

文章编号:1671-6833(2022)01-0001-06

基于深度学习的加油站销量预测与营销策略应用研究

卢晨辉¹, 冯 硕¹, 易爱华², 叶晓俊¹

(1.清华大学 软件学院,北京 海淀 100084; 2.中石化销售股份有限公司广东石油分公司,广东 广州 510620)

摘 要: 营销策略的制定是加油站业务的重要部分,而数据驱动的营销策略制定已成为加油站实现精准营销的迫切需求。为此提出了一种基于加油站历史数据、营销策略和关键特征的油品销量预测的深度学习模型和基于销量预测模型的营销策略制定方法。根据加油站历史数据特征,设计了一个多层次的神经网络结构处理不同类别特征的数据,并结合营销策略信息以执行油品的销量预测。另外,通过引入关键特征,提升了销量预测模型的准确度;通过输入营销策略信息的变更,实现了加油站营销策略的自动选择。在真实加油站数据构建的数据集上进行实验,结果显示:所提方法的销量预测模型相比其他主流方法具有更低的预测误差。

关键词: 销量预测; 数据驱动决策; 深度学习; 循环神经网络; 人工智能应用

中图分类号: TP181

文献标志码: A

doi: 10.13705/j.issn.1671-6833.2022.01.014

0 引言

加油站实施的营销策略对加油站的油品销量有着巨大的影响。传统的加油站营销策略的制定依赖于从业人员多年累积的经验,然而相关经验却难以通过直观的量化形式表现。因此,近年来人们试图通过机器学习方法对加油站历史数据进行训练,以应对市场快速变化情景下的不同地区、不同特征的加油站个性化的营销策略的制定和实施(简称精准营销)。数据驱动不同营销策略下的预测销量可以为从业人员的营销策略制定提供指导意见。换言之,数据驱动的加油站销量预测模型的建立对加油站的精准营销具有重要的现实意义。

当前,销量预测模型在很多领域得到了广泛的应用,例如房地产^[1]、餐饮业^[2]等。用于销量预测的历史数据通常存在时间序列形式。早期研究人员主要采用包括自回归滑动平均模型(ARMA)、整合移动平均自回归模型(ARIMA)等线性模型,而随着大数据技术的发展、企业业务信息化应用^[3]和深度学习技术的发展,循环神经网络(RNN)由于其在处理时间序列数据中表现出的优异性能,得到了业界广泛的认可。本文模型

的时序数据处理模块也是基于循环神经网络结构实现的,探讨的是一种端到端的基于加油站历史数据和关键特征的加油站销量预测模型,且提出的模型具有为加油站营销策略的制定提供参考的能力。

根据加油站原始数据多类别多维度的特性和端到端解决思路,本文针对性地设计了特征提取的网络结构,并结合提取到的不同类别的特征表示以及营销策略信息,使用全连接神经网络结构预测目标销量。此外,通过回归检验、关联分析、方差分析和回归模型预测选择4种方式,本文从原始的数据特征中筛选出若干影响销量较大的特征,并将其作为关键特征添加到预测网络的输出层之前的输入数据中,从而提供模型的营销策略敏感性。在真实数据集上训练并测试了本文提出的销量预测方法。实验结果显示,本文所提的销量预测模型和方法相比其他已知的销量预测方法更准确。

本文的主要创新点如下:

(1) 调研了加油站销量的影响因子并实现了基于长短记忆网络 and 全连接网络结构的多类别销量影响因素数据的特征抽取方法;

(2) 结合多类别销量影响因子特征表示和营

收稿日期:2021-07-12;修订日期:2021-09-22

基金项目:国家重点研发计划项目(2019QY1402)

通信作者:叶晓俊(1964—),男,江苏泰兴人,清华大学教授,博士,博士生导师,主要从事大数据安全、大数据分析技术、网络攻击智能发现技术、数据库测试技术等研究,E-mail:yexj@tsinghua.edu.cn。

销策略信息,实现了端到端的加油站销量预测模型,且模型具备指导营销策略制定的能力;

