

✓ In Depth: Naive Bayes Classification

The previous four sections have given a general overview of the concepts of machine learning. In this section and the ones that follow, we will be taking a closer look at several specific algorithms for supervised and unsupervised learning, starting here with naive Bayes classification.

Naive Bayes models are a group of extremely fast and simple classification algorithms that are often suitable for very high-dimensional datasets. Because they are so fast and have so few tunable parameters, they end up being very useful as a quick-and-dirty baseline for a classification problem. This section will focus on an intuitive explanation of how naive Bayes classifiers work, followed by a couple examples of them in action on some datasets.

✓ Bayesian Classification

Naive Bayes classifiers are built on Bayesian classification methods. These rely on Bayes's theorem, which is an equation describing the relationship of conditional probabilities of statistical quantities. In Bayesian classification, we're interested in finding the probability of a label given some observed features, which we can write as $P(L \mid \text{features})$. Bayes's theorem tells us how to express this in terms of quantities we can compute more directly:

$$P(L \mid \text{features}) = \frac{P(\text{features} \mid L)P(L)}{P(\text{features})}$$

If we are trying to decide between two labels—let's call them L_1 and L_2 —then one way to make this decision is to compute the ratio of the posterior probabilities for each label:

$$\frac{P(L_1 \mid \text{features})}{P(L_2 \mid \text{features})} = \frac{P(\text{features} \mid L_1) P(L_1)}{P(\text{features} \mid L_2) P(L_2)}$$

All we need now is some model by which we can compute $P(\text{features} \mid L_i)$ for each label. Such a model is called a *generative model* because it specifies the hypothetical random process that generates the data. Specifying this generative model for each label is the main piece of the training of such a Bayesian classifier. The general version of such a training step is a very difficult task, but we can make it simpler through the use of some simplifying assumptions about the form of this model.

This is where the "naive" in "naive Bayes" comes in: if we make very naive assumptions about the generative model for each label, we can find a rough approximation of the generative model for each class, and then proceed with the Bayesian classification. Different types of naive Bayes classifiers rest on different naive assumptions about the data, and we will examine a few of these in the following sections.

We begin with the standard imports:

```
%matplotlib inline
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns; sns.set()
from sklearn.datasets import make_blobs
from sklearn.naive_bayes import GaussianNB
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.naive_bayes import MultinomialNB
from sklearn.pipeline import make_pipeline
import pandas as pd
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
from sklearn.feature_selection import SelectKBest, chi2
```

✓ Gaussian Naive Bayes

Perhaps the easiest naive Bayes classifier to understand is Gaussian naive Bayes. In this classifier, the assumption is that *data from each label is drawn from a simple Gaussian distribution*. Imagine that you have the following data:

```
# from sklearn.datasets import make_blobs
X, y = make_blobs(100, 2, centers=2, random_state=2, cluster_std=1.5)
plt.scatter(X[:, 0], X[:, 1], c=y, s=50, cmap='RdBu');
```



One extremely fast way to create a simple model is to assume that the data is described by a Gaussian distribution with no covariance between dimensions. This model can be fit by simply finding the mean and standard deviation of the points within each label, which is all you need to define such a distribution. The result of this naive Gaussian assumption is shown in the following figure:

The ellipses here represent the Gaussian generative model for each label, with larger probability toward the center of the ellipses. With this generative model in place for each class, we have a simple recipe to compute the likelihood $P(\text{features} \mid L_1)$ for any data point, and thus we can quickly compute the posterior ratio and determine which label is the most probable for a given point.

This procedure is implemented in Scikit-Learn's `sklearn.naive_bayes.GaussianNB` estimator:

