

GPR:

原理:

高斯过程回归 (Gaussian Process Regression, GPR) 是一种基于高斯过程的非参数化贝叶斯回归方法, 其可以灵活适应不同类型数据, 用于建模和预测数据之间的复杂关系, 具有拟合能力强、泛化能力好等特点, 是一个植根于贝叶斯统计的较为流行的非参数机器学习技术, 因为它能够在小数据集上有出色的预测准确性而没有过拟合或欠拟合的问题, 并且能够给出关于模型输出的不确定性的信息。

GPR 有三个主要步骤:

- (1) 根据主观的先验知识选择合适的核函数并定义初始超参数;
- (2) 使用概率分布生成先验模型并对其进行训练, 即通过训练样本寻找最优超参数;
- (3) 对测试样本进行预测, 给出估计结果的均值和方差。

在其中, 如何获得最优超参是值得研究的问题, 常见如传统的基于梯度算法或现代智能优化算法的方法都可用于其中, 而负对数边际似然 (NLML) 值通常被设置为目标函数。

Codes:

```
import numpy as np
from sklearn.gaussian_process import GaussianProcessRegressor
from sklearn.gaussian_process.kernels import RBF
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error

# 生成模拟数据集
np.random.seed(0)
n_samples = 100
X = np.random.rand(n_samples, 1) # 一维输入特征
y = np.sin(10 * X) + np.cos(5 * X) + np.random.normal(scale=0.1, size=n_samples) # 添加噪声的目标变量

# 划分训练集和测试集
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# 定义高斯过程回归模型和核函数
kernel = RBF(length_scale=1.0, length_scale_bounds=(1e-1, 10.0)) # 使用 RBF 核函数, 并设置长度尺度的搜索范围
gpr = GaussianProcessRegressor(kernel=kernel, alpha=1e- 1/2 , optimizer='fmin_l_bfgs_b', n_restarts_optimizer=10)

# 训练模型
gpr.fit(X_train, y_train)

# 预测测试集结果
y_pred, y_pred_std = gpr.predict(X_test, return_std=True) # 获取预测均值和标准差 (方差的
```

平方根)

```
# 计算并打印预测误差
```

```
mse = mean_squared_error(y_test, y_pred)
```

```
print(f'Mean Squared Error (MSE): {mse:.2f}')
```

```
# 可视化结果
```

```
import matplotlib.pyplot as plt
```

```
plt.figure(figsize=(10, 6))
```

```
plt.plot(X_train, y_train, 'o', color='gray', label='Training Data')
```

```
plt.plot(X_test, y_test, '^', color='black', label='True Values')
```

```
plt.plot(X_test, y_pred, '-', color='red', label='Predictions')
```

```
plt.fill_between(X_test.ravel(),  
                 y_pred - y_pred_std,  
                 y_pred + y_pred_std,  
                 alpha=0.2, color='pink',  
                 label='Standard Deviation')
```

```
plt.title('Gaussian Process Regression with RBF Kernel')
```

```
plt.xlabel('Input Feature')
```

```
plt.ylabel('Target Variable')
```

```
plt.legend()
```

```
plt.show()
```

代码讲解：

1. 导入所需库：

numpy 用于生成模拟数据和数值计算。

sklearn.gaussian_process 包含高斯过程回归模型。

sklearn.gaussian_process.kernels 提供多种核函数。

sklearn.model_selection.train_test_split 用于划分训练集和测试集。

sklearn.metrics.mean_squared_error 计算均方误差 (MSE)。

2. 生成模拟数据集：

生成 100 个一维随机输入特征 X。

设置目标变量 y 为输入特征的复合非线性函数（包含噪声），模拟真实世界中的复杂关系。

3. 数据集划分：

使用 train_test_split 将数据集划分为 80% 作为训练集和 20% 作为测试集。

4. 定义高斯过程回归模型：

使用 GaussianProcessRegressor 类创建模型实例。

选择 RBF（径向基函数）核函数，并设置初始长度尺度为 1.0，同时指定长度尺度的搜索范围（用于超参数优化）。

设置 alpha 参数为先验噪声方差（这里取 0.5），用于平衡数据拟合与模型平滑度。

使用 fmin_l_bfgs_b 优化器进行超参数优化，并设置 n_restarts_optimizer 为 10，表示对优化过程进行 10 次重启以避免局部最优解。

5.训练模型:

调用 `gpr.fit()`方法, 使用训练集数据训练高斯过程回归模型。

6.模型预测:

调用 `gpr.predict()`方法, 传入测试集数据, 返回预测均值 `y_pred` 和标准差 `y_pred_std`。

7.计算预测误差:

使用 `mean_squared_error` 计算预测结果与真实值之间的均方误差。

8.可视化结果:

使用 `matplotlib` 绘制图表, 展示训练数据、真实值、预测值以及预测不确定性的区间 (由标准差确定)。

通过上述代码, 不仅实现了高斯过程回归算法, 还完成了模型训练、预测、误差计算以及结果可视化。这有助于直观地理解高斯过程回归如何拟合数据、预测新点以及量化预测不确定性。