

Classification and Feature Sets in the NLTK, Part 1

These examples and others appear in Chapter 6 of the NLTK book. In this chapter, we look at how to prepare data for classification. For each example, or instance, of the classification problem, we prepare a set of features that will represent that example to the machine learning algorithm.

Name Gender Classifier

We start with a simple problem that will illustrate the process of preparing text data for classification and training and using classifiers. This problem is based on the idea that male and female first names (in English) have distinctive characteristics. For example, names ending in a, e, and i are likely to be female, while names ending in k, o, r, s and t are likely to be male. We will build a classifier that will label any name with its gender.

For each item to be classified, in this case a single word, in NLTK we build the features of that item as a dictionary that maps each feature name to a value, which can be a Boolean, a number or a string. A feature set is the feature dictionary together with the label of the item to be classified, in this case the gender.

We first define a function that will extract or build the features for a single instance of the problem, in this case a single name. To start with, we will generate a single feature which consists of the last letter of the name. Note that the function returns a dictionary with a single item.

```
>>> def gender_features(word):  
    return{'last_letter': word[-1]}
```

We can apply this function to any name:

```
>>> gender_features('Shrek')
```

Now that we've defined our features, we need to construct the training data. This will be a list of first names, each of which will be labeled either male or female. The NLTK corpus contains a names corpus which has a list of male first names and another list of female first names, so we can use this data to create a list of all the first names, but where each is labeled with its gender.

```
>>> from nltk.corpus import names
```

The names corpus has a function words that will return either the names identified by the string 'male.txt' or 'female.txt'. Here we look at the first 20 male names.

```
>>> names.words('male.txt')[:20]
```

From the male and female names lists, we will create one long list with (name, gender) pairs to create the labeled data.

```
>>> namesgender = [(name, 'male') for name in names.words('male.txt')] +  
[(name, 'female') for name in names.words('female.txt')]
```

Take a look at this list with the first 20 names and the last 20 names.

```
>>> len(namesgender)
7944
>>> namesgender[:20]
>>> namesgender[7924:]
```

Now we create a random shuffle of the namesgender list so that we can easily split it into a training and test set.

```
>>> import random
>>> random.shuffle(namesgender)
>>> namesgender[:20]
```

Next we use the feature extractor function to create the list of instances of the problem that consists only of the features and the gender label. (This is the equivalent of the spread sheet for structured data in a classification problem.) But this example is unusual in that most of the actual text is ignored in the features!

```
>>> featuresets = [(gender_features(n), g) for (n,g) in namesgender]
>>> featuresets[:20]
```

We split this list into training and test sets and run the Naïve Bayes classifier algorithm to create a trained classifier. (The training set is the last 7444 examples and the test set is the first 500 examples.) What do you think will happen if some of the examples with last letter = “n” are female and some of them are male?

```
>>> train_set, test_set = featuresets[500:], featuresets[:500]
>>> classifier = nltk.NaiveBayesClassifier.train(train_set)
```

We can the compute the accuracy of the classifier on the test set.

```
>>> print nltk.classify.accuracy(classifier, test_set)
```

One of the things that we can use a classifier for is to label totally new instances of the problem, in this case, names that come from the future:

```
>>> classifier.classify(gender_features('Neo'))
>>> classifier.classify(gender_features('Trinity'))
```

Finally, the classifier class for Naïve Bayes has a function that shows the feature values that were most important in doing the classification.

```
>>> classifier.show_most_informative_features(20)
```

When the output of this function shows a line like:

```
last_letter = 'a'  female : male= 35.4 : 1.0
```

it means that when an example has the feature last_letter = ‘a’, it is 35.4 times more likely to be classified as female than male.

Choosing Good Features

Selecting relevant features can usually be the most important part of training a classifier. Often the approach is to throw in as many features as possible and then try to figure out which ones were important. For most machine learning algorithms, throwing in too many features can cause a problem known as “overfitting”, which is that the classifier is trained on so many of the exact details of the training set that it is not as good on new examples.

