基于大数据的用户行为预测模型研究

张红姣 贵州师范大学 贵州贵阳 550025

摘 要:随着互联网和信息技术的迅猛发展,用户在数字平台上的行为数据量急剧增加,推动了大数据技术在用户行为预测中的广泛应用。通过分析这些用户行为数据,企业可以更好地理解用户需求,从而提供更加精准的服务。本文主要研究了基于大数据的用户行为预测模型,讨论了传统模型和现代机器学习、深度学习模型的优缺点,并提出了一种基于深度学习和机器学习结合的用户行为预测新模型。通过实验数据验证,该模型在预测准确性和计算效率上都取得了显著提升,特别是在用户个性化推荐和精准营销领域具备广泛的应用前景。本文为未来用户行为预测领域的研究提供了理论基础和实践参考。

关键词:大数据;用户行为;预测模型;机器学习

随着互联网的普及,全球数据量呈爆发式增长,用 户在社交媒体、电商平台等途径的行为痕迹构成了庞大 的数据集。这些数据不仅反映了用户的兴趣和偏好,还 可用于预测未来的行为模式, 因此, 如何利用这些数据 进行有效预测成为学术界和工业界的热点问题。尤其在 电子商务和社交网络中,用户行为预测是个性化推荐、 精准营销和用户流失预警的核心技术。例如, 电商平台 可通过预测用户购买意图来精准推荐商品,提升销售转 化率。社交平台则可根据用户历史行为预测其未来偏好, 增强用户活跃度。用户行为预测技术对企业的经济效益 和用户体验提升有直接影响。本研究旨在通过结合大数 据技术、机器学习和深度学习,设计一种高效的用户行 为预测模型,能够处理大规模、多维度数据,在实际应 用中表现出高效性和准确性。该技术广泛应用于电子商 务、广告推荐等领域,不仅能提高转化率,还能减少客 户流失,优化企业资源分配。然而,传统方法难以应对 数据复杂性和规模的挑战,因此,通过大数据分析和深 度学习技术,构建高效预测模型将为企业智能决策提供 坚实技术基础。

一、用户行为预测的相关研究

1、传统用户行为预测方法

用户行为预测早期的研究大多集中在基于统计学和 数据挖掘的传统方法上。这些方法依赖于已有的数据和

作者简介: 张红姣(2004-), 女, 汉族, 贵州遵义, 本科, 研究方向: 大数据挖掘, 数据科学与大数据技术, 卷积神经网络, 机器学习。

经验进行建模。例如,基于回归分析的方法通过分析用户的行为与时间、地点、设备等变量之间的关系,预测未来的行为。这类方法简单易用,但往往假设数据之间存在线性关系,因而难以捕捉用户复杂的行为模式。另一类常用的传统方法是关联规则挖掘,尤其是市场篮分析(Market Basket Analysis)等技术。这类技术通过挖掘用户行为之间的关联性,找到频繁共现的行为模式。例如,在电商领域,关联规则可以揭示用户常购买的产品组合,从而为推荐系统提供支持。尽管传统方法具有一定的理论基础和应用价值,但在处理高维度、非线性、动态性强的数据时,预测效果不尽如人意。

2、基于机器学习的用户行为预测

在机器学习算法逐渐兴起的背景下,用户行为预测的研究方向发生了重大转变。相比于传统的统计学方法,机器学习算法能够在无需明确假设数据分布的前提下,自动从数据中学习模式,因而在面对复杂的用户行为数据时,表现出更高的灵活性和准确性。支持向量机和随机森林是用户行为预测中常用的机器学习算法。支持向量机是一种二分类模型,通过构建超平面来将不同类别的用户行为进行区分。在用户行为预测任务中,支持向量机可以用于分类任务,如用户是否会进行某种特定行为。随机森林是一种基于决策树的集成学习方法。其基本思想是通过构建多棵决策树,对每棵树的预测结果进行投票,从而得到最终的预测结果。随机森林在处理高维特征时表现良好,能够有效降低过拟合风险,并且具有较高的预测精度。此外,随机森林的优势还体现在对缺失数据的处理能力和计算效率的提升上。然而,机器



