

Article

# Machine-Learning-Based Approach for Anonymous Online Customer Purchase Intentions Using Clickstream Data

Zhanming Wen, Weizhen Lin \* and Hongwei Liu

School of Management, Guangdong University of Technology, Guangzhou 510520, China

\* Correspondence: 2112008011@mail2.gdut.edu.cn

**Abstract:** Since online shopping has become an important way for consumers to make purchases, consumers have signed up to e-commerce platforms to shop online. However, retailers are beginning to realise the critical role of predicting anonymous consumer purchase intent to improve purchase conversion rates and store profitability. Therefore, this study aims to investigate the prediction of anonymous consumer purchase intent. This research presents a machine learning model (MBT-POP) for predicting customer purchase behaviour based on multi-behavioural trendiness (MBT) and product popularity (POP) using 33,339,730 clicks generated from 445,336 sessions of real e-commerce customers. The results show that the MBT-POP model can effectively predict the purchase behaviour of anonymous customers ( $F1 = 0.9031$ ), and it achieves the best prediction result with a sliding window of 2 days. Compared to existing studies, the MBT-POP model not only improves the model performance, but also compresses the number of days required for accurate prediction. The present research has argued that product trendiness and popularity can significantly improve the predictive performance of the customer purchase behaviour model and can play an important role in predicting the purchase behaviour of anonymous customers.

**Keywords:** purchase intention; anonymous customer; clickstream; product tendency; product popularity



**Citation:** Wen, Z.; Lin, W.; Liu, H. Machine-Learning-Based Approach for Anonymous Online Customer Purchase Intentions Using Clickstream Data. *Systems* **2023**, *11*, 255. <https://doi.org/10.3390/systems11050255>

Academic Editor: William T. Scherer

Received: 20 April 2023

Revised: 11 May 2023

Accepted: 17 May 2023

Published: 18 May 2023



**Copyright:** © 2023 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

E-commerce has become a popular way to shop. The global e-commerce market is expected to exceed USD 5.7 trillion in 2022 and continue to grow in the coming years [1]. China has the world's largest group of digital shoppers (850 million people); online retail sales in China are even higher at CNY 13.79 trillion in 2022, while cross-border e-commerce imports and exports (including B2B) are CNY 2.11 trillion, both showing steady growth, according to the national online retail market development condition published by the Chinese Ministry of Commerce in 2023 [2]. As a result, online shopping has become a revolutionary way for customers to shop and an important alternative to the traditional market [3]. It is well known that the majority of customers who visit shopping websites tend to end their visit by simply browsing. This poses a challenge to retailers looking to increase their market share and profitability. As a result, there is a focus on improving purchase session rates based on insights into shopping site browsing behaviour. In particular, identifying which customer sessions result in purchases has become a key focus for improving conversion rates. Even a small increase in customer purchase conversion rates can be highly profitable for merchants [4–6].

Understanding online shopping behaviour and gaining insight into customers' decision-making processes can improve the customer experience and increase sales. With the rapid development of the e-commerce industry, it is now possible to record and obtain session logs and behavioural traces of customer groups on shopping websites. Clickstream datasets, which are considered to reflect customers' shopping preferences, have greatly improved in usability [7–9]. This makes it possible to analyse customers' shopping intentions and provides a new approach to understanding customers' decision-making behaviour. Previous studies have identified a phased approach to general customer shopping behaviour,

including an information-gathering stage, a consideration stage and a selection stage [10]. In reality, these stages generate massive clickstream datasets as customers repeatedly land on the shopping platform, capturing their browsing and clicking behaviour. Compared to other methods, clickstream datasets offer the advantage of data retention and can be used to predict customer purchase decision-making behaviour and infer intent.

Learning and analysing customers' historical consumption behaviour has been the focus of most studies on understanding customers' shopping behaviour [11–14]. However, these findings mainly focus on identifying and predicting the shopping behaviour of customer groups with consumption experience on the platform, based on their historical behaviour records. However, there are customer groups of customers who have not registered on the platform or have no historical purchase information, but who still play a critical role in driving revenue and sales. Unfortunately, due to a lack of historical consumption records and purchasing information, these groups have received less research attention. There is still much research to be carried out on inferring purchase intentions and predicting the behaviour of customer groups who have no historical purchase information. It has been found that predicting purchase intent using historical information is challenging for numbers of occasional online shoppers [15,16]. For the purposes of this paper, therefore, 'anonymous shoppers' are defined as those with no previous purchase record. Understanding the purchase intentions of anonymous visitors is crucial as they account for almost half of online purchases, including occasional and unidentified repeat shoppers. Previous studies have primarily focused on anonymous customers and have mined session clickstream datasets to identify known patterns, which are typical cases of frequent visitors with known purchase intentions [16]. This approach effectively boosts the purchase conversion rates and revenues for online retailers. However, there is still a need to explore effective ways to predict the purchase intent of anonymous visitors and improve the shopping experience for all customer groups.

Customer behaviour on e-commerce platforms can be categorised into two types: implicit feedback behaviour and explicit feedback behaviour. Implicit feedback behaviour includes actions such as clicking, favouriting and adding to the shopping basket, which do not directly reflect the customer's preferences for products. On the other hand, explicit feedback behaviour provides a direct indication of customers' preferences for products through activities such as giving praise or leaving bad reviews. These behaviours provide valuable information that e-commerce platforms can use to infer customers purchase intent and design personalised recommendation systems [17,18]. Group psychology also plays an important role in influencing customer behaviour. When people are in an unfamiliar environment, the views and information of social groups can have a strong influence on their behaviour patterns [19,20]. The feedback behaviour of customer groups on products reflects their preferences to some extent, which in turn influences product trends. Changing product trends can have a significant impact on top-selling products [21,22] as they reflect the popularity of the product on the website and dynamic changes in customer preferences [15,16].

The aim of this study is to construct a prediction model of customers' purchase intentions based on group feedback behaviour and to investigate the cumulative effect of the model on purchase intentions under different time windows. To achieve this, we used customer clickstream datasets from Jingdong Mall, a well-known comprehensive shopping platform in China. Various machine learning algorithms were applied during the session to test the predictive performance of the model. The ultimate goal was to predict the purchase behaviour of anonymous customers, which can help improve the purchase conversion rate.

## 2. Related Work

There are several recent studies that analyse customers' online behaviour, such as their historical purchases [11–13], page view history [23], page view rating and time, search behaviour [24], and psychological perceptions [25–27], to understand consumer behaviour and preferences. The advantage of these studies is that the experimental data

are more intuitive and easier to handle, but the disadvantage is that the generalisation and application of the study's findings are poor, and they do not provide a better perception and representation of what consumers really think inside. Clickstream data, which are naturally generated as customers browse online shopping sites, can provide direct or indirect feedback on their willingness and preferences. This has become a popular method for analysing customer behaviour. Currently, applications of clickstream datasets have focused primarily on customer profiling [28], customer segmentation [29,30], and the prediction of consumer behaviour [31,32]. However, these studies often fail to account for change over time and customer interest drift, even when recognising that product trends can influence product popularity. As a result, there is still a research gap in measuring product trends based on clickstream datasets and using them to predict customer purchase intentions. This research will use clickstream data to investigate the prediction of consumer purchase intent, taking into account product popularity, and thus fill the research gap.

