



基于机器学习的分布式系统故障诊断系统

系统治愈师团队 | System Healer Team

姓名	班级	学号	分工
张骁凯	软件2102班	U202117281	需求文档、代码编写、架构评审
陈德霆	软件2102班	U202117254	系统设计、架构设计

Github仓库地址: <https://github.com/dekrt/SystemHealer>

文档仓库地址: <https://github.com/dekrt/SoftwareArchitecture>

一、信息介绍

大数据时代, 分布式系统成为信息存储和处理的主流系统。相对于传统系统而言, 分布式系统更为庞大和复杂, 故障发生的平均几率比较高, 其运维的难度和复杂度大大提高。如何对分布式系统进行高效、准确的运维, 成为保障信息系统高效、可靠运行的关键问题。

为此, 我们团队设计了故障诊断模型, 可以高效地分析并识别故障类别, 实现分布式系统故障运维的智能化, 快速恢复故障的同时大大降低分布式系统运维工作的难度, 减少运维对人力资源的消耗。

在分布式系统中某个节点发生故障时, 故障会沿着分布式系统的拓扑结构进行传播, 造成自身节点及其邻接节点相关的 KPI 指标和发生大量日志异常。中兴通讯为我们提供分布式数据库的故障特征数据和标签数据, 其中特征数据是系统发生故障时的 KPI 指标数据, KPI指标包括由 `feature0`、`feature1` ... `feature106` 共107个指标, 标签数据为故障类别数据, 共6个类别, 用0、1、2、3、4、5分别表示6个故障, 根据这些数据, 我们借助机器学习、深度学习、web等技术搭建故障诊断系统, 该系统支持用户上传训练集对模型进行训练和模型下载, 同时支持用户上传单条或多条测试语句进行测试并可视化测试结果, 同支持测试结果与模型的下载。

我们的项目采用 `Django + BootStrap` 为框架进行web系统的开发, 模型采用 `多层感知机 (MLP) 模型`, 它是一种用于处理监督学习问题的神经网络, 是一种前馈神经网络, 在处理此类问题上具有不俗的表现。

二、Web系统介绍

2.1 概述



网页基于web3开发，主要使用html5+css3+js，使用bootstrap和jQuery框架构建交互操作，提高了web系统的美观性、易用性。

web系统分为头部栏、上传训练模块和结果分析模块，用户可以在用户栏访问到有关模型、算法的更多信息；上传训练模块给出了文件提示和上传训练数据、预测数据的接口；结果分析模块在未上传数据时给出了样例图标，上传数据训练和预测后会根据结果生成结果分析图标，供用户对模型进行评估。

用户初次进入web系统，网页处于初始状态，未选择上传的文件时，点击上传的按钮被禁用，同时下方结果分析模块处于默认状态；用户选择文件后，“点击上传”按钮处于激活状态，上传后，后端进行模型训练并返回结果，等待训练完毕后，用户可以在下方看到各类别所占比例的饼状图、分类结果可视化的tsne图、训练集和验证集损失随迭代次数变化的折线图，通过这些图用户可以对模型有基本的认识。

训练完毕后，用户可以上传需要预测的数据，等待模型预测后，同样的可以在下方看到预测结果的分析图标。在上传模块可以下载预测结果文件和模型文件，便于保存训练好的模型。

2.2 部署方法

2.2.1 使用脚本进行部署

在 `SystemHealer/` 目录运行 `setup.sh` 脚本，可以一键构建环境，安装有关依赖，并打开默认浏览器预览网页。

如果系统没有打开浏览器，可以在构建完成之后手动输入127.0.0.1:8000或者localhost:8000来访问web系统。

2.2.2 使用 Docker 进行部署

在 `SystemHealer/` 目录运行以下命令来构建 Docker 镜像：

```
1 | docker build -t systemhealer .
```

构建完成后，运行以下命令来启动 Docker 容器：

```
1 | docker run -p 8000:8000 systemhealer
```

现在，您可以在浏览器中输入127.0.0.1:8000或 `localhost:8000` 来访问我们的Web系统。

三、模型介绍

我们的项目采用 `多层感知机（MLP）模型`，它是一种用于处理监督学习问题的神经网络，是一种前馈神经网络，在处理此类问题上具有不俗的表现。下面是我们项目中MP的定义：

```
1 class MLP(nn.Module):
2     def __init__(self, input_size, hidden_size, num_classes):
3         super(MLP, self).__init__()
4         self.fc1 = nn.Linear(input_size, hidden_size)
5         self.fc2 = nn.Linear(hidden_size, hidden_size)
6         self.fc3 = nn.Linear(hidden_size, hidden_size)
7         self.fc4 = nn.Linear(hidden_size, num_classes)
8         self.activation = nn.LeakyReLU()
9
10    def forward(self, x):
11        out = self.fc1(x)
12        out = self.activation(out)
13        out = self.fc2(out)
14        out = self.activation(out)
15        out = self.fc3(out)
16        out = self.activation(out)
```

```
17         out = self.fc4(out)
18     return out
```

