

# 银瞳：基于自适应语义空间学习的中文金融多任务大模型

周宇航<sup>1</sup>, 李泽平<sup>1</sup>, 田思雨<sup>1</sup>, 倪雨琛<sup>2</sup>, 张健<sup>3</sup>, 刘响<sup>4</sup>, 叶广楠<sup>1</sup>, 吴杰<sup>1</sup>, 柴洪峰<sup>1</sup>

<sup>1</sup>复旦大学 计算机科学技术学院 金融科技研究院, 上海

<sup>2</sup>同济大学电子与信息工程学院, 上海

<sup>3</sup>达观数据, 上海

<sup>4</sup>纽约大学坦登工学院, 纽约

{yuhangzhou22, zepingli23, sytian23}@m.fudan.edu.cn

2230733@tongji.edu.cn; zhangjian@datagrand.com; xl493@nyu.edu

{yegn, jwu, hfchai}@fudan.edu.cn

## 摘要

大语言模型正逐渐被用于各种垂直领域, 利用其广泛的知识储备来赋能领域中的多种场景。然而, 各领域拥有多种待学习的特定任务, 且多源异构的领域数据容易引发模型进行任务迁移时的冲突。基于此, 本研究提出自适应语义空间学习框架, 利用对语义空间内数据的自适应重分布, 提升多专家模型的性能及选择效果, 并基于此框架训练了一个金融多任务大模型“银瞳”。研究结果表明, 我们的框架只需利用10%的数据就能达到接近全数据训练的效果, 并拥有较强的泛化表现。

**关键词:** 金融大模型; 自适应语义空间学习; 多任务学习

## SilverSight: A Multi-Task Chinese Financial Large Language Model Based on Adaptive Semantic Space Learning

Yuhang Zhou<sup>1,2</sup> Zeping Li<sup>1,2</sup> Siyu Tian<sup>1,2</sup> Yuchen Ni<sup>3</sup>  
Jian Zhang<sup>4</sup> Xiang Liu<sup>5</sup> Guangnan Ye<sup>1,2</sup> Jie Wu<sup>1,2</sup> Hongfeng Chai<sup>1,2</sup>

<sup>1</sup>Institute of FinTech, Fudan University

<sup>2</sup>School of Computer Science, Fudan University

<sup>3</sup>School of Electronics and Information Engineering, Tongji University

<sup>4</sup>DataGrand Inc.

<sup>5</sup>Tandon School of Engineering, New York University

## Abstract

Large language models (LLMs) are increasingly being applied across various specialized fields, leveraging their extensive knowledge to empower a multitude of scenarios within these domains. However, each field encompasses a variety of specific tasks that require learning, and the diverse, heterogeneous data across these domains can lead to conflicts during model task transfer. In response to this challenge, our study introduces an Adaptive Semantic Space Learning (ASSL) framework, which utilizes the adaptive reorganization of data distributions within the semantic space to enhance the performance and selection efficacy of multi-expert models. Utilizing this framework, we trained a financial multi-task LLM named "SilverSight". Our research findings demonstrate that our framework can achieve results close to those obtained with full data training using only 10% of the data, while also exhibiting strong generalization capabilities.

**Keywords:** Financial LLMs, Adaptive Semantic Space Learning, Multi-task Learning

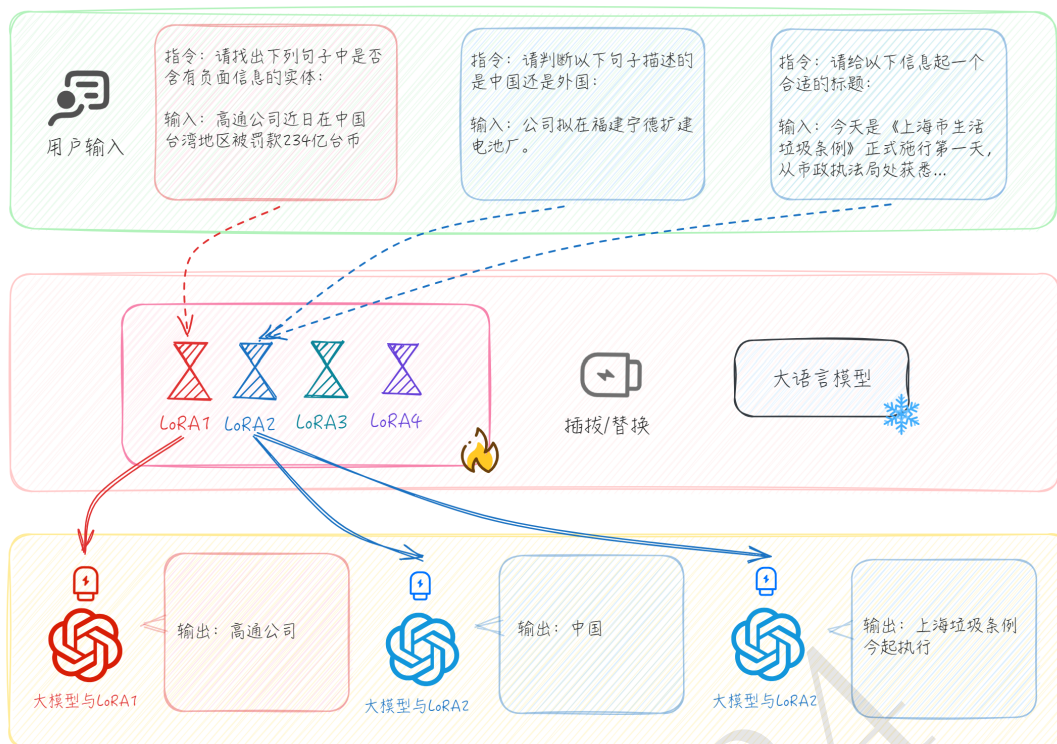


Figure 1: “银瞳”金融多任务大模型流程示意

## 1 引言

近期, 大型语言模型如GPT-4 (OpenAI, 2023)和LLaMA (Touvron et al., 2023)等, 在自然语言处理领域各项任务中展现了惊人的实力。这些模型凭借强大的知识存储和上下文理解能力, 其优异性能在金融 (Zhou et al., 2024; Chen et al., 2023)、法律 (Cui et al., 2023)等多个专业领域已得到广泛应用。在构建大语言模型的过程中, 指令微调技术扮演了核心角色, 它利用有监督的数据集对模型进行微调, 助力模型从简单的文本续写向复杂任务解答进行风格跃迁 (Qiu et al., 2020)。然而, 受制于特定领域的的数据量和计算资源限制, 指令微调往往需要借助参数高效的微调策略 (Houlsby et al., 2019), 这些策略仅需对基础模型进行少量参数的更新或增加, 便能显著提升模型对指令的响应能力。在此背景下, 低秩自适应(Low-Rank Adaptation, LoRA)方法 (Hu et al., 2021)应运而生, 通过为原始模型引入可分解的低参数量旁路矩阵, 实现了模块化改造及性能全面提升, 并具备了简便的“即插即用”特性。

当前, 自然语言处理领域的研究热点之一是采用混合专家模型(Mixture of Experts, MoE) (Jordan and Jacobs, 1994)来解决多任务处理的挑战。混合专家模型是一种集成学习策略, 它通过整合多个专家模型, 利用每个模型在子问题上的专长, 在处理复杂任务时提升整体性能 (Dou et al., 2023)。Zadouri (Zadouri et al., 2023)首次提出将LoRA方法与MoE框架相结合, 采用基于token级别的软路由机制, 对各个专家模型的输出进行加权融合, 有效提升了模型的性能。近期, 研究者进一步提出了LoRAMoE方法 (Dou et al., 2023), 该方法将专家模型划分为两组: 一组专注于通用世界知识的处理, 另一组则致力于在学习指令微调过程中遇到的新任务。这种设计旨在增强模型处理下游任务的能力, 同时保留大模型积累的世界知识。这些多专家策略在处理不同问题时展现出了优势, 但它们对于token级别的处理以及训练后专家和路由策略的固定性, 限制了模型快速适应新指令或场景的能力。为了解决这一问题, 研究人员提出了LoraRetriever方法 (Zhao et al., 2024), 该方法借鉴了检索思想, 将每个专家模型中的12条数据嵌入的平均值作为LoRA的自身嵌入, 并利用输入问题的嵌入来检索最匹配的LoRA, 从而有效弥补了多专家模型在动态增删方面的不足, 进一步提升了混合训练模型的性能和泛化能力。

尽管上述方法在提升LoRA等适配器的动态适应性方面表现出色, 但它们在训练LoRA时往往依赖于任务类型对训练数据进行手动划分, 并未充分考虑数据在语义空间中的内在联系。这

种处理方式可能导致LoRA专家的嵌入产生偏差，从而影响系统的整体性能。与此同时，来自不同来源的数据在多样性和质量上存在差异 (Touvron et al., 2023)，而如何合理配置不同数据量的比例在先前的研究中并未得到充分探讨。这进一步加剧了LoRA专家选择和数据选择的问题，使其成为一个亟待解决的关键挑战。因此，探索一种能够同时优化LoRA专家选择和数据选择策略的方法，以保持数据在语义空间中的连续性和关联性，成为了提升混合专家模型性能的重要研究方向。这要求我们不仅要关注LoRA的训练过程，还要重视数据的分布修正机制，以确保模型能够在面对复杂多变的任务时，依然能够保持高效和准确的性能表现。