The implicit feedback behaviours of customer groups have been used to describe product trends and predict the purchase intentions of anonymous customers [15,16]. For example, Bogina et al. [33] constructed product tendencies from the clicking behaviour of customer groups at different times to predict customers' purchase intentions during sessions. Mokryn et al. [16] used behavioural datasets on viewed and clicked products to differentiate product trend degree features and combined them with time features to predict customers' purchase intentions in the current session. Similarly, Esmeli et al. [15] used behavioural datasets to determine product popularity and predict the purchase intentions of anonymous customers in early sessions based on the minimum and maximum popularity of each session. These studies confirm the role of implicit feedback behaviours based on customer groups in predicting purchase intent. Unfortunately, these studies have the disadvantage of considering only a single type of feedback behaviour, such as browsing clicks, while ignoring other types of implicit feedback behaviours, such as adding to cart, following, commenting, etc. Our research will attempt to fill this gap by further considering multiple types of implicit feedback behaviours in the trend-measure construct, and by exploring the possibility of applying multiple behavioural trend measures to predict the purchase intention of anonymous customers.

In addition, the explicit feedback behaviour of customer groups can influence the electronic word-of-mouth of products and further influence customers' purchase behaviour or intentions [34–38]. Explicit feedback from customers can be positive (good reviews) or negative (bad reviews). When customers browse goods, positive and negative reviews are directly presented to them as a quality signal, which can influence their purchase intentions and decisions [39–41]. These show that the number of positive and negative reviews of a product can reflect changes in the electronic word of mouth about the product. It helps to capture the dynamic preferences of customer groups for products and, in turn, recommends customers' preferred products to facilitate purchase behaviour. However, these studies have not yet defined and explored the nature of the changing trends in positive and negative reviews and, in particular, they have not used review changes to construct variables that reflect electronic word-of-mouth trends. Therefore, our research will attempt to address this shortcoming by proposing to construct variables that characterise eWOM change trends based on review changes and define them as Product Popularity (POP).

Previous studies can be mainly classified as questionnaire surveys or psychological experiments, and are usually based on static data, such as a single type of explicit consumer behaviour characteristic, past shopping experiences, and online review texts, without considering the cumulative effect over time and implicit feedback behaviour data. This study focuses more on dynamic trends and cumulative effects over time, taking into account multiple types of implicit behaviours, and conducts research on predicting customers' purchase intentions based on popularity and the degree of trend.

### 3. Research Methodology

#### 3.1. Data Collection

Our datasets were obtained from JingDong, a leading comprehensive online shopping platform in China, covering millions of brands in 12 categories, including home appliances, digital communications, computers, home department stores, clothing and apparel, mother and child, books, food, and online travel. The anonymous clickstream datasets include click logs from 10,296 stores, 73 categories, 11,199 brands, and 126,441 products, with a total of 445,336 sessions and 3,339,730 clicks generated by anonymous customers between February and April in 2018. Specifically, Jingdong’s clickstream datasets include behaviour, review, and product logs—including anonymous user IDs, behaviour types and times, product categories, and stock keeping unit (SKU) IDs in the behaviour logs, and SKU IDs, total reviews, good and bad reviews, and review times in the review logs. We present the click events and purchases in Table 1, which shows that approximately 1.7% of sessions end in purchases.

**Table 1.** JingDong (JD) dataset general data statistics.

Name	Clicks	Buying Sessions	Non-Buying Sessions	Items
JingDong (JD)	3,339,730	7550	437,786	126,441

At the end of each session, we categorise sessions into two groups based on whether or not a purchase was made, and describe them by the number of clicks in the clickstream dataset, as shown in Table 1. When consumers click no more than 8 times in a session, more consumers do not purchase (76.40%) than do purchase (19.36%). Consumers are more likely to buy than not to buy when they have accumulated eight or more clicks in the current session, and this gap in likelihood will continue to grow. This trend can be seen in the line graph as a gradual increase in the relative likelihood of purchase in Figure 1. Of the non-purchase sessions, those with only one click account for the majority (24.27%), while sessions with two and three clicks ending without a purchase account for 15.26% and 10.74%, respectively. This means that non-purchase sessions with less than three clicks account for more than half (50.27%) of all non-purchase sessions. Conversely, sessions that end with a purchase after only one click represent the smallest proportion of all purchase sessions, at just 0.46%. Recent research suggests that this behaviour occurs when shoppers have gathered enough information about the intended product to make a purchase decision, particularly when discounts are offered during the waiting period. However, the small proportion of such behavioural sessions suggests that most customers’ purchasing behaviour results from the consideration of the match between products and needs. At the same time, it also reflects that there are relatively few ‘smart’ shoppers who buy immediately after waiting-out the promotion period.