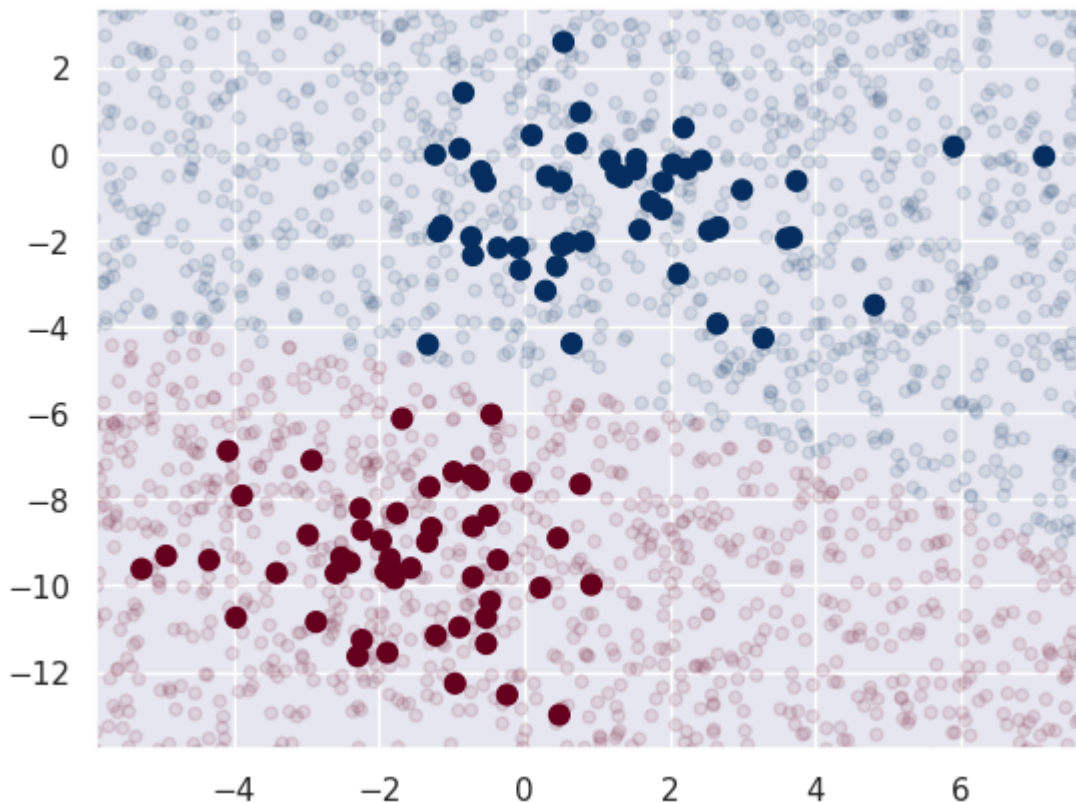
```
# from sklearn.naive_bayes import GaussianNB
model = GaussianNB()
model.fit(X, y);
```

Now let's generate some new data and predict the label:

```
rng = np.random.RandomState(0)
Xnew = [-6, -14] + [14, 18] * rng.rand(2000, 2)
ynew = model.predict(Xnew)
```

Now we can plot this new data to get an idea of where the decision boundary is:

```
plt.scatter(X[:, 0], X[:, 1], c=y, s=50, cmap='RdBu')
lim = plt.axis()
plt.scatter(Xnew[:, 0], Xnew[:, 1], c=ynew, s=20, cmap='RdBu', alpha=0.1)
plt.axis(lim);
```



We see a slightly curved boundary in the classifications—in general, the boundary in Gaussian naive Bayes is quadratic.

A nice piece of this Bayesian formalism is that it naturally allows for probabilistic classification, which we can compute using the `predict_proba` method:

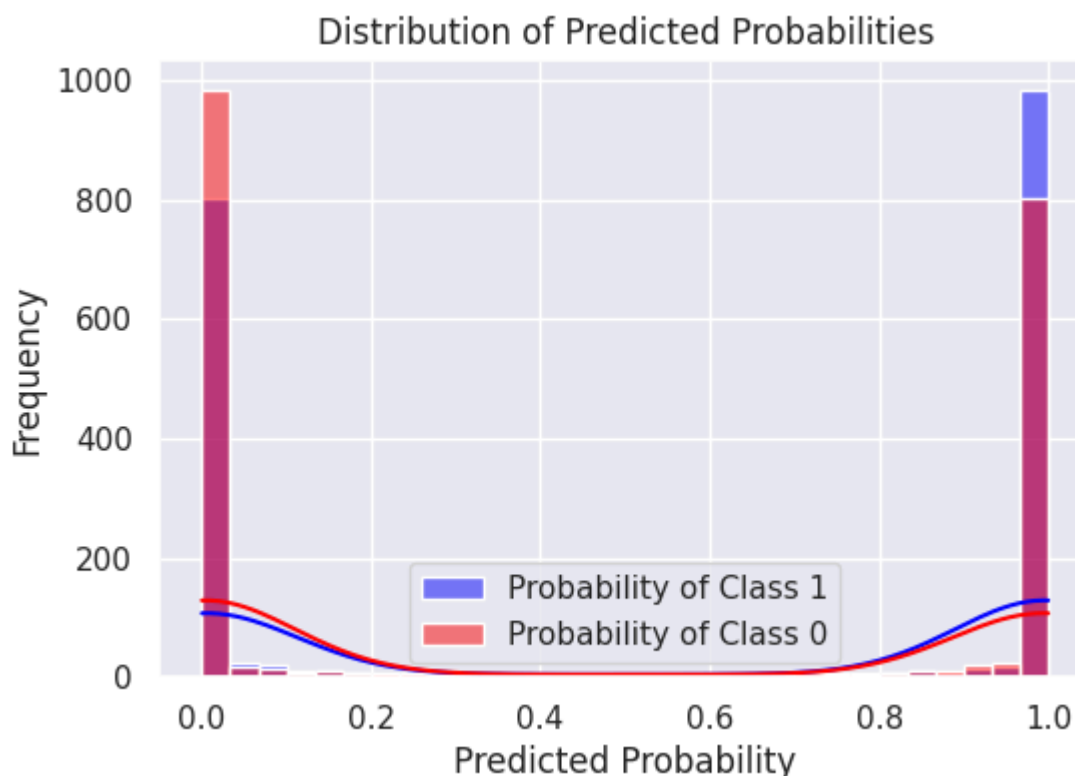
```
yprob = model.predict_proba(Xnew)
yprob[-8:].round(2)
```

```
array([[0.89, 0.11],
       [1.  , 0.  ],
       [1.  , 0.  ],
       [1.  , 0.  ],
       [1.  , 0.  ],
       [1.  , 0.  ],
       [0.  , 1.  ],
       [0.15, 0.85]])
```

