

面向用户

1. 项目目标（项目动机）

本企业工资预测项目的长期目标是为企业提供一个综合的、数据驱动的薪酬规划和预测工具，通过深度分析员工数据、市场趋势和公司财务情况，帮助企业实现更精准、合理和可持续的薪资决策。以下是项目目标和动机的详细阐述：

1.1 精准薪资预测

通过利用先进的数据科学和机器学习技术，项目致力于构建一个高度准确的工资预测模型，能够考虑多个因素，如员工职位、绩效历史、市场趋势等，以提供企业更具预测性的薪酬规划。

1.2 优化薪酬结构

项目的动机在于帮助企业优化薪酬结构，确保薪资与员工价值和市场竞争力相符，从而提高员工满意度、激励积极表现，并降低离职率。

1.3 成本效益管理

通过全面的薪酬分析，项目旨在帮助企业实现更有效的成本效益管理。企业可以更好地了解 and 预测薪资支出，从而更有针对性地进行预算规划和资源分配。

1.4 人才留存和吸引

通过准确的薪酬预测，企业可以更好地吸引和留住高素质员工。项目的目标是帮助企业打造具有竞争力的薪酬体系，吸引业界翘楚并提高员工忠诚度。

1.5 持续更新和优化

项目将定期更新模型，确保薪资预测与市场和公司内部变化保持同步。这有助于企业随时调整薪酬策略，以适应变化的业务环境和市场趋势。

1.6 数据驱动的决策支持

通过提供直观的可视化报告和数据解释，项目旨在帮助企业管理层更好地理解薪资预测结果，从而做出更为明智的决策。

通过实现上述目标，本项目追求帮助企业建立一个灵活、敏捷、能够适应变化的薪酬管理体系，为企业提供战略性的人力资源支持，最终实现员工和企业共赢的局面。

2.向最终用户展现模型

2.1 模型优点

(1) 可解释性强:

Logistic 回归模型的优势之一是其可解释性强。由于它是一种线性模型，模型的系数可以被直观地解释为每个特征对输出的影响程度。例如，如果某个特征的系数为正，说明该特征与输出正相关；反之，如果系数为负，则说明特征与输出负相关。这使得分析人员或领域专家能够更容易理解模型，从而更好地理解特征之间的关系和对预测结果的影响。

(2) 简单而有效:

Logistic 回归是一种简单而有效的分类算法，尤其适用于二分类问题。其简单性使得模型易于理解和实现，同时在相对较小的数据集上表现良好。通过最小化逻辑损失函数，模型能够从训练数据中学习出适用于分类的边界，使得在实际应用中具有出色的性能。

(3) 适用性广泛:

Logistic 回归在多个领域得到广泛应用，包括医学、社会科学和工业等。其广泛的应用证明了其在不同场景下的通用性和效果。由于其良好的性能和可解释性，Logistic 回归经常被选择用于解决各种分类问题，使其成为许多实际应用的首选算法之一。

(4) 对异常值不敏感:

Logistic 回归对一些异常值的存在相对不敏感。由于逻辑回归模型基于最大似然估计，而不是依赖于数据的均值或中位数，因此它对于异常值的影响相对较小。这使得模型更具鲁棒性，能够在存在一些异常值的情况下仍然表现出色。

(5) 容易实现和使用:

通过使用 Python 中的 statsmodels 库，Logistic 回归模型的实现变得简单而直观。该库提供了方便的接口，使得模型的拟合、系数提取和摘要获取等步骤轻松实现。此外，代码中的交叉验证等技术进一步提高了对模型性能的全面评估，为实际应用提供了更可靠的结果。

2.2 技术细节

(1) 数据预处理 (preprocess_data):

采用 pandas 库读取 CSV 文件，使用自定义的列名来解释数据。

处理缺失值：通过将 "?" 替换为 pandas 中的 pd.NA，然后根据缺失值比例判断是删除带有缺失值的行还是用均值填充。

将标签列进行二元分类编码 (">50K"和"<=50K")，以便模型训练。

对分类变量进行独热编码，将其转换为模型可以接受的数值型特征。

最后，将数据类型转换为 float 以满足模型的输入要求。

(2) 模型训练 (train_model):

使用 statsmodels 库中的 Logit 类构建逻辑回归模型。

在训练数据上调用 fit()方法，通过最大似然估计拟合模型。

添加常数项（截距）以适应模型。

返回训练好的逻辑回归模型，该模型包含了系数等关键信息。

(3) 提取模型系数 (extract_coefficients):

通过访问训练好的逻辑回归模型的 params 属性, 提取模型的系数。
通过遍历特征和系数, 输出每个特征对应的系数。

(4) 获取模型摘要 (get_model_summary):

使用 summary() 方法获取逻辑回归模型的详细摘要信息。
打印模型的统计学指标、系数、p 值等信息, 以支持对模型性能和特征重要性的全面理解。

(5) 评估和性能指标计算 (evaluate_model):

针对测试数据进行与训练数据相同的处理, 包括处理缺失值和独热编码。
使用训练好的模型进行预测, 根据设定的阈值进行二元分类。
计算准确率、召回率、F1 分数和 ROC AUC 分数等性能指标。
对 ROC AUC 分数的计算进行异常值处理, 以避免潜在的数值计算问题。

(6) 交叉验证 (perform_cross_validation):

使用 StratifiedKFold 来保持类别分布, 确保交叉验证的可靠性。
初始化逻辑回归模型, 并使用交叉验证计算准确率和召回率。
打印交叉验证的平均准确率和召回率。

3. 展现如何使用模型 (此处暂未开发, 为预计使用方法)

在使用模型时, 通常需要经历以下步骤, 以下是一个简单的示例:

3.1. 导入必要的库和加载训练好的模型

```
import statsmodels.api as sm
import pandas as pd

trained_model = sm.load('trained_logistic_model.pickle') # 请确保替换为实际保存的模型
```

3.2 准备待预测的数据

假设有一组新的数据, 命名为 new_data, 格式应与训练数据相同

```
new_data = pd.read_csv('new_data.csv') # 请确保替换为实际的新数据文件路径
```

3.3 进行与训练数据相同的数据预处理

```
# 使用预处理函数对新数据进行处理
new_data_processed = preprocess_data(new_data)
```

3.4 使用训练好的模型进行预测

```
# 添加常数项
new_data_processed = sm.add_constant(new_data_processed)

# 使用训练好的模型进行预测
predictions = trained_model.predict(new_data_processed)
```

3.5 设置阈值进行分类

```
# 假设阈值为 0.5
threshold = 0.5
predicted_labels = (predictions > threshold).astype(int)
```

3.6 查看预测结果

```
# 打印预测的标签
print("Predicted Labels:")
print(predicted_labels)
```

3.7 根据需要进一步分析

根据具体问题，可以根据模型的预测结果进行进一步的业务分析。
可以使用模型系数来解释每个特征对于预测结果的贡献。
针对实际场景，可以根据模型输出调整阈值，以平衡准确率和召回率。

3.8 总结和报告

汇总预测结果，形成最终的报告或结果输出。
根据业务需求，可能需要将结果可视化或以其他形式呈现。