



特征工程：特征工程在推荐系统中的应用

特征工程基础

1. 特征选择的重要性

在构建推荐系统时，特征选择是关键的一步。它涉及到从大量潜在特征中挑选出最相关、最有预测能力的特征，以提高模型的准确性和效率。特征选择不仅可以减少模型训练的时间和资源消耗，还能避免过拟合，使模型更加泛化。例如，在用户行为预测中，可能有成百上千个特征，如用户年龄、性别、历史购买记录、浏览时间等，但并非所有特征都对预测有同等贡献。通过特征选择，我们可以识别出哪些特征是最重要的，比如历史购买记录可能比用户的性别对预测下一次购买行为更有价值。

1.1 示例代码

假设我们有一个包含用户特征和购买行为的数据集，我们可以使用递归特征消除（RFE）方法来选择特征。

```
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression
import pandas as pd

# 加载数据
data = pd.read_csv('user_behavior.csv')
X = data.drop('purchased', axis=1) # 特征
y = data['purchased'] # 目标变量

# 创建模型
model = LogisticRegression()

# 特征选择
rfe = RFE(model, n_features_to_select=5)
fit = rfe.fit(X, y)

# 输出选择的特征
print("Selected Features: %s" % fit.support_)
print("Feature Ranking: %s" % fit.ranking_)
```

2. 特征编码技术

特征编码是将非数值特征转换为数值形式的过程，这对于机器学习模型至关重要，因为大多数模型只能处理数值输入。常见的特征编码技术包括独热编码（One-Hot Encoding）、标签编码

以上内容仅为本文档的试下载部分，为可阅读页数的一半内容。如要下载或阅读全文，请访问：<https://d.book118.com/938021047024006111>