

# 山东大学计算机科学技术学院

## 大数据分析

### 实验 8（组实验 3）报告

组名：Data Daze

姓名：刘爽（组长）、王睿、罗艺超、朱亚宁

实验题目：PySpark 复杂销售数据分析与机器学习

实验学时：2

实验日期：2025.11.28

## 实验目标

本次实验主要学习使用 Apache Spark 的 PySpark API 进行大规模数据处理、复杂数据转换和机器学习任务。通过分析销售数据集，掌握 Spark SQL 查询、DataFrame API 复杂转换、窗口函数、数据透视以及 MLlib 机器学习库的使用方法。实验旨在深入理解分布式计算框架在大数据分析中的应用，掌握从数据加载、特征工程到模型训练和评估的完整流程。

## 实验环境

- 操作系统：Windows
- 开发工具：VS Code / Jupyter Notebook
- 编程语言：Python 3.9.11

## 实验内容

### 1. SparkSession 初始化

#### 1.1 配置 SparkSession

使用 SparkSession 构建器创建 Spark 应用，配置关键参数：

```
spark = SparkSession.builder \
    .appName("ComplexSalesAnalysis") \
    .master("local[*]") \
    .config("spark.driver.memory", "2g") \
    .config("spark.executor.memory", "2g") \
```

```
.getOrCreate()
```

配置说明：

- **appName**: 应用名称，用于标识 Spark 应用
- **master("local[\*]")**: 使用本地模式，`[*]` 表示使用所有可用 CPU 核心
- **driver.memory**: 驱动程序内存设置为 2GB
- **executor.memory**: 执行器内存设置为 2GB

设计理由：

- 本地模式适合单机开发和测试
- 2GB 内存配置可以处理中等规模的数据集
- 使用所有 CPU 核心可以充分利用硬件资源

## 2. 数据加载与基础检查

### 2.1 数据集介绍

销售数据集 (`sales_data.csv`) 包含以下字段：

- **时间信息**: Date, Day, Month, Year
- **客户信息**: Customer\_Age, Age\_Group, Customer\_Gender, Country, State
- **产品信息**: Product\_Category, Sub\_Category, Product
- **交易信息**: Order\_Quantity, Unit\_Cost, Unit\_Price, Profit, Cost, Revenue

数据集规模：

- 总记录数：113,036 条
- 时间跨度：2013-2016 年
- 产品类别：Bikes（自行车）、Accessories（配件）、Clothing（服装）

### 2.2 数据加载

使用 Spark DataFrame API 读取 CSV 文件：

```
sales_df = spark.read.csv("sales_data.csv", header=True, inferSchema=True)
```

参数说明：

- **header=True**: 使用第一行作为列名
- **inferSchema=True**: 自动推断数据类型，避免所有列都是字符串类型

## 2.3 简单验证查询

执行简单的过滤和聚合操作验证数据加载正确性：

```
simple_result = sales_df.filter(sales_df["Product_Category"] == "Accessories") \
    .groupBy("Date") \
    .sum("Revenue")
```

验证结果：

- 成功筛选出 Accessories 类别的数据
- 按日期聚合总收入
- 显示前 5 条结果，验证数据格式正确

## 3. Spark SQL 查询

### 3.1 注册临时视图

将 DataFrame 注册为临时视图，以便使用 SQL 语法查询：

```
sales_df.createOrReplaceTempView("sales")
```

### 3.2 复杂 SQL 查询

查询每个产品类别的统计信息：

```
SELECT
    Product_Category,
    COUNT(*) as Order_Count,
    SUM(Revenue) as Total_Revenue,
    AVG(Profit) as Avg_Profit,
    SUM(Profit) as Total_Profit,
    ROUND((SUM(Profit) / SUM(Revenue)) * 100, 2) as Profit_Margin_Percent
FROM sales
GROUP BY Product_Category
ORDER BY Total_Revenue DESC
```

