

深度学习报告

July 5, 2025

摘要

本项目旨在介绍商业智能选址预测的深度学习方法。通过分析品牌历史门店的地理分布数据，建立能够预测品牌下一家门店最优选址的模型。项目以网格（Grid）为基本地理单元，利用品牌在城市中的历史开店网格序列及每个网格的地理属性特征，挖掘品牌扩张的空间模式和内在偏好。通过对历史数据的建模与学习，预测品牌未来最有可能选址的网格位置。模型由序列编码器、空间编码器和环境编码器组成，采用多模态融合设计。实验结果表明，该模型在选址预测任务中表现优异，能够为商业决策提供科学依据。

关键词

选址预测, 深度学习, 网格模型, 多模态融合, 序列编码器, 空间编码器, 环境编码器

目录

1	引言	3
2	相关工作	3
3	方法探索	3
3.1	数据预处理	3
3.1.1	数据加载	4
3.1.2	网格信息加载	4
3.1.3	特征归一化	4
3.1.4	基于密度的序列排序	4
3.1.5	滑动窗口样本生成	5
3.2	对于标签的处理尝试	5
3.3	数据增强和特征增强	5
3.4	多模态特征融合设计	5
3.4.1	序列特征编码	5
3.4.2	空间特征编码	5
3.4.3	环境特征编码	5
3.4.4	特征融合策略	5
3.5	模型最终设计	5
A	附录	6
A.1	项目信息	6
A.2	核心模型代码	6
A.2.1	神经网络模型实现 (model.py)	6
A.2.2	主程序入口 (main.py)	9
A.2.3	训练逻辑实现 (train.py)	10
A.2.4	数据预处理模块 (data_preprocessing.py)	11
A.2.5	数据特征增强 (测试数据文件.py)	14
A.2.6	模型评估模块 (evaluate.py)	15
A.3	数据样本展示	16
A.3.1	训练数据格式	16
A.3.2	网格坐标映射	17
A.3.3	测试数据格式	17
A.4	完整源代码文件结构	17

1 引言

在现代商业环境中，选址决策是品牌扩张和市场布局过程中至关重要的一环。一个优质的门店位置不仅能够提升品牌曝光度和客流量，还直接影响企业的经济效益和市场竞争能力。传统的选址方法主要依赖于专家经验和简单的人口统计分析，存在主观性强、难以量化、适应性差等局限。随着大数据、机器学习和深度学习技术的快速发展，基于数据驱动的选址预测方法为选址预测提供更多的思路启发，其能够通过对历史数据的深入分析，挖掘出潜在的空间模式和内在偏好，从而实现更科学、精准、有效的选址决策。

本课题聚焦于“商业智能选址预测”，旨在通过分析品牌历史门店的地理分布数据，建立能够预测品牌下一家门店最优选址的模型。具体而言，课题以网格（Grid）为基本地理单元，利用品牌在城市中的历史开店网格序列及每个网格的地理属性特征，挖掘品牌扩张的空间模式和内在偏好。通过对历史数据的建模与学习，预测品牌未来最有可能选址的网格位置，为企业提供科学、量化的决策依据。

本项目的数据集包含多个品牌的历史门店分布（训练集和测试集），以及覆盖研究区域的网格地理坐标信息。提供数据已划分为训练集和测试集，网格划分保证了空间分析的精度和一致性。研究区域的经纬度范围明确，便于空间特征的提取和建模。

2 相关工作

项目提供了多种基线（Baseline）模型或方法，包括简单模型和深度模型。在项目文档中，对比了多个基线方案的实验结果，展示了不同方法在选址预测任务中的表现，作为模型实现方案的参考和对比。

随机猜测（Random guess）方法，将所有未出现过的网格作为候选集，随机选取若干作为预测结果，作为最简单的基线 [1]。一般认为，一个模型至少有一点作用，就可以将它和随机猜测进行对比，如果这个模型居然结果还不如随机猜测，那么这个模型就没有意义了，一定存在很明显的缺陷，比如严重的过拟合、特征选择不当、数据标签出现错误、数据分布出现重大偏差等 [2]。

随着空间数据挖掘的发展，基于 Word2vec 的嵌入方法被用到选址任务中 [3]。这种方法，虽然最初被用于自然语言处理，但是其本质是一种时序性的序列预测方法 [4]，因此也可以被利用在这个任务中。该方法通过 skip-gram 模型训练每个网格的 embedding，测试时将品牌历史网格序列的 embedding 取平均，计算与候选网格 embedding 的相似度（如点积、余弦相似度等），并据此排序预测结果。

深度学习方法方面，LSTM+MLP 结构利用嵌入层将输入网格序列编码为向量，随后通过 LSTM 捕捉序列中的时序依赖，最后通过多层感知机（MLP）进行分类预测。Transformer+MLP 则用 Transformer 替代 LSTM，利用自注意力机制更好地建模序列中各网格之间的复杂关系。进一步地，基于注意力融合和对比学习的方法，通过对输入序列进行 self-attention，获得全局语义表达，并采用 InfoNCE 损失进行正负样本的对比训练，有效提升了模型的判别能力和泛化性能。

上述方法在本项目数据集上的实验结果如表所示。可以看出，深度学习模型（如 Transformer+MLP、对比学习方法）在准确率和平均排序等指标上均优于传统方法，尤其是对比学习方法在 acc@1 和 acc@10 等指标上表现突出，显示了其在空间选址预测任务中的潜力。此外，模型结构的层级设计（如嵌入层、序列建模层、融合层和输出层）对最终性能有显著影响，合理的结构选择和层级组合能够更好地捕捉空间和序列特征，从而提升预测效果。

3 方法探索

3.1 数据预处理

数据集中有已经分配好的训练集和测试集。数据的内容主要是品牌历史门店的选址信息。具体分析数据结构如表1所示。

字段	类型	描述
brand_name	string	品牌名称
brand_type	string(csv)	品牌类型（如餐饮、零售等）
longitude_list	list(float)	网格经度列表
latitude_list	list(float)	网格纬度列表
grid_id_list	list(int)	网格 ID 列表

表 1: 数据集字段说明

此外还通过了网格数据，给出对应网格 ID 在真实地理世界中的经纬度坐标。网格数据的结构如表2所示。

字段	类型	描述
grid_id	int	网格 ID
grind_lon_min	float	网格左下角经度
grind_lon_max	float	网格右上角经度
grind_lat_min	float	网格左下角纬度
grind_lat_max	float	网格右上角纬度