Changed.

```
# Plot probability distributions
plt.figure(figsize=(6,4))
sns.histplot(yprob[:, 1], bins=30, kde=True, color='blue', label="Probability of Class 1")
sns.histplot(yprob[:, 0], bins=30, kde=True, color='red', label="Probability of Class 0")
plt.xlabel("Predicted Probability")
plt.ylabel("Frequency")
plt.legend()
plt.title("Distribution of Predicted Probabilities")
```

```
plt.show()
```



The columns give the posterior probabilities of the first and second label, respectively. If you are looking for estimates of uncertainty in your classification, Bayesian approaches like this can be a useful approach.

Of course, the final classification will only be as good as the model assumptions that lead to it, which is why Gaussian naive Bayes often does not produce very good results. Still, in many cases—especially as the number of features becomes large—this assumption is not detrimental enough to prevent Gaussian naive Bayes from being a useful method.

Changed.

```
# from sklearn.model_selection import train_test_split
```

```
X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.4, random_state=42)
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_state=42)
```

Changed.

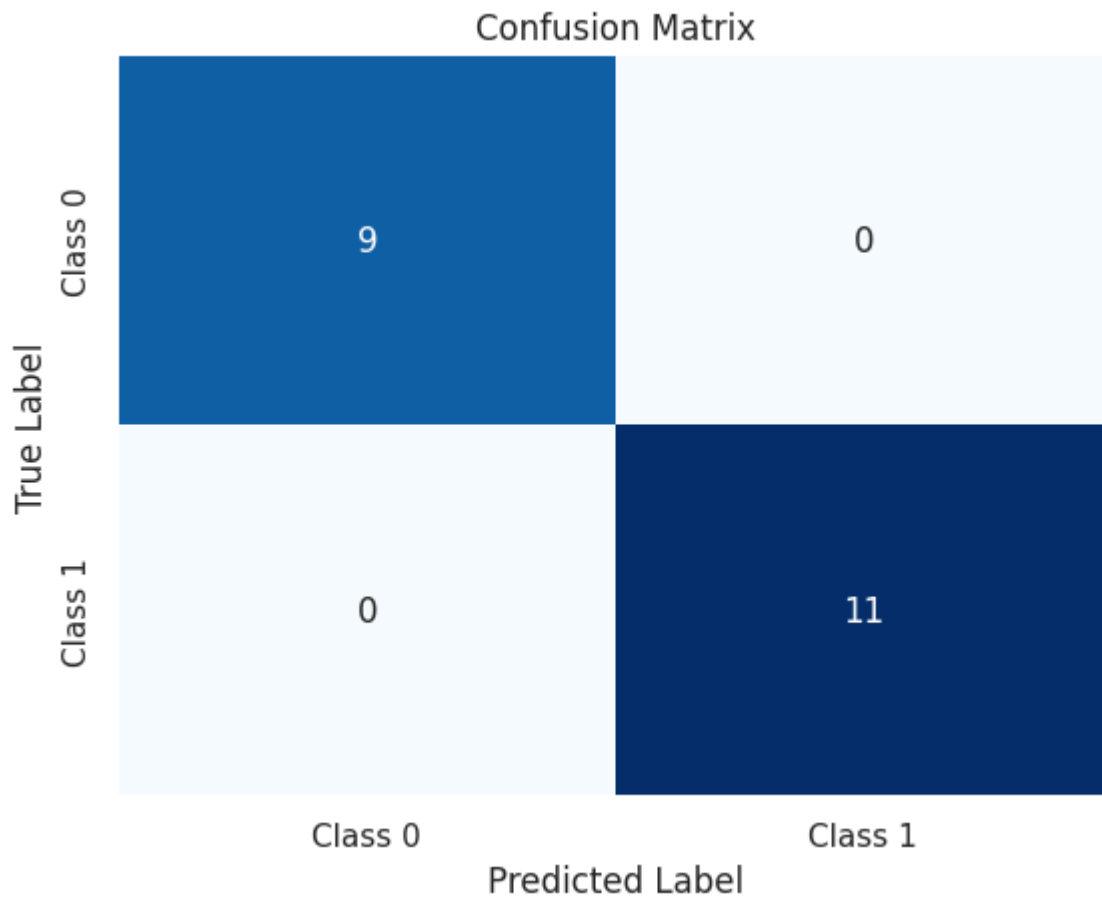
```
model = GaussianNB()
model.fit(X_train, y_train)
```

```
y_val_pred = model.predict(X_val)
```