为了填补现有研究的空白，我们提出了自适应语义空间学习框架(Adaptive Semantic Space Learning, ASSL)，并通过金融多任务数据集训练了一个名为“银瞳”(SilverSight)的金融多任务大语言模型。ASSL框架利用多任务训练数据在语义空间中的相似性进行聚类，从而摒弃了对训练数据任务类型的预定义。实验结果表明，与基于预设任务类型的分类方法相比，我们的方法具有较为显著的优势，确保了每个专家模型能够被分配到与其最相关的下游任务中。同时，我们发现不同来源和格式的训练任务可能引发不可避免的任务冲突，影响模型的指令遵循能力。此外，通过利用语义空间的相似性，我们的方法能够聚合互补任务的多种数据，从而在相关任务上提升模型性能。在ASSL框架下，我们在每个聚类中采用基于模型自我进化的数据重分布策略，进行自适应的数据挑选。通过两阶段的数据自适应筛选，我们确保了用于训练每个LoRA专家的少量数据具备高质量、高覆盖率以及高必要性。在原始数据质量及数量分布不均的情况下，ASSL框架通过考虑语义空间分布密度和模型自反馈机制，对每个LoRA的训练数据进行自适应筛选。这种方法使得模型能够对长尾分布中的头部相似性数据和尾部小样本数据进行更均衡的拟合，从而提高了模型在多样任务上的泛化能力和性能表现。

我们从金融领域的23个不同来源收集了22万条中文金融微调数据，并将这些数据分类为情感分析、金融问答、文本生成、金融选择题等多个任务类型。基于这些实验数据，我们得出以下主要发现：(1) 通过利用语义空间中的相似性进行聚类，我们能够识别出相互促进和存在冲突的训练任务。采用多个专家模型分别针对特定领域任务进行学习，可以使得每个专家专注于其擅长的领域，实现“各司其职”的效果。(2) 结合数据在语义空间中的密度分布和模型自身对训练数据的需求，我们能够有效地对数据进行语义平滑重分布。这种方法使得整个系统能够在仅使用10%的数据进行微调的情况下，达到与使用全量数据进行微调相似的效果。(3) 通过对聚类内部的数据进行分布平滑处理，我们利用聚类内数据嵌入的质心作为LoRA专家的嵌入，从而实现了LoRA选择的最优化。

这些发现不仅验证了ASSL框架在金融多任务学习中的有效性，而且为如何高效利用有限数据资源、提高模型性能提供了新的视角和方法。这些成果对于推动领域自然语言处理技术的发展具有重要意义。我们的主要贡献如下：

1. 我们创新性地提出了自适应语义空间学习(ASSL)框架，通过分析数据在语义空间的分布特性，实现了多任务场景下LoRA专家与数据的自适应选择机制。此举有效调整了专家与训练数据的分布平衡，借助语义空间信息对训练数据进行平滑处理，从而避免了专家过拟合或训练过程中的冲突现象。
2. 我们在金融领域对ASSL框架的有效性进行了实证验证，并成功训练了一个中文金融多任务大语言模型“银瞳”。在两个评测基准上，我们对模型的泛化能力及多任务表现进行了全面测试，充分展现了该框架的优越性和应用潜力。
3. 我们对模型不同能力的来源进行了深入的多维度分析，为未来专家与数据自适应选择领域的研究提供了新的思考方向。

## 2 准备工作

### 2.1 LoRA

对大语言模型进行全参数微调需要耗费大量的数据和计算资源，并且可能会导致灾难性遗忘的问题 (Luo et al., 2023)。为了解决这一问题，近期研究人员开发了多种高效的微调方法，其中LoRA (Hu et al., 2021)因其有效性而被广泛应用。LoRA方法利用低秩矩阵分解的思想来减少微调过程中需要调节的参数数量，通过在网络中引入一个附加的旁路结构 $W_0^{m \times n}$ 来实现。具

体来说, LoRA通过增加 $W_0^{m \times n}$ 到原始权重矩阵 $W^{m \times n}$ 上, 来修改网络中的权重, 即:

$$W'^{m \times n} = W^{m \times n} + W_0^{m \times n} = W^{m \times n} + A^{m \times r} B^{r \times n} \quad (1)$$

其中,  $A^{m \times r}$ 和 $B^{r \times n}$ 是从 $W_0^{m \times n}$ 通过低秩分解得到的两个矩阵,  $r$ 为矩阵的秩。当输入模型的是token表征 $x^{n \times 1}$ 时, 模型的输出 $y^{m \times 1}$ 可以表示为:

$$y = W'x = (W + AB)x = Wx + ABx \quad (2)$$

这里,  $x^{n \times 1}$ 代表模型输入的token表征, 而 $y^{m \times 1}$ 代表模型输出。

## 2.2 任务定义

在ASSL框架中, 我们利用数据在语义空间的分布情况, 设计了一个多专家系统。主要关注两个核心任务: 专家自适应选择和数据自适应选择。这两个任务旨在提高系统在处理多任务环境下的灵活性和效率。

专家自适应选择任务侧重于在给定的输入情况下, 如何自适应地挑选最合适的LoRA专家来处理特定的任务。考虑到本框架的多任务系统是一个由 $n$ 种不同任务组成的集合 $T = \{t_1, t_2, \dots, t_n\}$ , 以及 $m$ 个不同的LoRA专家构成的集合 $E = \{e_1, e_2, \dots, e_m\}$ 。其中每个专家 $e_i$ 都经过若干特定任务的训练, 以优化其性能。对于任意输入 $x$ , 我们的目标是找到一个映射函数 $f: X \rightarrow E$ , 使得 $f(x)$ 能够选择出最适合当前输入 $x$ 的专家 $e$ 。这可以通过最大化输入 $x$ 与专家 $e$ 之间相关度的方式来实现, 即:

$$e^* = f(x) = \arg \max_{e \in E} \text{Relevance}(x, e) \quad (3)$$

其中, 函数 $\text{Relevance}(x, e)$ 衡量了输入 $x$ 与专家 $e$ 之间的相关性。

另一方面, 数据自适应选择任务专注于处理来自多个源的异构数据, 并解决数据中存在的长尾分布问题。我们有一个数据集 $D = \{d_1, d_2, \dots, d_l\}$ , 每个 $d_i$ 代表不同源的数据。在实际应用中, 这些数据往往遵循长尾分布, 即大量数据集中在少数类别中, 而大部分类别仅包含少量数据。我们的目标是通过一个转换函数 $g: D \rightarrow D'$ 来处理原始数据集 $D$ , 将其转换为一个更平滑、分布更均匀的数据集 $D'$ , 从而减轻长尾分布的影响。

## 3 自适应语义空间学习框架

我们提出了一种名为ASSL (自适应语义空间学习) 的框架, 该框架旨在实现对LoRA专家及其数据自适应选择的功能。框架示意如 Figure 2 所示, 它确保了每个专家能够发挥最佳性能, 并使整个系统在多任务集合上拥有卓越的表现。在本章中, 我们将详细阐述框架的两个核心组成部分, 即如何通过数据在语义空间的重新分布来统一实现专家以及数据之间的自适应选择。

### 3.1 LoRA专家的自适应选择

在处理LoRA专家的自适应选择时, 我们的目标是优化数据划分以避免任务间的潜在冲突, 并确保对于每个输入, 都能匹配到最擅长处理该类问题的专家。这一目标主要分为两个方面: (1) 如何更有效地训练多任务指令; (2) 如何根据用户输入选择最合适的LoRA专家。

首先, 为了增强系统对多样化指令的泛化能力, 我们对任务集中的每个子任务进行了指令扩充。具体而言, 我们为每个任务人工编写了30条具有相似语义但表述不同的指令。然后, 使用句子编码器 $\text{Emb}(\cdot)$ 对每条指令及其输入数据的拼接 $\text{Ins} \oplus \text{Inp}$ 进行编码, 获得所有数据在同一语义空间中的嵌入向量。我们采用K-Means聚类方法 (MacQueen and others, 1967) 将语义空间中相近的句子聚集成 $K$ 类, 以此来优化任务的数据划分。聚类过程可以通过以下公式表示:

$$\text{Cluster}_k = \arg \min_{x \in X} \sum_{i=1}^K \sum_{x \in C_i} \|x - \mu_i\|^2 \quad (4)$$

