

## 小k交叉验证过程

具体详细到每一步地讲解一下小k的交叉验证的这个两层嵌套循环的具体过程，每个步骤拆分

直接开拆，不拐弯。

你看到的“小k交叉验证的两层嵌套循环”，**本质就是：**

- 外层循环：枚举一种情况（最常见是“一个超参数组合 / 一个模型设置”）
- 内层循环：在这 k 个折叠（fold）上轮流做“训练 + 验证”

我用一个很具体的例子来讲，每一步都拆开。

### 0. 先把场景说清楚

假设：

- 有 100 条样本数据
- 想做 5 折交叉验证（ $k=5$ ）
- 有 2 个超参数组合要比较：
  - 组合 A：learning\_rate=0.01
  - 组合 B：learning\_rate=0.1

那这个“两层嵌套循环”的逻辑就是：

- 外层：先用组合 A 跑完 5 折，再用组合 B 跑 5 折
- 内层：每一折都要做「划分训练/验证 → 训练模型 → 评估 → 记录得分」

### 1. 数据准备 + 划分成 k 折（在循环外做一次）

步骤 1：整理原始数据

- 有特征矩阵  $X$  ( $100 \times d$ )
- 有标签向量  $y$  ( $100 \times 1$ )

步骤 2：打乱顺序

- 为了避免数据有顺序偏差（比如前一半是一类，后一半是另一类），先随机打乱索引：
  - 原索引：[0, 1, 2, ..., 99]
  - 打乱后：例如 [57, 2, 88, 13, ...]

步骤 3：平均切成 k 份

- $k=5$ ，所以每折大约  $100/5=20$  个样本
- 形成 5 个折（fold）：
  - fold\_1：索引 [打乱后的前 20 个]
  - fold\_2：索引 [接下来的 20 个]
  - ...
  - fold\_5：索引 [最后的 20 个]

注意：这一步只做一次，不在循环里反复切，否则就乱套了。

### 2. 外层循环：枚举超参数 / 模型设置

伪代码长这样：

```
python

for param in param_grid:          # 外层循环，遍历超参数组合
    fold_scores = []               # 存这一组参数在 5 折上的得分
    for i in range(k):            # 内层循环：k 折交叉验证
        ...
```

先把逻辑拆开说一遍再配代码。

步骤 2.1：从参数集合里拿出“当前这一组”

假设  $\text{param\_grid} = [\text{A}, \text{B}]$ ：

- 第一次外层循环： $\text{param} = \text{A}$
- 第二次外层循环： $\text{param} = \text{B}$

步骤 2.2：为这一组参数准备一个“折得分列表”

- 初始化一个空列表： $\text{fold_scores} = []$
- 用来存 [ $\text{fold1\_score}$ ,  $\text{fold2\_score}$ , ...,  $\text{fold5\_score}$ ]

### 3. 内层循环：在 k 个折上轮流做训练和验证

伪代码：

```
python

for i in range(k):                  # i = 0, 1, 2, 3, 4 对应第 1~5 折
    val_idx = folds[i]              # 当前折作验证集
    train_idx = 所有其他折的索引拼在一起

    X_train = X[train_idx]
    y_train = y[train_idx]
    X_val = X[val_idx]
    y_val = y[val_idx]

    model = Model(**param)          # 用当前这组超参数 新建一个模型
    model.fit(X_train, y_train)      # 在训练集上训练

    y_pred = model.predict(X_val)    # 在验证集上预测
    score = metric(y_val, y_pred)    # 算一个评价指标，例如准确率 / MSE 等

    fold_scores.append(score)        # 把这一折的成绩存起来
```

我再用“非常啰嗦的中文”复述一遍每一步：

步骤 3.1：选当前哪一折做验证集

- 内层循环第  $i$  次：
  - 第  $i$  个 fold 为 验证集：
  - 剩下的  $k-1$  个 fold 合成 训练集：
  - $\text{train\_idx} = \text{folds}[0] + \dots + \text{folds}[i-1] + \text{folds}[i+1] + \dots + \text{folds}[k-1]$

步骤 3.2：根据索引切出训练/验证数据

- $\text{X\_train} = \text{X}[\text{train\_idx}]$
- $\text{y\_train} = \text{y}[\text{train\_idx}]$
- $\text{X\_val} = \text{X}[\text{val\_idx}]$
- $\text{y\_val} = \text{y}[\text{val\_idx}]$

本质：同一批原始数据，只是每一折轮流谁当验证，谁当训练。

步骤 3.3：构造模型（带上当前外层的超参数设置）

- 用当前  $\text{param}$  创建模型：
  - 比如  $\text{param} = \{\text{lr': 0.01, 'max\_depth': 5}\}$
  - $\text{model} = \text{Model}(\text{lr}=0.01, \text{max\_depth}=5)$

关键点：

每一折都要重新建一个“干净”的模型实例，不能复用上一次训练过的模型，否则信息就泄漏了。

步骤 3.4：在训练集上拟合

- $\text{model}.fit(\text{X\_train}, \text{y\_train})$

模型只看到训练集的数据，不允许用验证集来训练。

步骤 3.5：在验证集上评估

- 预测：
  - $\text{y\_pred} = \text{model}.predict(\text{X\_val})$
- 计算指标：
  - 分类：准确率、F1、AUC 等
  - 回归：MSE、MAE、R<sup>2</sup> 等
- 得出一个数字： $\text{score}$

