracy: 0.978

Test set accuracy score for best params: 0.971

Classifier with best test set accuracy: Support Vector Machine

Saved Support Vector Machine grid search pipeline to file:  
best\_grid\_search\_pipeline.pkl

Now as i have compared Logistic Regression, Random Forest and SVM in which i could definitely see that SVM is the best model with an accuracy of 0.978 .we also obtained the best parameters from the Grid Search cross validation.

So i have taken accuracy as a scoring parameter. It is completely based on the problem statement.You can take scoring parameter as recall or f1-score or precision which are completely derived from confusion matrix

Finally, Grid search builds a model for every combination of hyper parameters specified and evaluates each model. Another efficient technique for hyper parameter

tuning is the Randomized search — where random combinations of the hyper parameters are used to find the best solution.

The Confusion matrix and the scoring parameters can be understood from the below images.

		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

Confusion Matrix(image by author)

## Accuracy, Precision, Recall, F-measure

- $\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$
- $\text{Precision} = \frac{TP}{TP + FP}$
- $\text{Recall} = \frac{TP}{TP + FN}$
- $\text{F-measure} = 2 * \left( \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \right)$

Scoring Parameters(image by author)

You can define your own custom transformers which should definitely contain fit and transform methods in it. if you had noticed i had left the column transformer for you to try on your own.

Voila, there you go. Now you can try building your own machine learning pipeline and do not forget to try a custom transformer. Also change the parameters in the grid search and run experiments. The best way to learn anything is by experimenting with it.

I hope this article empowers your knowledge. Keep supporting and Happy Learning.

Nikhil Pentapalli - Data Scientist || Machine Learning Engineer - Jio Platforms Limited | LinkedIn

Experience in Real world AI product building. ->Experienced in transferring real problem into requirements and solution...

[www.linkedin.com](https://www.linkedin.com)

---

## Sign up for Analytics Vidhya News Bytes

By Analytics Vidhya

Latest news from Analytics Vidhya on our Hackathons and some of our best articles! [Take a look.](#)

Your email

---

Get this newsletter

By signing up, you will create a Medium account if you don't already have one. Review our [Privacy Policy](#) for more information about our privacy practices.

[Gridsearchcv](#)

[Scikit Learn](#)

[Machine Learning](#)

[Data Science](#)

[Model Tuning](#)

[About](#) [Help](#) [Legal](#)

Get the Medium app