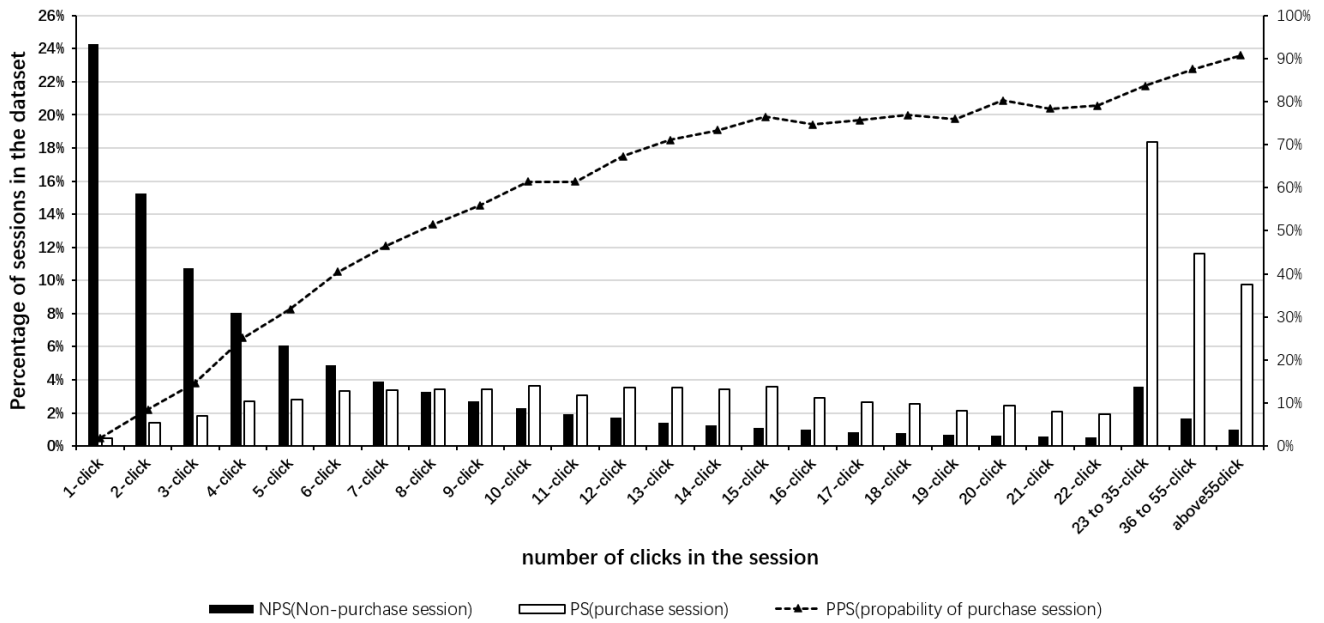


Figure 1. Distribution of the number of clicks per sessions.

### 3.2. Modelling the Trendiness of Products and Sessions

We divide the number of implicit feedback behaviours ( $P_i^n$ ) generated by the customer browsing product  $i$  within  $n$  days into two categories: the number of behaviours received in purchase-ending sessions (PS) and the number of behaviours received in non-purchase-ending sessions (NPS).

$$P_i^n(t) = \sum_{j=n-1}^{t-1} (PS(j, i) + NPS(j, i))$$

Multi-behaviour product tendency ( $MBT_i^n$ ) refers to the ratio of the number of implicit feedback behaviours received by product  $i$  in sessions ending with purchases to the total number of implicit feedback behaviours received by product  $i$  in all viewed sessions over the past  $n$  days. In addition, we classify clicks into four types based on different types of click behaviours on customer pages: attention (SC), add to shopping cart (GWC), review (PL), and browsing (LL):

$$MBT_i^n(t) = \frac{\sum_{j=n-1}^{t-1} PS(j, i)}{P_i^n(t)} = \frac{\sum_{j=n-1}^{t-1} SC(j, i) + GWC(j, i) + PL(j, i) + LL(j, i)}{P_i^n(t)}$$

Multi-behavioural session tendency ( $MBT_s$ ): Among all the sessions  $S$  occurring for product  $i$  on day  $t$ , the session with the maximum product tendency value is defined as the session tendency.

$$MBT_S(t) = \max_i MBT_i^n(t)$$

In addition, we use product tendency based on the number of clicks to obtain the click-tendency of the product (click-tendency of product, CT) [16,33]. The number of clicks in all sessions generated by customers browsing product  $i$  in time  $n$  days ( $CP_i^n$ ) is divided into two categories: clicks ending in purchase sessions (PS) and clicks ending with non-purchase sessions (NPS).

$$CP_i^n(t) = \sum_{j=n-1}^{t-1} (PS^C(j, i) + NPS^C(j, i))$$

Product tendency ( $CT_i^n$ ): In the last  $n$  days, the ratio of the number of clicks received by product  $i$  in sessions ending in purchase to the total number of clicks received by product  $i$  in all sessions that are viewed.

$$CT_i^n(t) = \frac{\sum_{j=n-1}^{t-1} PS^C(j, i)}{CP_i^n(t)} = \frac{\sum_{j=n-1}^{t-1} LL(j, i)}{CP_i^n(t)}$$

Session tendency based on click behaviour ( $CT_s$ ): Among all sessions  $S$  of product  $i$  on day  $t$ , the session with the maximum product tendency value is defined as the session tendency.

$$CT_S(t) = \max_i CT_i^n(t)$$

### 3.3. Modelling the Popularity of Products and Sessions

Product popularity ( $POP_i^n$ ): The value of the number of positive reviews ( $G_i^n$ ) minus the number of negative reviews ( $B_i^n$ ) received for product  $i$  in the last  $n$  days.

$$POP_i^n(t) = \sum_{j=n-1}^{t-1} (G(j, i) - B(j, i))$$

Session popularity ( $POP_s$ ): Among all the sessions  $s$  on day  $t$  of product  $i$ , the maximum difference between the number of positive reviews and the number of negative reviews is defined as session popularity.

$$POP_{S_i}^n(t) = \max \sum_{j=n-1}^{t-1} POP_i^n(t)$$

### 3.4. Modelling the Temporal and Clickstream Characteristics of a Session

This study describes the time and clickstream characteristics of customer-generated sessions, as shown in Table 2.

**Table 2.** Definition of temporal features and clickstream features.

Feature	Symbol	Definition
Month	M	The month in which the customer has a session, with values in the range (2, 3, 4)
Weekday	WKD	Sunday to Saturday
Festival	FES	Whether it is a traditional holiday
Dwell time	DT	The time spent by customers browsing product pages during current session
Clicks counts	CN	The total number of clicks generated by the customer during current session

## 4. Data Analysis and Results

The aim of this study is to develop a model for predicting customer purchase intentions based on product tendency and popularity in the form of group feedback. We also seek to investigate the cumulative effect of the model on purchase intention over different time windows. Inspired by Mokryn et al. [16], we conduct three sets of comparative experiments to confirm its validity and applicability. For the classification task of predicting the purchase intention of anonymous visitors, we train a set of machine learning classifiers to evaluate the effects of different dynamic features, as shown in Tables 3–5.

The time window for the dynamic in the experiments is set from 2 to 6 days. The three groups of experiments include four main types of features: product tendency, product popularity, session time features, and session clickstream features. Product tendency includes the CT variable and the MBT variable, while product popularity includes the POP variable. Session time features include month (M), festival (FES), and session dwell time of a session (DT). Experimental group 1 consists of three experiments: control group, with the CT variable, and with the MBT variable, as shown in Table 1. Group 2 consists of two experiments: control group and with the POP variable, as shown in Table 2. Finally, group



3 consists of three experiments: control group, with the CT variable and POP variable, and with the MBT variable and POP variable, as shown in Table 3.