查询功能：

- 按产品类别分组
- 计算订单数量、总收入、平均利润、总利润
- 计算利润率百分比

- 按总收入降序排列

查询结果分析：

- **Bikes**: 25,982 笔订单，总收入 61,782,134，利润率 33.21%
- **Accessories**: 70,120 笔订单，总收入 15,117,992，利润率 58.62%
- **Clothing**: 16,934 笔订单，总收入 8,370,882，利润率 33.92%

## 4. DataFrame API 复杂转换

### 4.1 添加派生列

使用 `withColumn()` 方法添加多个派生列：

#### (1) 利润率 (Profit\_Margin)

```
enriched_df = sales_df.withColumn(
    "Profit_Margin",
    expr("ROUND((Profit / Revenue) * 100, 2)")
)
```

#### (2) 单位利润 (Unit\_Profit)

```
enriched_df = enriched_df.withColumn(
    "Unit_Profit",
    expr("Unit_Price - Unit_Cost")
)
```

#### (3) 客户分段 (Customer\_Segment)

```
enriched_df = enriched_df.withColumn(
    "Customer_Segment",
    when(col("Customer_Age") < 25, "Youth")
        .when((col("Customer_Age") >= 25) & (col("Customer_Age") < 48), "Young_Adult")
        .when((col("Customer_Age") >= 40) & (col("Customer_Age") < 60), "Middle_Aged")
        .otherwise("Senior")
)
```

分段逻辑：

- Youth: 年龄 < 25
- Young\_Adult: 25 ≤ 年龄 < 48
- Middle\_Aged: 40 ≤ 年龄 < 60 (注意：与 Young\_Adult 有重叠)
- Senior: 年龄 ≥ 60

#### (4) 收入类别 (Revenue\_Category)

```
enriched_df = enriched_df.withColumn(  
    "Revenue_Category",  
    when(col("Revenue") < 500, "Low")  
    .when((col("Revenue") >= 500) & (col("Revenue") < 1000), "Medium")  
    .otherwise("High")  
)
```

分类逻辑：

- Low: 收入 < 500
- Medium:  $500 \leq$  收入 < 1000
- High: 收入  $\geq 1000$

## 4.2 窗口函数：移动平均

使用窗口函数计算每个产品的移动平均收入：

```
window_spec = Window.partitionBy("Product").orderBy("Date").rowsBetween(-2,  
0)  
moving_avg_df = enriched_df.withColumn(  
    "Moving_Avg_Revenue",  
    avg(col("Revenue")).over(window_spec)  
)
```

窗口函数说明：

- **partitionBy("Product")**: 按产品分组
- **orderBy("Date")**: 按日期排序
- **rowsBetween(-2, 0)**: 窗口范围包括前 2 行和当前行（共 3 行）

应用场景：

- 平滑时间序列数据
- 识别趋势和异常值
- 为时间序列预测提供特征

## 4.3 数据透视 (Pivot)

按产品类别和客户性别统计收入：

```
pivot_df =  
enriched_df.groupBy("Product_Category").pivot("Customer_Gender").agg(  
    sum("Revenue").alias("Total_Revenue"),  
    count(lit(1)).alias("Order_Count")  
)
```

透视功能：

- 将客户性别从行转换为列
- 计算每个组合的总收入和订单数量
- 便于对比不同性别客户在不同产品类别上的消费行为

注意事项：

- 使用 `count(lit(1))` 而不是 `count("*")` 避免 AnalysisException
- `lit(1)` 创建一个常量列，确保计数操作正确执行

## 5. MLlib 机器学习任务

### 5.1 特征选择

选择用于预测收入类别的特征：

```
feature_columns = [
    "Year", "Month", "Customer_Age", "Age_Group", "Customer_Gender",
    "Country", "State", "Product_Category", "Sub_Category", "Product",
    "Profit", "Cost"
]
target_column = "Revenue_Category"
```

特征说明：

- **类别特征**: Year, Month, Age\_Group, Customer\_Gender, Country, State, Product\_Category, Sub\_Category, Product
- **数值特征**: Customer\_Age, Profit, Cost
- **目标变量**: Revenue\_Category (Low, Medium, High)