表 2: 网格数据字段说明

3.1.1 数据加载

数据加载是数据预处理的第一步。我们需要从提供的训练集和测试集文件中读取品牌历史门店的选址信息，并将其转换为适合模型输入的格式，提交给神经网络。

3.1.2 网格信息加载

“坐标”含有的信息过于的稀少，难以直接被模型“理解”，从而掌握其中内在的一些联系。需要对数据增加有关的信息，以扩展其表达能力，否则模型表达能力的上限就会受到限制 [5]。为了得出更加合理的结论，必须在数据预处理阶段强化并且增加特征。

网格（Grid）是商业选址预测任务中的基本地理单元。在原始数据集里，每个网格由其唯一的网格 ID 和对应的地理坐标（经纬度）组成。这个体现了网格在空间上的位置信息。但是这个信息并不能很有效地描述网格的全部特性。为了增强网格的信息，需要从数据集中重新挖掘出可能存在的信息。兴趣点（POI, Point of Interest）是地理信息系统（GIS）和位置服务中非常核心的概念，它指的是地图上具有特定意义或吸引力的地点。它代表一个具有空间位置和语义信息的地理目标，比如餐厅、医院、公交站、商场、景点等，在经纬度以外，其还通常包括有名称、地址、类别等信息，并经常以向量的形式表示。POI 可以将三维世界中的复杂实体抽象为一个零维点，从而便于管理、分析和计算 [6]。

因此，可以从全部的门店数据中，统计每个网格的相关信息，再将其关联到网格 ID 上，构成更加立体、具有丰富特征的高维网格对象。具体而言，再项目实现上，从网格坐标文件中提取每个网格的地理信息，包括经纬度坐标和 10 种类型的 POI 特征（医疗、住宿、摩托、体育、餐饮、公司、购物、生活、科教、汽车）。这种类型的统计信息，虽然引入了一些和其他若干特征高度相关的特征，但是会在大型复杂模型的训练上甚至有性能提升的作用 [7]。

3.1.3 特征归一化

为了消除不同特征间的尺度差异，我们对坐标和 POI 特征分别进行归一化处理：对于坐标特征，我们计算网格的中心点坐标，然后将其归一化到 [0,1] 区间：

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min} + \epsilon} \quad (1)$$

对于 POI 特征，我们同样采用归一化方法，确保不同类型 POI 特征的数值在同一尺度上。

3.1.4 基于密度的序列排序

考虑到商业扩张通常遵循“从高密度区域向外扩展”的规律，我们设计了基于密度的序列排序算法。该算法通过计算每个网格到其 K 个最近邻网格的平均距离来衡量密度，距离越小表示密度越大：

$$density_score_i = \frac{1}{K} \sum_{j=1}^K d(grid_i, neighbor_j) \quad (2)$$

其中 $d(grid_i, neighbor_j)$ 表示网格 i 到其第 j 个最近邻的欧氏距离。通过这种排序方式，我们能够构建符合商业扩张规律的序列，提高模型的学习效果。

3.1.5 滑动窗口样本生成

为了最大化利用有限的训练数据，我们采用滑动窗口方法从每个品牌的店铺序列中生成多个训练样本。例如，对于序列 [A,B,C,D]，我们生成以下样本：

- 已有 [A]，预测 B
- 已有 [A,B]，预测 C
- 已有 [A,B,C]，预测 D

这种处理方法不仅增加了训练样本数量，还使模型能够学习不同长度序列的预测模式，更好地模拟品牌实际扩张过程中的决策链。

3.2 对于标签的处理尝试

3.3 数据增强和特征增强

3.4 多模态特征融合设计

针对商业选址预测任务的特点，我们设计了一个多模态融合的神经网络架构，该架构能够同时处理和融合三种不同类型的信息：

3.4.1 序列特征编码

历史选址序列包含了品牌扩张的时序信息和空间模式。我们使用嵌入层将网格 ID 映射到低维向量空间，然后通过 LSTM 网络捕获序列中的时序依赖关系：

$$h_t = LSTM(embed(grid_id_t), h_{t-1}) \quad (3)$$

LSTM 的长短期记忆特性使其能够有效学习品牌在不同时期的选址偏好变化。

3.4.2 空间特征编码

地理坐标信息反映了选址的空间连续性。我们通过多层感知机 (MLP) 对归一化后的坐标特征进行编码，提取空间分布模式。

3.4.3 环境特征编码

POI 特征反映了网格周边的商业环境和设施分布。同样使用 MLP 对 POI 特征向量进行编码，捕获环境因素对选址决策的影响。

3.4.4 特征融合策略

我们采用特征拼接的方式将三种编码后的特征向量组合，然后通过多层融合网络进一步提取联合特征：

$$f_{fused} = MLP_{fusion}(concat(f_{seq}, f_{coord}, f_{poi})) \quad (4)$$

最终通过线性分类器预测下一个网格的概率分布。

3.5 模型最终设计

考虑到数据集中特征的特点，模型采用时序的方案进行建模。整体架构采用多模态融合设计，具体包括：

- 序列编码器：基于嵌入层和 LSTM 的序列特征提取
- 空间编码器：基于 MLP 的坐标特征编码
- 环境编码器：基于 MLP 的 POI 特征编码

- 融合网络：多层感知机进行特征融合
- 分类器：线性层输出概率分布

该设计充分考虑了商业选址的多因素特性，通过端到端的学习方式自动发现不同模式特征间的内在关联，为精准的选址预测提供了有力支撑。

参考文献

- [1] T. G. Dietterich, “Approximate statistical tests for comparing supervised classification learning algorithms,” *Neural Computation*, vol. 10, no. 7, pp. 1895–1923, 1998.
- [2] P. Domingos, “A few useful things to know about machine learning,” *Commun. ACM*, vol. 55, no. 10, p. 78–87, Oct. 2012. [Online]. Available: <https://doi.org/10.1145/2347736.2347755>
- [3] T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” 2013. [Online]. Available: <https://arxiv.org/abs/1301.3781>
- [4] A. Mohan and K. V. Pramod, “Link prediction in dynamic networks using time-aware network embedding and time series forecasting,” *Journal of Ambient Intelligence and Humanized Computing*, vol. 12, no. 2, pp. 1981–1993, Feb 2021. [Online]. Available: <https://doi.org/10.1007/s12652-020-02289-0>
- [5] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. USA: Cambridge University Press, 2014. [Online]. Available: <https://www.cs.huji.ac.il/~shais/UnderstandingMachineLearning/understanding-machine-learning-theory-algorithms.pdf>
- [6] A. Psyllidis, S. Gao, Y. Hu, E.-K. Kim, G. McKenzie, R. Purves, M. Yuan, and C. Andris, “Points of interest (POI): a commentary on the state of the art, challenges, and prospects for the future,” *Computational Urban Science*, vol. 2, no. 1, p. 20, 2022. [Online]. Available: <https://doi.org/10.1007/s43762-022-00047-w>
- [7] J. Heaton, “An empirical analysis of feature engineering for predictive modeling,” in *SoutheastCon 2016*. IEEE, Mar. 2016, p. 1–6. [Online]. Available: <http://dx.doi.org/10.1109/SECON.2016.7506650>