**Table 3.** Experiment about measuring the impact of MBT on purchase intention predictions.

Group-1			
Features	Control Group	with CT	with MBT
Temporal	Trendiness		-
	Month	Month	Month
	Weekday	Weekday	Weekday
	Festival	Festival	Festival
	Dwell Time	Dwell Time	Dwell Time
Clicks	Number of Clicks	Number of Clicks	Number of Clicks
Window size		2–6 Days	

**Table 4.** Experiment about measuring the impact of POP on purchase intention predictions.

Group-2		
Features	Control Group	with POP
Popularity	-	POP
Temporal	Month	Month
	Weekday	Weekday
	Festival	Festival
	Dwell Time	Dwell Time
Clicks	Number of Clicks	Number of Clicks
Window size		2–6 Days

**Table 5.** Experiment about measuring the synergy of POP and CT/MBT on purchase intention predictions.

Group-3			
Features	Control Group	with CT-POP	with MBT-POP
Trendiness	-	CT	MBT
Popularity	-	POP	POP
Temporal	Month	Month	Month
	Weekday	Weekday	Weekday
	Festival	Festival	Festival
	Dwell Time	Dwell Time	Dwell Time
Clicks	Number of Clicks	Number of Clicks	Number of Clicks
Window sizes		2–6 Days	

Figure 2 shows our research process, which is mainly divided into two parts: feature engineering and comparative experiments. In the first part, we use 60% of the data to learn the dynamic features of the product and construct the lookup table of product tendency (CT, MBT) and product popularity (POP) under different dates and time windows. The remaining 40% is then used to generate the anonymous visitor session dataset. Finally, we calculate the dynamic variable (CT, MBT and POP) based on the lookup table of product tendency and product popularity. In the second part, we use 80% of the session dataset as the training set and 20% as the test set. In both datasets, the number of sessions ending with a purchase is similar to the number of sessions ending without a purchase. Due to the relatively large proportion of sessions ending with no purchase in the training set, we randomly down-sample them to reduce data imbalance and train a better model. We sequentially train 120 models based on different time windows (ranging from 2 to 6 days) and different experimental groups (6 models in total) using logistic regression (LR), random forest (RF), histogram-based gradient boosting (HGBDT) and XGBoost (XGB) classifiers. Finally, we evaluate and compare the models based on their F1 values on the test set.

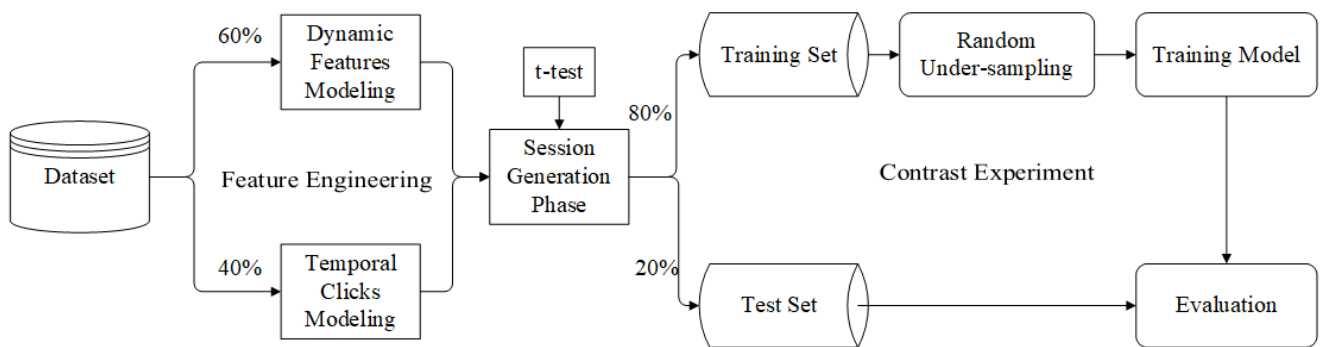


Figure 2. Research framework.

#### 4.1. The Effect of Product Trendiness over Different Window Sizes

In experimental group 1, we use the purchase intention prediction model based on session time and clickstream characteristics as a control group. We evaluate the predictive quality of the model without product tendency and with the CT variable and MBT variable included under different time windows. Specifically, we evaluate the predictive performance of the model with the CT variable and MBT variable replacing product tendency over different time windows. For each scenario, the prediction performance is compared and the four different classifiers are compared separately.

The F1 value is used for the classification performance of the model. Table 6 shows the prediction performance of different classifiers for anonymous customers’ purchase intentions across different time windows and product tendencies. Among the three classifiers (HGBT, RF, and XGB), both the CT model and MBT model significantly improve the prediction performance of customers’ purchase intentions compared to the control group in each time window. Similarly, the predictive quality of the control groups’ integrated classifiers decreases as the time window increases. In addition, the tendency of multi-type behaviour is better than just browsing tendency alone. The above results show that the MBT model, based on the implicit feedback behaviours of multiple groups, is more useful for predicting purchase intentions than the CT model alone, which is based on product browsing and clicking behaviour.

Table 6. The quality of prediction (F1) of the product trendiness model over different time windows.

Days	Features	HGBDT			LR			RF			XGB		
		F1	Precision	Accuracy	F1	Precision	Accuracy	F1	Precision	Accuracy	F1	Precision	Accuracy
2	Control group	0.7898	0.7088	0.6927	0.7861	0.6477	0.6478	0.8344	0.7701	0.7660	0.7729	0.7004	0.6720
	With CT	0.8033	0.7609	0.7303	0.7863	0.6478	0.6481	0.8680	0.8093	0.8157	0.8004	0.7351	0.7163
	With MBT	0.8178	0.8021	0.7593	0.7765	0.6440	0.6356	0.8881	0.8449	0.8473	0.7960	0.8014	0.7376
3	Control group	0.7408	0.6617	0.6499	0.7460	0.5949	0.5950	0.8099	0.7365	0.7489	0.7334	0.6557	0.6403
	With CT	0.7799	0.7315	0.7197	0.7458	0.5947	0.5947	0.8624	0.7957	0.8214	0.7610	0.7161	0.6967
	With MBT	0.7938	0.7637	0.7447	0.7389	0.5920	0.5869	0.8875	0.8329	0.8568	0.7747	0.7676	0.7295
4	Control group	0.7040	0.6426	0.6361	0.7146	0.5560	0.5560	0.7954	0.7223	0.7469	0.7027	0.6305	0.6267
	With CT	0.7425	0.7083	0.6991	0.7144	0.5558	0.5557	0.8495	0.7858	0.8179	0.7226	0.6860	0.6741
	With MBT	0.7708	0.7659	0.7435	0.7146	0.5560	0.5560	0.8659	0.8175	0.8415	0.7460	0.7275	0.7102
5	Control group	0.6942	0.6385	0.6495	0.6029	0.6067	0.5871	0.8019	0.7352	0.7721	0.6924	0.6253	0.6394
	With CT	0.7514	0.7300	0.7320	0.5425	0.5517	0.5291	0.8612	0.8017	0.8431	0.7255	0.6876	0.6961
	With MBT	0.7551	0.7633	0.7465	0.5874	0.5747	0.5586	0.8798	0.8272	0.8657	0.7537	0.7386	0.7369
6	Control group	0.6546	0.6291	0.6383	0.4554	0.6245	0.5694	0.7960	0.7184	0.7702	0.6545	0.6205	0.6327
	With CT	0.7230	0.7087	0.7160	0.4522	0.5865	0.5521	0.8528	0.7881	0.8388	0.7061	0.6750	0.6904
	With MBT	0.7378	0.7376	0.7364	0.4724	0.6102	0.5676	0.8706	0.8199	0.8614	0.7321	0.7237	0.7277

#### 4.2. The Effect of Product Popularity over Different Window Sizes

Table 7 illustrates the predictive performance of the different classifiers for anonymous customer purchase intentions across different time windows and product popularity features. For each time window with different integrated classifiers, the POP variable group outperforms the control group in terms of prediction. Product popularity significantly improves the classifier’s prediction quality of anonymous customers’ purchase intentions.