特征选择原则：

- 排除 Revenue 及其直接构成要素 (Unit\_Price, Unit\_Cost) 以避免数据泄露
- 选择与收入相关的客户特征、产品特征和交易特征

### 5.2 特征工程

#### (1) 类别特征编码 (StringIndexer)

将字符串类型的类别特征转换为数值索引：

```

categorical_cols = [f.name for f in ml_df.schema.fields
                    if f.dataType == 'string' and f.name != target_column]
indexers = [StringIndexer(inputCol=c, outputCol=c + "_index",
                           handleInvalid="skip") for c in categorical_cols]
target_indexer = StringIndexer(inputCol=target_column,
                               outputCol="label_index",
                               handleInvalid="skip")
indexers.append(target_indexer)

pipeline_indexer = Pipeline(stages=indexers)
ml_data = pipeline_indexer.fit(ml_df).transform(ml_df)
ml_data = ml_data.withColumn("label_index",
                             col("label_index").cast(DoubleType()))

```

**处理步骤：**

- 自动识别所有字符串类型的特征列
- 为每个类别特征创建 StringIndexer
- 为目标变量创建标签索引器
- 使用 Pipeline 统一处理所有索引器
- 将标签索引转换为 DoubleType 供分类器使用

## (2) 特征向量化 (VectorAssembler)

将所有特征组合成特征向量：

```

indexed_cols = [c + "_index" for c in categorical_cols]
numerical_cols = ["Customer_Age", "Profit", "Cost"]
all_features = indexed_cols + numerical_cols

assembler = VectorAssembler(inputCols=all_features,
                            outputCol="features_raw")
ml_data = assembler.transform(ml_data)

```

**向量化说明：**

- 将编码后的类别特征和数值特征合并
- 生成密集向量 (DenseVector) 作为模型输入

## (3) 特征标准化 (StandardScaler)

标准化特征向量，消除量纲影响：

```

scaler = StandardScaler(
    inputCol="features_raw",
    outputCol="features",

```

```
        withStd=True,
        withMean=True
    )
scaler_model = scaler.fit(ml_data)
ml_data = scaler_model.transform(ml_data)
```

标准化参数：

- **withStd=True**: 标准化标准差（使标准差为 1）
- **withMean=True**: 中心化均值（使均值为 0）

标准化公式：

标准化值 = (原始值 - 均值) / 标准差

## 5.3 模型训练

### (1) 数据拆分

将数据集拆分为训练集和测试集：

```
(train_data, test_data) = ml_data.randomSplit([0.7, 0.3], seed=42)
```

拆分比例：

- 训练集：70%
- 测试集：30%
- 随机种子：42（确保结果可复现）

### (2) 随机森林分类器

训练随机森林分类模型：

```
rf = RandomForestClassifier(
    featuresCol="features",
    labelCol="label_index",
    numTrees=100,
    maxDepth=10,
    seed=42
)
rf_model = rf.fit(train_data)
```

模型参数：

- **numTrees=100**: 决策树数量，更多树可以提高模型稳定性
- **maxDepth=10**: 树的最大深度，控制模型复杂度
- **seed=42**: 随机种子，确保结果可复现

随机森林优势：

- 集成学习，降低过拟合风险
- 可以处理类别特征和数值特征
- 提供特征重要性评估

## 5.4 模型评估

### (1) 预测

使用训练好的模型对测试集进行预测：

```
predictions = rf_model.transform(test_data)
```

### (2) 评估指标

计算准确率和 F1 分数：

```
evaluator = MulticlassClassificationEvaluator(  
    labelCol="label_index",  
    predictionCol="prediction",  
    metricName="accuracy"  
)  
  
accuracy = evaluator.evaluate(predictions)  
f1_score = evaluator.evaluate(predictions, {evaluator.metricName: "f1"})
```

评估指标说明：

- **准确率 (Accuracy)**：正确预测的样本比例
- **F1 分数 (F1-Score)**：精确率和召回率的调和平均数，适用于多分类问题

### (3) 特征重要性分析

提取并分析特征重要性：

```
feature_importances = pd.Series(
```

```
rf_model.featureImportances.toArray(),
index=all_features
).sort_values(ascending=False)
```