Changed.

```
cm = confusion_matrix(y_val_pred, y_val)
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False, xticklabels=["Class 0", "C
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.title("Confusion Matrix")
```

```
➞ Text(0.5, 1.0, 'Confusion Matrix')
```

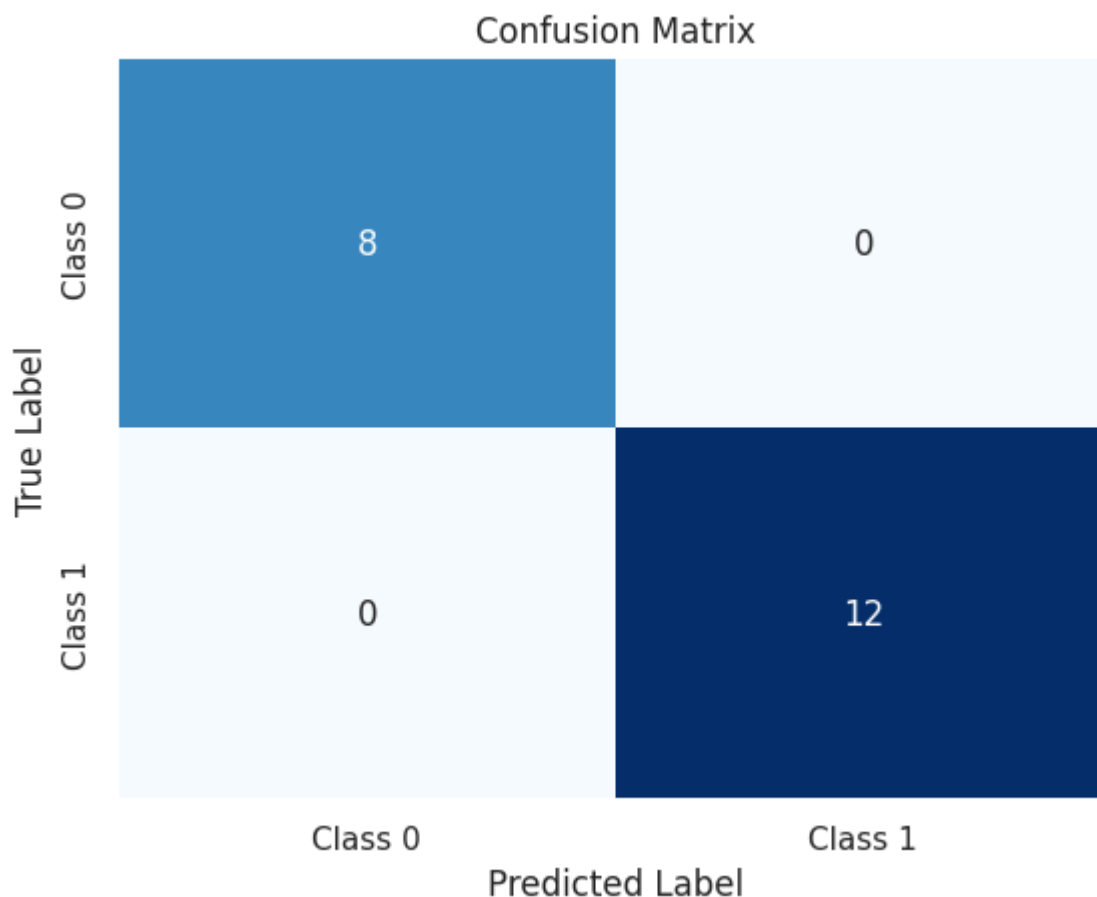


Changed.

```
# from sklearn.metrics import confusion_matrix
y_pred = model.predict(X_test)

cm = confusion_matrix(y_test, y_pred)
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False, xticklabels=["Class 0", "C
plt.xlabel("Predicted Label")
plt.ylabel("True Label")
plt.title("Confusion Matrix")
```

```
Text(0.5, 1.0, 'Confusion Matrix')
```



✓ Multinomial Naive Bayes

The Gaussian assumption just described is by no means the only simple assumption that could be used to specify the generative distribution for each label. Another useful example is multinomial naive Bayes, where the features are assumed to be generated from a simple multinomial distribution. The multinomial distribution describes the probability of observing counts among a number of categories, and thus multinomial naive Bayes is most appropriate for features that represent counts or count rates.

The idea is precisely the same as before, except that instead of modeling the data distribution with the best-fit Gaussian, we model the data distribution with a best-fit multinomial distribution.

✓ Example: Classifying Text

One place where multinomial naive Bayes is often used is in text classification, where the features are related to word counts or frequencies within the documents to be classified. We discussed the extraction of such features from text in [Feature Engineering](#); here we will use the sparse word count features from the 20 Newsgroups corpus to show how we might classify these short documents into categories.