在这个模型中，我们有一个输入层，两个隐藏层，以及一个输出层。首先，我们定义了一个名为 `MLP` 的类，它继承了 `PyTorch` 的 `nn.Module` 类。`nn.Module` 是所有神经网络模块的基类，我们自定义的模型也继承这个类。

在 `MLP` 类的构造函数中，我们定义了四个全连接层（也称为线性层）。每个全连接层都使用 `nn.Linear` 类创建，该类在输入数据上应用线性变换。第一个全连接层（`fc1`）的输入大小为 `input_size`，输出大小为 `hidden_size`。接下来的两个全连接层（`fc2` 和 `fc3`）的输入和输出大小都是 `hidden_size`。最后一个全连接层（`fc4`）的输入大小为 `hidden_size`，输出大小为 `num_classes`，这表示我们的模型将输出 `num_classes` 个类别的预测。

我们还定义了一个激活函数，使用了 `nn.LeakyReLU` 类。`LeakyReLU` 是修正线性单元（`ReLU`）的一个变种，它解决了 `ReLU` 的死亡 `ReLU` 问题，即一些神经元可能会停止学习，因为它们进入了一个永远不会被激活的状态。

在前向传播函数（`forward`）中，我们将输入数据 `x` 通过每个全连接层和激活函数，最后输出结果。注意，我们没有在最后一个全连接层后应用激活函数，这是因为我们通常会在计算损失函数时使用包含激活函数的损失函数，例如 `nn.CrossEntropyLoss`。

使用多层感知机模型有以下优点和特点：

- 1. 处理监督学习问题：** 多层感知机（`MLP`）是一种用于处理监督学习问题的神经网络模型。监督学习是一种机器学习任务，其中模型从带有标签的训练数据中学习，并预测未标记数据的标签。`MLP` 通过前向传播算法和反向传播算法来学习输入数据与标签之间的关系。
- 2. 前馈神经网络：** `MLP` 是一种前馈神经网络，这意味着数据在网络中沿一个方向流动，不会形成循环连接。这种结构使得 `MLP` 可以对输入数据进行逐层处理，每一层都产生中间的表示，最终输出预测结果。
- 3. 多层结构：** `MLP` 包含多个层次的神经元，通常包括输入层、隐藏层和输出层。隐藏层是位于输入层和输出层之间的层次，通过这些层次可以学习更加复杂和抽象的特征表示。在给定足够的隐藏单元和适当的参数设置下，`MLP` 可以学习到非线性关系，从而具备更强的模型表达能力。
- 4. 全连接层：** `MLP` 的每个层都是全连接层，也称为线性层。全连接层中的神经元与上一层的所有神经元都有连接，这样可以实现输入数据和权重的线性组合，并通过激活函数引入非线性。
- 5. 激活函数：** 在 `MLP` 中，我们使用了 `LeakyReLU` 作为激活函数。激活函数的作用是引入非线性，使得模型可以学习非线性关系。`LeakyReLU` 是修正线性单元（`ReLU`）的一个改进版本，当输入小于零时，它引入一个小的斜率，避免了 `ReLU` 的“死亡神经元”问题。
- 6. 模型训练：** `MLP` 通过前向传播和反向传播算法进行模型训练。前向传播是将输入数据通过网络的每一层，得到最终的预测结果。反向传播通过计算预测结果与真实标签之间的误差，并沿着网络的方向更新权重，使得模型能够逐渐减小误差，提高预测的准确性。

7. 模型应用： `MLP` 模型在分类和回归等监督学习任务中表现出色。它可以处理各种类型的输入数据，例如图像、文本、声音等。`MLP` 的强大表达能力和灵活性使得它成为机器学习和深度学习领域的重要工具之一。

总之，`MLP` 模型是一种灵活且强大的神经网络模型，适用于处理各种监督学习问题。它的多层结构和全连接层使其能够学习非线性模型，并通过适当的调整隐藏层的数量和大小来控制模型的复杂性。在我们的项目中，采用`MLP`模型可以有效处理分布式故障诊断等任务，并获得优秀的结果。

四、实现思路

4.1 数据预处理

在我们的故障诊断系统中，数据预处理是一个重要的步骤。我们需要确保输入到模型的数据是完整的，没有缺失值，并且格式正确。下面是我们进行数据预处理的主要步骤：

4.1.1 去除重复值

首先，我们使用 `drop_duplicates` 函数去除数据中的重复值。这是一个重要的步骤，因为重复的数据点可能会导致模型过拟合，即模型过于复杂，以至于它开始记忆训练数据，而不是从中学习模式。当我们的数据集中存在重复的样本时，模型在训练过程中可能会过度关注这些重复样本，导致对它们的预测过于自信，而对其他样本的泛化能力下降。