(3)通过多种相关性分析方法选择得到影响销量较大的关键特征,并通过预测网络引入关键特征的方法降低了销量预测的误差,提升了销量预测模型的营销策略敏感度。

1 相关工作

基于本文的目标,回顾了若干文献作为参考和对照。其中,与本文工作主要相关的包括销量预测与时间序列模型。

1.1 销量预测问题

为了具有足够的竞争力和获取更多的利益,商业组织在持续地研究关于销量预测和关键信息维护的更好的模型和技术^[4]。鉴于销量预测任务使用的数据通常为时间序列格式,ARIMA 等时间序列模型可被用于销量预测^[5]。机器学习方法例如广义线性模型、决策树、梯度提升树等也在销量预测任务中得到了广泛的应用^[6]。对于销量预测中的历史数据,研究者可以通过时间窗口的形式提取特征并通过机器学习算法进行预测^[7]。销量预测的结果可以为销量影响因子的确立和销售策略的制定提供指导^[8]。此外,深度学习同样也应用于销量预测任务中。有研究者使用循环神经网络与注意力机制处理相关时间序列

以解决销量预测问题^[9]。除了销量预测任务本身,辅助的分类任务也可用于增强销量预测的表现^[10]。

1.2 时间序列模型

针对加油站精准营销这样的需求,用于销量预测的各种数据一般含有时序形式。支持向量机(SVM)^[11]、线性回归(LR)、决策树、Adaboost^[12]等学习模型在处理时序数据时,由于时序特性引起的高维度,存在消耗资源大、效率低等问题。为此,研究者提出了循环神经网络(RNN)以及其两个知名变种:长短记忆网络(LSTM)^[13]和门控循环单元(GRU)^[14]。在处理时序数据上,RNN 及其变种有着优异的效果。此外,注意力机制^[15]开始用于为时序数据赋予权值,使得深度网络更加关注重要的部分。此后,不依赖 RNN 的自注意力机制模型被提出,同样在时序数据的处理中取得了很好的效果^[16]。

2 加油站销量预测模型与方法

在本节中,本文先描述要解决的问题,包括该问题相关的数据预处理相关工作。在此基础上,介绍本文提出的一种基于深度学习的加油站销量预测框架(如图 1 所示),包括特征表示和销量预测。

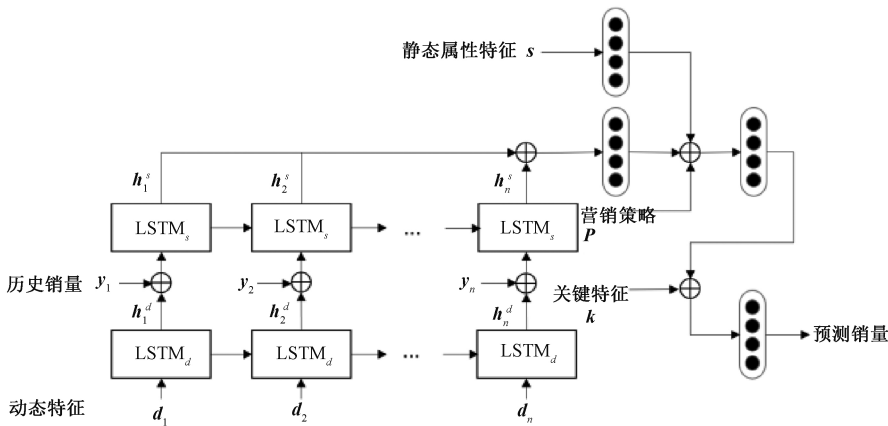


图 1 基于深度学习的油品销量预测模型框架

Figure 1 Framework of the gas sales prediction model

2.1 问题描述

本文研究目标是基于加油站的历史数据构建加油站销量预测模型,使其具备指导相关从业人员制定加油站营销方案的作用。在此应用背景下,模型的输入应当包括加油站的历史数据以及加油站计划采用的营销策略。此外,鉴于不同加油站之间存在的差异性,本文引入了加油站的静

态属性特征作为输入数据的一部分。由于模型中还使用了影响油品销量的一些重要特征,因此模型的输入中还包括采用特征工程分析得到的对预测当天销量重要的关键特征信息。在描述本文提出的销量预测模型前,本文先给出问题形式化的符号说明和含义。