Let's download the data and take a look at the target names:

```
# from sklearn.datasets import fetch_20newsgroups
```

```
data = fetch_20newsgroups()
data.target_names
```

```

↳ ['alt.atheism',
   'comp.graphics',
   'comp.os.ms-windows.misc',
   'comp.sys.ibm.pc.hardware',
   'comp.sys.mac.hardware',
   'comp.windows.x',
   'misc.forsale',
   'rec.autos',
   'rec.motorcycles',
   'rec.sport.baseball',
   'rec.sport.hockey',
   'sci.crypt',
   'sci.electronics',
   'sci.med',
   'sci.space',
   'soc.religion.christian',
   'talk.politics.guns',
   'talk.politics.mideast',
   'talk.politics.misc',
   'talk.religion.misc']

```

For simplicity here, we will select just a few of these categories, and download the training and testing set:

```
categories = ['talk.religion.misc', 'soc.religion.christian',
              'sci.space', 'comp.graphics']
train = fetch_20newsgroups(subset='train', categories=categories)
test = fetch_20newsgroups(subset='test', categories=categories)
```

Here is a representative entry from the data:

```
print(train.data[5])
```

```

↳ From: dmcgee@uluhe.soest.hawaii.edu (Don McGee)
   Subject: Federal Hearing
   Originator: dmcgee@uluhe
   Organization: School of Ocean and Earth Science and Technology
   Distribution: usa
   Lines: 10

```

Fact or rumor....? Madalyn Murray O'Hare an atheist who eliminated the use of the bible reading and prayer in public schools 15 years ago is now going to appear before the FCC with a petition to stop the reading of the Gospel on the airways of America. And she is also campaigning to remove Christmas programs, songs, etc from the public schools. If it is true then mail to Federal Communications Commission 1919 H Street Washington DC 20054 expressing your opposition to her request. Reference Petition number

2493.

In order to use this data for machine learning, we need to be able to convert the content of each string into a vector of numbers. For this we will use the TF-IDF vectorizer (discussed in [Feature Engineering](#)), and create a pipeline that attaches it to a multinomial naive Bayes classifier:

```
#from sklearn.feature_extraction.text import TfidfVectorizer
#from sklearn.naive_bayes import MultinomialNB
#from sklearn.pipeline import make_pipeline

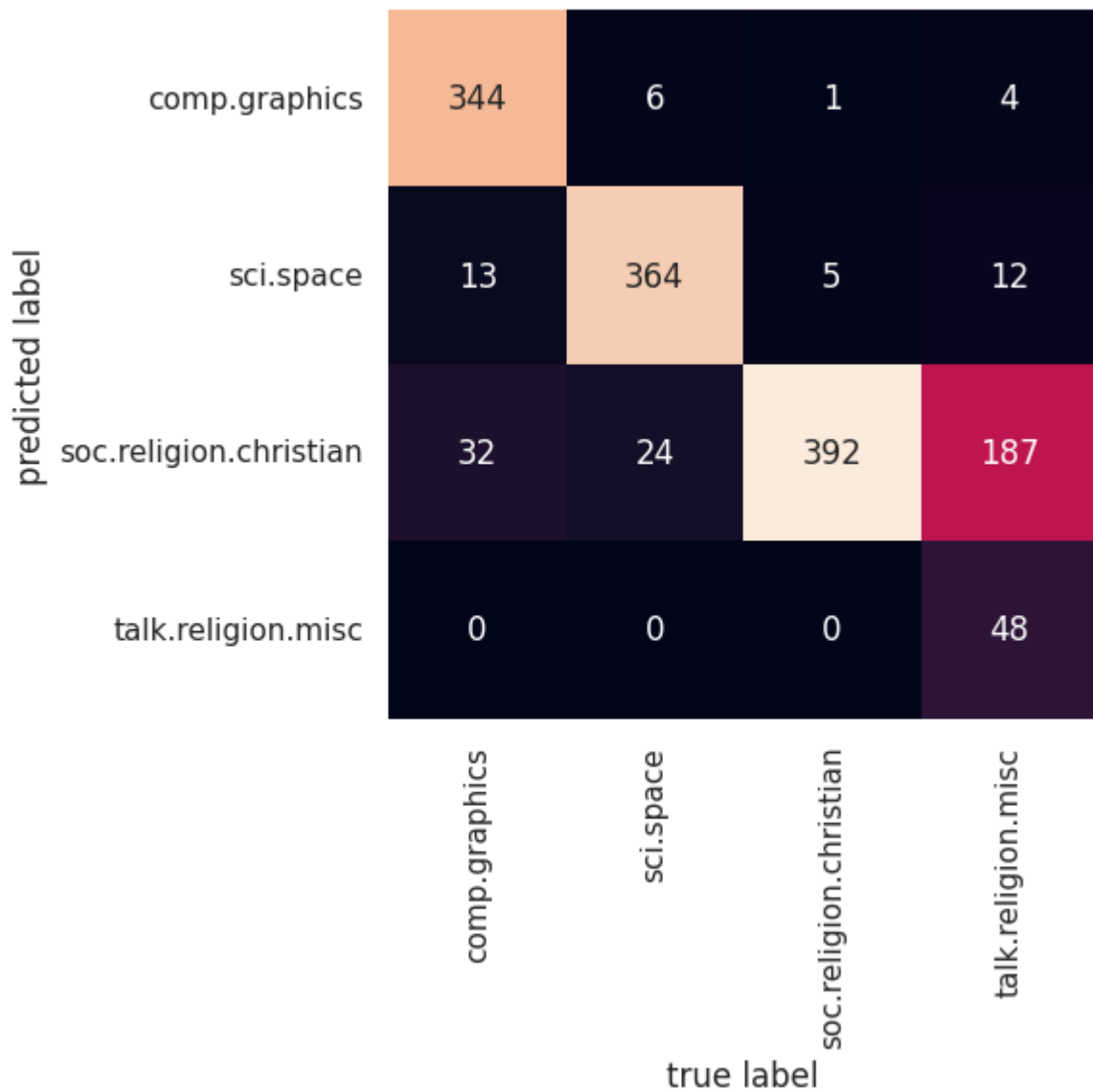
model = make_pipeline(TfidfVectorizer(), MultinomialNB())
```