重复数据的存在可能是由于多种原因引起的，例如数据收集过程中的误操作、系统错误或者数据来源的特殊性。无论是哪种情况，重复数据都可能对我们的分析和建模过程产生负面影响，因此需要及时处理。

通过调用 `drop_duplicates` 函数并将参数 `inplace` 设置为 `True`，我们可以直接在原始数据上进行修改，而无需创建新的数据副本。这样做不仅可以减少内存的使用，还可以提高代码的执行效率。

`drop_duplicates` 函数会比较数据集中的每个数据点，并检查其与其他数据点是否完全相同。如果找到了重复的数据点，函数会将其删除，只保留一个副本。这样，我们就能确保数据集中的每个样本都是唯一的，避免了冗余数据的干扰。

通过去除重复数据，我们可以消除数据集中的冗余信息，提高模型的训练效果。当模型面临更清晰、更准确的数据时，它可以更好地学习数据中的模式和规律。这将有助于改善模型的泛化能力，使其能够更好地应用于未见过的数据，从而提高预测和决策的准确性。

总之，去除重复数据是数据预处理的重要一步，它有助于净化数据集，提高模型的性能和鲁棒性。在进行任何进一步的分析和建模之前，我们应该始终注意并处理数据中的重复值，以确保我们所使用的数据是可靠、准确且唯一的。

```
1 def preprocess(data):
2     # 去掉重复值
3     data.drop_duplicates(inplace=True)
```

4.1.2 缺失值处理

然后，我们使用 `IterativeImputer` 类来填充数据中的缺失值。`IterativeImputer` 是一种高级的填充方法，它通过迭代的方式使用其他特征的值来预测缺失值，并在每次迭代中更新缺失值的预测结果。这种迭代的过程可以提高填充的准确性和可靠性。

在代码中，我们创建了一个名为 `imputer` 的 `IterativeImputer` 对象，并设置了 `random_state` 为0，以确保每次运行时得到相同的结果。接下来，我们使用 `imputer.fit_transform(data)` 方法对数据进行拟合和转换操作，其中 `data` 是待填充的数据集。

在拟合的过程中，`IterativeImputer` 会首先对数据集中的特征进行分析，并建立一个预测模型来预测缺失值。然后，它使用已有的非缺失值特征作为输入，利用建立的模型来预测缺失值。这个过程会不断迭代，直到达到收敛条件或者达到预先设定的最大迭代次数。

转换的结果是一个新的数据集 `data_imputed`，它是一个 `DataFrame` 对象，具有与原始数据相同的列名。这个新数据集中的缺失值已经被填充，不再存在。为了保持一致性，我们将 `data_imputed` 的列名设置为与原始数据相同。

接着，我们通过 `data_imputed.dropna()` 方法去除仍然存在缺失值的行，以保证数据的完整性和一致性。然后，我们将数据集中的目标变量（标签）提取出来，存储在 `y` 变量中，以备后续的建模和分析使用。

最后，根据需求，我们定义了一个列表 `columns_to_drop`，其中包含需要从数据集中删除的特征列名。通过调用 `data_imputed.drop(columns_to_drop, axis=1)`，我们可以删除这些特征列，以满足后续分析的需要。经过这一步操作后，`data_imputed` 将只包含用于建模的特征列，不再包含标签列。

总结而言，我们进行了一次完整的数据预处理流程，包括使用 `IterativeImpute` 类进行缺失值填充，删除含有缺失值的行，提取标签列，并根据需求删除特定的特征列。通过这些处理步骤，我们能够准备好适用于机器学习模型训练和分析的数据集。

```
1 imputer = IterativeImputer(random_state=0)
2 data_imputed = imputer.fit_transform(data)
3 data_imputed = pd.DataFrame(data_imputed, columns=data.columns)
4 data_imputed = data_imputed.dropna()
5 y = data_imputed['label']
6 columns_to_drop = ['label']
7 data_imputed = data_imputed.drop(columns_to_drop, axis=1)
```


4.1.3 标准化处理

最后，我们使用 `StandardScaler` 类对数据进行标准化。标准化是一种常见的预处理步骤，它可以将所有特征的值调整到同一尺度，通常是均值为0，标准差为1。这个过程涉及两个主要步骤：计算每个特征的均值和标准差，然后将每个特征的值减去均值并除以标准差。

首先，我们创建了一个名为 `scaler` 的 `StandardScaler` 对象。这个对象将用于计算均值和标准差，并进行标准化转换。然后，我们使用 `fit_transform` 方法将数据集 `data_imputed` 传递给 `scaler` 对象进行处理。`fit_transform` 方法执行了两个步骤：首先，它通过计算每个特征的均值和标准差来学习数据集的统计信息，然后将数据集进行标准化转换。

标准化的目的是消除不同特征之间的量纲差异，确保它们都在相同的尺度上进行比较。在机器学习中，很多算法都对特征的尺度敏感。如果特征具有不同的尺度，某些特征的值范围可能会主导算法的训练过程，导致其他特征的影响力被忽略。通过将所有特征标准化到相同的尺度，我们可以确保每个特征对算法的影响程度是均等的。