其中,  $X$ 表示所有数据点的集合,  $C_i$ 是第 $i$ 个聚类中的数据点集合, 而 $\mu_i$ 是第 $i$ 个聚类的质心。

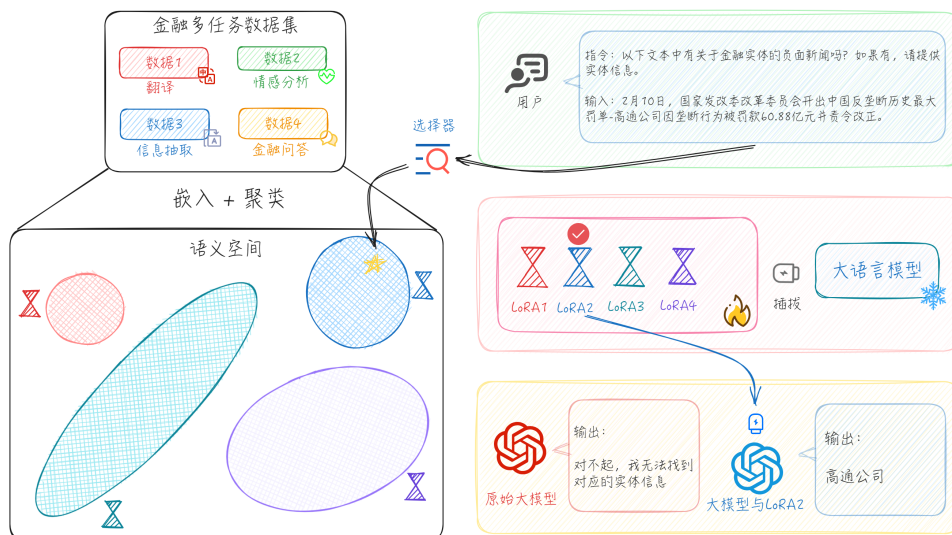


Figure 2: 自适应语义空间学习框架

通过 4.5 实验证明，与混合任务训练以及按预定义标签进行划分训练的方式相比，这种基于语义的聚类训练法能够显著提升系统性能。针对每个LoRA专家，我们选择其聚类的质心作为该专家的语义嵌入，每个聚类的质心 $\mu_i$ 是聚类中所有点语义嵌入的平均位置，计算公式如下：

$$\text{Emb}(e_i) = \mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x \quad (5)$$

其中， $C_i$  是第  $i$  个聚类中的数据点集合，而  $|C_i|$  表示集合  $C_i$  中元素的数量。 $\mu_i$  代表聚类  $i$  中所有点的均值，即该聚类的质心。每当有用户输入时，系统通过以下公式寻找与用户输入语义嵌入最接近的专家进行响应：

$$e^* = \arg \min_{e \in E} \|\text{Emb}(x) - \text{Emb}(e)\| \quad (6)$$

这里， $\text{Emb}(x)$  是用户输入的嵌入向量， $\text{Emb}(e)$  是专家的语义嵌入向量，而  $e^*$  则是被选中的专家。通过此种匹配方式，系统能够在语义空间中找到与训练任务最匹配的专家。

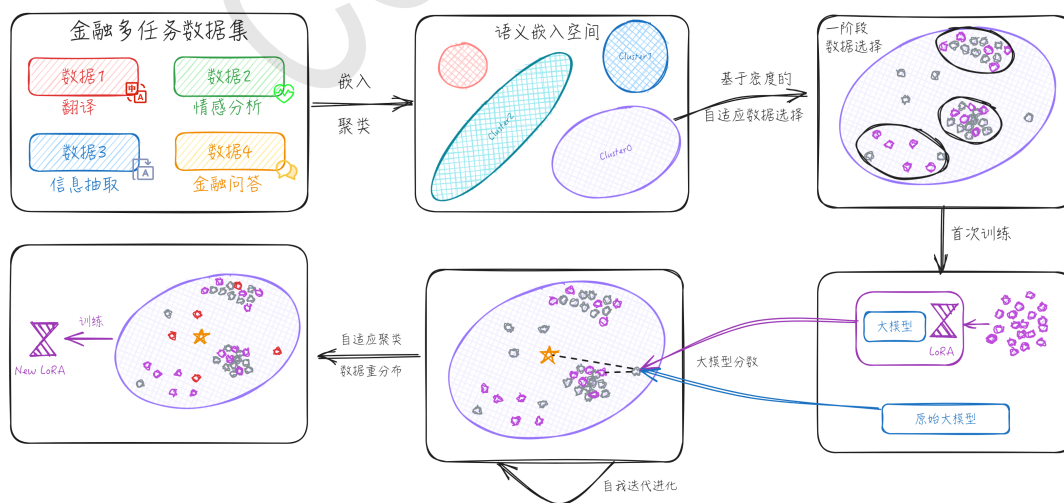


Figure 3: 自适应语义空间数据重分布流程

### 3.2 多源异构数据的自适应选择

在本章中，我们讨论如何对多任务有监督数据进行聚类，以隔离相互冲突的任务并聚集相互增强的任务。虽然这种方法有效地解决了任务间的冲突，但它也引入了新的挑战，如数据配比失衡和质量不一，特别是长尾数据分布的问题。为了解决这些问题，我们为每个聚类设计了一个两阶段的数据重分布操作，目的是在保持少量数据集上进行有效微调的同时，达到比拟全数据集微调的效果，具体流程见 Figure 3。

在第一阶段，针对聚类簇中数据失衡的问题，我们基于DBSCAN算法 (Schubert et al., 2017)设计了一种自适应调整的A-DBSCAN算法，对每个聚类中的数据进行嵌套聚类，具体算法流程见 A.1。该算法可以根据数据密度在不同区域中动态调整连通数量。具体实现步骤如下：首先，算法通过K-最近邻算法(K-Nearest Neighbor, KNN) (Peterson, 2009)距离计算框架来评估数据点在语义空间中的局部密度。每个数据点的局部密度为其至最近的 $k$ 个邻居的平均距离的倒数，其数学表达式为：

$$\rho_i = \frac{1}{\frac{1}{k} \sum_{j=1}^k d(x_i, x_{ij})} \quad (7)$$

其中， $d(x_i, x_{ij})$ 表示数据点 $x_i$ 与其第 $j$ 个最近邻居 $x_{ij}$ 之间的距离。

随后，算法根据计算得到的局部密度值将数据点进行排序，形成一个优先队列，优先处理局部密度较高的数据点。在每次迭代中，选取队列中局部密度最高的点作为起始点，围绕此点进行簇的生成。

在自适应过程中，算法利用队列中所有数据点KNN距离的中位数定义邻域半径 $\epsilon$ 。同时，启发式地将邻域节点数 $MinPts$ 的初始值设置为：

$$MinPts_{init} = \frac{\epsilon \times \rho_{max}}{2} \quad (8)$$

其中， $\rho_{max}$ 表示全局最大局部密度，即初始优先队列中第一个数据点的局部密度。在每个簇形成后， $MinPts$ 根据以下公式进行更新，以适应当前的局部密度环境：

$$MinPts_{update} = \max \left( 2, \frac{\rho_{current}}{\rho_{max}} \times MinPts \right) \quad (9)$$

此处， $\rho_{current}$ 表示当前优先队列中第一个点的局部密度。这种动态调整策略使得算法能够更灵活地适应不同密度的数据分布，提高聚类的准确性和效率。同时，这一策略允许我们在每个子簇中平均选取数据量，同时过滤掉不连通的噪声点，有效地对高密度数据进行下采样，对低密度数据进行上采样，从而避免过拟合和欠拟合。

在第二阶段，我们利用第一阶段筛选出的少量数据对模型进行初步的LoRA微调。参考先前的研究，大语言模型在微调阶段主要学习新的语言风格分布，而难以获得新的领域知识 (Ren et al., 2024)。因此，我们考虑数据与大模型自身知识的冲突情况，根据初步训练前后大模型对未选数据的得分差异，设计了两种得分机制来评估每条数据对当前模型的价值。差异得分和比例得分分别定义为：

$$Score_{diff}(x) = LLM_{Raw}(x) - LLM_{LoRA}(x) \quad (10)$$

$$Score_{prop}(x) = -\frac{LLM_{LoRA}(x)}{LLM_{Raw}(x) + 1} \quad (11)$$

$$Score_{llm}(x) = Score_{diff}(x) + Score_{prop}(x) \quad (12)$$

其中， $Score_{Raw}(x)$ 以及 $Score_{LoRA}(x)$ 通过Rouge方法 (Lin, 2004)计算。此外，为了确保新选数据的质量和对整个聚类簇的覆盖度，受启发于MMR公式 (Carbonell and Goldstein, 1998)，我们设计了训练数据的效用函数：

$$U(d_{new}) = \lambda_1 \cdot \text{sim}(\mu, d_{new}) - \lambda_2 \cdot \max_{d \in D_{selected}} \text{sim}(d, d_{new}) + \lambda_3 \cdot Score_{llm}(d_{new}) \quad (13)$$

其中,  $\mu$  代表聚类中心点,  $d_{new}$  代表待加入的数据点,  $D_{selected}$  表示已经被选中的数据点集合,  $\lambda_1$ ,  $\lambda_2$  和  $\lambda_3$  是三个可设定的权重参数, 用于调整相似度、多样性和模型得分对最终数据点效用值的贡献。

经过两阶段的数据筛选, 每个聚类中的数据将在语义空间中进行重分布, 从而促使模型学习那些罕见但有益的数据, 避免在常见数据集上产生过拟合现象。

## 4 实验

本章中, 我们将利用金融领域的公开数据集, 训练一个中文金融多任务大模型“银瞳”, 利用两个中文金融的评测分析来验证本研究提出的自适应语义空间学习算法的有效性。

### 4.1 数据介绍

我们从金融领域23个不同的来源收集了约22万条中文金融微调数据, 并将这些数据分类为情感分析、金融问答、金融问答等7类任务, 具体介绍见 A.2。同时, 我们使用CFLEB (Lu et al., 2023)和FinEval (Zhang et al., 2023)评测集作为评估工具, 旨在评估大语言模型在金融领域的知识储备、指令跟随及任务执行能力。

CFLEB评测集是利用公开的研报及新闻等项目构建的高质量、高实用性的评测基准, 包含了六种自然语言处理任务, 衡量金融领域模型在情感分析、问答、摘要生成、信息抽取、语义匹配等方面理解与生成的全面能力。