With this pipeline, we can apply the model to the training data, and predict labels for the test data:

```
model.fit(train.data, train.target)
labels = model.predict(test.data)
```

Now that we have predicted the labels for the test data, we can evaluate them to learn about the performance of the estimator. For example, here is the confusion matrix between the true and predicted labels for the test data:

```
# from sklearn.metrics import confusion_matrix
mat = confusion_matrix(test.target, labels)
sns.heatmap(mat.T, square=True, annot=True, fmt='d', cbar=False,
            xticklabels=train.target_names, yticklabels=train.target_names)
plt.xlabel('true label')
plt.ylabel('predicted label');
```



Changed.

```
# from sklearn.metrics import accuracy_score
print("Accuracy:", accuracy_score(test.target, labels))
```



Accuracy: 0.8016759776536313

Changed.

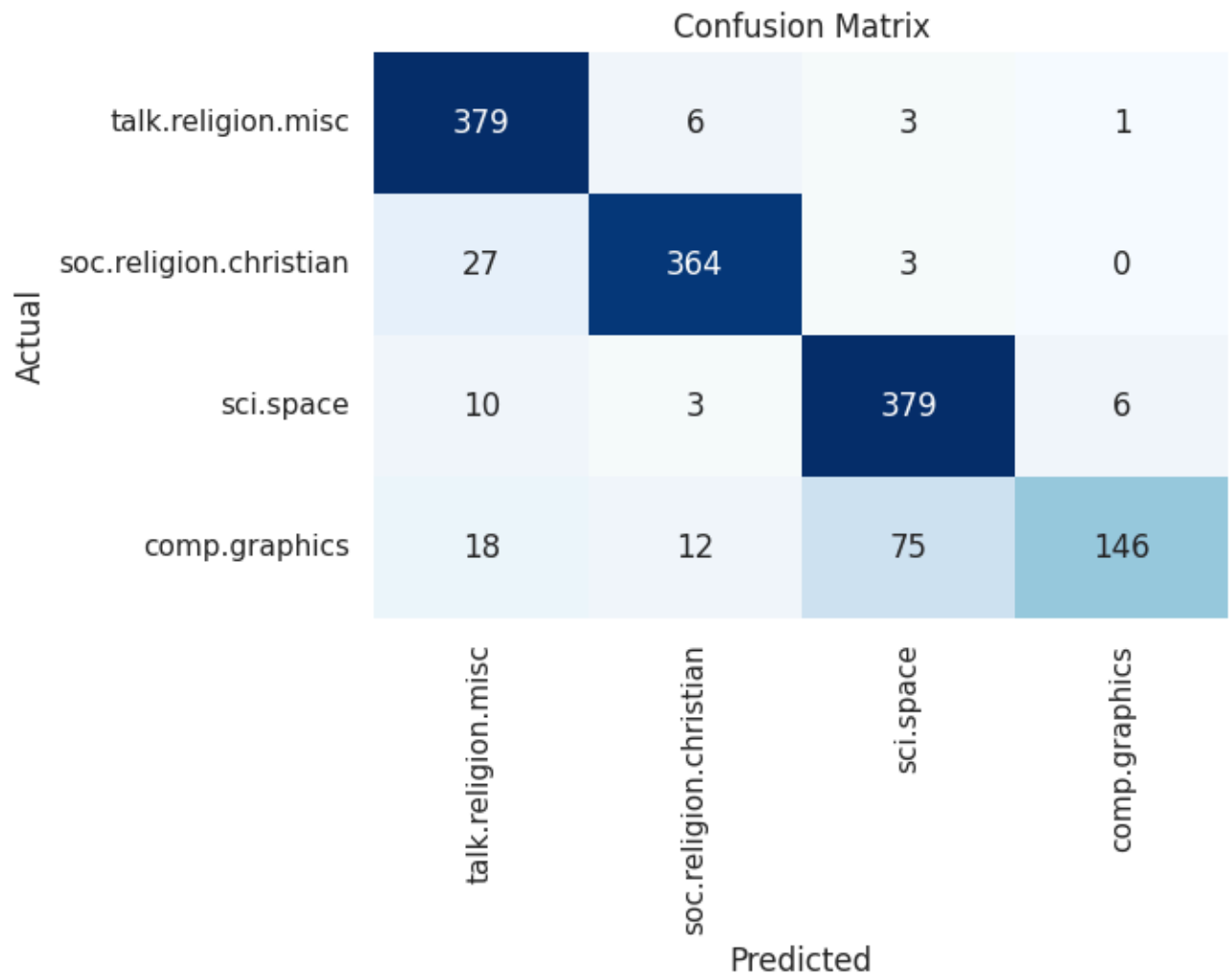
```
model1 = make_pipeline(TfidfVectorizer(stop_words='english', ngram_range=(1,2), max_featu
model1.fit(train.data, train.target)

labels1 = model1.predict(test.data)

cm = confusion_matrix(test.target, labels1)
plt.figure(figsize=(6, 4))
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False, xticklabels=set(train.targ
plt.xlabel("Predicted")
plt.ylabel("Actual")
```

```
plt.title("Confusion Matrix")
plt.show()
```

```
print("Accuracy:", accuracy_score(test.target, labels1))
```