These results demonstrate that the POP variable is useful in determining anonymous shoppers' purchase intentions by taking into account the dynamic changes in the electronic word-of-mouth of products.

**Table 7.** The quality of prediction (F1) of the product popularity model over different time windows.

Days	Features	HGBDT			LR			RF			XGB		
		F1	Precision	Accuracy	F1	Precision	Accuracy	F1	Precision	Accuracy	F1	Precision	Accuracy
2	Control group	0.7898	0.7088	0.6927	0.7861	0.6477	0.6478	0.8344	0.7701	0.7660	0.7729	0.7004	0.6720
	With POP	0.7950	0.7274	0.7074	0.7863	0.6478	0.6481	0.8424	0.7723	0.7756	0.7869	0.7146	0.6930
3	Control group	0.7408	0.6617	0.6499	0.7460	0.5949	0.5950	0.8099	0.7365	0.7489	0.7334	0.6557	0.6403
	With POP	0.7557	0.6969	0.6827	0.7458	0.5947	0.5947	0.8314	0.7617	0.7793	0.7272	0.6756	0.6487
4	Control group	0.7040	0.6426	0.6361	0.7146	0.5560	0.5560	0.7954	0.7223	0.7469	0.7027	0.6305	0.6267
	With POP	0.7194	0.6699	0.6631	0.7146	0.5560	0.5560	0.8196	0.7503	0.7790	0.6985	0.6449	0.6344
5	Control group	0.6942	0.6385	0.6495	0.6029	0.6067	0.5871	0.8019	0.7352	0.7721	0.6924	0.6253	0.6394
	With POP	0.7146	0.6734	0.6819	0.5766	0.5856	0.5637	0.8343	0.7677	0.8102	0.6822	0.6642	0.6582
6	Control group	0.6546	0.6291	0.6383	0.4554	0.6245	0.5694	0.7960	0.7184	0.7702	0.6545	0.6205	0.6327
	With POP	0.6867	0.6507	0.6668	0.4528	0.5983	0.5577	0.8307	0.7638	0.8136	0.6621	0.6392	0.6479

#### 4.3. The Synergistic Promotion Effect of Product Trendiness and Popularity

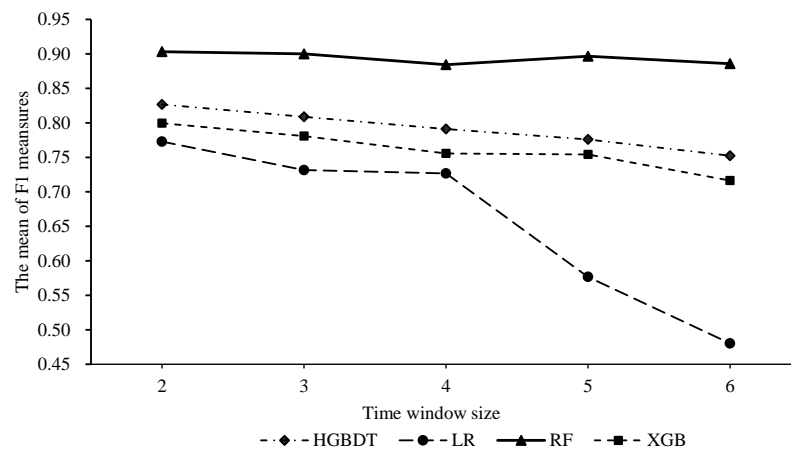
We combined the POP variable with the CT variable and the MBT variable to evaluate the predictive performance of the classifiers on anonymous users' purchase intentions. The results show that the POP variable synergistically improves the prediction of anonymous users' purchase intentions in both the CT variable and the MBT variable (see Table 8). Furthermore, the predictive quality of the MBT-POP model is better than that of the CT-POP in all time windows, indicating that the MBT-POP model has the best predictive performance among all groups (F1 = 0.9031). Compared to similar studies by Mokryn et al. (2019), the MBT-POP model performs better across all time windows and achieves optimal prediction accuracy with only a 2-day time window. These results highlight the importance of considering multi-behaviour with the MBT variable and the POP variable when predicting the purchase intentions of anonymous visitors.

**Table 8.** The quality of prediction (F1) of the MBT-POP model over different time windows.

Days	Features	HGBDT			LR			RF			XGB		
		F1	Precision	Accuracy	F1	Precision	Accuracy	F1	Precision	Accuracy	F1	Precision	Accuracy
2	Control group	0.7898	0.7088	0.6927	0.7861	0.6477	0.6478	0.8344	0.7701	0.7660	0.7729	0.7004	0.6720
	With CT-POP	0.8044	0.7656	0.7332	0.7863	0.6478	0.6481	0.8767	0.8185	0.8282	0.7836	0.7509	0.7070
3	With MBT-POP	0.8267	0.8352	0.7778	0.7728	0.6465	0.6344	0.9031	0.8700	0.8696	0.7995	0.8388	0.7520
	Control group	0.7408	0.6617	0.6499	0.7460	0.5949	0.5950	0.8099	0.7365	0.7489	0.7334	0.6557	0.6403
4	With CT-POP	0.7774	0.7399	0.7211	0.7458	0.5947	0.5947	0.8662	0.8069	0.8282	0.7617	0.7185	0.6985
	With MBT-POP	0.8088	0.8082	0.7724	0.7315	0.5960	0.5866	0.9000	0.8625	0.8756	0.7808	0.7900	0.7423
5	Control group	0.7040	0.6426	0.6361	0.7146	0.5560	0.5560	0.7954	0.7223	0.7469	0.7027	0.6305	0.6267
	With CT-POP	0.7474	0.7167	0.7065	0.7144	0.5558	0.5557	0.8559	0.7969	0.8270	0.7274	0.6863	0.6776
6	With MBT-POP	0.7911	0.7937	0.7685	0.7267	0.5947	0.6094	0.8843	0.8433	0.8648	0.7557	0.7357	0.7207
	Control group	0.6942	0.6385	0.6495	0.6029	0.6067	0.5871	0.8019	0.7352	0.7721	0.6924	0.6253	0.6394
5	With CT-POP	0.7516	0.7332	0.7334	0.5464	0.5571	0.5343	0.8714	0.8172	0.8559	0.7349	0.6892	0.7029
	With MBT-POP	0.7760	0.7887	0.7693	0.5766	0.5867	0.5645	0.8966	0.8582	0.8867	0.7542	0.7474	0.7405
6	Control group	0.6546	0.6291	0.6383	0.4554	0.6245	0.5694	0.7960	0.7184	0.7702	0.6545	0.6205	0.6327
	With CT-POP	0.7284	0.7195	0.7237	0.4833	0.5796	0.5548	0.8597	0.8008	0.8479	0.7037	0.6783	0.6907
6	With MBT-POP	0.7524	0.7679	0.7561	0.4802	0.6038	0.5665	0.8858	0.8473	0.8798	0.7164	0.7612	0.7309