FinEval评测集是针对金融领域知识评估而设计的一套全面数据集, 包含了共4661道选择题。这些题目覆盖了34个不同的学术领域, 如金融学、保险学、宏观经济学和税法等, 主要分为四大类别: 金融、经济、会计以及资格证考试, 以全面测试大型模型对金融领域通用知识的掌握程度, 评估其在金融专业领域内的先进知识与实际应用能力。

### 4.2 实验设置

我们使用了具有代表性的中文主导模型Qwen1.5-7B-Chat (Bai et al., 2023)作为基座模型, 在2张A800显卡上进行实验。对于每一个专家, 都为其训练一个LoRA适配器, 其中 $r$ 设置为16,  $\alpha$ 设计为32。对于所有的微调, 我们设置学习率为 $1e-4$ , 经过10%步的预热, 再利用余弦衰减方式将其减小至 $1e-5$ , 总共经过3轮的训练。在K-Means聚类算法中, 经过我们的多次尝试, 并利用误差平方和 (Sum of Squared Errors, SSE) 和轮廓系数 (Silhouette Coefficient, SC) 对 $K=2$ 至 $K=15$ 的所有数据进行测量, 如Figure 4所示。最终设置 $K=6$ 。同时在自适应的A-DBSCAN算法中, 我们选取了 $k=20$ 作为最初的KNN计算参数。对于修正的MMR公式, 我们将 $\lambda_1$ ,  $\lambda_2$  和  $\lambda_3$  分别设定为0.2, 0.2以及0.6, 以平衡每个数据与聚类中心相似度、多样性以及模型得分对最终效用函数的贡献。

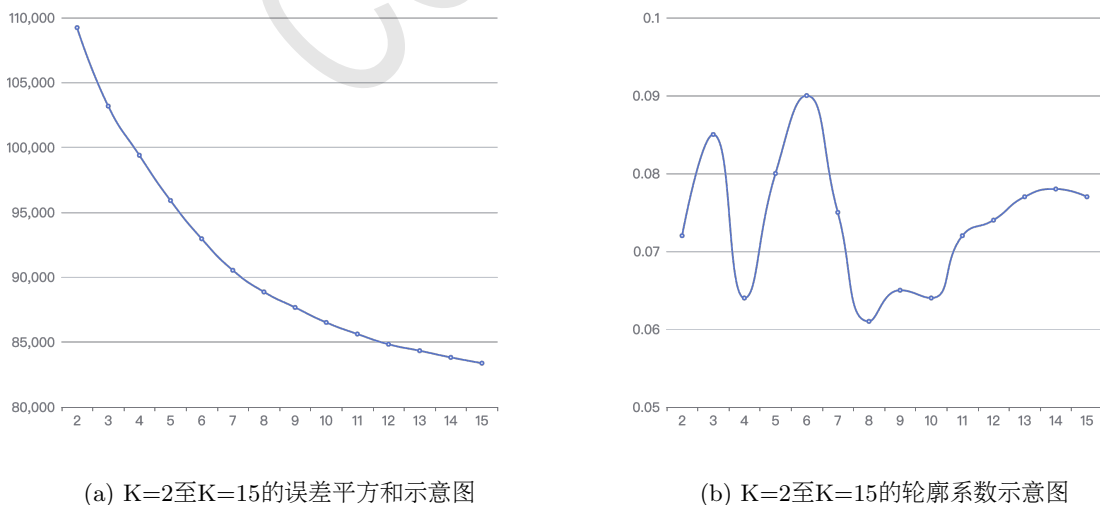


Figure 4: K-Means算法参数K数值挑选图

### 4.3 “银瞳”金融多任务大模型

在本研究中，我们利用设计的ASSL框架对数据进行处理，并训练出了针对中文金融领域的多任务大语言模型“银瞳”。首先我们采用了K-Means算法对公开金融领域的数据进行了语义空间聚类，最终形成六个类别，多种类别个数结果的数据分布见 A.4。接下来，我们对每个聚类中的数据进行平滑分布操作：①第一阶段根据自适应的密度聚类算法A-DBSCAN进行数据筛选，在每个聚类中对高密度区域和低密度区域分别进行下采样和上采样，每个聚类选取近2000条微调数据，将处于极低密度区域的离群数据作为噪声。②第二阶段的重分布筛选旨在通过训练前的原始模型以及第一阶段数据训练后的模型，自适应补充在第一阶段中未被选中的必要数据点，使得微调数据分布更加平滑，并增强数据选择的多样性。在每个聚类中总共筛选出4000条数据，合计约占总数据量的10%。我们利用这六类经过自适应语义分布平滑后的数据，分别训练了六个不同的LoRA专家模型，以适应不同的金融领域任务。当处理特定的金融问题时，通过计算问题在同一个语义空间中的表征与六个LoRA专家表征的相似度，自动选择最适合的LoRA专家进行回答，实现了一个高效、自适应且分工明确的金融多任务大模型系统。

Table 1: CFLEB的评估结果：All-data表示使用我们收集的所有数据训练的单个模型，Cluster-data表示使用所有聚类数据训练的多专家系统，SilverSight-mix表示使用所有过滤数据训练的单一模型。在评估结果中，我们使用10%数据训练的SilverSight模型获得了与All-data相似的结果，而Cluster-data系统获得了最佳性能。

| Model                               | FE<br>(ACC) | QA<br>(F1)  | NA<br>(Rouge) | RE<br>(F1)  | NSP1<br>(ACC) | NSP2<br>(F1) | NL<br>(F1)  | AVG         |
|-------------------------------------|-------------|-------------|---------------|-------------|---------------|--------------|-------------|-------------|
| Qwen-1.5-7B-Chat                    | 65.3        | 31.6        | 30.9          | 10.4        | 90.8          | 9.1          | 33.4        | 38.8        |
| All-data                            | <b>72.3</b> | <b>87.1</b> | <b>53.8</b>   | <b>35.8</b> | 93.2          | 31.1         | <b>90.1</b> | 66.2        |
| Cluster-data                        | 70.9        | 83.4        | 53.1          | 35.3        | <b>93.5</b>   | <b>56.2</b>  | <b>90.1</b> | <b>68.9</b> |
| PIXIU                               | 44.7        | 25.4        | 16.1          | 6.5         | 0             | 26.5         | <b>26.1</b> | 20.8        |
| FinGPT                              | 50.6        | 14.8        | 25.9          | 4.3         | 38.6          | 4.7          | 7.1         | 20.9        |
| DISC-FinLLM                         | <b>69.3</b> | <b>42.4</b> | <b>31.0</b>   | <b>10.1</b> | <b>84.6</b>   | <b>57.9</b>  | 23.4        | <b>45.5</b> |
| <i>SilverSight<sub>LoRA-0</sub></i> | 68.2        | 60.8        | 32.5          | 10.9        | 83.8          | 10.6         | 45.7        | 44.6        |
| <i>SilverSight<sub>LoRA-1</sub></i> | <b>70.2</b> | <b>74.5</b> | 49.5          | <b>31.4</b> | 87.8          | <b>64.8</b>  | 62.9        | <b>63.0</b> |
| <i>SilverSight<sub>LoRA-2</sub></i> | 65.4        | 61.9        | 45.2          | 11.9        | 11.8          | 45.4         | 25.9        | 38.2        |
| <i>SilverSight<sub>LoRA-3</sub></i> | 67.3        | 53          | 37.7          | 12.4        | <b>92</b>     | 9.1          | 43.8        | 45.0        |
| <i>SilverSight<sub>LoRA-4</sub></i> | 57.5        | 72.9        | <b>49.6</b>   | 14.9        | 89.2          | 18           | <b>80.5</b> | 54.7        |
| <i>SilverSight<sub>LoRA-5</sub></i> | 66.6        | 61.8        | 38.4          | 14          | 73.2          | 18.5         | 49.3        | 46.0        |
| SilverSight(our, 10%)               | <b>66.6</b> | 73.3        | <b>49.2</b>   | <b>34.3</b> | <b>93.0</b>   | <b>57.1</b>  | <b>82.4</b> | <b>65.1</b> |
| SilverSight-mix                     | 65.7        | <b>75.4</b> | 48.3          | 30.4        | 89.2          | 31.1         | 82.3        | 60.3        |

### 4.4 主要结果

Table 1与 Table 2展示了我们对“银瞳”大型模型进行的评测实验结果。在针对CFLEB数据集的评测中，我们注意到FinFE、FinQA、FinNA、FinRE的指标与FinNSP1、FinNSP2的指标之间存在一定的矛盾，往往表现出此消彼长的趋势。然而，我们采用的多专家方法有效地缓解了这种偏差，几乎所有指标都超越了采用混合数据训练的SilverSight-mix模型。此外，我们的方法仅使用总量10%的数据，就能达到与全数据微调模型相媲美的测试成绩，在平均分数上只有1%的差距，在FinNSP2任务上的F1指标甚至为全数据微调模型的2倍。

同时，我们发现对每个聚类进行全数据微调的模型，在CFLEB数据集上的表现相较于全数据微调模型有所提升，这充分证明了ASSL框架的有效性。为了检验模型的泛化能力以及对金融领域知识的掌握，我们在未经过训练的FinEval评测基准上进行了测试。结果显示，训练后的模型在数学计算和资格证考试的选择题上的正确率下降了不到2%，但在经济和金融领域的选择题上的正确率提升了4%，FinEval的平均分数也相较于原始模型有所提高。