Accuracy: 0.8854748603351955

Changed.

```
# from sklearn.model_selection import GridSearchCV

# Hyperparameter tuning (alpha tuning for MultinomialNB)
param_grid = {'multinomialnb__alpha': [0.01, 0.1, 1, 10]}
grid = GridSearchCV(model1, param_grid, cv=5)
grid.fit(train.data, train.target)

print("Best alpha:", grid.best_params_['multinomialnb__alpha'])
```



Best alpha: 0.01

Changed.

```
#best model
```

```
# Train the final model with best alpha
```

```
best_model = make_pipeline(TfidfVectorizer(stop_words='english', ngram_range=(1,2), max_f
```

```
best_model.fit(train.data, train.target)
```

```
predictions = best_model.predict(test.data)
```

Changed.

```
print("Accuracy:", accuracy_score(test.target, predictions))
```

```
➡ Accuracy: 0.9120111731843575
```

Changed.

```
cm = confusion_matrix(test.target, predictions)
```

```
plt.figure(figsize=(6, 4))
```

```
sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False, xticklabels=set(train.targ
```

```
plt.xlabel("Predicted")
```

```
plt.ylabel("Actual")
```

```
plt.title("Confusion Matrix")
```

```
plt.show()
```



Confusion Matrix

Actual	talk.religion.misc	373	8	3	5
	soc.religion.christian	24	368	1	1
	sci.space	9	3	375	11
	comp.graphics	9	10	42	190
		talk.religion.misc	soc.religion.christian	sci.space	comp.graphics
		Predicted			

Evidently, even this very simple classifier can successfully separate space talk from computer talk, but it gets confused between talk about religion and talk about Christianity. This is perhaps an expected area of confusion!

The very cool thing here is that we now have the tools to determine the category for *any* string, using the `predict()` method of this pipeline. Here's a quick utility function that will return the prediction for a single string:

```
def predict_category(s, train=train, model=model):
    pred = model.predict([s])
    return train.target_names[pred[0]]
```

Let's try it out:

```
predict_category('sending a payload to the ISS')
```



```
'sci.space'
```

```
predict_category('discussing islam vs atheism')
```

→ 'soc.religion.christian'

`predict_category('determining the screen resolution')`

→ 'comp.graphics'

Changed.

```
def predict_category_best(s, train=train, model=best_model):
    pred_best = best_model.predict([s])
    return train.target_names[pred_best[0]]
```

`predict_category_best('sending a payload to the ISS')`

→ 'sci.space'

`predict_category_best('discussing islam vs atheism')`

→ 'soc.religion.christian'

`predict_category_best('determining the screen resolution')`

→ 'comp.graphics'

Changed.

```
def predict_category_with_prob(s, train=train, model=best_model):
    pred = best_model.predict([s])[0]
    prob = best_model.predict_proba([s])[0]
    category = train.target_names[pred]

    # Show top 3 most probable categories
    top_categories = sorted(
        zip(train.target_names, prob), key=lambda x: x[1], reverse=True
    )[:3]

    print(f"Predicted Category: {category}")
    print("Top 3 Predictions:")
    for cat, p in top_categories:
        print(f"- {cat}: {p:.4f}")

    return category
```