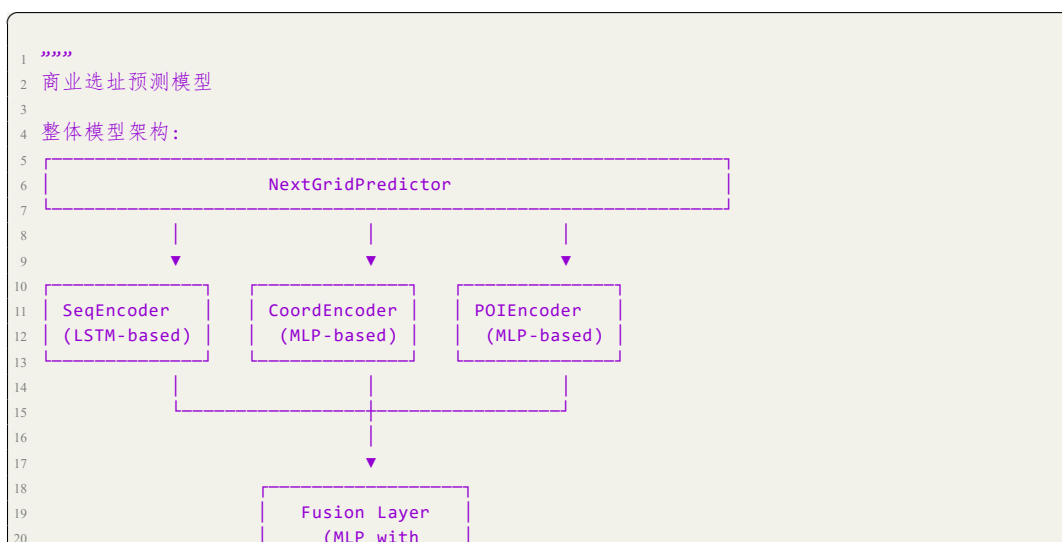
A 附录

A.1 项目信息

GitHub 仓库地址: <https://github.com/szw0407/DL-project-2025>

A.2 核心模型代码

A.2.1 神经网络模型实现 (model.py)



```

21         Dropout)
22     )
23     ↓
24     ▼
25     Classifier
26     (Linear Layer)
27     )
28     ↓
29     ▼
30     预测下一个网格的概率分布
31
32 该模型采用多模态融合架构，综合利用三种信息进行商业选址预测：
33 1. 历史选址序列信息（通过LSTM编码）
34 2. 地理坐标信息（通过MLP编码）
35 3. 兴趣点(POI)特征信息（通过MLP编码）
36
37 这三种信息经过各自的编码器处理后，在特征层面融合，
38 再通过多层感知机进行进一步特征提取，最终用于预测下一个最佳选址位置。
39 """
40
41 import torch
42 import torch.nn as nn
43
44 class SeqEncoder(nn.Module):
45     """
46     序列编码器：用于编码历史选址序列的特征
47
48     该编码器首先将网格ID映射到低维嵌入空间，然后通过LSTM网络捕获序列中的
49     时序依赖关系和选址模式。使用LSTM的优势在于能够记忆长期依赖，
50     适合捕捉商业选址中的空间扩张规律。
51
52     架构：
53     输入序列 → 嵌入层 → LSTM → 最后时刻的隐藏状态
54
55     参数：
56         vocab_size: 词汇表大小，即网格总数
57         embed_dim: 嵌入维度
58         lstm_hidden: LSTM隐藏层维度
59         lstm_layers: LSTM层数
60         dropout: Dropout比率，用于防止过拟合
61     """
62     def __init__(self, vocab_size, embed_dim, lstm_hidden, lstm_layers, dropout):
63         super().__init__()
64         self.embed = nn.Embedding(vocab_size, embed_dim)
65         self.lstm = nn.LSTM(embed_dim, lstm_hidden, lstm_layers,
66                             batch_first=True, dropout=dropout if lstm_layers > 1 else 0)
67
68     def forward(self, seq_ids):
69         """
70         前向传播
71
72         参数：
73             seq_ids: 网格ID序列，形状为(batch_size, seq_len)
74
75         返回：
76             序列的最终表示，形状为(batch_size, lstm_hidden)
77         """
78         emb = self.embed(seq_ids) # (batch_size, seq_len, embed_dim)
79         out, (h, _) = self.lstm(emb) # out: (batch_size, seq_len, lstm_hidden)
80         return out[:, -1, :] # 返回最后一个时间步的输出作为序列表示
81
82 class MLPEncoder(nn.Module):
83     """
84     多层感知机编码器：用于编码坐标和POI特征
85
86     这是一个简单的两层感知机，通过非线性变换将输入特征映射到固定维度的
87     输出空间。适用于处理非序列性的数值特征，如坐标和POI统计数据。
88
89     架构：
90     输入特征 → 线性层 → ReLU → 线性层 → ReLU
91
92     参数：
93         in_dim: 输入特征维度
94         out_dim: 输出特征维度
95         hidden_dim: 隐藏层维度，默认为32
96     """