Figure 3 shows the average prediction quality (mean F1) of all experiments conducted in different time windows for the four classifiers. Specifically, the mean F1 score of the  $i$ th classifier in the  $j$ th time window is shown:

$$\text{Mean F1}(i, j) = \frac{1}{N} \sum_N f1(i, j)$$

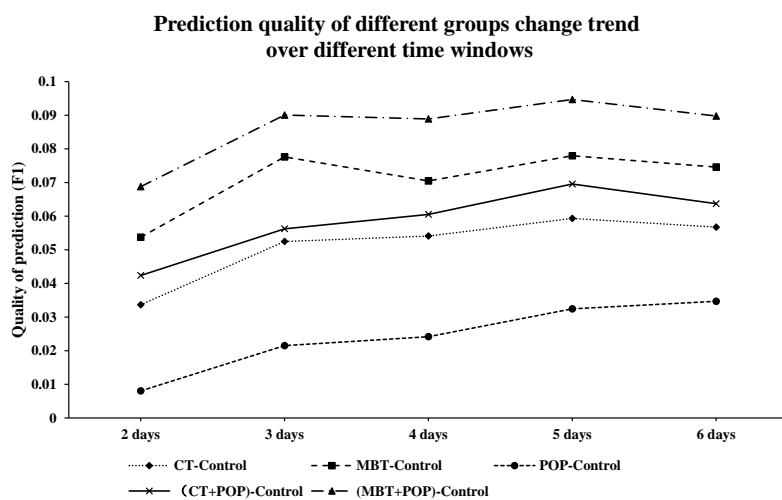


**Figure 3.** The mean performance of classifiers over different window sizes.

The random forest algorithm shows the best overall prediction performance of all classifiers, suggesting its suitability for predicting the purchase intent of anonymous customers.

Comparing the experimental results of the three groups, we obtain the best performance with the random forest classifier of the MBT-POP model and a sliding time window of 2 days (F1 value of 0.9031). It is worth noting that as the sliding time window increases, the F1 value gradually decreases, indicating the importance of recent information in predicting anonymous customers’ purchase intentions.

To investigate how the model improves the quality of purchase intention predictions, we compare the prediction performance of the tendency and popularity models under the random forest classifier with control groups. Figure 4 illustrates how well the best-performing random forest classifier improves the prediction performance of purchase intentions under different feature combinations. The vertical axis in Figure 4 reflects the difference in F1 values between the tendency or popularity model and the control group. Of all the models, the MBT-POP model shows the most significant improvement in prediction quality across all selected time windows. Despite the small number of sessions containing trend products within a time window of 5 to 6 days [16,33], the MBT-POP model achieves an F1 value of 0.8966 for predicting purchase intentions, which is 9.47% higher than that of the control group. This result suggests that the MBT-POP model still shows excellent predictive performance even with small sample data.



**Figure 4.** Enhancement degree of different features on the prediction quality of purchase intent over different window sizes.

## 5. Discussion and Implication

### 5.1. Discussion

The aim of this study was to investigate the factors that influence the purchase intentions of anonymous customers. Previous research on online consumer behaviour has primarily focused on inferring preferences from historical and repurchased customers, which is not applicable to determining the purchase intentions of anonymous visitors with limited historical purchase records, occasional online shoppers, and cold-start users. By identifying the factors that influence the purchase intent of anonymous visitors, we can gain a better understanding of occasional visitors, who account for nearly half of all online purchases and are critical to increasing conversion rates and revenues for online retailers. Using data-driven empirical research, we find that factors such as product tendency, popularity, temporal characteristics and clickstream characteristics based on group feedback behaviours significantly influence the purchase intent of anonymous visitors. Through exploratory analysis (Figure 1), we observe that when the cumulative number of clicks in the current session does not exceed 8, the cumulative proportion of non-purchase sessions reaches 76.40%, while only 19.36% of sessions end in purchases. However, when the number of clicks exceeds 8, the number of sessions ending with a purchase exceeds the number of sessions ending without a purchase, and the probability of sessions ending with a purchase continues to increase with the number of clicks. Before making a purchase decision, customers often evaluate the degree of fit between perceived product value and their own needs, which takes time and effort. We therefore use the number of clicks and time spent to predict the purchase intent of anonymous visitors at the session level.

From the perspective of the dynamic features of products and sessions, we conduct further analysis to explore the impact of product tendency, based on the implicit feedback behaviour of groups, and product popularity, based on dynamic changes in electronic word-of-mouth, on the prediction of anonymous visitors' purchase intentions. Using session time and clickstream features as controls, we conduct three sets of comparative experiments to address the three sub-questions related to our overall research objective, as outlined in Section 4. These experiments have confirmed the effectiveness and applicability of our research.

### 5.2. Theoretical and Practical Implications

First, we show that the predictive performance of the MBT model, which considers multiple types of behaviour, is significantly superior to that of the CT model, which only considers browsing and clicking behaviour. This highlights the importance of product tendency based on consumer information from various implicit feedback behaviours in reflecting the dynamic preferences of anonymous customers. During the consumer-decision process, received information can stimulate preference drift, leading to dynamic changes in consumer preferences, which are directly reflected in different types of clicking behaviour during the shopping journey. By considering multiple types of behaviours, our model can better capture the dynamic changes in consumer preferences, ultimately improving the accuracy of predicting customer behaviour.

Second, our experiment examines the impact of product popularity on the prediction of purchase intentions. We find that including the POP variable of popularity in any time window significantly improves the predictive accuracy of anonymous users' purchase intentions. By tracking the changes in positive and negative reviews, product popularity reflects the identification trends of customer groups for products, which essentially indicates the dynamic trend of electronic word-of-mouth. As the number of positive reviews gradually increases over time, it indicates that customers' appreciation of the products is growing, providing insight into their willingness to purchase based on the level of electronic word-of-mouth and customer recognition.

In conclusion, our research highlights the synergistic role of product tendency and popularity based on group feedback behaviour in predicting anonymous customers' purchase intentions. As mentioned earlier, product popularity reflects customer identification trends

and electronic word-of-mouth, while product tendency reveals customer behavioural patterns and dynamic preferences from a behavioural flow perspective. By combining these two factors, we obtain a comprehensive purchase signal from customers' implicit and explicit feedback. Therefore, considering product tendency and popularity together is more effective and synergistic than focusing on either factor alone.