Remember that this is nothing more sophisticated than a simple probability model for the (weighted) frequency of each word in the string; nevertheless, the result is striking. Even a very naive algorithm, when used carefully and trained on a large set of high-dimensional data, can be surprisingly effective.

```
predict_category_with_prob('sending a payload to the ISS')
```

```
➞ Predicted Category: sci.space  
Top 3 Predictions:  
- sci.space: 0.6983  
- comp.graphics: 0.1344  
- soc.religion.christian: 0.0992  
'sci.space'
```

```
predict_category_with_prob('discussing islam vs atheism')
```

```
➞ Predicted Category: soc.religion.christian  
Top 3 Predictions:  
- soc.religion.christian: 0.6039  
- talk.religion.misc: 0.3080  
- comp.graphics: 0.0800  
'soc.religion.christian'
```

```
predict_category_with_prob('determining the screen resolution')
```

```
➞ Predicted Category: comp.graphics  
Top 3 Predictions:  
- comp.graphics: 0.9513  
- sci.space: 0.0311  
- talk.religion.misc: 0.0090  
'comp.graphics'
```

✓ When to Use Naive Bayes

Because naive Bayesian classifiers make such stringent assumptions about data, they will generally not perform as well as a more complicated model. That said, they have several advantages:

- They are extremely fast for both training and prediction
- They provide straightforward probabilistic prediction
- They are often very easily interpretable
- They have very few (if any) tunable parameters

These advantages mean a naive Bayesian classifier is often a good choice as an initial baseline classification. If it performs suitably, then congratulations: you have a very fast, very interpretable classifier for your problem. If it does not perform well, then you can begin exploring more sophisticated models, with some baseline knowledge of how well they should perform.

Naive Bayes classifiers tend to perform especially well in one of the following situations:

- When the naive assumptions actually match the data (very rare in practice)
- For very well-separated categories, when model complexity is less important
- For very high-dimensional data, when model complexity is less important

The last two points seem distinct, but they actually are related: as the dimension of a dataset grows, it is much less likely for any two points to be found close together (after all, they must be close in *every single dimension* to be close overall). This means that clusters in high dimensions tend to be more separated, on average, than clusters in low dimensions, assuming the new dimensions actually add information. For this reason, simplistic classifiers like naive Bayes tend to work as well or better than more complicated classifiers as the dimensionality grows: once you have enough data, even a simple model can be very powerful.

Changed.

Application of Naive Bayes - Breast cancer dataset

Unsupported cell type. Double-click to inspect/edit the content.

```
# Load the Breast Cancer dataset
data = load_breast_cancer()

X, y = data.data, data.target

# Convert to DataFrame
df = pd.DataFrame(X, columns=data.feature_names)
df['target'] = y

# Check for missing values
print("Missing Values:\n", df.isnull().sum())
```

```
➡ Missing Values:
  mean radius      0
  mean texture     0
  mean perimeter   0
  mean area        0
  mean smoothness  0
  mean compactness 0
  mean concavity    0
  mean concave points 0
  mean symmetry     0
  mean fractal dimension 0
  radius error      0
  texture error     0
  perimeter error   0
  area error        0
  smoothness error  0
  compactness error 0
  concavity error   0
  concave points error 0
  symmetry error    0
  fractal dimension error 0
  worst radius      0
  worst texture     0
  worst perimeter   0
  worst area        0
```



```

worst smoothness      0
worst compactness     0
worst concavity        0
worst concave points  0
worst symmetry         0
worst fractal dimension 0
target                0
dtype: int64

```

```
# Check class balance
```

```
print("Class Distribution:\n", df['target'].value_counts())
```

```

↪ Class Distribution:
   target
1      357
0      212
Name: count, dtype: int64

```

```
# Split the dataset into training and testing sets
```

```
X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.4, random_state=42)
```

```
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_sta
```

```
# Initialize and train the Naïve Bayes model
```

```
model = GaussianNB()
```

```
model.fit(X_train, y_train)
```

```
# Make predictions
```

```
y_val_pred = model.predict(X_val)
```

```
# Evaluate the model
```

```
accuracy = accuracy_score(y_val, y_val_pred)
```

```
print("Accuracy:", accuracy)
```

```
print("Classification Report:\n", classification_report(y_val, y_val_pred))
```

```
print("Confusion Matrix:\n", confusion_matrix(y_val, y_val_pred))
```

```

↪ Accuracy: 0.9736842105263158
Classification Report:
              precision    recall  f1-score   support

     0       0.98         0.95         0.96         42
     1       0.97         0.99         0.98         72

   accuracy          0.97              114
  macro avg          0.97              114
weighted avg          0.97              114

Confusion Matrix:
[[40  2]
 [ 1 71]]

```

```
# Feature Selection using SelectKBest (Chi-Square Test)
```

```
selector = SelectKBest(chi2, k=10)
```

```
X_selected = selector.fit_transform(X, y)
```

```
selected_features = np.array(data.feature_names)[selector.get_support()]
```