输入的数据包括:加油站的静态属性特征 s ,

例如占地面积、油枪数等;加油站过去 n d 的动态特征 $D = \{d_1, d_2, \dots, d_n\}$, d_i 包括历史营销策略、节假日、天气信息等;加油站过去 n d 的销量数据 $Y = \{y_1, y_2, \dots, y_n\}$; 加油站在第 j 天计划使用的营销策略 p_j ; 第 j 天的关键特征信息 k_j 。输出的数据为预测的油站油品销量 \hat{y}_j , 计算表达式如下:

$$\hat{y}_j = F(s, D, Y, p_j, k_j), \quad j \geq n + 1. \quad (1)$$

式中: $F(\cdot)$ 为预测模型的数学表示, 是预测模型所对应的非线性预测函数。

2.2 数据预处理

销售业务中的加油站营销相关的原始数据包含了多种格式(例如文本数据), 而模型训练与测试过程中需要的是数值类型的数据。鉴于加油业务原始数据中的文本数据多为枚举类型, 在数据预处理过程中, 制定了每个枚举类型特征的取值编码表, 将其转化为数值类型的编码。对于不同取值间未存在大小关系的特征, 使用独热编码对原始特征进行编码。而对于取值类型为连续型数值变量的原始特征, 为消除不同特征存在的不同规模取值范围的影响, 采用了归一化的方法把原始数值转化为 $[0, 1]$ 中的数值。

2.3 数据特征表示

在本文研究的应用场景下, 模型的输入数据包含静态数据、历史动态数据、历史销量数据、营销策略数据以及关键特征数据。特征表示模块主要的目的是设计合适的网络处理以得到静态数据、历史动态数据以及历史销量数据的隐藏层表示。由于历史动态数据和历史销量数据属于时间序列数据而静态数据不随着时间的推移发生变化, 因此设计了不同类型的数据特征提取网络用于得到加油站业务系统汇聚得到的原始特征集成和降维后的特征表示。

对于历史动态数据和历史销量数据, 参考文献[8]中的不同类别时间序列数据分别处理的思路, 使用了一个双 LSTM 网络的结构用于处理该时序特征。令 $LSTM_{dyn}$ 代表处理历史动态数据的网络, 给定输入 d_i , 每一个循环单元均会产生一个隐藏输出 h_i^d ; 令 $LSTM_{sal}$ 代表处理历史销量数据的网络, 以每一个循环单元接受拼接后的 h_i^d 和 y_i 作为输入, 输出集成了历史动态信息和历史销量消息的隐藏层表示 h_i^s 。

$$h_i^d = LSTM_{dyn}(h_{i-1}^d, d_i); \quad (2)$$

$$h_i^s = LSTM_{sal}(h_{i-1}^s, h_i^d, y_i). \quad (3)$$

对于加油站的各种静态数据, 本文通过一个

全连接层获取其特征表示 s_l 。

最终, 通过特征提取模块, 本文获取了加油站静态数据、历史动态数据、历史销量数据的隐藏层表示 $s_l, \{h_i^s\}_{1 \leq i \leq n}$, 并将其输入到后续的销量预测网络中。

2.4 基于关键特征的销量预测

使用一个双层的全连接神经网络进行销量预测任务。在未使用策略敏感性等关键特征信息的情形下, 销量预测网络的输入包括特征提取模块获得的隐藏层表示 $s_l, \{h_i^s\}_{1 \leq i \leq n}$, 以及待预测当天计划执行的营销策略 p_j 。由于 $\{h_i^s\}_{1 \leq i \leq n}$ 拼接后的张量维度较高, 设计一个额外使用了单个的全连接网络, 以便对前面隐藏层表示 $\{h_i^s\}_{1 \leq i \leq n}$ 进行降维:

$$h_l = \text{ReLU}(W_h[h_1^s, h_2^s, \dots, h_n^s] + b_h). \quad (4)$$

式中: h_l 为降维后的特征表示; ReLU 函数为线性整流函数; W_h 和 b_h 为待训练的全连接层权值和偏置值。此后, 拼接 s_l, h_l, p_j 作为预测网络的输入, 通过两层全连接网络输出预测的销量值。

由于加油站实际的销量受当天的部分动态特征影响较大, 例如节假日期间的汽油销量相比非节假日有明显的上升。因此, 在预测模型的设计过程中需要考虑到影响销量的一些重要特征。为此通过回归检验、关联分析、方差分析和回归模型预测选择 4 种方式对每一维度的动态特征数据序列和销量数据序列做相关性分析, 并且从中筛选出与销量数据关联性较大的部分动态特征作为关键特征。目前本文筛选出的关键特征包括节假日、天气、油品挂牌价等动态特征。关键特征作为所提模型的最后一个全连接层的部分输入数据。销量预测网络的计算过程如下:

$$h_a = \text{ReLU}(W_a[s_l, h_l, p_j] + b_a); \quad (5)$$

$$\hat{y}_j = W_o[h_a, k_j] + b_o. \quad (6)$$

式(5)代表了第 1 个全连接层的计算过程, 式(6)代表了第 2 个全连接层的计算过程。其中, W_a, b_a 为待训练的第 1 层全连接层的权值和偏置值; W_o, b_o 为待训练的第 2 层全连接层的权值和偏置值; h_a 为第 1 层全连接层输出的隐藏状态。