```

```

97 def __init__(self, in_dim, out_dim, hidden_dim=32):
98     super().__init__()
99     self.model = nn.Sequential(
100         nn.Linear(in_dim, hidden_dim),
101         nn.ReLU(),
102         nn.Linear(hidden_dim, out_dim),
103         nn.ReLU()
104     )
105
106 def forward(self, x):
107     """
108     前向传播
109
110     参数:
111         x: 输入特征, 形状为(batch_size, in_dim)
112
113     返回:
114         编码后的特征, 形状为(batch_size, out_dim)
115     """
116     return self.model(x)
117
118 class NextGridPredictor(nn.Module):
119     """
120     下一个网格预测器: 预测品牌下一个最佳选址位置
121
122     这是整个模型的主体部分, 采用多模态融合方法, 将三种不同类型的信息
123     (历史选址序列、地理坐标、POI特征)整合起来, 共同预测下一个最佳选址位置。
124
125     模型处理流程:
126     1. 通过各自的编码器分别编码三种信息
127     2. 将编码后的特征向量拼接起来
128     3. 通过多层感知机进行特征融合和提取
129     4. 最终通过线性分类器预测下一个网格的概率分布
130
131     多模态融合的优势:
132     - 序列信息: 捕捉品牌扩张的时序模式和依赖关系
133     - 坐标信息: 考虑地理位置的连续性和空间分布
134     - POI信息: 考虑周边设施和商业环境的影响
135
136     参数:
137         num_classes: 类别数量, 即网格总数
138         embed_dim: 网格ID的嵌入维度, 默认为32
139         lstm_hidden: LSTM隐藏层维度, 默认为64
140         lstm_layers: LSTM层数, 默认为1
141         coord_dim: 坐标特征维度, 默认为2 (经度、纬度)
142         poi_dim: POI特征维度, 默认为10 (10种POI类型)
143         coord_out_dim: 坐标编码后的维度, 默认为16
144         poi_out_dim: POI编码后的维度, 默认为16
145         fusion_dim: 特征融合后的维度, 默认为64
146         dropout: Dropout比率, 用于防止过拟合, 默认为0.1
147     """
148     def __init__(self, num_classes, embed_dim=32, lstm_hidden=64, lstm_layers=1,
149                 coord_dim=2, poi_dim=10, coord_out_dim=16, poi_out_dim=16, fusion_dim=64,
150                 dropout=0.1):
151         super().__init__()
152         # 序列编码器: 处理历史选址序列
153         self.seq_encoder = SeqEncoder(num_classes, embed_dim, lstm_hidden, lstm_layers,
154                                     dropout)
155         # 坐标编码器: 处理地理坐标信息
156         self.coord_encoder = MLPDecoder(coord_dim, coord_out_dim)
157         # POI编码器: 处理兴趣点特征信息
158         self.poi_encoder = MLPDecoder(poi_dim, poi_out_dim)
159         # 特征融合网络: 整合三种编码后的特征
160         self.fusion = nn.Sequential(
161             nn.Linear(lstm_hidden + coord_out_dim + poi_out_dim, fusion_dim),
162             nn.ReLU(),
163             nn.Dropout(dropout),
164             nn.Linear(fusion_dim, fusion_dim),
165             nn.ReLU(),
166             nn.Dropout(dropout)
167         )
168         # 分类器: 预测下一个网格的概率分布
169         self.classifier = nn.Linear(fusion_dim, num_classes)
170     def forward(self, seq_ids, seq_coords, seq_poi):
171         """
172         前向传播

```



```

171     数据流程:
172     1. 各编码器独立处理对应的输入特征
173     2. 对坐标和POI特征取序列平均值, 降维处理
174     3. 拼接所有特征向量
175     4. 通过融合网络进一步提取联合特征
176     5. 通过分类器预测下一个网格的概率分布
177
178     参数:
179     seq_ids: 网格ID序列, 形状为(batch_size, seq_len)
180     seq_coords: 坐标序列, 形状为(batch_size, seq_len, 2)
181     seq_poi: POI特征序列, 形状为(batch_size, seq_len, 10)
182
183     返回:
184     logits: 下一个网格的预测概率分布, 形状为(batch_size, num_classes)
185 """
186 # 输入: (batch, seq_len, dim)
187 seq_out = self.seq_encoder(seq_ids) # (batch, lstm_hidden)
188
189 # 取坐标序列的平均值, 简化处理同时保留整体分布信息
190 coords_out = self.coord_encoder(seq_coords.mean(dim=1)) # (batch, coord_out_dim)
191
192 # 取POI特征序列的平均值, 同样简化处理
193 poi_out = self.poi_encoder(seq_poi.mean(dim=1)) # (batch, poi_out_dim)
194
195 # 特征拼接: 将三种编码后的特征向量拼接起来
196 x = torch.cat([seq_out, coords_out, poi_out], dim=-1)
197
198 # 特征融合: 通过MLP进一步提取联合特征
199 f = self.fusion(x)
200
201 # 分类预测: 预测下一个网格的概率分布
202 logits = self.classifier(f)
203 return logits
204

```

Listing 1: 多模态神经网络模型实现

A.2.2 主程序入口 (main.py)

```

1 import torch
2 from data_preprocessing import load_all_data
3 from model import NextGridPredictor
4 from train import train_model
5 from evaluate import evaluate_model
6
7 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
8
9 train_csv = 'data/train_data.csv'
10 test_csv = 'data/test_data.csv'
11 grid_csv = 'data/grid_coordinates-2.csv'
12 if __name__ == '__main__':
13     print("加载并处理数据...")
14     train_set, val_set, test_set, num_classes, grid2idx = load_all_data(
15         train_csv, test_csv, grid_csv, val_size=0.2
16     )
17
18     print(f"训练样本数: {len(train_set)}, 验证样本数: {len(val_set)}, 测试样本数: {len(test_set)}")
19     model = NextGridPredictor(num_classes=num_classes)
20
21     print("开始训练...")
22     model = train_model(model, train_set, val_set, device, num_epochs=40, batch_size=32, lr=1e-3, patience=5)
23
24     print("在测试集上评估...")
25     acc_k, mrr = evaluate_model(model, test_set, device)
26     print(f"Test_MRR: {mrr:.4f}")
27     for k in [1, 5, 10]:
28         print(f"Test_Acc@{k}: {acc_k[k]:.3f}")