Furthermore, considering the perspective of machine learning classifiers in our study, we identify the random forest algorithm as a top performing model. The predictive power of anonymous customers' purchase intentions increases with more recent information. In the random forest classifier of the MBT-POP model, we achieve an optimal prediction quality (F1 value) of 0.9031, which corresponds to an optimal time window of 2 days. In addition, we investigate how the combination of different features can improve the quality of purchase intention prediction. As shown in Section 4.3, the MBT-POP model shows the greatest improvement in model prediction performance. This highlights the effectiveness of using two dynamic features based on multi-type behaviour to determine the purchase intentions of anonymous visitors. In particular, our results indicate that customers' purchase intentions are influenced by fashion products and public reputation, suggesting a social influence on purchasing behaviour.

Our research has important practical implications. First, we suggest that product tendency and popularity can be used to design recommendation algorithms that target customers who are infrequent or first-time visitors to e-commerce sites. Second, we recommend that merchants emphasise the social attributes of their products and focus on strengthening their word-of-mouth marketing strategies.

## 6. Conclusions

We present a novel approach to predicting the purchase intent of anonymous customer groups using clickstream datasets to construct customer behaviour tendency and product popularity. By considering these factors, we introduce a new perspective to identify changes in behavioural tendency and product popularity, which helps to detect customer intention and preference drift. In addition, we discover the optimal prediction time scale and machine learning methods, which we apply to independent datasets from the real e-commerce industry to efficiently identify the behavioural intentions of anonymous customer groups. The MBT-POP model achieves the best prediction accuracy by considering product tendency and popularity under the random forest algorithm with a time window of 2 days. These findings provide practical opportunities for structuring a real-time recommendation system based on predicting anonymous customers' purchase intentions. However, there are some limitations to this study. Due to the limitations of the dataset, our research cannot obtain more information related to product reputation, such as retailer reputation and brand effects, which could allow for us to have a more comprehensive definition of product popularity. In future studies, we will continue to explore the predictive performance and applicability of customer feedback information for different customer types and behaviours. We will also refine the analysis of how different implicit feedback behaviours (browsing, favouriting, commenting, adding to cart) affect the prediction of anonymous customer purchase intentions, and how product tendency affects these predictions.

**Author Contributions:** Conceptualization, H.L.; methodology, Z.W. and H.L.; software, W.L.; writing—original draft, Z.W. and W.L.; writing—review and editing, Z.W. and W.L. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by National Natural Science Foundation of China (Grant No.: 71671048). National Education Science Planning Youth Project of the Ministry of Education (Grant No.: EIA210424). The 14th Five-Year Plan of Philosophy and Social Sciences of Guangdong Province 2022 Regular Projects (Grant No.: GD22YJY13).

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** Not applicable.

**Conflicts of Interest:** The authors declare that there is no conflict of interest.

## References

1. Statista. E-Commerce Worldwide—Statistics & Facts. 2023. Available online: <https://www.statista.com/topics/871/online-shopping/> (accessed on 27 February 2023).
2. Research, People's Republic of China Ministry of Commerce. China e-Tailing Market Development in 2022. 2023. Available online: <http://www.mofcom.gov.cn/article/syxwfb/202301/20230103380919.shtml> (accessed on 30 January 2023).
3. Yassein, M.B.; Alomari, O. Detecting the Online Shopping Factors Using the Arab Tweets on Media Technology. *Int. J. Commun. Antenna Propag. (IRECAP)* **2020**, *10*, 206. [\[CrossRef\]](#)
4. Tong, T.; Xu, X.; Yan, N.; Xu, J. Impact of different platform promotions on online sales and conversion rate: The role of business model and product line length. *Decis. Support Syst.* **2022**, *156*, 113746. [\[CrossRef\]](#)
5. Zimmermann, R.; Auinger, A. Developing a conversion rate optimization framework for digital retailers—Case study. *J. Mark. Anal.* **2022**, *11*, 233–243. [\[CrossRef\]](#)
6. Koehn, D.; Lessmann, S.; Schaal, M. Predicting online shopping behaviour from clickstream data using deep learning. *Expert Syst. Appl.* **2020**, *150*, 113342. [\[CrossRef\]](#)
7. Blasco-Arcas, L.; Lee HH, M.; Kastanakis, M.N.; Alcañiz, M.; Reyes-Menendez, A. The role of consumer data in marketing: A research agenda. *J. Bus. Res.* **2022**, *146*, 436–452. [\[CrossRef\]](#)
8. Kukar-Kinney, M.; Scheinbaum, A.C.; Orimoloye, L.O.; Carlson, J.R.; He, H. A model of online shopping cart abandonment: Evidence from e-tail clickstream data. *J. Acad. Mark. Sci.* **2022**, *50*, 961–980. [\[CrossRef\]](#)
9. Gao, X.S.; Currim, I.S.; Dewan, S. Validation of the information processing theory of consumer choice: Evidence from travel search engine clickstream data. *Eur. J. Mark.* **2022**, *56*, 2250–2280. [\[CrossRef\]](#)
10. Jobber, D.; Ellis-Chadwick, F. *EBOOK: Principles and Practice of Marketing, 9e*; McGraw Hill: London, UK, 2019.
11. Zhang, C.; Qiu, J.; Yang, Y.; Zhao, J. Residential customers-oriented customized electricity retail pricing design. *Int. J. Electr. Power Energy Syst.* **2023**, *146*, 108766. [\[CrossRef\]](#)
12. Liao, S.-H.; Widowati, R.; Hsieh, Y.-C. Investigating online social media users' behaviors for social commerce recommendations. *Technol. Soc.* **2021**, *66*, 101655. [\[CrossRef\]](#)
13. Li, Y.; Jia, X.; Wang, R.; Qi, J.; Jin, H.; Chu, X.; Mu, W. A new oversampling method and improved radial basis function classifier for customer consumption behavior prediction. *Expert Syst. Appl.* **2022**, *199*, 116982. [\[CrossRef\]](#)
14. Wang, S.Y.; Qiu, J.T. A deep neural network model for fashion collocation recommendation using side information in e-commerce. *Appl. Soft Comput.* **2021**, *110*, 107753. [\[CrossRef\]](#)
15. Esmeli, R.; Bader-El-Den, M.; Abdullahi, H. Towards early purchase intention prediction in online session based retailing systems. *Electron. Mark.* **2021**, *31*, 697–715. [\[CrossRef\]](#)
16. Mokryn, O.; Bogina, V.; Kuflik, T. Will this session end with a purchase? Inferring current purchase intent of anonymous visitors. *Electron. Commer. Res. Appl.* **2019**, *34*, 100836. [\[CrossRef\]](#)
17. Ko, H.; Lee, S.; Park, Y.; Choi, A. A survey of recommendation systems: Recommendation models, techniques, and application fields. *Electronics* **2022**, *11*, 141. [\[CrossRef\]](#)
18. Roy, D.; Dutta, M. A systematic review and research perspective on recommender systems. *J. Big Data* **2022**, *9*, 59. [\[CrossRef\]](#)
19. Fei, T.; Liu, X. Herding and market volatility. *Int. Rev. Financ. Anal.* **2021**, *78*, 101880. [\[CrossRef\]](#)
20. Loxton, M.; Truskett, R.; Scarf, B.; Sindone, L.; Baldry, G.; Zhao, Y. Consumer behaviour during crises: Preliminary research on how coronavirus has manifested consumer panic buying, herd mentality, changing discretionary spending and the role of the media in influencing behaviour. *J. Risk Financ. Manag.* **2020**, *13*, 166. [\[CrossRef\]](#)
21. Yalcin, E.; Bilge, A. Investigating and counteracting popularity bias in group recommendations. *Inf. Process. Manag.* **2021**, *58*, 102608. [\[CrossRef\]](#)
22. Yi, S.; Kim, D.; Ju, J. Recommendation technologies and consumption diversity: An experimental study on product recommendations, consumer search, and sales diversity. *Technol. Forecast. Soc. Chang.* **2022**, *178*, 121486. [\[CrossRef\]](#)
23. Lu, X.; He, S.; Lian, S.; Ba, S.; Wu, J. Is user-generated content always helpful? The effects of online forum browsing on consumers' travel purchase decisions. *Decis. Support Syst.* **2020**, *137*, 113368. [\[CrossRef\]](#)
24. Huang, S.-L.; Lin, Y.-H. Exploring consumer online purchase and search behavior: An FCB grid perspective. *Asia Pac. Manag. Rev.* **2021**, *27*, 245–256. [\[CrossRef\]](#)
25. Dong, X.; Jiang, B.; Zeng, H.; Kassoh, F.S. Impact of trust and knowledge in the food chain on motivation-behavior gap in green consumption. *J. Retail. Consum. Serv.* **2022**, *66*, 102955. [\[CrossRef\]](#)
26. Klein, A.; Sharma, V.M. Consumer decision-making styles, involvement, and the intention to participate in online group buying. *J. Retail. Consum. Serv.* **2022**, *64*, 102808. [\[CrossRef\]](#)
27. Zhou, Y.; Huang, W. The influence of network anchor traits on shopping intentions in a live streaming marketing context: The mediating role of value perception and the moderating role of consumer involvement. *Econ. Anal. Policy* **2023**, *78*, 332–342. [\[CrossRef\]](#)