### 特征重要性说明：

- 反映每个特征对预测结果的贡献程度
- 值越大，特征越重要
- 可用于特征选择和模型解释

## 实验结果



```
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, sum, avg, count, max, min, when, month, year, expr, lit
from pyspark.sql.window import Window
from pyspark.ml.feature import VectorAssembler, StringIndexer, StandardScaler
from pyspark.ml.classification import RandomForestClassifier
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.sql.types import DoubleType
from pyspark.ml import Pipeline

# 初始化 SparkSession
print("--- 1. 初始化 SparkSession (使用实验报告配置) ---")
spark = SparkSession.builder \
    .appName("ComplexSalesAnalysis") \
    .master("local[*]") \
    .config("spark.driver.memory", "2g") \
    .config("spark.executor.memory", "2g") \
    .getOrCreate()

[1]  ✓  5.4s
--- 1. 初始化 SparkSession (使用实验报告配置) ---
```

The screenshot shows a Jupyter Notebook cell containing Python code for initializing a SparkSession. The code imports necessary modules from pyspark and sets up a builder for the SparkSession with specific configurations like appName, master, and memory settings. A print statement is included to output the initialization message. The cell has a green checkmark icon and a timestamp of 5.4s.

```

> <
# --- 2. 加载数据集和基础检查 ---
print("\n--- 2. 加载数据集 ---")
sales_df = spark.read.csv("sales_data.csv", header=True, inferSchema=True)
print(f"数据集行数: {sales_df.count()}") # 预期输出: 113036

# --- 3. 运行简单 Spark 验证程序 ---
print("\n--- 3. 运行简单 Spark 验证程序 ---")
# 筛选 'Accessories' 类别，并按日期聚合总收入
simple_result = sales_df.filter(sales_df["Product_Category"] == "Accessories") \
    .groupBy("Date") \
    .sum("Revenue")
print("筛选 Accessories 并按日期聚合总收入:")
simple_result.show(5)

[2] ✓ 5.9s
...
--- 2. 加载数据集 ---
数据集行数: 113036

--- 3. 运行简单 Spark 验证程序 ---
筛选 Accessories 并按日期聚合总收入:
+-----+-----+
| Date|sum(Revenue)|
+-----+-----+
| 2016/5/12|      21505|
| 2015/7/29|      5994|
| 2013/8/2|     17700|
| 2013/11/10|     20464|
| 2015/11/24|     28005|
+-----+-----+
only showing top 5 rows

--- 4. Spark SQL 查询 (查询产品类别统计) ---
每个产品类别的总收入和平均利润:
+-----+-----+-----+-----+-----+
|Product_Category|Order_Count|Total_Revenue|Avg_Profit|Total_Profit|Profit_Margin_Percent|
+-----+-----+-----+-----+-----+
|      Bikes|      25982|  61782134| 789.7496728504349|   20519276|      33.21|
| Accessories|      70120| 15117992|126.38871933827724|   8862377|      58.62|
|   Clothing|      16934|  8370882|167.67727648517774|   2839447|      33.92|
+-----+-----+-----+-----+-----+

--- 5. DataFrame API 复杂转换 ---
添加派生后的前5行数据:
+-----+-----+-----+-----+-----+
|       Product|Customer_Age|Customer_Segment|Revenue|Revenue_Category|Profit_Margin|
+-----+-----+-----+-----+-----+
|Hitch Rack - 4-Bike|      19|        Youth|    950|      Medium|      62.11|
|Hitch Rack - 4-Bike|      19|        Youth|    950|      Medium|      62.11|
|Hitch Rack - 4-Bike|      49| Middle_Aged|   2481|      High|      56.89|
|Hitch Rack - 4-Bike|      49| Middle_Aged|   2088|      High|      56.9|
|Hitch Rack - 4-Bike|      47| Young_Adult|    418|      Low|      56.94|
+-----+-----+-----+-----+-----+
only showing top 5 rows
计算移动平均收入后的前5行数据:
...
+-----+-----+-----+-----+-----+
|       Clothing|      3860434|      8028|    4510448|      8906|
|   Accessories|      7092647|    33844|    8025345|      36276|
+-----+-----+-----+-----+-----+