为验证该LoRA自适应挑选算法的有效性，我们对聚类后的六个LoRA专家逐一在CFLEB数据集和FinEval数据集上进行测试，如 Table 1与 Table 2所示。实验结果显



示，LoRA自适应挑选算法在每个任务上的表现和在该任务上表现最出色的单个LoRA专家相近，证明大多数时候这种LoRA自适应挑选算法均能挑选到最合适的LoRA专家回答问题。同时，Figure 5不仅证明了LoRA自适应挑选算法的有效性，也证明了同时使用聚类方法和LoRA自适应挑选算法能够保证语义空间的一致性与连贯性，最终提升系统的综合性能。

Table 2: FinEval的评估结果

| FinEval                             | Accounting  | Certificate | Economy     | Finance     | AVG         |
|-------------------------------------|-------------|-------------|-------------|-------------|-------------|
| Qwen-7B-Chat                        | 44.5        | 53.6        | 52.1        | 51.5        | 50.5        |
| ChatGPT                             | 45.2        | 55.1        | 61.6        | 59.3        | 55.0        |
| Qwen-1.5-7B-Chat                    | <b>69.5</b> | <b>71.3</b> | 62.8        | 65.6        | 67.8        |
| GPT-4                               | 59.3        | 70.4        | <b>74.5</b> | <b>71.0</b> | <b>68.6</b> |
| PIXIU                               | 36.1        | 34.7        | 37.2        | 34.1        | 35.5        |
| FinGPT                              | <b>50.8</b> | <b>52.1</b> | 42.5        | 50.5        | <b>49.6</b> |
| DISC-FinLLM                         | 42.0        | 49.1        | <b>47.9</b> | <b>63.8</b> | 48.7        |
| <i>SilverSight<sub>LoRA-0</sub></i> | <b>68.9</b> | <b>72.2</b> | 64.7        | 65.9        | <b>68.3</b> |
| <i>SilverSight<sub>LoRA-1</sub></i> | 66.9        | 70.7        | 59.4        | 62          | 65.3        |
| <i>SilverSight<sub>LoRA-2</sub></i> | 62          | 65.6        | 58.5        | 62.6        | 62.6        |
| <i>SilverSight<sub>LoRA-3</sub></i> | 67.5        | 69.8        | <b>67.1</b> | 67.2        | 68          |
| <i>SilverSight<sub>LoRA-4</sub></i> | 68.5        | 70          | 63.8        | 66.2        | 67.5        |
| <i>SilverSight<sub>LoRA-5</sub></i> | 68.2        | 71.3        | 63.3        | 65.6        | 67.5        |
| SilverSight(our, 10%)               | 67.9        | 70          | 66.7        | <b>67.9</b> | <b>68.3</b> |

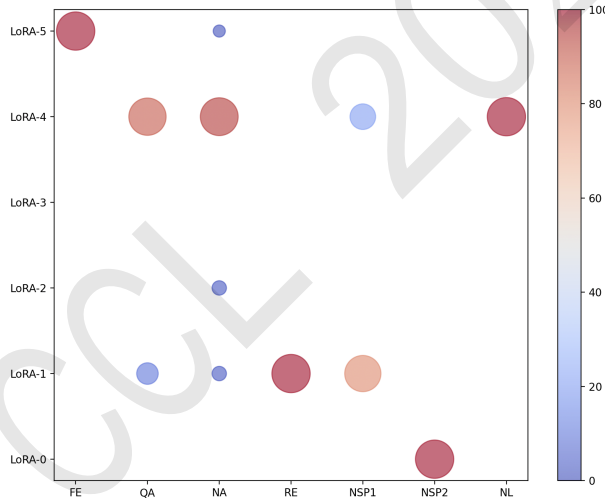


Figure 5: CFLEB数据集每类任务对专家的选择情况

#### 4.5 消融实验

为了证明ASSL框架中算法的有效性，我们对以下三个问题进行探究：

1. 问题1：在语义空间进行聚类学习相较于人工任务分类是否存在优势？
2. 问题2：进行第一阶段数据重分布是否具有意义？
3. 问题3：进行第二阶段的数据筛选补充是否具有意义？

为了回答问题1，我们设计了一组对照实验，两个模型系统均采用LoRA专家自适应选择的方式。对于第一个模型系统，在数据预处理阶段，根据自然语言处理任务定义，人工将数据分为七类，每一类数据训练一个LoRA专家，组成多专家系统。同时，对于第二个模型系统，我

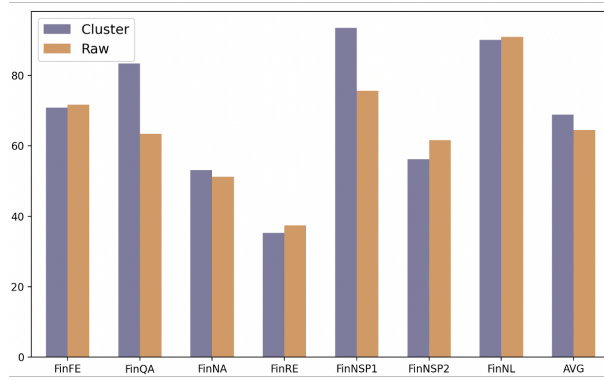


Figure 6: CFLEB数据集对于聚类 and 预定义模型系统的评测结果

们在语义空间中将数据重新聚类为六类，旨在验证语义空间中的距离能反映任务间的互补程度，相较于人工预定义类别具有优势。为此，我们对采用聚类算法得到的六类数据和原始的七类数据进行了详尽的对比分析，实验结果如 Figure 6 所示，预定义人工分类的方法在CFLEB上平均表现不如语义空间聚类的方法，在FinQA和FinNSP1任务上显著低于语义空间聚类方法。由此证明了语义空间聚类学习的方法相较于预定义任务类型方法的优越性，证明了利用语义空间相似性的方法能够聚拢互补任务的多种数据，起到促进模型在这些任务上的性能的作用。

为了回答问题2，我们将两种数据筛选算法作为基线：随机筛选以及K-Center贪婪算法，对ASSL框架中第一部分数据重分布的A-DBSCAN算法进行对比测试。三种数据筛选算法的侧重点不同，Random算法旨在实现从各聚类中按照原始数据分布情况均匀地抽取数据，K-Center贪婪算法注重选择数据的多样性，A-DBSCAN算法则对原始数据进行了重分布，根据密度的稀疏性对数据进行采样，同时也过滤了大量的离群噪声数据。实验结果如 Figure 7所示，在CFLEB数据集上，A-DBSCAN和Random算法表现突出；而在FinEval数据集上，K-Center贪婪和A-DBSCAN算法则展现了较好的性能。在两个数据集上都证实了A-DBSCAN算法的鲁棒性，表明这种通过对语义空间密度进行平滑分布方法的优越性，但由于其过滤了大量噪声，还需进一步进行数据补充。

为了回答问题3，并进一步确认增强数据平滑分布的重要性，我们将利用效用函数筛选数据与Random算法筛选数据的实验结果进行对比，为保证公平性，两个实验类型均在每个聚类中选择4000条数据进行训练。实验结果如 Figure 8 所示，通过第二阶段的数据扩充，模型的综合性能相比于仅经过第一阶段筛选有了显著的提升，同时优于Random算法，这进一步验证了第二阶段进行数据扩充的必要性和有效性，增强了对于原始数据长尾分布的平滑性。

## 5 讨论

我们的研究基于设计ASSL框架训练了一个中文金融领域的多任务大模型，并在多个评测集上取得优异的性能。通过ASSL框架，研究人员可以深入了解基于语义空间的自适应学习优势，利用语义空间的特性，将多专家系统的优势注入语义空间中，实现专家自适应选择、数据自适应选择，从而得到更好的模型性能。基于此，我们可以进一步讨论这种基于语义空间自适应学习的互补特性。

每个领域都拥有众多不同的自然语言处理或领域特定任务，要求领域大模型具备全面而综合的能力，但目前的多源异构数据容易引发模型进行任务风格迁移时的冲突。进行多任务学习时，利用所有语料对一个大模型进行训练极易影响模型指令遵循的能力，进而引发性能下降、幻觉产生等现象。基于此，制定一个多专家系统，并根据所有数据在同一个语义空间中的分布情况，为模型制定自适应的任务分类、数据筛选的方案变得尤为重要。根据 Figure 5，利用语义空间数据点的相似性能够有效的聚拢相互促进模型能力提升的任务数据，分隔互斥的任务数据，提高微调数据的均衡性。利用语义空间相似性不仅优化了数据分类的过程，也提供了一种高效的LoRA专家模型的挑选机制。通过比较问题嵌入与各个LoRA专家模型嵌入之间的相似性，系统能够精准地选取与特定问题最为匹配的LoRA专家模型。这种方法能够确保数据、任务与模型在语义空间上的连贯性与一致性，最终提高模型的综合能力。实验结果进一步验证了该方法的有效性和鲁棒性。相较于随机选择数据进行模型微调，采用数据筛选算法能够更加精

