





This is a valuable diagnostic, because it gives us a visual depiction of how our model responds to increasing training data. In particular, when your learning curve has already converged (i.e., when the training and validation curves are already close to each other) *adding more training data will not significantly improve the fit!* This situation is seen in the left panel, with the learning curve for the degree-2 model.

这是一项非常有价值的分析，因为它为我们提供了模型随着训练数据增加发生性能变化的可视化展示。而且当你的学习曲线已经收敛时（例如当训练和验证曲线已经非常接近的情况下），*增加更多的训练数据不会显著的提升拟合度*。这个结论很容易从左图二阶模型的学习曲线中获得。

The only way to increase the converged score is to use a different (usually more complicated) model. We see this in the right panel: by moving to a much more complicated model, we increase the score of convergence (indicated by the dashed line), but at the expense of higher model variance (indicated by the difference between the training and validation scores). If we were to add even more data points, the learning curve for the more complicated model would eventually converge.

要提升已经收敛的学习曲线的性能唯一方法就是使用一个不同的（通常更复杂的）模型。我们可以从右图中看到：当使用了复杂的多的模型后，我们将收敛的分数值（使用虚线表示）提升了，付出的代价是更高的模型方差（图中训练曲线和验证曲线的间距）。如果我们继续增加更多的样本，更复杂模型的学习曲线最终也会收敛。

Plotting a learning curve for your particular choice of model and dataset can help you to make this type of decision about how to move forward in improving your analysis.

绘制模型和数据集的学习曲线能帮助你作出进一步改善性能的决定。

## Validation in Practice: Grid Search

### 验证实践：网格搜索

The preceding discussion is meant to give you some intuition into the trade-off between bias and variance, and its dependence on model complexity and training set size. In practice, models generally have more than one knob to turn, and thus plots of validation and learning curves change from lines to multi-dimensional surfaces. In these cases, such visualizations are difficult and we would rather simply find the particular model that maximizes the validation score.

前面的讨论意在为你提供直观的偏差和方差权衡的知识，它取决于模型复杂度和训练集规模。在实践中，模型通常有多于一个开关进行调节，因此前面关于验证曲线和学习曲线的二维线条就会变成多维平面。在这些情况下，要将它可视化出来是很困难的，并且我们更希望简单的找到特定模型能最大化验证分数。

Scikit-Learn provides automated tools to do this in the grid search module. Here is an example of using grid search to find the optimal polynomial model. We will explore a three-dimensional grid of model features; namely the polynomial degree, the flag telling us whether to fit the intercept, and the flag telling us whether to normalize the problem. This can be set up using Scikit-Learn's `GridSearchCV` meta-estimator:

Scikit-Learn提供了自动化的工具来完成这项任务，它们在网格搜索模块中。下面是一个使用网格搜索找到最优多项式模型的例子。我们会探索模型特征的一个三维网格：包括多项式阶数，一个是否拟合截距的标志和一个是否归一化问题的标志。这可以通过Scikit-Learn的`GridSearchCV`元评估器来设置：

```
In [18]: from sklearn.model_selection import GridSearchCV

param_grid = {'polynomialfeatures__degree': np.arange(21),
              'linearregression__fit_intercept': [True, False],
              'linearregression__normalize': [True, False]}

grid = GridSearchCV(PolynomialRegression(), param_grid, cv=7)
```

Notice that like a normal estimator, this has not yet been applied to any data. Calling the `fit()` method will fit the model at each grid point, keeping track of the scores along the way.

网格搜索模型和普通模型一样，实例化后还未应用到任何数据集上。通过调用`fit()`方法会将模型的每个网格点拟合到数据集上，同时过程中保存了验证的分数：

```
In [19]: grid.fit(X, y);
```

Now that this is fit, we can ask for the best parameters as follows:

拟合完后，我们可以使用下面代码来获得最佳参数：

```
In [20]: grid.best_params_

Out[20]: {'linearregression__fit_intercept': False,
          'linearregression__normalize': True,
          'polynomialfeatures__degree': 4}
```

Finally, if we wish, we can use the best model and show the fit to our data using code from before:

最终，需要的话，我们可以使用代码将最佳模型、数据及它们的拟合情况绘制出来：

```
In [21]: model = grid.best_estimator_

plt.scatter(X.ravel(), y)
lim = plt.axis()
y_fit = model.fit(X, y).predict(X_test)
plt.plot(X_test.ravel(), y_test);
plt.axis(lim);
```



The grid search provides many more options, including the ability to specify a custom scoring function, to parallelize the computations, to do randomized searches, and more. For information, see the examples in [10-Depend\\_Kernel\\_Density\\_Estimation](#) and [Feature\\_Engineering\\_Working\\_with\\_Images](#), or refer to Scikit-Learn's [grid search documentation](#).

网格搜索提供很多其他参数，包括指定自定义的评分函数，并行化计算和执行随机搜索等等。需要更多信息，参见[深入：核密度估计和特征工程](#)，或者参考Scikit-Learn的[网格搜索在线文档](#)。

## Summary

### 总结

In this section, we have begun to explore the concept of model validation and hyperparameter optimization, focusing on intuitive aspects of the bias–variance trade-off and how it comes into play when fitting models to data. In particular, we found that the use of a validation set or cross-validation approach is vital when tuning parameters in order to avoid over-fitting for more complex/flexible models.

在本节中，我们开始探讨模型验证和超参数优化的概念，聚焦在偏差方差权衡的直观概念和它在模型拟合数据时扮演的角色。特别是，我们强调使用测试集验证和交叉验证方法的重要性，当在复杂灵活模型中调节参数时要避免过拟合。

In later sections, we will discuss the details of particularly useful models, and throughout will talk about what tuning is available for these models and how these free parameters affect model complexity. Keep the lessons of this section in mind as you read on and learn about these machine learning approaches!

在后续章节中，我们会讨论每种模型的细节，并在过程中介绍这些模型可以调节哪些参数以及这些参数如何影响模型复杂度。请将本节的内容牢记，当你在后面继续学习机器学习方法的时候，本节内容会提供重要的帮助。

< [Scikit-Learn简介](#) | [目录](#) | [继续工程](#) >

[Open in Colab](#)