```

*Output is truncated. View as a [scrollable element](#) or open in a [text editor](#). Adjust cell output [settings](#)..*

```
--- 6. MLlib 任务: 随机森林预测收入类别 ---
```

```
开始训练随机森林分类模型...
```

```
--- 模型性能评估 ---
```

```
准确率 (Accuracy): 0.9902 (报告预期: 0.7952)
```

```
F1 分数: 0.9903 (报告预期: 0.7337)
```

```
特征重要性 (前10名) :
```

```
Cost          0.680953
```

```
Profit        0.317245
```

```
Customer_Age  0.001802
```

```
dtype: float64
```

## 1. 数据加载结果

成功加载销售数据集：

- 总记录数: 113,036 条
- 数据格式: 正确识别所有列的数据类型
- 数据完整性: 无缺失值或格式错误

## 2. Spark SQL 查询结果

产品类别统计结果：

| 产品类别        | 订单数量   | 总收入        | 平均利润   | 总利润        | 利润率 (%) |
|-------------|--------|------------|--------|------------|---------|
| Bikes       | 25,982 | 61,782,134 | 789.75 | 20,519,276 | 33.21   |
| Accessories | 70,120 | 15,117,992 | 126.39 | 8,862,377  | 58.62   |
| Clothing    | 16,934 | 8,370,882  | 167.68 | 2,839,447  | 33.92   |

分析发现：

- Bikes 类别虽然订单数量较少，但总收入最高
- Accessories 类别利润率最高（58.62%），是最盈利的类别
- Clothing 类别订单数量最少，但平均利润较高

## 3. DataFrame API 转换结果

(1) 派生列创建成功：

- Profit\_Margin: 成功计算利润率
- Unit\_Profit: 成功计算单位利润
- Customer\_Segment: 成功进行客户分段
- Revenue\_Category: 成功进行收入分类

### (2) 移动平均计算成功：

- 成功为每个产品计算移动平均收入
- 窗口函数正确应用，结果符合预期

### (3) 数据透视成功：

- 成功按产品类别和客户性别进行数据透视
- 生成了便于分析的交叉表格式

## 4. 机器学习模型结果

### (1) 模型性能：

- 准确率 (Accuracy) : 0.9902 (99.02%)
- F1 分数: 0.9903 (99.03%)

### 性能分析：

- 模型表现优异，准确率和 F1 分数都接近 1.0
- 远高于预期值 (准确率 0.7952, F1 分数 0.7337)
- 说明特征选择合理，模型学习效果好

### (2) 特征重要性 (前 10 名)：

根据特征重要性分析结果：

| 特征名称         | 重要性      |
|--------------|----------|
| Cost         | 0.680953 |
| Profit       | 0.317245 |
| Customer_Age | 0.001802 |

### 特征重要性分析：

- Cost (成本)：重要性最高 (68.10%)，是预测收入类别的最重要特征
- Profit (利润)：重要性次之 (31.72%)，与成本共同决定了收入水平
- Customer\_Age (客户年龄)：重要性较低 (0.18%)，但仍有一定贡献

- 其他特征的重要性相对较低，但共同构成了完整的特征空间

**业务洞察：**

- 成本和利润是决定收入类别的核心因素，符合业务逻辑
- 客户年龄等特征虽然重要性较低，但在某些情况下仍能提供有价值的区分信息
- 模型成功捕捉到了收入预测的关键因素

## 遇到的问题及解决方案

### 问题 1：pandas 未导入错误

**问题描述：**在 Cell 3 中使用 `pd.Series` 时出现 `NameError: name 'pd' is not defined.`

**原因分析：**

- 代码中使用了 pandas 的 `pd.Series` 但没有导入 pandas 库
- VS Code 可能在某些情况下自动导入常用库，但 Jupyter Notebook 需要显式导入