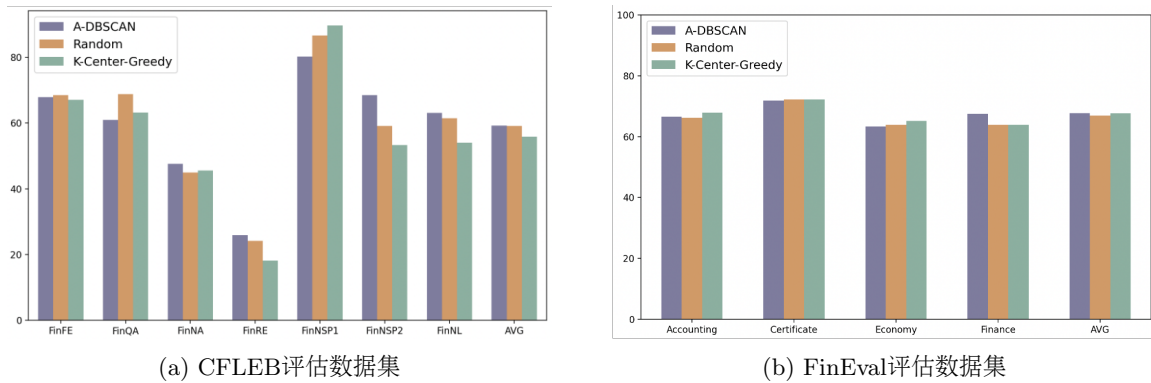


Figure 7: 消融实验2:第一阶段数据筛选算法对比

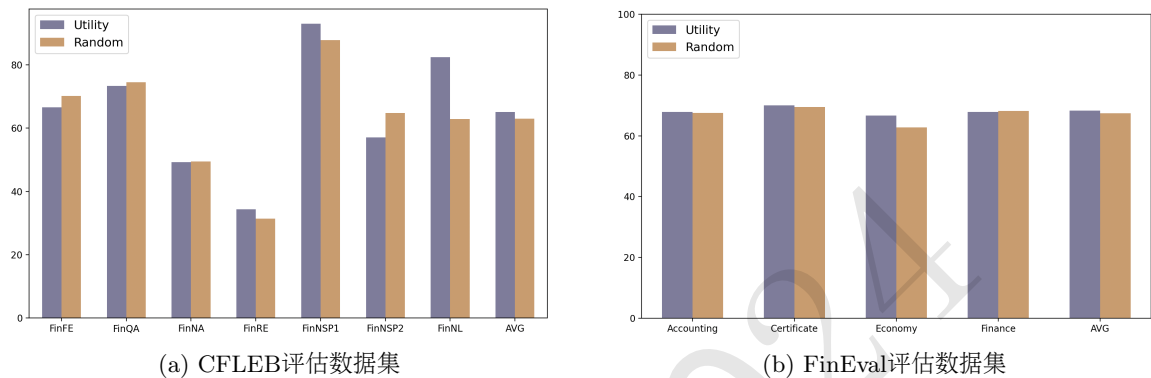


Figure 8: 消融实验3:第二阶段数据筛选算法对比

确地匹配合适的LoRA专家模型，显著提高了模型处理不同领域问题时的整体表现。

## 6 结论

本研究中我们主要探索了基于语义空间的自适应学习对于多专家大语言模型系统的性能影响，目标是根据语义空间对互补以及互斥的任务数据进行划分，利用其中的嵌入位置信息对多专家进行自适应的选择，并利用模型本身以及语义空间的密度分布对训练数据分布进行两阶段调整，该框架使得模型在进行多任务学习时具有较好的性能及泛化性。我们基于自适应语义空间学习框架，利用金融领域的公开数据集训练了“银瞳”大语言模型，并在基准上对“银瞳”模型系统进行了评测，显示出了其优异的表现。这一结果不仅证明了我们方法的可行性和有效性，也为未来多专家系统的大语言模型发展开辟了新的视野。

## 致谢

本研究得到了国家重点研发计划（2023YFC3304800）、中国工程院战略研究咨询项目“数字化转型背景下的金融风险监测与预警系统战略研究”（2023-XY-43）和上海市自然科学基金（23ZR1404900）的支持。

## 参考文献

- Janze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Jaime Carbonell and Jade Goldstein. 1998. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 335–336.
- Wei Chen, Qiushi Wang, Zefei Long, Xianyin Zhang, Zhongtian Lu, Bingxuan Li, Siyuan Wang, Jiarong Xu, Xiang Bai, Xuanjing Huang, et al. 2023. Disc-finllm: A chinese financial large language model based on multiple experts fine-tuning. *arXiv preprint arXiv:2310.15205*.
- Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. 2023. Chatlaw: Open-source legal large language model with integrated external knowledge bases. *arXiv preprint arXiv:2306.16092*.
- Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Jun Zhao, Wei Shen, Yuhao Zhou, Zhiheng Xi, Xiao Wang, Xiaoran Fan, et al. 2023. Lora-moe: Revolutionizing mixture of experts for maintaining world knowledge in language model alignment. *arXiv preprint arXiv:2312.09979*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Michael I Jordan and Robert A Jacobs. 1994. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 6(2):181–214.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Dakuan Lu, Hengkui Wu, Jiaqing Liang, Yipei Xu, Qianyu He, Yipeng Geng, Mengkun Han, Yingsi Xin, and Yanghua Xiao. 2023. Bbt-fin: Comprehensive construction of chinese financial domain pre-trained language model, corpus and benchmark. *arXiv preprint arXiv:2302.09432*.
- Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2023. An empirical study of catastrophic forgetting in large language models during continual fine-tuning. *arXiv preprint arXiv:2308.08747*.
- James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA.
- OpenAI. 2023. Gpt-4 technical report.
- Leif E Peterson. 2009. K-nearest neighbor. *Scholarpedia*, 4(2):1883.
- Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. 2020. Pre-trained models for natural language processing: A survey. *Science China Technological Sciences*, 63(10):1872–1897.
- Mengjie Ren, Boxi Cao, Hongyu Lin, Cao Liu, Xianpei Han, Ke Zeng, Guanglu Wan, Xunliang Cai, and Le Sun. 2024. Learning or self-aligning? rethinking instruction fine-tuning.
- Erich Schubert, Jörg Sander, Martin Ester, Hans Peter Kriegel, and Xiaowei Xu. 2017. Dbscan revisited, revisited: why and how you should (still) use dbscan. *ACM Transactions on Database Systems (TODS)*, 42(3):1–21.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Ted Zadouri, Ahmet Üstün, Arash Ahmadian, Beyza Ermiş, Acyr Locatelli, and Sara Hooker. 2023. Pushing mixture of experts to the limit: Extremely parameter efficient moe for instruction tuning. *arXiv preprint arXiv:2309.05444*.

Liwen Zhang, Weige Cai, Zhaowei Liu, Zhi Yang, Wei Dai, Yujie Liao, Qianru Qin, Yifei Li, Xingyu Liu, Zhiqiang Liu, et al. 2023. Fineval: A chinese financial domain knowledge evaluation benchmark for large language models. *arXiv preprint arXiv:2308.09975*.

Ziyu Zhao, Leilei Gan, Guoyin Wang, Wangchunshu Zhou, Hongxia Yang, Kun Kuang, and Fei Wu. 2024. Loraretriever: Input-aware lora retrieval and composition for mixed tasks in the wild. *arXiv preprint arXiv:2402.09997*.

Yuhang Zhou, Yuchen Ni, Xiang Liu, Jian Zhang, Sen Liu, Guangnan Ye, and Hongfeng Chai. 2024. Are large language models rational investors? *arXiv preprint arXiv:2402.12713*.

## A 附录

## A.1 A-DBSCAN

---

**Algorithm 1** 简化且包含子簇处理的基于密度自适应A-DBSCAN算法
 

---

**Data:** 数据集 $D$ , 邻域大小 $k$ **Result:** 形成的簇集合 $Clusters$ 初始化簇集合 $Clusters = \emptyset$  计算每个数据点 $p$ 的局部密度 $\rho_p$ :

$$\rho_p = \frac{1}{\frac{1}{k} \sum_{j=1}^k d(p, p_j)}$$

根据 $\rho$ 对 $D$ 中的点进行排序, 形成优先队列 $Q$  初始化 $\varepsilon$ 和 $MinPts_{init}$ :

$$\varepsilon = \text{中位数}\{d(p, p_k)\} \quad \forall p \in D$$

$$MinPts_{init} = \frac{\varepsilon \times \rho_{max}}{2}$$

**while** 从 $Q$ 中选取队列首个点 $p$  **do**  **if** 未访问过 $p$  **then**    标记 $p$ 为已访问 获取 $p$ 的 $\varepsilon$ -邻域中的点 $N_p$  **if**  $|N_p| \geq MinPts_{init}$  **then**      初始化一个新簇 $cluster$ 并添加 $p$  **for** 每个点 $q \in N_p$  **do**        **if**  $q$ 未被访问 **then**          标记 $q$ 为已访问 获取 $q$ 的 $\varepsilon$ -邻域中的点 $N_q$  **if**  $|N_q| \geq MinPts_{init}$  **then**          | 将 $N_q$ 中的点添加到 $N_p$ 中          **end**        **end**        **if**  $q$ 不属于任何簇 **then**          | 将 $q$ 添加到 $cluster$         **end**      **end**       $Clusters = Clusters \cup \{cluster\}$  更新 $MinPts$ :

$$MinPts = \max\left(2, \frac{\rho_p}{\rho_{max}} \times MinPts_{init}\right)$$

**end**  **end****end**处理子簇( $Clusters$ ) **return**  $Clusters$ **Procedure** 处理子簇 $Clusters$ :  计算平均量 $N_{avg} = \frac{\sum_i N_{C_i}}{|Clusters|}$ , 其中 $N_{C_i}$ 是子簇 $C_i$ 的数据量  **foreach** 子簇 $C_i \in Clusters$  **do**    **if**  $N_{C_i} > N_{avg}$  **then**      | 采样: 从 $C_i$ 中随机选择 $N_{avg}$ 个数据点    **end****return**