28. Pernot, D. Internet shopping for Everyday Consumer Goods: An examination of the purchasing and travel practices of click and pickup outlet customers. *Res. Transp. Econ.* **2021**, *87*, 100817. [\[CrossRef\]](#)
29. Miller, K.; Rosenberg, J.; Pickard, O.; Hawrusik, R.; Karlage, A.; Weintraub, R. Segmenting Clinicians' Usage Patterns of a Digital Health Tool in Resource-Limited Settings: Clickstream Data Analysis and Survey Study. *JMIR Form. Res.* **2022**, *6*, e30320. [\[CrossRef\]](#)
30. Zavali, M.; Lacka, E.; De Smedt, J. Shopping hard or hardly shopping: Revealing consumer segments using clickstream data. *IEEE Trans. Eng. Manag.* **2021**, *70*, 1353–1364. [\[CrossRef\]](#)
31. Ozyurt, Y.; Hatt, T.; Zhang, C.; Feuerriegel, S. A deep Markov model for clickstream analytics in online shopping. In Proceedings of the ACM Web Conference 2022, Lyon, France, 25–29 April 2022; pp. 3071–3081.
32. Gadepally, K.C.; Dhal, S.B.; Kalafatis, S.; Nowka, K. Privacy First Path Analysis Using Clickstream Data. *Preprints.org* **2023**, 2023040904. [\[CrossRef\]](#)
33. Bogina, V.; Kuflik, T.; Mokryn, O. Learning item temporal dynamics for predicting buying sessions. In Proceedings of the 21st International Conference on Intelligent User Interfaces, Sonoma, CA, USA, 7–10 March 2016.
34. González-Rodríguez, M.R.; Díaz-Fernández, M.C.; Bilgihan, A.; Okumus, F.; Shi, F. The impact of eWOM source credibility on destination visit intention and online involvement: A case of Chinese tourists. *J. Hosp. Tour. Technol.* **2022**, *13*, 855–874. [\[CrossRef\]](#)
35. Rahaman, M.A.; Hassan, H.K.; Asheq, A.A.; Islam, K.A. The interplay between eWOM information and purchase intention on social media: Through the lens of IAM and TAM theory. *PLoS ONE* **2022**, *17*, e0272926. [\[CrossRef\]](#)
36. Kurdi, B.; Alshurideh, M.; Akour, I.; Alzoubi, H.; Obeidat, B.; Alhamad, A. The role of digital marketing channels on consumer buying decisions through eWOM in the Jordanian markets. *Int. J. Data Netw. Sci.* **2022**, *6*, 1175–1186. [\[CrossRef\]](#)
37. Nofal, R.; Bayram, P.; Emeagwali, O.L.; Al-Mu'ani, L.A. The Effect of eWOM Source on Purchase Intention: The Moderation Role of Weak-Tie eWOM. *Sustainability* **2022**, *14*, 9959. [\[CrossRef\]](#)
38. Rahayu, S.; Utomo, B.; Kustiningsih, N. The Impact Of Electronic Word Of Mouth (Ewom), Ease Of Use, Trust, And Brand Images To Purchase Intention On Tokopedia: Evidence From Indonesia. *Int. J. Eng. Technol. Manag. Res.* **2022**, *9*, 77–89. [\[CrossRef\]](#)
39. Majali, T.; Alsoud, M.; Yaseen, H.; Almajali, R.; Barkat, S. The effect of digital review credibility on Jordanian online purchase intention. *Int. J. Data Netw. Sci.* **2022**, *6*, 973–982. [\[CrossRef\]](#)
40. Al-Abadi, L.; Bader, D.; Mohammad, A.; Al-Quran, A.; Aldaihani, F.; Al-Hawary, S.; Alathamneh, F. The effect of online consumer reviews on purchasing intention through product mental image. *Int. J. Data Netw. Sci.* **2022**, *6*, 1519–1530. [\[CrossRef\]](#)
41. Duan, Y.; Liu, T.; Mao, Z. How online reviews and coupons affect sales and pricing: An empirical study based on e-commerce platform. *J. Retail. Consum. Serv.* **2022**, *65*, 102846. [\[CrossRef\]](#)

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.