步骤 3.6：把这一折的分数存起来

- $\text{fold_scores.append(score)}$

这样跑完  $k$  次之后，就有一个长度为  $k$  的分数列表，比如：

$\text{fold_scores} = [0.81, 0.79, 0.83, 0.80, 0.82]$

### 4. 内层结束：对这一组参数做折平均

内层循环结束后，回到外层：

```
python

mean_score = np.mean(fold_scores)
results.append((param, mean_score))
```

步骤 4.1：算平均分

- $\text{mean\_score} = (0.81 + 0.79 + 0.83 + 0.80 + 0.82) / 5$
- 得到某个平均表现，例如 0.81

步骤 4.2：把这个参数组合 + 它的平均分记下来

比如  $\text{results}$  里放一堆这样的东西：

text

```
[ (A, 0.81),
  (B, 0.78),
  (C, 0.84),
  ...
 ]
```

### 5. 外层继续：换下一组参数，再重复整个 k 折

- 外层循环跳到下一次：
  - param = B
  - 又重新  $\text{fold_scores} = []$
  - 再跑一遍 完整的 k 折内循环
- 对每一组参数都重复上面“3 + 4”的过程

最后，你就有了每组参数在  $k$  折下的平均成绩。

### 6. 选最优参数 + 训练最终模型（经常在循环外做）

通常最后的步骤是：

1. 从  $\text{results}$  中找出平均分最高的一组参数  $\text{param\_best}$

2. 用  $\text{param\_best}$  在  $\text{全量数据 X, y}$  上再训练一次模型：

```
python

from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import cross_val_score, KFold

# 假设你已经有特征和标签
X = ...; y = ...
model = KNeighborsClassifier(n_neighbors=param_best)
```

3. 这个  $\text{best\_model}$  才是你真正拿去上线 / 做预测的模型

这样跑完  $k$  次之后，就有一个长度为  $k$  的分数列表，比如：

$\text{fold_scores} = [0.81, 0.79, 0.83, 0.80, 0.82]$

步骤 4.6：对这一组参数做折平均

```
python

mean_score = np.mean(fold_scores)
results.append((param, mean_score))
```

关键点：

每一折都要重新建一个“干净”的模型实例，不能复用上一次训练过的模型，否则信息就泄漏了。

步骤 4.7：在训练集上拟合

- $\text{model}.fit(\text{X\_train}, \text{y\_train})$

模型只看到训练集的数据，不允许用验证集来训练。

步骤 4.8：在验证集上评估

- 预测：
  - $\text{y\_pred} = \text{model}.predict(\text{X\_val})$
- 计算指标：
  - 分类：准确率、F1、AUC 等
  - 回归：MSE、MAE、R<sup>2</sup> 等
- 得出一个数字： $\text{score}$

步骤 4.9：把这一折的分数存起来

- $\text{fold_scores.append(score)}$

这样跑完  $k$  次之后，就有一个长度为  $k$  的分数列表，比如：

$\text{fold_scores} = [0.81, 0.79, 0.83, 0.80, 0.82]$

步骤 4.10：对这一组参数做折平均

```
python

mean_score = np.mean(fold_scores)
results.append((param, mean_score))
```

关键点：

每一折都要重新建一个“干净”的模型实例，不能复用上一次训练过的模型，否则信息就泄漏了。

步骤 4.11：在训练集上拟合

- $\text{model}.fit(\text{X\_train}, \text{y\_train})$

模型只看到训练集的数据，不允许用验证集来训练。

步骤 4.12：在验证集上评估

- 预测：
  - $\text{y\_pred} = \text{model}.predict(\text{X\_val})$
- 计算指标：
  - 分类：准确率、F1、AUC 等
  - 回归：MSE、MAE、R<sup>2</sup> 等
- 得出一个数字： $\text{score}$

步骤 4.13：把这一折的分数存起来

- $\text{fold_scores.append(score)}$

这样跑完  $k$  次之后，就有一个长度为  $k$  的分数列表，比如：

$\text{fold_scores} = [0.81, 0.79, 0.83, 0.80, 0.82]$

步骤 4.14：对这一组参数做折平均

```
python

mean_score = np.mean(fold_scores)
results.append((param, mean_score))
```

关键点：

每一折都要重新建一个“干净”的模型实例，不能复用上一次训练过的模型，否则信息就泄漏了。

步骤 4.15：在训练集上拟合

- $\text{model}.fit(\text{X\_train}, \text{y\_train})$

模型只看到训练集的数据，不允许用验证集来训练。

步骤 4.16：在验证集上评估

- 预测：
  - $\text{y\_pred} = \text{model}.predict(\text{X\_val})$
- 计算指标：
  - 分类：准确率、F1、AUC 等
  - 回归：MSE、MAE、R<sup>2</sup> 等
- 得出一个数字： $\text{score}$

步骤 4.17：把这一折的分数存起来

- $\text{fold_scores.append(score)}$

这样跑完  $k$  次之后，就有一个长度为  $k$  的分数列表，比如：

$\text{fold_scores} = [0.81, 0.79, 0.83, 0.80, 0.82]$

步骤 4.18：对这一组参数做折平均

```
python

mean_score = np.mean(fold_scores)
results.append((param, mean_score))
```

关键点：

每一折都要重新建一个“干净”的模型实例，不能复用上一次训练过的模型，否则信息就泄漏了。

步骤 4.19：在训练集上拟合

- $\text{model}.fit(\text{X\_train}, \text{y\_train})$

模型只看到训练集的数据，不允许用验证集来训练。</p