## A.2 金融NLP任务

在这一部分, 我们将会对用于构造指令的金融NLP数据集进行更为详细的介绍。数据集涵盖了7类任务, 包括情感分析、信息抽取、文本分类、文本生成、语义匹配、金融问答和金融考

试选择题，具体统计数据如 Table 3所示。

Table 3: 数据集统计

| 数据集                 | 主要任务类型  | 细分任务类型       | 数量 (条)      |
|---------------------|---------|--------------|-------------|
| flare-zh-fe         | 情感分析    | 情感分析         | 2020        |
| flare-zh-stocka     | 情感分析    | 情感分析         | 1477        |
| flare-zh-stockb     | 情感分析    | 情感分析         | 1962        |
| CFSC-ABSA           | 情感分析    | 情感分析         | 33184       |
| ficuge-finnl        | 文本分类    | 金融新闻分类       | 7071        |
| flare-zh-nl         | 文本分类    | 金融新闻分类       | 884         |
| flare-zh-nl2        | 文本分类    | 金融新闻分类       | 884         |
| flare-zh-nsp        | 文本分类    | 金融负面消息及其主体判定 | 500         |
| fincuge-finna       | 文本生成    | 新闻摘要         | 28799       |
| flare-zh-na         | 文本生成    | 新闻摘要         | 3600        |
| flare-zh-19ccks     | 信息抽取    | 事件类型与实体识别    | 2936        |
| flare-zh-20ccks     | 信息抽取    | 事件类型与实体识别    | 9159        |
| flare-zh-21ccks     | 信息抽取    | 事件类型与实体识别    | 1400        |
| flare-zh-22ccks     | 信息抽取    | 事件类型与实体识别    | 11829       |
| fincuge-finre       | 信息抽取    | 关系抽取         | 13486       |
| flare-zh-ner        | 信息抽取    | 命名实体识别       | 337         |
| flare-zh-corpus     | 语义匹配    | 语义匹配         | 10000       |
| flare-zh-afqmc      | 语义匹配    | 语义匹配         | 4316        |
| fincuge-finqa       | 金融问答    | 问答           | 19906       |
| fincuge-fincqa      | 金融问答    | 因果问答         | 21965       |
| Duxiaoman/FinanceIQ | 金融考试选择题 | 单项选择题        | 3573        |
| fingpt-fineval      | 金融考试选择题 | 单项选择题        | 1056        |
| Duxiaoman/FinCorpus | 金融考试选择题 | 多项选择题        | 40000 (过滤后) |

**情感分析** 情感分析能够反映金融新闻的情感趋向，影响决策。我们采集了4个数据集，包括flare-zh-fe、flare-zh-stocka、flare-zh-stockb和CFSC-ABSA。flare-zh-stocka分析市场数据和公司相关的公告，根据新闻对股票数据的影响判断公司的股票运动走势是跑赢大盘、中性还是表现不佳。其余数据集是判断所给的金融新闻中句子的情绪，给出积极、消极或中性的回答。

**文本分类** 对于金融新闻分类，选取ficuge-finnl、flare-zh-nl、flare-zh-nl2这3个数据集，根据所给金融文本，输出2-3个类别关键词。对于金融负面消息及其主体判定，flare-zh-nsp数据集根据金融新闻和实体，判断该实体是否含有负面消息。

**文本生成** 文本生成任务主要是根据金融新闻长文本生成摘要，要求简洁但能包含关键信息，使用的数据集有fincuge-finna和flare-zh-na。

**信息抽取** 信息抽取任务主要包括事件类型与实体识别、关系抽取和命名实体识别，使用的数据集有flare-zh-19ccks、flare-zh-20ccks、flare-zh-21ccks、flare-zh-22ccks、fincuge-finre和flare-zh-ner。其中前四个数据集属于事件类型与实体识别类型，分析给定金融文本，判断其中所有的事件类型及其对应主体；关系抽取数据集fincuge-finre根据金融文本及实体对，从44个关系中选择出能正确表示该实体对的关系；命名实体识别数据集flare-zh-ner是对金融文本中存在实体的抽取。

**语义匹配** 对于语义匹配任务，我们选取了两个数据集：flare-zh-corpus和flare-zh-afqmc，判断两个金融文本表达的语义是否一致。

**金融问答** 针对普通的金融问答，使用fincuge-finqa数据集，从金融文本中提取事件信息，并结合上下文回答问题。而金融因果数据集fincuge-fincqa还关注金融文本中的因果关系。

**金融考试选择题** 包括来自Duxiaoman和FinGPT的金融各领域的单项选择题，以及Duxiaoman金融预料中的多项选择题，涵盖了银行、证券、会计、基金、期货等多领域的知识和计算。

### A.3 金融NLP任务示例

将所有训练数据统一为三元组{instruction, input, output}格式，具体各类数据集的示例如图 Figure 9 - Figure 13所示。

```
{
  "instruction": "针对金融股民论坛的这段话，请评估它的情绪倾向：是积极、消极，还是处于中立态度？欧股涨幅扩大，德国DAX指数、欧洲斯托克50指数均涨超1%，法国CAC40指数涨0.82%，英国富时100指数涨0.71%。\\n\\n答案中只需包含\\\"积极\\\"、\\\"中性\\\"或\\\"消极\\\"",
  "input": "欧股涨幅扩大，德国DAX指数、欧洲斯托克50指数均涨超1%，法国CAC40指数涨0.82%，英国富时100指数涨0.71%。\\n",
  "output": "积极"
}
```

Figure 9: 情感分析示例

```
{
  "instruction": "阅读以下金融声明，并根据其内容，从中国、外国、国际、公司、行业、大盘、经济、政策、政治、期货、债券、房地产、外汇、虚拟货币、新冠、能源中挑选两至三个最相关的类别，答案用空格隔开。\\n国家外汇管理局：加强风险防控和监管能力建设方面，完善外汇市场“宏观审慎+微观监管”两位一体管理框架，加强跨境资金流动风险监测预警和逆周期调节，严厉打击外汇领域违法违规活动，维护外汇市场稳定和健康秩序。\\n\\n",
  "input": "国家外汇管理局：加强风险防控和监管能力建设方面，完善外汇市场“宏观审慎+微观监管”两位一体管理框架，加强跨境资金流动风险监测预警和逆周期调节，严厉打击外汇领域违法违规活动，维护外汇市场稳定和健康秩序。",
  "output": "中国政策"
}
```

Figure 10: 文本分类示例

```
{
  "instruction": "\\n 央行金融消费者权益保护局副局长尹优平日前表示，人民银行将统筹推进违法金融广告和营销宣传行为治理，科学运用各种技术手段，有效抓取非法金融广告线索；加强金融矛盾源头预防、前端化解，为人民群众提供低成本、高效率、一站式、法制化的金融纠纷化解渠道；提高金融监管的灵敏度、精准度和及时性，推动治理能力现代化建设。\\n\\n阅读后，请简述所给资讯的核心要点，摘要长度限制在一到两句话之内。",
  "input": "央行金融消费者权益保护局副局长尹优平日前表示，人民银行将统筹推进违法金融广告和营销宣传行为治理，科学运用各种技术手段，有效抓取非法金融广告线索；加强金融矛盾源头预防、前端化解，为人民群众提供低成本、高效率、一站式、法制化的金融纠纷化解渠道；提高金融监管的灵敏度、精准度和及时性，推动治理能力现代化建设。",
  "output": "央行：将统筹推进违法金融广告和营销宣传行为治理"
}
```

Figure 11: 文本生成示例

```
{
  "instruction": "\\n此后外资对国内OTC重磅品种的并购沉寂数年,直到日前拜耳宣布全资收购滇虹药业,OTC并购战火才重新燃起。头实体：拜耳，尾实体：滇虹药业。\\n\\n在财经金融相关的文本和问题中，只从提供的选项中挑选出一个正确答案，直接给出这个答案，不包含其他任何信息。你可以选择的关系类别包括：分析、借壳、合作、转让、买资、入股、商讨、被拟收购、被成立、重组、自己、被注资、被分析、被帮助、交易、被持股、被入股、注资、成立、被买资、被借壳、增持、被拥有、发行、订单、拥有、纠纷、被增持、被转让、合资、减持、欠款、其他、被减持、签约、拟收购、被收购、合并、帮助、被发行、被欠款、持股、收购、竞争\\n请注意，如果头尾实体不含以上关系，请输出其他。",
  "input": "此后外资对国内OTC重磅品种的并购沉寂数年,直到日前拜耳宣布全资收购滇虹药业,OTC并购战火才重新燃起。头实体：拜耳，尾实体：滇虹药业。",
  "output": "收购"
}
```