```

Listing 2: 主程序实现

A.2.3 训练逻辑实现 (train.py)

```
1 """
2 模型训练模块
3 该模块提供了模型训练的相关函数，包括训练数据批处理和模型训练流程。
4 实现了带有早停机制的模型训练，以及基于验证集性能的模型选择。
5 """
6 import torch
7 import torch.nn as nn
8 import torch.optim as optim
9 import numpy as np
10 from evaluate import evaluate_model
11
12 def batch_iter(samples, batch_size=32, shuffle=True):
13     """
14     训练数据批次迭代器
15
16     与evaluate模块中的batcher函数类似，但增加了数据打乱功能，
17     适用于训练过程中需要随机打乱数据的场景。
18
19     参数：
20         samples: 样本列表，每个样本是(前缀索引，前缀坐标，前缀POI特征，目标索引)的元组
21         batch_size: 批次大小，默认为32
22         shuffle: 是否打乱数据顺序，默认为True
23
24     生成：
25         每次生成一个批次的数据，包含：
26         - seq_ids: 序列ID张量，形状为(batch_size, max_seq_len)
27         - seq_coords: 序列坐标张量，形状为(batch_size, max_seq_len, 2)
28         - seq_poi: 序列POI特征张量，形状为(batch_size, max_seq_len, 10)
29         - targets: 目标索引张量，形状为(batch_size,)
30     """
31     idxs = np.arange(len(samples))
32     if shuffle:
33         np.random.shuffle(idxs) # 随机打乱索引，增加训练随机性
34     for i in range(0, len(samples), batch_size):
35         batch = [samples[j] for j in idxs[i:i+batch_size]]
36         maxlen = max(len(x[0]) for x in batch)
37         # pad to maxlen
38         seq_ids = np.zeros((len(batch), maxlen), dtype=np.int64)
39         seq_coords = np.zeros((len(batch), maxlen, 2), dtype=np.float32)
40         seq_poi = np.zeros((len(batch), maxlen, 10), dtype=np.float32)
41         for i, (pid, pcoord, ppoi, _) in enumerate(batch):
42             L = len(pid)
43             # 右对齐填充：将数据放在数组的右侧，左侧用0填充
44             seq_ids[i, -L:] = pid
45             seq_coords[i, -L:, :] = pcoord
46             seq_poi[i, -L:, :] = ppoi
47         targets = np.array([x[3] for x in batch], dtype=np.int64)
48         yield (
49             torch.from_numpy(seq_ids),
50             torch.from_numpy(seq_coords),
51             torch.from_numpy(seq_poi),
52             torch.from_numpy(targets)
53         )
54
55 def train_model(model, train_set, val_set, device, num_epochs=40, batch_size=32, lr=1e-3,
56                 patience=5):
57     """
58     训练模型的主函数
59
60     实现了完整的模型训练流程，包括：
61     1. 模型训练与参数优化
62     2. 验证集评估
63     3. 早停机制
64     4. 最佳模型保存
65
66     使用MRR(平均倒数排名)作为模型选择的指标，
67     当验证集上的MRR不再提升时，触发早停机制。
68
69     参数：
70         model: 待训练的模型
71         train_set: 训练数据集
72         val_set: 验证数据集
73         device: 计算设备(CPU/GPU)
```

```

73     num_epochs: 最大训练轮数, 默认为40
74     batch_size: 批次大小, 默认为32
75     lr: 学习率, 默认为1e-3
76     patience: 早停耐心值, 当验证指标连续多少轮未改善时停止训练, 默认为5
77
78     返回:
79         训练完成的模型 (已加载最佳状态)
80     """
81     model = model.to(device) # 将模型移至指定设备
82     opt = optim.Adam(model.parameters(), lr=lr) # Adam优化器
83     criterion = nn.CrossEntropyLoss() # 交叉熵损失函数, 适用于多分类问题
84     best_mrr = -1 # 最佳MRR值
85     best_state = None # 最佳模型状态
86     no_improve = 0 # 未改善轮数计数器
87
88     for epoch in range(1, num_epochs+1):
89         # 训练阶段
90         model.train() # 设置模型为训练模式
91         losses = [] # 记录每个批次的损失
92
93         for seq_ids, seq_coords, seq_poi, targets in batch_iter(train_set, batch_size):
94             # 将数据移至指定设备
95             seq_ids, seq_coords, seq_poi, targets = (
96                 seq_ids.to(device), seq_coords.to(device), seq_poi.to(device), targets.to(
97                     device)
98             )
99             opt.zero_grad() # 梯度清零
100             logits = model(seq_ids, seq_coords, seq_poi) # 前向传播
101             loss = criterion(logits, targets) # 计算损失
102             loss.backward() # 反向传播
103             opt.step() # 参数更新
104             losses.append(loss.item()) # 记录损失值
105
106         # 验证阶段
107         val_acc_k, val_mrr = evaluate_model(model, val_set, device)
108
109         # 打印当前训练状态
110         print(f"Epoch {epoch} | loss={np.mean(losses):.4f} | Val_MRR={val_mrr:.4f} | Acc@1={
111             val_acc_k[1]:.3f} Acc@5={val_acc_k[5]:.3f} Acc@10={val_acc_k[10]:.3f}")
112
113         # 模型选择与早停机制
114         if val_mrr > best_mrr: # 如果验证集MRR有提升
115             best_mrr = val_mrr # 更新最佳MRR
116             best_state = model.state_dict() # 保存当前模型状态
117             no_improve = 0 # 重置未改善计数器
118         else:
119             no_improve += 1 # 未改善计数器加1
120             if no_improve >= patience: # 如果连续patience轮未改善
121                 print("Early stop triggered.") # 触发早停
122                 break
123
124         # 加载最佳模型状态
125         if best_state is not None:
126             model.load_state_dict(best_state)
127
128     return model

```

Listing 3: 模型训练实现

A.2.4 数据预处理模块 (data_preprocessing.py)

```

1     """
2     数据预处理模块
3     该模块用于处理位置预测任务的数据, 包括网格信息加载、样本生成和数据集划分。
4     主要实现了基于密度的店铺序列排序和序列到预测样本的转换。
5     """
6     import pandas as pd
7     import numpy as np
8     import ast
9     from sklearn.model_selection import train_test_split
10
11     def parse_list(s):
12         """