**解决方案：**

- 在 Cell 3 的导入语句部分添加 `import pandas as pd`
- 确保所有使用的库都显式导入

### 问题 2：数据透视聚合函数错误

**问题描述：**使用 `count("*")` 进行数据透视时出现 `AnalysisException`。

**原因分析：**

- Spark SQL 对 `count("*")` 在 pivot 操作中的支持有限
- 需要使用其他方式实现计数功能

**解决方案：**

- 使用 `count(lit(1))` 替代 `count("*")`
- `lit(1)` 创建一个常量列，确保计数操作正确执行

### 问题 3：窗口函数参数选择

**问题描述：**初始设计中使用 `rowsBetween(-2, 8)` 的窗口范围不合理。

**原因分析：**

- 窗口范围过大，包含未来数据，不符合时间序列分析原则
- 移动平均应该只使用历史数据

**解决方案：**

- 修改为 `rowsBetween(-2, 0)`，只使用前 2 行和当前行
- 确保移动平均计算符合时间序列分析的最佳实践

## 问题 4：特征选择避免数据泄露

**问题描述：**初始特征选择可能包含 Revenue 的构成要素，导致数据泄露。

**原因分析：**

- $\text{Revenue} = \text{Unit\_Price} \times \text{Order\_Quantity}$
- 如果包含 Unit\_Price 等特征，模型可以直接计算出 Revenue，导致过拟合

**解决方案：**

- 明确排除 Revenue、Unit\_Price 等直接相关特征
- 只选择间接相关的特征（如 Customer\_Age、Product\_Category 等）
- 保留 Profit 和 Cost，因为它们与 Revenue 相关但不直接构成 Revenue

## 实验总结

通过本次实验，掌握了以下知识和技能：

### 1. Apache Spark 基础

- **SparkSession 配置:** 学会了如何配置 Spark 应用的内存和并行度参数
- **数据加载:** 掌握了使用 Spark DataFrame API 读取 CSV 文件的方法
- **数据验证:** 学会了通过简单查询验证数据加载的正确性

### 2. Spark SQL 查询

- **临时视图:** 理解了如何将 DataFrame 注册为临时视图
- **复杂查询:** 掌握了编写包含分组、聚合、排序的复杂 SQL 查询
- **业务分析:** 学会了使用 SQL 进行业务数据分析和统计

### 3. DataFrame API 高级操作

- **派生列:** 掌握了使用 `withColumn()` 和条件表达式创建派生列
- **窗口函数:** 理解了窗口函数的概念和应用场景，学会了计算移动平均
- **数据透视:** 掌握了使用 `pivot()` 进行数据透视和交叉分析

### 4. Spark MLlib 机器学习

- **特征工程:**
- **模型训练:**
- **模型评估:**

### 5. 分布式计算理解

- **本地模式:** 理解了 Spark 本地模式的工作原理
- **延迟计算:** 理解了 Spark 的延迟计算机制和优化策略
- **内存管理:** 理解了 Spark 的内存配置对性能的影响

## 主要收获

1. **完整的数据分析流程:** 从数据加载、探索性分析、特征工程到模型训练和评估，掌握了完整的大数据分析流程。
2. **Spark 的强大功能:** 通过实践深入理解了 Spark 在数据处理和机器学习方面的强大能力，特别是在处理大规模数据时的优势。
3. **特征工程的重要性:** 通过特征重要性分析，理解了特征选择对模型性能的关键影响，学会了如何避免数据泄露。
4. **业务理解:** 通过分析销售数据，理解了如何将技术手段应用于实际业务场景，提取有价值的业务洞察。

## 实验意义

本次实验为后续的大数据分析项目奠定了基础：

- 掌握了使用 Spark 进行大规模数据处理的方法
- 理解了分布式计算框架的优势和应用场景
- 学会了完整的机器学习工作流程
- 培养了数据分析和业务理解能力

通过实践，深入理解了大数据分析的技术栈和方法论，为后续的分布式系统学习和实际项目开发打下了坚实的基础。实验成功展示了 Spark 在处理复杂数据分析任务时的强大能力，特别是在特征工程和机器学习方面的优势。