预测网络使用的损失函数为平均绝对误差损失函数, 计算表达式如下:

$$MAE = \frac{1}{m} \sum_{k=1}^m |\hat{y}_j^k - y_j^k|. \quad (7)$$

式中: \hat{y}_j^k 和 y_j^k 分别为第 k 个样本的预测销量和实际销量。

2.5 营销策略制定

依据目标加油站的静态数据及过去 n d 的动态数据、销量数据以及预测当天的关键特征信息,观察以不同的预设营销策略作为模型输入的情形下,模型输出的预测销量数值。此后,基于预先设定的营销目标,例如销量优先、效益优先、量效平衡等,计算不同方案的评分,加油站就可筛选出最优的若干营销方案。

算法1 营销策略选择算法。

输入:销量预测模型 F 、静态属性特征 s 、历史动态数据 D 、历史销量数据 Y 、候选营销策略 P 、关键特征 k ;

输出:最佳营销方案 $best_p$ 。

```
① PromotionSelection( $F, s, D, Y, P, k$ );
② for  $p$  in  $P$ :
③    $sales \leftarrow F(s, D, Y, P, k)$ ;
④    $score \leftarrow ComputeScore(sales)$ ; /* 根据营销
目标计算当前销量下的评分 */
⑤   if  $score > best\_score$ 
⑥      $best\_score \leftarrow score$ 
⑦      $best\_p \leftarrow p$ 
⑧   end if
⑨ end for
```

3 实验与结果

本节中,本文使用真实加油站销售数据构建了训练数据集、验证数据集和测试数据集,并在构建的这些数据集上对本文的模型进行了实际测试和验证。

3.1 数据集

本文收集了 201 个加油站从 2020 年 1 月 1 日到 12 月 31 日的影响销量的业务数据,并且基于收集到的数据构建了本文使用的销量预测数据集。数据集的每条记录包含了加油站的静态属性特征、动态特征、营销策略数据以及销量数据。在实验过程中,考虑到真实场景下的销量预测是基于每个站点的历史数据预测未来销量,本文划分数据集的方式为将 3 月 1 日到 10 月 31 日的数据作为训练数据集(考虑到新冠疫情影响,本文去除了 2020 年前两个月的业务数据),11 月 1 日到 12 月 31 日的数据作为测试数据集。实验过程中,加油站历史数据的天数设置为 30 d,且分别训练模型用于预测未来第 1 至第 7 天的销量。

3.2 评价指标

考虑到实际应用的需求,本文使用了两种评

价指标用于评价模型的预测效果。一种直观的评价指标是平均绝对误差(MAE),计算表达式如式(7)所示。此外,考虑到不同站点的不同销量规模,另一种评价指标为平均相对误差(MRE),计算表达式如下:

MRE = 1/m * sum_{k=1}^m (|y_j^k - y_j^k| / y_j^k). (8)

在实际应用过程中,平均绝对误差和平均相对误差相比于均方误差,可以更加直观地体现模型的预测效果。

3.3 单日销量预测实验

本文使用不同的方法,在构建的销量预测数据集上基于过去 30 d 的加油站数据预测未来第 1 天的销量数据,实验结果如表 1 所示。

表 1 预测未来第 1 天销量的实验结果

方法	MAE	MRE
K 近邻	0.789 4	0.128 3
极端随机树	0.811 6	0.135 7
Adaboost	0.857 6	0.132 8
GradientBoosting	0.817 1	0.132 2
基于关键特征的模型	0.641 8	0.106 4

实验结果显示,相比于其他常用的预测算法,本文的模型在销量预测上取得了最低的预测误差。本文的销量预测解决方案是端到端的,预测模型通过分别处理不同类型的原始特征、网络的训练过程调节不同特征对销量预测的影响程度,以及引入关键特征,降低了销量预测的误差。

3.4 单周销量预测实验

实际的业务场景下,销量预测任务需要具备未来一定周期内的预测能力,而对于模型的销量预测能力的评估是基于模型预测给定周期内累计销量的预测误差。因此,本文分别训练模型以预测未来 7 d 的销量进而获取第 1 周的累计预测销量,并与真实数据对比。预测误差及站点误差分布比例如表 2 所示。