标准化还有助于提高算法的收敛速度和稳定性。某些优化算法（如梯度下降）在处理标准化的数据时更容易找到最优解。此外，标准化还可以提高特征的解释性。当特征处于相同的尺度上时，我们可以更直观地理解它们对目标变量的影响。

需要注意的是，标准化是在训练数据上进行的，并且使用同一套参数对测试数据进行相同的标准化处理。这是为了保持训练集和测试集之间的一致性。在实际应用中，我们通常会将整个数据集划分为训练集和测试集，并且只使用训练集的统计信息来进行标准化，以避免测试集信息泄露。

总之，通过使用 `StandardScaler` 类对数据进行标准化，我们可以消除特征之间的尺度差异，提高机器学习算法的效果，并增强特征的解释能力。这是一个在数据预处理中常用且重要的步骤，为我们提供了更可靠和一致的分析结果。

```
1 | scaler = StandardScaler()
2 | data_standerd = scaler.fit_transform(data_imputed)
3 | return data_standerd, y
```

4.2 过采样处理

然后，我们使用 `KMeansSMOTE` 类进行过采样。`KMeansSMOTE` 是一种用于处理不平衡数据集的技术，它结合了 `K-Means` 聚类和 `SMOTE`（合成少数类过采样技术）来生成新的少数类样本。`K-Means` 聚类是一种常见的无监督学习算法，通过将数据分为 `K` 个簇来寻找数据的内在结构。它根据数据的特征相似性将样本分组，将相似的样本归为同一簇，从而在数据中发现隐藏的模式和结构。

在我们的情况中，`KMeansSMOTE` 利用 `K-Means` 聚类算法对数据集进行聚类，将样本分成不同的簇。然后，针对少数类样本较少的簇，我们使用 `SMOTE` 技术生成合成样本。`SMOTE`（`Synthetic Minority Over-sampling Technique`）是一种经典的过采样技术，它通过插值生成新的合成少数类样本，以增加少数类的样本数量。

`KMeansSMOTE` 的主要优势在于它结合了两方法的优点：`K-Means` 聚类可以帮助我们理解数据集的结构和模式，而 `SMOTE` 可以增加样本数量来平衡类别不平衡问题。通过将这两种技术结合起来，`KMeansSMOTE` 可以生成更加合理和真实的合成样本，有助于改善模型的性能和泛化能力。

通过使用 `KMeansSMOTE` 进行过采样，我们能够解决数据集中存在的类别不平衡问题。类别不平衡可能导致模型对多数类样本进行过度拟合，从而忽略少数类样本。通过生成合成样本，我们可以增加少数类样本的数量，使得模型能够更好地学习并识别少数类的特征和模式。

总结而言，`KMeansSMOTE` 是一种结合了 `K-Means` 聚类和 `SMOTE` 技术的过采样方法，适用于处理不平衡数据集。它通过聚类分析和合成样本生成，帮助我们处理类别不平衡问题，提高模型性能和泛化能力。这种技术在许多实际应用中被广泛使用，特别是在医疗诊断、金融风险预测和欺诈检测等领域，以处理不平衡数据集并改善模型的预测能力。

```
1 | kmeans_smote = KMeansSMOTE(random_state=42)
2 | X_train, y_train = kmeans_smote.fit_resample(X_train, y_train)
```

4.3 损失函数和优化器

然后，我们定义了 `criterion` 损失函数和 `optimizer` 优化器。损失函数是 `nn.CrossEntropyLoss()`，这是一种常用于多分类问题的损失函数。交叉熵损失函数用于衡量模型的预测结果与实际标签之间的差异。通过最小化交叉熵损失，我们可以优化模型的性能。

而优化器是 `optim.Adam(model.parameters(), lr=0.0001)`，它是一种基于梯度的优化算法。Adam优化器结合了动量（momentum）和自适应学习率的特性，能够高效地更新模型的参数。它根据参数的梯度信息自适应地调整学习率，从而在训练过程中实现更快速的收敛和更好的性能。

在初始化Adam优化器时，我们传入了模型的参数 `model.parameters()` 和学习率 `lr=0.0001`。模型的参数是需要被优化的权重和偏置项，在训练过程中通过反向传播算法来更新这些参数，使得损失函数最小化。学习率决定了每次参数更新的步长大小，较小的学习率可以使得模型更稳定地收敛，但训练速度可能较慢，而较大的学习率可能导致训练过程不稳定或错过最优解。

综上所述，通过定义交叉熵损失函数和Adam优化器，我们为模型训练提供了损失的度量和参数更新的策略。这将帮助我们在训练过程中逐步优化模型，提高其对多分类任务的性能。

```
1 | criterion = nn.CrossEntropyLoss()
2 | optimizer = optim.Adam(model.parameters(), lr=0.0001)
```