```

```

13 将字符串形式的列表解析为Python列表对象
14
15 参数:
16     s: 字符串表示的列表, 如 "[1, 2, 3]"
17
18 返回:
19     解析后的Python列表, 若解析失败则返回空列表
20 """
21 try:
22     return ast.literal_eval(s)
23 except Exception:
24     return []
25
26 def load_grid_info(grid_csv_path):
27     """
28     加载网格信息, 包括坐标和POI特征, 并进行归一化处理
29
30     归一化的目的是将不同尺度的特征转换到同一尺度, 提高模型训练效果。
31     坐标归一化到[0,1]区间, 使模型对不同区域的预测更加公平。
32     POI特征归一化避免某些数值较大的特征主导模型训练。
33
34     参数:
35         grid_csv_path: 网格信息CSV文件路径
36
37     返回:
38         coords_map: 网格ID到归一化坐标的映射字典
39         poi_feat_map: 网格ID到归一化POI特征的映射字典
40     """
41     grid_df = pd.read_csv(grid_csv_path, encoding='gbk')
42     coords_map = {}
43     poi_feat_map = {}
44     poi_columns = ['医疗', '住宿', '摩托', '体育', '餐饮', '公司', '购物', '生活', '科教', '汽车']
45     for _, row in grid_df.iterrows():
46         gid = int(row["grid_id"])
47         x = (row["grid_lon_min"] + row["grid_lon_max"]) / 2.0
48         y = (row["grid_lat_min"] + row["grid_lat_max"]) / 2.0
49         coords_map[gid] = (x, y)
50         poi_feat_map[gid] = row[poi_columns].values.astype(float)
51     # 归一化空间和poi
52     xs, ys = zip(*coords_map.values())
53     x_min, x_max = min(xs), max(xs)
54     y_min, y_max = min(ys), max(ys)
55     for gid in coords_map:
56         x, y = coords_map[gid]
57         coords_map[gid] = [(x - x_min) / (x_max - x_min + 1e-8), (y - y_min) / (y_max - y_min + 1e-8)]
58     all_poi = np.stack(list(poi_feat_map.values()))
59     poi_min, poi_max = all_poi.min(axis=0), all_poi.max(axis=0)
60     for gid in poi_feat_map:
61         poi_feat_map[gid] = (poi_feat_map[gid] - poi_min) / (poi_max - poi_min + 1e-8)
62     return coords_map, poi_feat_map
63
64 def sort_by_density(gid_list, coords_map):
65     """
66     基于密度对网格ID进行排序
67
68     算法原理:
69     计算每个点到其他k个最近邻点的平均距离, 距离越小表示该点周围密度越大。
70     通过这种排序方式, 我们可以从高密度区域到低密度区域构建序列,
71     这种顺序更符合商业扩张的实际规律(先在核心区域布局, 再扩展到周边)。
72
73     参数:
74         gid_list: 网格ID列表
75         coords_map: 网格ID到坐标的映射字典
76
77     返回:
78         按照密度排序后的网格ID列表
79     """
80     if len(gid_list) <= 1: return gid_list[:]
81     k = 3 # 考虑最近的3个邻居点
82     locs = [coords_map[g] for g in gid_list]
83     scores = []
84     for i, g in enumerate(gid_list):
85         xi, yi = locs[i]

```

```

86     dists = [float(np.linalg.norm([xi-xj, yi-yj])) for j, (xj, yj) in enumerate(locs) if
87     i != j]
88     # 计算到k个最近邻点的平均距离, 如果点不足k个, 则取所有点
89     # 若没有其他点, 则设置一个极大值表示最低密度
90     avg = np.mean(sorted(dists)[:min(k, len(dists))]) if dists else 1e5
91     scores.append((avg, g))
92     scores.sort() # 按平均距离排序, 距离小的(密度大的)排前面
93     return [g for _, g in scores]
94
95 def make_samples(data_csv_path, coords_map, poi_feat_map, grid2idx, max_seq_len=10):
96     """
97     构建序列预测样本
98
99     采用滑动窗口方法, 从每个品牌的店铺序列中生成多个训练样本:
100     例如序列[A,B,C,D]会生成样本:
101     - 已有[A], 预测B
102     - 已有[A,B], 预测C
103     - 已有[A,B,C], 预测D
104
105     这种处理方法能够:
106     1. 最大化利用有限的数
107     2. 学习不同长度序列的预测模式
108     3. 模拟品牌实际扩张过程中的决策链
109
110     参数:
111     data_csv_path: 数据CSV文件路径
112     coords_map: 网格ID到坐标的映射字典
113     poi_feat_map: 网格ID到POI特征的映射字典
114     grid2idx: 网格ID到索引的映射字典
115     max_seq_len: 最大序列长度, 若超过则截断
116
117     返回:
118     brand_samples: 样本列表, 每个样本是(前缀索引, 前缀坐标, 前缀POI特征, 目标索引)的元组
119     """
120     data_df = pd.read_csv(data_csv_path)
121     brand_samples = []
122     for _, row in data_df.iterrows():
123         brand = row['brand_name']
124         gid_list = parse_list(row['grid_id_list'])
125         seq = sort_by_density(gid_list, coords_map)
126         if len(seq) < 2: continue # 至少需要两个点才能形成序列预测样本
127         for l in range(1, len(seq)):
128             prefix = seq[:l]
129             target = seq[l]
130             if len(prefix) > max_seq_len:
131                 prefix = prefix[-max_seq_len:] # 保留最近的max_seq_len个点
132             prefix_idx = [grid2idx[g] for g in prefix]
133             prefix_coords = [coords_map[g] for g in prefix]
134             prefix_poi = [poi_feat_map[g] for g in prefix]
135             target_idx = grid2idx[target]
136             brand_samples.append((prefix_idx, prefix_coords, prefix_poi, target_idx))
137     return brand_samples
138
139 def load_all_data(train_csv, test_csv, grid_csv, val_size=0.2):
140     """
141     加载并处理所有数据, 划分为训练、验证和测试集
142
143     该函数是数据处理的主入口, 完成以下步骤:
144     1. 加载网格信息(坐标和POI特征)
145     2. 构建全局网格ID到索引的映射
146     3. 生成训练和测试样本
147     4. 从训练样本中再划分出验证集
148
149     数据划分的合理性:
150     - 使用固定的随机种子确保实验可复现
151     - 验证集从训练集中划分, 保证测试集的纯净性
152     - 按品牌-网格对划分, 避免信息泄露
153
154     参数:
155     train_csv: 训练数据CSV文件路径
156     test_csv: 测试数据CSV文件路径
157     grid_csv: 网格信息CSV文件路径
158     val_size: 验证集比例, 默认为0.2
159
160     返回:
161     train_set: 训练集样本

```

```

161     val_set: 验证集样本
162     test_samples: 测试集样本
163     num_classes: 类别数量 (网格总数)
164     grid2idx: 网格ID到索引的映射字典
165 """
166 coords_map, poi_feat_map = load_grid_info(grid_csv)
167 # 构造全网格字典
168 all_grids = set()
169 for csvf in [train_csv, test_csv]:
170     df = pd.read_csv(csvf)
171     for _, row in df.iterrows():
172         all_grids.update(parse_list(row['grid_id_list']))
173 grid2idx = {gid: idx for idx, gid in enumerate(sorted(all_grids))}
174 num_classes = len(grid2idx)
175 train_samples = make_samples(train_csv, coords_map, poi_feat_map, grid2idx)
176 test_samples = make_samples(test_csv, coords_map, poi_feat_map, grid2idx)
177 # 再在train_samples中拆分出val
178 train_idx, val_idx = train_test_split(np.arange(len(train_samples)), test_size=val_size,
179                                       random_state=42)
179 train_set = [train_samples[i] for i in train_idx]
180 val_set = [train_samples[i] for i in val_idx]
181 return train_set, val_set, test_samples, num_classes, grid2idx