Figure 12: 信息抽取示例

### A.4 数据混合比例分布



```
{
  "instruction": "分析两个金融语句的意思是否一致，你只需要回答是或否。请分析下面所给语句语义是否一致：\n",
  "input": "使用花呗人脸识别不了怎么办\n花呗要人脸识别 识别不了怎么回事",
  "output": "是"
}
```

Figure 13: 语义匹配示例

```
{
  "instruction": "\n\n上汽集团股份回购事件对应的回购股份数量是什么？上汽集团公告，公司计划回购5842万股-1.17亿股股份，占公司总股本的0.5%-1%，拟回购资金总额为不超过人民币30.34亿元，回购实施期限自股东大会审议通过本次回购方案之日起6个月内。针对本段上下文和提出的问题，提取信息并作答。如不能回答，请回答无相应参数。",
  "input": "上汽集团股份回购事件对应的回购股份数量是什么？上汽集团公告，公司计划回购5842万股-1.17亿股股份，占公司总股本的0.5%-1%，拟回购资金总额为不超过人民币30.34亿元，回购实施期限自股东大会审议通过本次回购方案之日起6个月内。",
  "output": "5842万股-1.17亿"
}
```

Figure 14: 金融问答示例

```
{
  "instruction": "请根据上下文题目，不论是单选题还是多选题，列出所有正确答案并进行分析。在个人/家庭资产负债表中，下列属于资产项目的有（）。I.现金与现金等价物II.股票、债券等金融资产III.车辆、房产等实物资产IV.信用卡透支\nA、I、II、IV\nB、I、II、III\nC、II、III\nD、I、III、IV",
  "input": "在个人/家庭资产负债表中，下列属于资产项目的有（）。I.现金与现金等价物II.股票、债券等金融资产III.车辆、房产等实物资产IV.信用卡透支\nA、I、II、IV\nB、I、II、III\nC、II、III\nD、I、III、IV",
  "output": "B\n解析：在个人/家庭资产负债表中，资产项目包括现金与现金等价物、股票、债券等金融资产、车辆、房产等实物资产。信用卡透支属于负债，因此不在资产项目中。选项B包含了正确的资产项目。"
}
```

Figure 15: 金融考试选择题示例

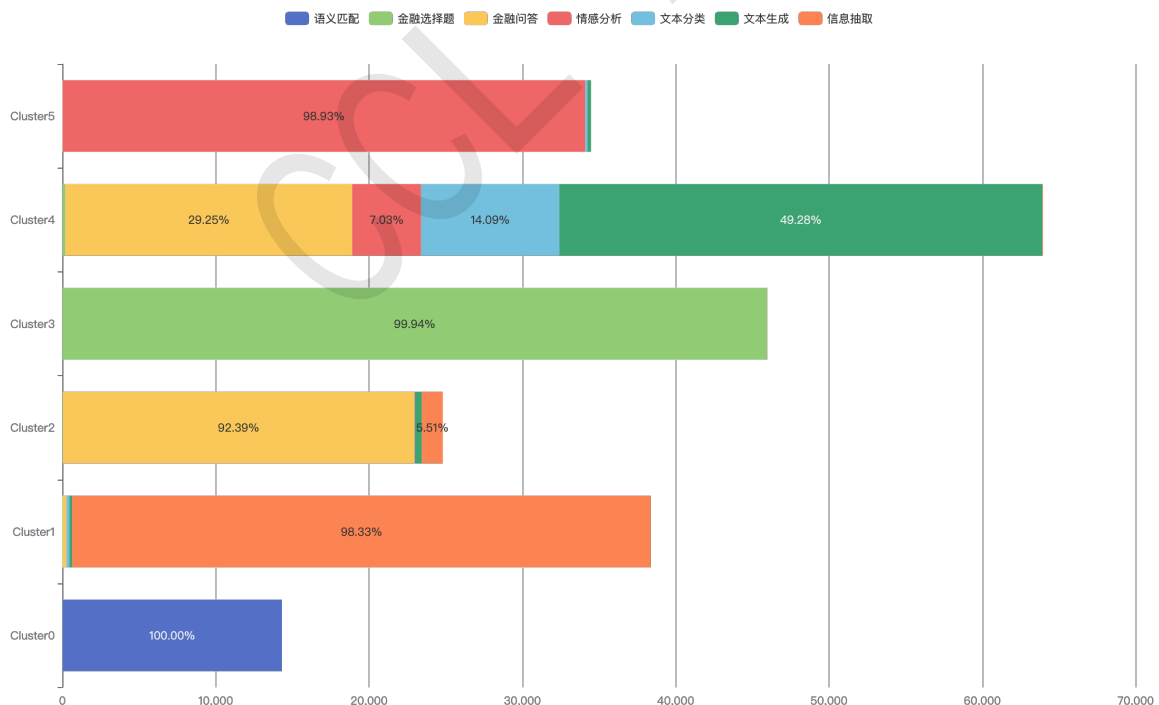


Figure 16: 经过聚类后每个类别的数据混合比例

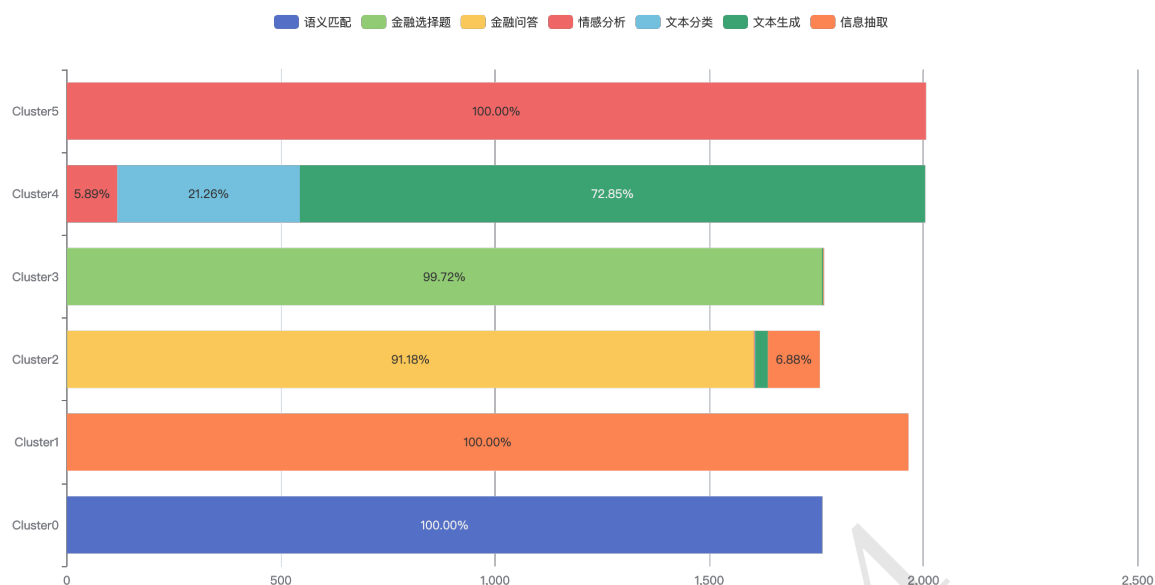


Figure 17: 经过一阶段数据重分布后，每个类别的数据混合比例

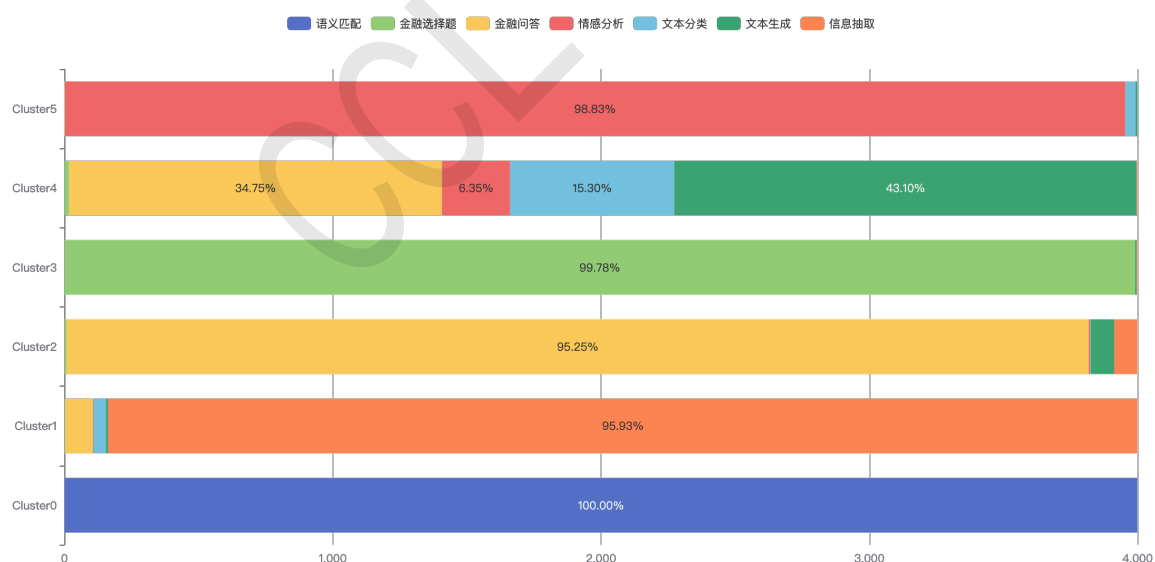


Figure 18: 经过二阶段数据重分布后，每个类别的数据混合比例