表 2 预测未来第 1 周销量的实验结果

误差阈值	站点比例/%	MRE
0.05	79.52	0.042 4
0.10	95.58	0.042 4
0.15	97.99	0.042 4

实验结果显示,在使用本文的模型预测未来 1 周累计销量的情形下,模型的平均相对误差为

0.042 4。其中 95% 以上的加油站站点销量预测误差在 10% 以内。这也说明本文的方法在实际业务场景下的未来周期内预测具有优秀的表现。

3.5 模型结构对销量影响实验

为验证模型不同模块设计的有效性,采用了 3 个不同版本的销量预测模型:第 1 个版本去除了模型中的关键特征部分;第 2 个版本在第 1 个版本的基础上,仅使用一个 LSTM 网络处理拼接后的销量数据和动态特征数据;第 3 个版本在第一个版本的基础上,仅使用 LSTM 的最后一个循环单元的输出而不是所有的循环单元的输出作为预测网络的输入。3 个不同版本模型和原始模型(无关键特征)实验结果如表 3 所示。

表 3 不同版本模型销量预测的实验结果

Table 3 Predicting results of different versions

模型版本	MRE
无关键特征	0.186 7
无关键特征和单 LSTM	0.222 8
无关键特征和单个单元	0.355 4
基于关键特征的模型	0.106 4

实验结果表明,模型设计中的关键特征部分以及特征提取网络的结构和输出形式对销量预测的误差降低具有较好的作用。

3.6 销量预测可视化实验

本文随机挑选了数据集中 3 个加油站,基于其预测销量和实际销量数据绘制曲线。销量预测的可视化结果如图 2 所示。

实验结果显示,预测销量与真实销量相差较小,且预测销量曲线的走势与实际销量曲线基本一致。因此,本文的模型具有较高的预测准确率且具备捕获销量变化趋势的能力。

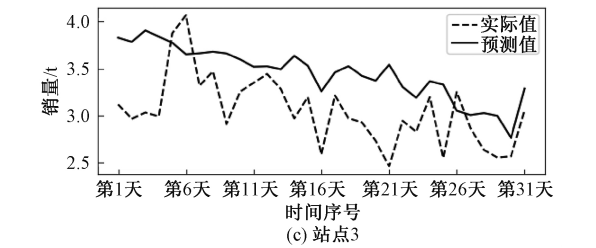
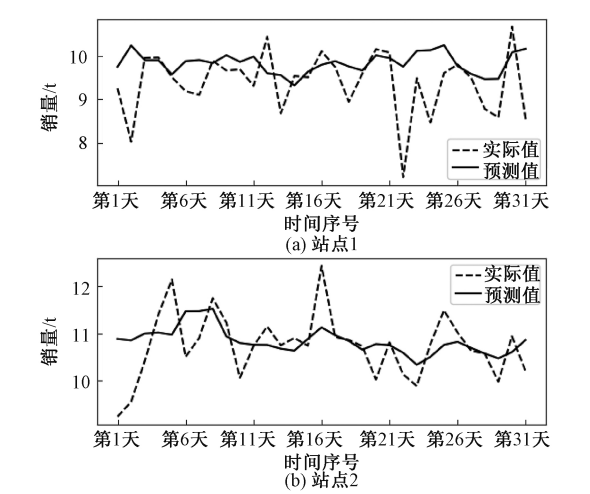


图 2 基于关键特征的销量预测模型结构图
Figure 2 Curves of predicted sales and actual sales

4 结论

本文针对销量预测模型指导下的精准营销的需求,提出了一种基于深度学习的加油站销量预测模型和营销策略选择方法,以满足端到端加油站营销业务需求。采用了针对性的特征提取网络结构来处理不同类别的加油站历史数据中的销量影响因子,并将获得的特征表示与营销策略结合以执行销量预测任务。为应对市场快速变化和营销策略敏感性等需求,本文采用了 4 种相关性分析方法筛选出与销量关联程度较高的关键特征并将其引入预测网络,有效地降低了模型的预测误差。通过观测输入不同营销策略下的销量输出,从业人员可以选择收效更好的营销方案。实验结果说明了本文提出的销量预测模型的预测准确性。未来将在实际应用场景下进一步验证本文提出的销量预测模型和营销策略选择方法的有效性,并针对实际应用下潜在的问题(如节假日、突发事件)及可能影响销量预测的因子数据进行集成,并依据模型应用实践对油站油品的销量预测模型做持续性的优化工作。

参考文献:

[1] CASTELLI M, DOBREVA M, HENRIQUES R, et al. Predicting days on market to optimize real estate sales strategy[J].Complexity,2020,2020:1-22.