```

Listing 4: 数据预处理实现

A.2.5 数据特征增强 (测试数据文件.py)

```

1 import pandas as pd
2 import ast
3 from collections import defaultdict
4
5 # 读取 train_data.csv 文件
6 df_train = pd.read_csv('data/train_data.csv')
7
8 # 初始化一个默认字典来存储 grid_id 和 brand_type 的计数
9 grid_brand_counts = defaultdict(lambda: defaultdict(int))
10
11 # 已知的 brand_type 前两个字符集合
12 brand_types = {'住宿', '摩托', '公司', '餐饮', '体育', '购物', '生活', '汽车', '医疗', '科教'}
13
14 # 遍历 train_data.csv 的每一行
15 for _, row in df_train.iterrows():
16     brand = row['brand_type'][:2] # 取 brand_type 的前两个字符
17     # 将 grid_id_list 从字符串解析为列表
18     grid_ids = ast.literal_eval(row['grid_id_list'])
19
20     # 为每个 grid_id 统计 brand_type 的数量
21     for grid_id in grid_ids:
22         grid_brand_counts[grid_id][brand] += 1
23
24 df_grid = pd.read_csv('data/grid_coordinates.csv')
25
26 # 为每个 brand_type 添加一行, 初始化为 0
27 for brand in brand_types:
28     df_grid[brand] = 0
29
30 # 更新 df_grid 中每个 grid_id 对应的 brand_type 数量
31 for grid_id in grid_brand_counts:
32     # 确保 grid_id 存在于 df_grid 中
33     if grid_id in df_grid['grid_id'].values:
34         for brand, count in grid_brand_counts[grid_id].items():
35             df_grid.loc[df_grid['grid_id'] == grid_id, brand] = count
36
37 # 保存更新后的表格到新的 Excel 文件
38 df_grid.to_csv(
39     "data/grid_coordinates-2.csv", encoding='gbk'
40 )

```

Listing 5: 数据特征增强实现

A.2.6 模型评估模块 (evaluate.py)

```
1 """
2 模型评估模块
3 该模块提供了评估模型性能的函数和用于批量处理数据的工具函数。
4 主要实现了模型在验证集/测试集上的评估指标计算，包括Top-K准确率和MRR。
5 """
6
7 import torch
8 import numpy as np
9
10 @torch.no_grad()
11 def evaluate_model(model, dataset, device, k_list=[1, 5, 10]):
12     """
13     评估模型在给定数据集上的性能
14
15     该函数计算两个关键指标：
16     1. Top-K准确率：预测结果的前K个选项中包含正确答案的比例
17     2. 平均倒数排名(MRR)：正确答案排名的倒数的平均值
18
19     使用@torch.no_grad()装饰器可以在评估时禁用梯度计算，节省内存并加速评估过程
20
21     参数：
22     model: 要评估的模型
23     dataset: 评估数据集
24     device: 计算设备(CPU/GPU)
25     k_list: Top-K准确率中K值的列表，默认为[1, 5, 10]
26
27     返回：
28     acc_k: 包含各个K值对应准确率的字典
29     mrr: 平均倒数排名(Mean Reciprocal Rank)
30     """
31     model.eval() # 设置模型为评估模式
32     acc_k = {k: 0 for k in k_list} # 初始化各K值的准确计数
33     mrr_sum = 0 # MRR累加值
34     total = 0 # 样本总数
35     for seq_ids, seq_coords, seq_poi, targets in batcher(dataset, 64):
36         # 将数据移至指定设备
37         seq_ids, seq_coords, seq_poi, targets = (
38             seq_ids.to(device), seq_coords.to(device), seq_poi.to(device), targets.to(device)
39         )
40         # 获取模型输出(预测分数)
41         logits = model(seq_ids, seq_coords, seq_poi)
42         # 获取预测分数最高的K个位置的索引
43         topk = torch.topk(logits, max(k_list), dim=1).indices.cpu().numpy()
44         targets_np = targets.cpu().numpy()
45
46         for i, target in enumerate(targets_np):
47             # 找出目标位置在预测的topk结果中的排名
48             rank = np.where(topk[i] == target)[0]
49             if len(rank) > 0: # 如果目标在topk中
50                 rank = rank[0] + 1 # 将索引转换为排名(从1开始)
51                 mrr_sum += 1.0 / rank # 计算倒数排名并累加
52                 for k in k_list:
53                     if rank <= k: # 如果排名在k以内，对应的Top-K准确率计数加1
54                         acc_k[k] += 1
55             total += 1
56
57     # 计算平均值
58     mrr = mrr_sum / total if total else 0
59     acc_k = {k: acc_k[k]/total for k in k_list}
60     return acc_k, mrr
61
62 def batcher(samples, batch_size=64):
63     """
64     将样本数据转换为批次形式的生成器函数
65
66     该函数实现了以下功能：
67     1. 将样本数据按批次划分
68     2. 对每个批次内的序列进行填充，使其长度一致
69     3. 构建张量格式的数据，用于模型输入
70
71     填充策略是从序列右侧对齐，这种方式在处理序列数据时更为合理，
72     因为序列的最新部分(右侧)通常包含更重要的信息。
73     """
```

```

73
74 参数:
75     samples: 样本列表, 每个样本是(前缀索引, 前缀坐标, 前缀POI特征, 目标索引)的元组
76     batch_size: 批次大小, 默认为64
77
78 生成:
79     每次生成一个批次的数据, 包含:
80     - seq_ids: 序列ID张量, 形状为(batch_size, max_seq_len)
81     - seq_coords: 序列坐标张量, 形状为(batch_size, max_seq_len, 2)
82     - seq_poi: 序列POI特征张量, 形状为(batch_size, max_seq_len, 10)
83     - targets: 目标索引张量, 形状为(batch_size,)
84
85 """
86 for i in range(0, len(samples), batch_size):
87     batch = samples[i:i+batch_size]
88     # 找出当前批次中最长的序列长度
89     maxlen = max(len(x[0]) for x in batch)
90     # 初始化批次数据数组
91     seq_ids = np.zeros((len(batch), maxlen), dtype=np.int64)
92     seq_coords = np.zeros((len(batch), maxlen, 2), dtype=np.float32)
93     seq_poi = np.zeros((len(batch), maxlen, 10), dtype=np.float32)
94
95     for j, (pid, pcoord, ppoi, _) in enumerate(batch):
96         L = len(pid)
97         # 右对齐填充: 将数据放在数组的右侧, 左侧用0填充
98         seq_ids[j, -L:] = pid
99         seq_coords[j, -L:, :] = pcoord
100        seq_poi[j, -L:, :] = ppoi
101
102    # 提取目标值
103    targets = np.array([x[3] for x in batch], dtype=np.int64)
104
105    # 将NumPy数组转换为PyTorch张量并返回
106    yield (
107        torch.from_numpy(seq_ids),
108        torch.from_numpy(seq_coords),
109        torch.from_numpy(seq_poi),
110        torch.from_numpy(targets)
111    )