学习模型尽管在一定程度上解决了非线性问题,但在面对大量时序数据和长依赖行为时,模型的表现仍然有限。

3、深度学习在用户行为预测中的应用

深度学习作为近年来人工智能领域的重大突破,已 经在图像处理、自然语言处理等领域取得了显著的成绩。 而在用户行为预测中,深度学习的应用也逐渐增多,尤 其是在处理时序数据方面,深度学习模型展现出了其独 特的优势。循环神经网络是处理时序数据的典型深度学 习模型。RNN通过引入循环结构,能够将之前的输入信 息保留在网络中,形成对序列数据的长期记忆。基于这 一特点, RNN特别适合用于用户行为预测。例如, 在电 商平台中, 用户的点击、浏览、购买等行为具有显著的 时序特性, RNN可以通过对这些历史行为的分析, 预测 用户的下一步行为。相比于传统RNN, 长短时记忆网 络(LSTM)在解决长时序依赖问题上表现得更加出色。 LSTM通过引入记忆单元和门控机制,能够有效避免梯度 消失问题,从而在处理长时间跨度的用户行为数据时具 有较强的表现。研究表明,LSTM模型在用户行为预测 任务中的表现显著优于传统的机器学习方法, 尤其是在 长序列数据上,能够捕捉到用户行为的潜在规律。此外, 卷积神经网络(CNN)也逐渐被引入用户行为预测领域。 尽管 CNN 最初是为处理图像数据而设计的, 但其在提取 高阶特征方面的能力同样适用于用户行为预测。例如, 用户的行为数据可以被视为时序信号,通过卷积操作提 取局部特征,再通过池化层降低维度,最终进行预测。

二、基于大数据的用户行为预测模型设计

1、数据预处理

在用户行为预测中,数据预处理是模型构建的关键环节之一。由于用户行为数据通常来源广泛,且数据质量参差不齐,常常包含缺失值、异常值和噪声数据。因此,进行数据清洗是保证预测模型效果的前提。针对缺失值问题,常见的处理方法包括删除缺失数据或使用插值法进行填补。例如,用户行为记录中的某些字段可能存在空值,这些空值可能是由于用户未完成某些操作或数据采集系统的局限性所导致的。针对这些情况,可以根据其他类似用户的行为数据进行合理填补,以保证数据的完整性。另外,异常值的处理也是数据预处理的重要步骤。用户行为数据中可能存在一些不合理的异常行为,例如在短时间内大量点击、短时间内频繁切换商品等。这些异常行为可能是由机器人或恶意用户引起的,因此需要通过设定合理的阈值或使用机器学习算法进行自动检测,并对异常数据进行剔除或修正。

2、特征提取

特征提取是用户行为预测中至关重要的一步。特征 是机器学习或深度学习模型进行预测的依据,因此,提 取到高质量的特征对模型的准确性有着直接影响。

在用户行为预测中,常见的特征提取方法包括时间 特征、频率特征和行为序列特征等。时间特征包括用户 的活跃时间段、行为间隔时间等;频率特征包括用户的 点击频次、购买频次等;而行为序列特征则反映了用户 行为的时间序列模式。例如,一个用户的浏览历史可以 作为序列特征,反映其对某类商品的持续关注程度。

3、模型设计

本文提出了一种结合深度学习和机器学习的用户行为预测模型。该模型包括特征提取模块、预测模块和评估模块。特征提取模块主要通过深度学习提取用户行为的高级特征,预测模块则通过机器学习方法(如随机森林、GBDT)对用户未来行为进行预测,评估模块用于评估预测效果并对模型进行优化。在特征提取阶段,模型采用了LSTM网络对用户的时序行为数据进行处理,从而捕捉用户行为的长时依赖关系。在预测阶段,本文采用随机森林算法来预测用户的未来行为,这样可以充分利用随机森林对非线性数据的处理能力。

通过实验验证,这种深度学习与机器学习结合的模型在用户行为预测中的表现优于传统模型,尤其在处理复杂的时序数据时效果显著。

三、实验与分析

1、实验数据

为了验证本文提出的用户行为预测模型,实验数据 来源于某大型电商平台,涵盖了多种用户行为信息,包 括点击、浏览、购买等关键操作。该数据集包含50万条 用户行为记录, 涉及10,000名活跃用户, 这些用户的 行为反映了电商平台上典型的用户互动模式。为了确保 实验结果的公平性与科学性,数据集被划分为训练集和 测试集。训练集包含80%的数据,主要用于模型的训练, 即通过学习用户的历史行为模式,模型逐渐调整内部参 数,从而提高预测的准确性;剩余20%的数据作为测试 集,用于验证模型的泛化能力和预测效果。这样,训练 集和测试集在时间维度上保持一致性, 避免了未来信息 泄漏,同时保证了评估过程的客观性。此外,实验中还 确保了不同用户行为类别在训练集和测试集中的分布均 衡,以避免数据偏斜导致的预测误差。通过对训练集的 充分学习,模型能够捕捉用户的行为模式,而测试集的 使用则能准确评估模型在新数据上的表现, 从而判断其 在实际应用场景中的有效性和可推广性。