[2] POSCH K,TRUDEN C,HUNGERLÄNDER P,et al.A Bayesian approach for predicting food and beverage sales in staff canteens and restaurants[EB/OL].(2020-05-26)[2021-07-06].<https://arxiv.org/abs/2005.12647>.

[3] SAYLI A, OZTURK I, USTUNEL M. Brand loyalty analysis system using K-means algorithm[J].Journal of engineering technology and applied sciences,2016,1(3):107-126.

[4] SASTRY S, BABU M S P. Analysis & prediction of sales data in SAP-ERP system using clustering algorithms[EB/OL].(2013-12-10)[2021-07-06].

- <https://arxiv.org/abs/1312.2678>.
- [5] CHIRU C G, POSEA V V. Time series analysis for sales prediction [C]//International Conference on Artificial Intelligence: Methodology, Systems, and Applications. Cham: Springer, 2018: 163–172.
- [6] CHERIYAN S, IBRAHIM S, MOHANAN S, et al. Intelligent sales prediction using machine learning techniques [C]//2018 International Conference on Computing, Electronics & Communications Engineering. Piscataway: IEEE, 2018: 53–58.
- [7] XIA Z C, XUE S, WU L B, et al. ForeXGBoost: passenger car sales prediction based on XGBoost[J]. Distributed and parallel databases, 2020, 38(3): 713–738.
- [8] BRANDA F, MAROZZO F, TALIA D. Ticket sales prediction and dynamic pricing strategies in public transport[J]. Big data and cognitive computing, 2020, 4(4): 1–17.
- [9] CHEN T, YIN H Z, CHEN H X, et al. Online sales prediction via trend alignment-based multitask recurrent neural networks [J]. Knowledge and information systems, 2020, 62(6): 2139–2167.
- [10] XIN S, ESTER M, BU J J, et al. Multi-task based sales predictions for online promotions[C]//Proceedings of the 28th ACM International Conference on Information and Knowledge Management. New York: ACM, 2019: 2823–2831.
- [11] CORTES C, VAPNIK V. Support-vector networks[J]. Machine learning, 1995, 20(3): 273–297.
- [12] FREUND Y, SCHAPIRE R E. A decision-theoretic generalization of on-line learning and an application to boosting[J]. Journal of computer and system sciences, 1997, 55(1): 119–139.
- [13] HOCHREITER S, SCHMIDHUBER J. Long short-term memory [J]. Neural computation, 1997, 9(8): 1735–1780.
- [14] CHUNG J, GULCEHRE C, CHO K, et al. Empirical evaluation of gated recurrent neural networks on sequence modeling [EB/OL]. (2014–12–11) [2021–07–06]. <https://arxiv.org/abs/1412.3555>.
- [15] XU K, BA J, KIRO S R, et al. Show, attend and tell: neural image caption generation with visual attention [EB/OL]. (2015–02–10) [2021–07–06]. <https://arxiv.org/abs/1502.03044v2>.
- [16] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need [EB/OL]. (2017–01–12) [2021–07–06]. <https://arxiv.org/abs/1706.03762>.

Gasoline Station Sales Prediction Method Based on Deep Learning and Its Application of Promotion Strategy

LU Chenhui¹, FENG Shuo¹, YI Aihua², YE Xiaojun¹

(1. School of Software, Tsinghua University, Beijing 100084, China; 2. Sinopec Sales Co., Ltd., Guangdong Branch, Guangzhou 510620, China)

Abstract: Promotion strategy is an important part of gas station business, and data-driven promotion strategy has become an urgent demand for gas stations to achieve precise marketing. A deep learning model was proposed for forecasting gasoline sales based on historical gas station data, promotion strategies and key features, and a promotion strategy formulation method based on sales forecasting models. Due to the historical data characteristics of gas stations, a multi-level network structure was designed to process data of different types, and combine promotion strategy information to perform oil sales forecasts. In addition, by introducing key features, the accuracy of the sales forecast model was improved; by inputting different promotion strategies, the automatic selection of gas station marketing strategies was realized. The results of experiments conducted on a data set constructed from real gas station data showed that the sales forecast model proposed had lower forecast errors than other mainstream methods.

Keywords: sales prediction; data-driven decision; deep learning; recurrent neural network; artificial intelligence applications