For the name gender problem, we will try a second feature extraction function that has the first letter, the last letter, a count of each letter, and the individual letters of the name.

(Note that the string “count(%)” uses the string formatting features of Python. In this case, it is inserting the value of the variable *letter* into the string. For more details, see <http://docs.python.org/library/stdtypes.html#string-formatting-operations> .)

```
>>> def gender_features2(name):
    features = {}
    features["firstletter"] = name[0].lower()
    features["lastletter"] = name[-1].lower()
    for letter in 'abcdefghijklmnopqrstuvwxyz':
        features["count(%) " % letter] = name.lower().count(letter)
        features["has(%) " % letter] = (letter in name.lower())
    return features

>>> features = gender_features2('Shrek')
>>> len(features)
>>> features
```

Create new feature sets for all our names in the namesgender list.

```
>>> featuresets2 = [(gender_features2(n), g) for (n,g) in namesgender]
```

We can just look at the feature sets, but it’s more informative to print the name as well:

```
>>> for (n,g) in namesgender[:5]:
    print n, gender_features2(n), '\n'
```

We create a new training and test set based on these features. According to the NLTK book, they got a lower accuracy, but I get slightly higher (of course, we each have our own random shuffle of the dataset, so we cannot expect the same results):

```
>>> train_set, test_set = featuresets2[500:], featuresets2[:500]
>>> classifier = nltk.NaiveBayesClassifier.train(train_set)
>>> print nltk.classify.accuracy(classifier, test_set)
```

In developing a classifier for a classification problem, we will want to do some error analysis of the test set and then perhaps change our features and retrain the classifier. This is known as the development process, and it is important then to keep a separate test set that was not used in the error analysis for our final evaluation. So we actually divide our labeled data into a training set, a development test set, and a test set.

```
>>> train_names = namesgender[1500:]
>>> devtest_names = namesgender[500:1500]
>>> test_names = namesgender[:500]
```

Going back to using our features with just the last letter of each name, we generate the features for each name and train the classifier.

```
>>> train_set = [(gender_features(n), g) for (n,g) in train_names]
>>> devtest_set = [(gender_features(n), g) for (n,g) in devtest_names]
>>> test_set = [(gender_features(n), g) for (n,g) in test_names]
>>> classifier = nltk.NaiveBayesClassifier.train(train_set)
>>> print nltk.classify.accuracy(classifier, devtest_set)
```

(Save this classifier accuracy number on the development test set to use for comparison in today's exercise.)

Next we define a function that will get a list of errors by running the classifier on the development test names and comparing it with the original name gender labels (comparing the classifier labels with the gold standard labels).

```
>>> def geterrors(devtest):
    errors = []
    for (name, tag) in devtest:
        guess = classifier.classify(gender_features(name))
        if guess != tag:
            errors.append( (tag, guess, name) )
    return errors

>>> errors = geterrors(devtest_names)
>>> len(errors)
```

Then we define a function to print all the errors, sorted by the correct labels, so that we can look at the differences (with even more string formatting).

```
>>> def printerrors(errors):
    for (tag, guess, name) in sorted(errors):
        print 'correct=%-8s guess=%-8s name=%-30s' % (tag, guess, name)

>>> printerrors(errors)
```

Looking through the list of errors, we observe cases where using the last two letters of each name might be more informative. As observed in the book, “For example, names ending in *yn* appear to be predominantly female, despite the fact that names ending in *n* tend to be male; and names ending in *ch* are usually male, even though names that end in *h* tend to be female.”

Exercise:

Define a new feature extraction function that includes features for two-letter suffixes, such as the one here:

```
>>> def gender_features3(word):  
    return {'suffix1': word[-1],  
            'suffix2': word[-2]}
```

Using the same training, development test, and test data already defined in the variables `train_names`, `devtest_names`, and `test_names`, make new `train_set`, `devtest_set` and `test_set` variables. Carry out the classification and look at the errors in the development test set.

Is this classification more accurate? Can you see error examples that you could use to make new features to improve?

Will this work for names in your mother tongue? How would you set about trying it out?