4.4 EarlyStopping | 早停策略

EarlyStopping（早停策略）是一种在训练过程中用于自动确定何时停止训练的技术。它基于模型的性能指标，例如验证集上的损失函数值，来判断模型是否已经达到了最佳状态或者是否进入了过拟合状态。当模型的性能不再改善或开始恶化时，**EarlyStopping** 能够及时中断训练，从而避免浪费计算资源和时间。

我们的代码中定义了一个名为 **EarlyStopping** 的类，用于实现提前停止的功能。它具有以下属性：

- **patience**（耐心）：指定在性能不再改善之前要等待的训练轮次数。如果在等待的轮次数内模型性能没有改善，则触发提前停止。
- **verbose**（冗长模式）：控制是否打印出提前停止的计数信息。
- **delta**（变化阈值）：用于定义性能改善的最小阈值。只有当性能指标的变化超过此阈值时，才被认为是真正的改善。

在初始化过程中，我们设置了一些初始变量。**counter**（计数器）用于跟踪在等待期间模型性能没有改善的训练轮次数。**best_score**（最佳分数）用于存储当前最佳的性能指标。**early_stop**（提前停止标志）用于指示是否应该提前停止训练。**val_loss_min**（最小验证损失）初始化为正无穷大，用于跟踪最佳验证损失的值。

__call__ 方法是该类的主要方法，用于每次验证集上计算损失函数后的回调。它接收验证集上的损失值 **val_loss** 和模型实例 **model** 作为输入参数。在每次调用时，它首先将 **val_loss** 转换为分数（负数，因为我们的目标是最小化损失函数）。然后，它根据以下几种情况来判断是否需要提前停止训练：

1. 如果 **best_score** 为 **None**（即第一次调用），则将当前分数设置为 **best_score**，并保存模型的检查点（即模型的当前状态）。
2. 否则，如果当前分数小于等于 **best_score + delta**，则说明模型性能没有改善。此时，**counter** 增加1，并打印出当前的提前停止计数信息。如果 **counter** 达到或超过了 **patience**，则将 **early_stop** 设置为 **True**，表示应该提前停止训练。
3. 否则，如果当前分数大于 **best_score + delta**，说明模型性能有所改善。我们更新 **best_score** 为当前分数，并保存模型的检查点。同时，将 **counter** 重置为0，以便在下一次等待期间重新计数。

通过使用 **EarlyStopping**，我们可以在模型开始过拟合之前自动停止训练，从而提高模型的泛化能力。这对于大型模型或训练时间较长的任务尤为重要，因为它可以避免不必要的计算开销和时间浪费。代码片段中的 **np.Inf** 表示正无穷大，是 **numpy** 库中的一个特殊值，用于初始化最小验证损失。

```
1 class EarlyStopping:
2     def __init__(self, patience=7, verbose=False, delta=0):
3         self.patience = patience
4         self.verbose = verbose
```

```

5         self.counter = 0
6         self.best_score = None
7         self.early_stop = False
8         self.val_loss_min = np.Inf
9         self.delta = delta
10
11     def __call__(self, val_loss, model):
12         score = -val_loss
13
14         if self.best_score is None:
15             self.best_score = score
16             self.save_checkpoint(val_loss, model)
17         elif score < self.best_score + self.delta:
18             self.counter += 1
19             print(f'EarlyStopping counter: {self.counter} out of
{self.patience}')
20             if self.counter ≥ self.patience:
21                 self.early_stop = True
22         else:
23             self.best_score = score
24             self.save_checkpoint(val_loss, model)
25             self.counter = 0

```

4.5 训练、保存模型

代码中涉及到 **t-SNE** 算法的部分，将在 [六、结果分析](#) 的分析进行详解。

在我们的代码中，模型的训练和保存流程如下：

首先，我们定义了一些训练参数，包括训练的轮数（`num_epochs`）、每批次的样本数（`batch_size`）以及早停策略（`EarlyStopping`）。

接着，我们通过循环迭代每个训练轮次（`epoch`）。在每个轮次中，我们使用随机排列（`permutation`）的方式对训练数据集（`X_train`）进行处理，以获取随机的批次索引（`indices`）。然后，根据这些索引获取相应的训练样本（`batch_X`）和标签（`batch_y`）。

在获取了训练样本和标签后，我们进行前向传播，将批次样本输入模型（`model`），得到模型的输出（`outputs`）。然后，我们计算损失函数（`loss`）来评估模型的预测结果与真实标签之间的差异。

在计算了损失后，我们进行反向传播和优化。首先，我们将梯度归零（`optimizer.zero_grad()`），然后根据损失进行反向传播（`loss.backward()`），最后执行优化器的一步更新（`optimizer.step()`）来更新模型的参数。

在每个轮次结束后，我们打印训练损失（`loss.item()`）和验证集损失（`loss_test.item()`），并将它们分别添加到训练损失列表（`train_loss_list`）和测试损失列表（`test_loss_list`）中。