2、实验设计

实验的设计过程分为三个关键阶段:数据预处理、 特征提取和模型训练与预测。在数据预处理阶段,由于 用户行为数据通常存在缺失值、异常值和噪声数据,首 先进行了数据清洗和修正。对于缺失值,采取了适当的 填补策略, 例如使用平均值、中位数填补或基于相似用 户行为的推测填补方法,确保数据的完整性。对于异常 值,采用了基于统计规则的异常值检测技术,剔除了极 端行为数据,避免其对模型造成干扰。在清洗数据后, 进行了数据的标准化和归一化处理,消除了不同特征在 量纲上的差异, 使得各特征在模型中能够同等对待, 避 免数值差异过大导致模型训练偏向某些特征。在特征提 取阶段,本文聚焦于挖掘影响用户行为的关键特征。首 先, 提取了用户的时序行为特征, 即用户行为的时间序 列模式, 例如用户在不同时间段的购买频率、活跃时间 点以及行为间隔时间。这些时序特征有助于模型捕捉用 户行为的规律性。其次,提取了用户的频率特征,诸如 用户在一定时间窗口内的点击频率、浏览频率和购买频率 等,这些频率特征能够反映用户对商品的兴趣和参与度。

3、实验结果

实验结果表明,本文提出的深度学习与机器学习结 合的用户行为预测模型在准确率和表现上显著优于传统 方法。在测试集上的预测准确率达到了87.5%,相比传 统的回归分析模型提升了12个百分点,体现了模型在处 理复杂用户行为数据时的卓越性能。具体来说,LSTM模 型在长时序数据的处理上表现尤为优异。LSTM模型通过 其内部的记忆单元和门控机制, 能够有效捕捉到用户行 为数据中的长短期依赖关系。这意味着该模型能够理解 用户的长期行为模式,并根据用户的过去行为做出更为 精准的未来行为预测。例如, 电商平台的用户可能在几 周或几个月内展现出某种持续的兴趣, LSTM 可以识别这 些长期趋势,从而做出更为精准的购买或点击预测。与 此同时, 随机森林模型在处理用户的高维特征时展现出 了强大的鲁棒性。用户行为数据往往包含多种维度的信 息,例如用户的频次特征、时间特征以及社交网络互动 等。面对这些多维度的数据,随机森林通过构建多棵决 策树来处理复杂的非线性关系,并通过多次投票得出最 终的预测结果。这使得随机森林能够在保持模型稳定性 的同时,提升预测的准确性。在实验中,随机森林模型 有效地处理了用户的行为频率、浏览习惯等高维数据, 补充了LSTM在处理非时序数据时的不足,进一步提升 了整体模型的性能。

4、误差分析与优化

尽管本文的模型在预测准确率上表现优异,但仍然 存在一定的误差。这些误差可能源于数据集的不平衡性, 某些行为模式的样本数量较少,导致模型在预测时表现 不佳。为了解决这一问题,未来可以考虑采用样本平衡 技术,如过采样、欠采样等,来提高模型在小样本情况 下的预测能力。此外,模型的超参数选择对结果也有一 定影响,未来可以通过自动化超参数搜索方法,如网格 搜索和贝叶斯优化,进一步提升模型性能。

结束语

本文研究了大数据环境下的用户行为预测问题,提出了一种结合深度学习与机器学习的混合预测模型。通过对用户历史行为数据的分析,该模型能够准确预测用户的未来行为,并在实验中表现出较高的准确率。未来研究可以进一步扩展到多模态数据融合、多任务学习和强化学习等技术领域,以进一步提高预测模型的精度和适用性。

用户行为预测是个性化推荐、精准营销和客户关系管理中的关键技术,随着大数据技术的发展,其重要性日益凸显。同时,企业在应用这些技术时,应高度重视用户数据的隐私保护,以确保技术发展与社会责任的平衡。

参考文献

[1]朱佳燕.基于数据挖掘的电商广告对用户购买行 为影响研究[D].内蒙古财经大学,2024.

[2]于赛赛.基于互联网用户的电影个性化行为预测及推荐研究[D].临沂大学,2024.

[3]许亚如,陶俊宇,梁蕊,等.机器学习在建筑垃圾处理领域的应用与现状[J].环境卫生工程,2024,32(02):10-19.

[4] 李特.基于机器学习的个性化推荐算法研究[D]. 上海应用技术大学,2023.

[5] 丁婷婷.基于机器学习的国家公园游客景观偏好研究[]].城市建筑,2023,20(24):114-118.