```

Listing 6: 模型评估实现

A.3 数据样本展示

A.3.1 训练数据格式

以下直接展示训练数据文件的前 10 行内容:

```

1 brand_name,brand_type,longitude_list,latitude_list,grid_id_list
2 四季啤酒屋,体育休闲服务;娱乐场所;酒吧,"[116.928274525344, 117.043320008084, 116.9422963244636, 117.0778221078006]", "[36.69774403270743,
3 36.69430636057923, 36.69940378336072, 36.65480559942927]", "[71, 77, 72, 25]"
4 携程旅游,生活服务;生活服务场所;生活服务场所,"[117.0099652551294, 117.0083278100292, 116.9029874347601, 117.1258233788835, 117.1318264911793,
5 117.130692205233, 117.0276045368049, 116.9643505011395, 117.0700358356303, 116.9762270603227, 117.0415532841417]", "[36.615152171727929,
6 36.6453498096663, 36.67805564190155, 36.65345336291342, 36.66021335891504, 36.676643613408, 36.65479701387262, 36.656813123805596,
7 36.64246444489299, 36.61303461884545, 36.68037991342308]", "[4, 21, 51, 28, 45, 64, 22, 19, 25, 3, 59]"
8 幸福时光KTV,体育休闲服务;娱乐场所;KTV,"[116.9786056408031, 116.9167728180969, 117.0098096139525, 116.9590453043534,
9 117.04950909466829]", "[36.66161180954602, 36.65290255903322, 36.6710547488655, 36.67043566438539, 36.60606090605174]", "[37, 16, 38, 36, 41]"
10 小电,生活服务;共享设备;充电宝,"[116.9469669662845, 117.0165404020267, 117.0442775933112, 117.0251578749105, 117.0251360192659, 117.0703663502994,
11 116.9679460276465, 117.1241885471427, 117.0454510621802, 116.9908858098853, 117.0164047908445, 117.0595249194539, 117.1492552466137,
12 117.0387305279752, 117.0026990920988, 116.9833750721236, 117.0259239931365, 117.1239399242847, 116.9688828178708, 117.0047431014325,
13 116.9807510358491, 117.0817309886335]", "[36.70479258447448, 36.66309812601193, 36.64531065644391, 36.66396584217707, 36.66344022296101,
14 36.67629346501127, 36.7221597591807, 36.65502138143246, 36.67928044322549, 36.66917117386832, 36.65574600192463, 36.68491232846421,
15 36.70537523707548, 36.65513173060003, 36.61779391991816, 36.68115049015188, 36.66329622415145, 36.69306744935307, 36.6628680099412,
16 36.61928189507721, 36.66546105430051, 36.70016712172751]", "[72, 39, 23, 39, 39, 61, 89, 28, 59, 37, 22, 60, 83, 23, 4, 56, 39, 64, 36, 4,
17 37, 79]"
18 康明眼镜,购物服务;专卖店;眼镜店,"[117.1474925587893, 116.9899725033711, 117.0496668536213, 117.0009462913962, 116.9162896279549,
19 117.1410070518251]", "[36.72060927039779, 36.69860146269658, 36.67619517995269, 36.67480000224907, 36.65303092587671,
20 36.73127828863626]", "[98, 74, 60, 38, 16, 105]"
21 中财管道,购物服务;家居建材市场;建材五金市场,"[117.03345333680044, 116.9703254004627, 116.9880174777508, 116.9558548711589, 116.992987253108,
22 116.9364190460066, 117.143078283607, 116.9750794367394, 117.0050266022032, 117.0778553647915, 117.031448155798, 117.0102340663978,
23 116.9578983721978, 117.0242489307601, 117.039477183328, 117.0215823229305, 116.9612181682283, 116.9550518132942, 116.9431431740838,
24 116.9606835793541, 117.0478167565158]", "[36.68642073855862, 36.71318820869229, 36.66762686110791, 36.71761468549576, 36.68152706900885,
25 36.6848510759297, 36.7221657010541, 36.716158675651, 36.70174365839906, 36.71288009941338, 36.70789597294158, 36.69540520391379,
26 36.69043405380769, 36.69456713613525, 36.68870724027503, 36.71320839471499, 36.70611753046003, 36.66261143837973, 36.71354397295513,
27 36.69362291074891, 36.70676831386894]", "[59, 89, 37, 88, 56, 53, 98, 89, 75, 94, 77, 75, 55, 76, 59, 91, 73, 35, 88, 73, 78]"
28 美利达自行车,购物服务;专卖店;自行车专卖店,"[116.9524783955336, 116.9154449123776, 117.0606450771626, 117.2299536114407, 116.9908572113357,
29 117.075930872932, 117.2280146835699, 117.1371797256954, 116.956813342363, 117.038657173431]", "[36.67667890606948, 36.65417998169556,
30 36.73888578931914, 36.68444825883368, 36.65774427341513, 36.68036777295877, 36.72384491010683, 36.65061320002995, 36.64512170736571,
31 36.6779178612031]", "[54, 16, 102, 68, 37, 61, 101, 29, 18, 59]"
32 酒快到,购物服务;专卖店;烟酒专卖店,"[117.1316969616346, 117.3016293700698, 117.0343191810293, 116.9759775287286,
33 117.1013689834804]", "[36.65471352253594, 36.67943315359312, 36.67821879365496, 36.64535937215268, 36.685721602042]", "[28, 70, 59, 20, 63]"
34 望京烧烤,餐饮服务;中餐;地方风味餐厅,"[116.9464233558445, 117.0335236445724, 117.0455703904242, 117.218120082778, 116.9148782739088,
35 117.1755609416811]", "[36.61671256380603, 36.68151265323061, 36.71683623281009, 36.70442907328825, 36.67787789073074,
36 36.67421657284811]", "[1, 59, 92, 85, 52, 48]"

```

Listing 7: 训练数据样本 (train_data.csv)

- data/train_data.csv - 训练数据
- data/test_data.csv - 测试数据
- data/grid_coordinates-2.csv - 网格坐标映射

详细的代码实现请参考 GitHub 仓库: <https://github.com/szw0407/DL-project-2025>