接下来，我们进行模型在测试集上的评估，通过将测试集样本（`X_test`）输入模型，得到预测结果（`outputs`）。然后，我们计算预测的准确率（`accuracy`）并打印出来。

我们还计算模型在验证集上的损失（`val_loss`）。然后，我们运行早停策略（`early_stopping`）来判断是否提前停止训练，如果满足早停的条件（`early_stopping.early_stop`），则打印"early stop"并跳出训练循环。

在所有训练轮次完成后，在测试集上进行最终评估，计算模型的准确率并打印出来。

最后，我们调用 `torch.save` 函数对模型进行保存，并更新模型保存的路径。这样，我们就完成了模型的训练和保存流程。

```
1  def train(input_path, file_basic_path, output_path):
2      global input_size_global, model_path
3      train_loss_list = []
4      test_loss_list = []
5      data = pd.read_csv(input_path)
6      X_train, y_train = preprocess(data)
7      kmeans_smote = KMeansSMOTE(cluster_balance_threshold=0.064,
random_state=42)
8      X_train, y_train = kmeans_smote.fit_resample(X_train, y_train)
9      X_train, X_test, y_train, y_test = train_test_split(X_train, y_train,
test_size=0.2, random_state=42)
10     # -----MLP-----#
11     # 转换为张量
12     X_train = torch.Tensor(X_train)
13     X_test = torch.Tensor(X_test)
14     y_train = torch.LongTensor(y_train.values)
15     y_test = torch.LongTensor(y_test.values)
16     input_size = X_train.shape[1] # 输入特征的维度
17     input_size_global = X_train.shape[1]
18     hidden_size = 200 # 隐藏层的大小
19     num_classes = 6 # 类别数量
20     model = MLP(input_size, hidden_size, num_classes)
21
22     # 定义损失函数和优化器
23     criterion = nn.CrossEntropyLoss()
24     optimizer = optim.Adam(model.parameters(), lr=0.0001)
25     # 训练模型
26     num_epochs = 100
27     batch_size = 30
28     early_stopping = EarlyStopping(patience=5, verbose=True)
```

```

29
30     for epoch in range(num_epochs):
31         permutation = torch.randperm(X_train.size()[0])
32         for i in range(0, X_train.size()[0], batch_size):
33             indices = permutation[i:i + batch_size]
34             batch_X, batch_y = X_train[indices], y_train[indices]
35
36             # 前向传播
37             outputs = model(batch_X)
38             loss = criterion(outputs, batch_y)
39             predict = model(X_test)
40             loss_test = criterion(predict, y_test)
41
42             # 反向传播和优化
43             optimizer.zero_grad()
44             loss.backward()
45             optimizer.step()
46
47
48             # 每个epoch打印损失
49             print(f"Epoch [{epoch + 1}/{num_epochs}], Loss: {loss.item():.4f},
Validate Loss: {loss_test.item():.4f}")
50             train_loss_list.append(loss.item())
51             test_loss_list.append(loss_test.item())
52
53             # early_stopping(loss.item(), model, val_features,
'checkpoint.pt')
54             # 在测试集上进行评估
55             with torch.no_grad():
56                 outputs = model(X_test)
57                 _, predicted = torch.max(outputs.data, 1)
58                 accuracy = (predicted == y_test).sum().item() / y_test.size(0)
59                 print(f"Test Accuracy: {accuracy:.4f}")
60                 val_loss = criterion(outputs, y_test)
61                 early_stopping(val_loss, model)
62                 if early_stopping.early_stop:
63                     print("early stop")
64                     break
65
66             # 在测试集上进行评估
67             with torch.no_grad():
68                 outputs = model(X_test)
69                 _, predicted = torch.max(outputs.data, 1)
70                 accuracy = (predicted == y_test).sum().item() / y_test.size(0)
71                 print(f"Test Accuracy: {accuracy:.4f}")
72             label_count_dict = [0, 0, 0, 0, 0, 0]

```

```

73     for value in data['label']:
74         label_count_dict[value] += 1
75     val_features = model(X_test, return_features=True)
76     val_features_tsne = TSNE(n_components=2, random_state=33, init='pca',
77
learning_rate='auto').fit_transform(
78         val_features)
79     font = {"color": "darkred",
80            "size": 13,
81            "family": "serif"}
82
83     plt.style.use("dark_background")
84     plt.figure(figsize=(9, 8))
85
86     plt.scatter(val_features_tsne[:, 0], val_features_tsne[:, 1],
c=y_test.cpu().numpy(), alpha=0.6,
87                 cmap=plt.cm.get_cmap('rainbow', num_classes))
88     plt.title("t-SNE", fontdict=font)
89     cbar = plt.colorbar(ticks=range(num_classes))
90     cbar.set_label(label='Class label', fontdict=font)
91     plt.clim(-0.5, num_classes - 0.5)
92     plt.tight_layout()
93     plt.savefig(file_basic_path + f'_tsne.png', dpi=300)
94     result = {"train_losses": train_loss_list, "val_losses":
test_loss_list, "label_counts": label_count_dict}
95     torch.save(model.state_dict(), file_basic_path + '_model.pth')
96     model_path = file_basic_path + '_model.pth'
97     return result

