

Using our custom estimator

使用我们的自定义评估器

Let's try this custom estimator on a problem we have seen before: the classification of hand-written digits. Here we will load the digits, and compute the cross-validation score for a range of candidate bandwidths using the `GridSearchCV` meta-estimator (refer back to [Hyperparameters and Model Validation](#)).

下面我们试一下这个自定义评估器，使用前面我们研究过的问题：手写数字分类。我们载入手写数字数据，然后针对一定范围的带宽值使用 `GridSearchCV` 元评估器计算交叉验证结果（参见[超参数和模型验证](#)）：

译者注：下面代码做了修改以适应新版本Scikit-Learn，包括`GridSearchCV`从属的包，参数`cv`和结果中使用`cv_result_`字典取分值。

```
In [17]: from sklearn.datasets import load_digits
         from sklearn.model_selection import GridSearchCV

digits = load_digits()

bandwidths = 10 ** np.linspace(0, 2, 100)
grid = GridSearchCV(KernelDensity(), {'bandwidth': bandwidths}, cv=5)
grid.fit(digits.data, digits.target)

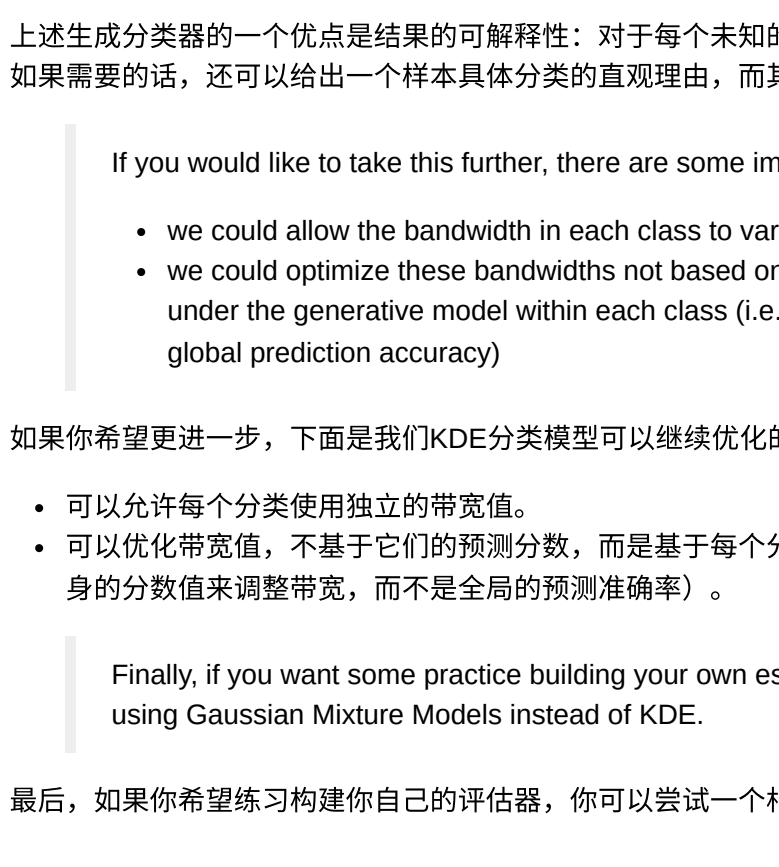
scores = grid.cv_results_['mean_test_score']
```

Next we can plot the cross-validation score as a function of bandwidth:

接下来我们可以绘制交叉验证分值与带宽之间的函数图像：

```
In [18]: plt.semilogx(bandwidths, scores)
         plt.xlabel('bandwidth')
         plt.ylabel('accuracy')
         plt.title('KDE Model Performance')
         print(grid.best_params_)
         print("accuracy = ", grid.best_score_)

{'bandwidth': 6.135867273413174}
accuracy = 0.9677298659319276
```



We see that this not-so-naive Bayesian classifier reaches a cross-validation accuracy of just over 96%; this is compared to around 80% for the naive Bayesian classification.

我们看到这个不那么朴素的贝叶斯分类器达到了交叉验证准确率超过96%；而朴素贝叶斯分类只有大约80%：

```
In [19]: from sklearn.naive_bayes import GaussianNB
         from sklearn.model_selection import cross_val_score
         cross_val_score(GaussianNB(), digits.data, digits.target, cv=5).mean()

Out[19]: 0.8669281956658759
```

One benefit of such a generative classifier is interpretability of results: for each unknown sample, we not only get a probabilistic classification, but a *full model* of the distribution of points we are comparing it to! If desired, this offers an intuitive window into the reasons for a particular classification that algorithms like SVMs and random forests tend to obscure.

上述生成分类器的一个优点是结果的可解释性：对于每个未知的样本，我们不但得到了概率分类，还获得了数据点分布情况的完整模型。如果需要的话，还可以给出一个样本具体分类的直观理由，而其他算法像SVM和随机森林在这点上通常是模糊的。

If you would like to take this further, there are some improvements that could be made to our KDE classifier model:

- we could allow the bandwidth in each class to vary independently
- we could optimize these bandwidths not based on their prediction score, but on the likelihood of the training data under the generative model within each class (i.e. use the scores from `KernelDensity` itself rather than the global prediction accuracy)

如果你希望更进一步，下面是我们KDE分类模型可以继续优化的一些建议：

- 可以允许每个分类使用独立的带宽值。
- 可以优化带宽值，不基于它们的预测分数，而是基于每个分类在生成模型下拟合训练数据的似然值（也就是使用 `KernelDensity` 本身的分数值来调整带宽，而不是全局的预测准确率）。

Finally, if you want some practice building your own estimator, you might tackle building a similar Bayesian classifier using Gaussian Mixture Models instead of KDE.

最后，如果你希望练习构建你自己的评估器，你可以尝试一个相似的贝叶斯分类器，使用高斯混合模型而不是KDE。

< [深入：高斯混合模型](#) | [目录](#) | [应用：脸部识别验证](#) >

[Open in Colab](#)