```

4.6 调用模型进行预测

代码中涉及到 **t-SNE** 算法的部分，将在 [六、结果分析](#) 的分析进行详解。

在我们的代码中，调用模型进行预测的流程如下：

首先，它使用pandas的 `read_csv` 函数从 `input_path` 读取数据，然后对数据进行预处理。预处理的具体步骤在这段代码中并未给出，但通常包括数据清洗、缺失值处理、数据标准化等步骤。

接着，它定义了一个多层感知器（**MLP**）模型。这个模型有一个输入层，一个隐藏层和一个输出层。输入层的大小由全局变量 `input_size_global` 决定，隐藏层的大小为200，输出层的大小为6，这意味着模型可以预测6个不同的类别。

然后，它从 `model_path` 加载预训练的模型参数，并使用这些参数对数据进行预测。预测的结果是一个概率分布，表示数据属于每个类别的概率。`torch.max` 函数用于找到概率最大的类别作为预测结果。

之后，它统计了每个类别的预测数量，并将预测结果保存到一个json文件中。这个文件的路径由 `file_basic_path` 和固定的字符串 `'_predicted.json'` 拼接而成。

接下来，它使用 `t-SNE` 算法对模型的特征进行降维，并将降维后的特征可视化。在这个过程中，它首先获取模型的特征，然后使用 `t-SNE` 算法将特征降到2维，最后使用 `matplotlib` 库将降维后的特征绘制成散点图。在这个散点图中，每个点的颜色表示它的类别。

最后，它返回一个字典，其中包含每个类别的预测数量。

这段代码的主要作用是对输入的数据进行预测，并将预测结果以及一些相关信息保存到文件中。同时，它还对模型的特征进行了可视化，这有助于我们理解模型的预测结果。

```
1  def test(input_path, file_basic_path, output_path):
2      global input_size_global, model_path
3      data = pd.read_csv(input_path)
4      data = test_preprocess(data)
5      input_size = input_size_global # 输入特征的维度
6      hidden_size = 200 # 隐藏层的大小
7      num_classes = 6 # 类别数量
8      model = MLP(input_size, hidden_size, num_classes)
9      model.load_state_dict(torch.load(model_path))
10     with torch.no_grad():
11         data = torch.Tensor(data)
12         outputs = model(data)
13         _, predicted = torch.max(outputs.data, 1)
14
15     label_count_list = [0, 0, 0, 0, 0, 0]
16     for value in predicted:
17         label_count_list[value] += 1
18
19     result_dict = {}
20     predicted_new = predicted.tolist()
21     for i, value in enumerate(predicted_new, start=1):
22         result_dict[str(i)] = value
23     with open(file_basic_path + '_predicted.json', 'w') as json_file:
24         json.dump(result_dict, json_file)
25     result = {"label_counts": label_count_list}
26     val_features = model(data, return_features=True)
27     val_features_tsne = TSNE(n_components=2, random_state=33, init='pca',
28
29     learning_rate='auto').fit_transform(
30         val_features)
```

```

30     font = {"color": "darkred",
31             "size": 13,
32             "family": "serif"}
33
34     plt.style.use("dark_background")
35     plt.figure(figsize=(9, 8))
36
37     plt.scatter(val_features_tsne[:, 0], val_features_tsne[:, 1],
38               c=predicted.cpu().numpy(), alpha=0.6,
39               cmap=plt.cm.get_cmap('rainbow', num_classes))
40     plt.title("t-SNE", fontdict=font)
41     cbar = plt.colorbar(ticks=range(num_classes))
42     cbar.set_label(label='Class label', fontdict=font)
43     plt.clim(-0.5, num_classes - 0.5)
44     plt.tight_layout()
45     plt.savefig(file_basic_path + f'_tsne.png', dpi=300)
46     return result

```

##

五、原理分析

5.1 深度学习与神经网络

深度学习是机器学习的一个分支，它试图模拟人脑的工作原理，通过训练大量的数据来学习数据的内在规律和表示。在深度学习中，最基本的单元是神经元，这是受到生物神经元的启发而设计的。在人工神经网络中，神经元通过连接和权重将输入信号转化为输出信号。一个深度神经网络由多个这样的神经元层叠加而成。

5.2 多层感知机 (MLP)

在一些任务中，最常用的深度学习模型类型是多层感知机（MLP）。MLP是一种前向结构的人工神经网络，映射一组输入向量到一组输出向量。MLP由多个层级组成，每个层级都完全连接到下一个层级。在MLP中，最后一个层级是输出层，而除输入层外的其它层级都是隐藏层。

5.3 激活函数

ReLU（Rectified Linear Unit）是深度学习中常用的一种激活函数。它的公式是：

$$f(x) = \max(0, x)$$

ReLU 函数有一个很好的特性，那就是它能够在训练深度神经网络时帮助缓解梯度消失问题。这是因为 **ReLU** 函数在输入大于0时，其梯度始终为1。因此，对于深度模型，ReLU函数可以保证反向传播时梯度不会消失，从而有效地更新模型的参数。

5.4 损失函数和优化器

在神经网络的训练过程中，我们需要有一种方式来度量模型的表现。这就是所谓的损失函数，它可以量化模型的预测结果与真实结果之间的差距。在这个项目中，我们使用的是交叉熵损失函数，它是分类问题中常用的损失函数。

优化器是用来更新和调整模型参数以最小化损失函数的工具。在这个项目中，我们使用的是Adam优化器。Adam是一种自适应的学习率优化算法，它结合了动量优化和RMSProp的思想。

5.5 数据预处理

在训练MLP模型时，对输入数据进行预处理是非常重要的。常见的预处理步骤包括标准化（使输入数据的每一特征都有相同的尺度）、缺失值处理（填充或删除缺失值）以及类别编码（将类别变量转化为数值变量）。预处理步骤可以帮助提高模型的性能和稳定性。

六、结果分析

6.1 t-SNE分析

t-SNE（t-Distributed Stochastic Neighbor Embedding）是一种用于数据可视化的机器学习算法。它是一种降维技术，特别适合于高维数据的可视化。t-SNE的主要目标是在低维空间（通常是二维或三维）中保留高维空间中的相似性结构。

以下是对t-SNE的基本介绍：

- 工作原理：** t-SNE首先在高维空间中计算数据点之间的相似度，然后在低维空间中创建一个概率分布，使得低维空间中的点的相似度尽可能接近高维空间中的相似度。这样，高维空间中相似的点在低维空间中也会相似。
- 优点：** t-SNE相比于其他降维方法（如PCA）的优点在于，它能更好地保留局部结构。这意味着在高维空间中相互靠近的点在低维空间中也会靠近。这使得t-SNE非常适合用于数据可视化。
- 缺点：** t-SNE的一个主要缺点是它可能不会很好地保留全局结构。也就是说，高维空间中远离的点在低维空间中可能不会远离。此外，t-SNE的结果也可能对初始参数（如随机种子）敏感。
- 应用：** t-SNE广泛应用于各种领域，包括机器学习、数据挖掘、信息检索和生物信息学等。它常常用于探索性数据分析，以帮助研究人员理解高维数据的结构。
- 算法步骤：**
 - 在高维空间中计算点之间的条件概率，使得相似的对象有更高的概率被选中。
 - 在低维空间中定义相似的条件概率，并最小化两个条件概率分布的KL散度。
 - 通过梯度下降等优化方法，找到最佳的低维空间表示。

以上是对t-SNE的基本介绍，如果你需要更详细的信息或者具体的使用方法，我可以提供更多的帮助。

6.2 模型性能分析

- 模型大小：我们的模型大小为49kb，这是一个非常轻量级的模型，可以在各种设备上快速运行。
- 性能分析：我们通过F1评估指标对模型的性能进行了详细的分析。从F1评估指标可以看出，我们的模型在各个故障类别上的表现情况。例如，故障类别0的F1评估指标为0.80，准确率为0.72，召回率为0.90，这表明模型在识别故障类别0时具有较好的性能。然而，对于故障类别1和2，模型的性能相对较差，可能需要进一步优化。

6.3 误差分析

- 误识别情况：我们通过F1评估指标对模型的误识别情况进行了详细的分析。从F1评估指标可以看出，我们的模型在识别某些故障类别时存在误识别和漏识别的问题。例如，故障类别1的F1评估指标较低，可能是由于准确率和召回率都不高，表明模型在识别故障类别1时容易出现误识别和漏识别的问题。另外，故障类别2的F1评估指标也相对较低，主要是由于准确率较低，说明模型在识别故障类别2时容易产生误识别的问题。

6.4 结论与建议

综合考虑模型的性能分析和误差分析，我们得出以下结论和建议：

1. 故障类别0、3、4和5的性能表现较好，模型能够较为准确地识别这些故障类型，但仍有优化的空间。
2. 故障类别1和2的性能相对较差，可能需要增加样本量或进行数据增强，以提高模型对这些故障的识别能力。
3. 我们需要进一步优化模型的特征选择和超参数调整，以提高模型的整体性能和泛化能力。
4. 我们需要继续收集更多的真实故障数据，以进一步改进模型的训练效果。
5. 对于特定故障类型的误识别和漏识别问题，我们可以考虑引入更多的相关特征，或者尝试其他的机器学习或深度学习模型。

通过不断的优化和改进，我们有望提高模型的准确性和可靠性，进一步支持分布式系统的高效运维，提供更好的故障诊断服务。