

# AI-Driven Customer Journey Stage Prediction for Real-Time Personalization in E-Commerce

Praveen Kumar

Omni Channel, Discount Tire

Phoenix, Arizona, USA

\*\*\*

**Abstract** - Large amounts of behavioral data are collected by e-commerce platforms, yet most of them are based on personalization pipelines that only take the action after a session has ended. The user experience is worse, and the conversion is restricted by this delay. This paper presents an AI-driven platform that uses session-level behavioral signals, including session duration, click count, page visits, cart activity, and inter-click idle time, to predict the customer's current journey stage in real time: Explorer, Evaluator, Intent Buyer, or Drop-Off Risk. Ground-truth stage labels are generated by a rule-based labeling scheme, and the raw signals are enhanced by three engineering features (engagement score, cart ratio, and idle time). Four classifiers are trained, compared, and thoroughly examined: logistic regression, decision tree, random forest, and gradient boosting. A personalization engine consumes each prediction and initiates stage-appropriate actions: Recovery prompts for drop-off candidates, contextual discounts for intent buyers, comparative tools for evaluators, and popular products for explorers. Random Forest and Gradient Boosting tied at 92.7% test accuracy with macro F1-scores of 0.92 on a behavioral dataset of 2,400 simulated sessions with 5% injected label noise, significantly surpassing the Logistic Regression baseline (71.5% accuracy). Sample efficiency is confirmed by per-class one-vs-rest AUC values exceeding 0.92 for all four phases and the model maintaining accuracy above 90% with as little as 60% of the training data. With a calculated prediction latency of less than 5 ms per session on commodity hardware, the system combines actionable personalization and real-time intent classification in a single low-latency pipeline appropriate for production deployment.

**Key Words:** Customer Journey, Real-Time Personalization, Machine Learning, Random Forest, Gradient Boosting, Behavioural Analytics, E-Commerce, Intent Prediction, Recommendation Systems, Click stream Mining.

## 1. INTRODUCTION

### 1.1 Background and Motivation

The global e-commerce industry has undergone unprecedented expansion over the past decade, with worldwide online retail sales projected to surpass USD 7 trillion by 2027 [1]. Behind every product page view, click, scroll, and cart addition lies a continuous stream of behavioral signals that, when interpreted correctly, reveal the user's intent and the current stage of the purchase

funnel [2], [3]. Despite the volume and richness of this data, most platforms still rely on personalization pipelines driven by aggregated historical profiles and offline batch scoring, which fail to capture the dynamic, session-level shifts in intent that occur within a single visit [4], [5].

The customer journey, originally a marketing construct describing the sequence of touchpoints a buyer encounters between awareness and post-purchase engagement [1], has evolved into a computational object: a time-ordered sequence of behavioral events that can be modeled, predicted, and acted upon. In e-commerce, this journey is compressed: a user may begin a session as a casual Explorer browsing categories, transition into an Evaluator while comparing alternatives and reading reviews, become an Intent Buyer after adding items to a cart, and — if the experience falters — turn into a Drop-Off Risk before completing checkout. Each of these stages calls for a fundamentally different personalization strategy [6], [7]. Showing a discount banner to an Explorer is wasteful and may even feel intrusive; recommending trending categories to an Intent Buyer is irrelevant; failing to recover a Drop-Off candidate is a missed conversion. The cost of mistimed personalization is therefore both reduced revenue and degraded user experience.

### 1.2 Limitations of Conventional Approaches

Conventional approaches to user segmentation in e-commerce, such as post-session funnel analysis [8], A/B testing [9], and static rule-based segmentation [10], operate on a delay that ranges from minutes to days. By the time the analytics pipeline identifies a high-intent user, that user has often already left the platform. Even modern recommendation systems, which excel at next-item prediction [11], [12], typically optimize a single ranking objective and produce no interpretable label that downstream marketing automation can act upon. Real-time prediction of user intent has therefore emerged as a critical research and engineering problem at the intersection of machine learning, web analytics, and human-computer interaction [13], [14], [15].

A second limitation is interpretability. Most published models for purchase prediction are deep neural networks [16], [17] whose decisions are difficult to audit and update. Marketing teams need transparent rules and confidence scores they can reason about, especially when the predictions trigger user-facing interventions such as discounts. There is therefore demand for interpretable,

lightweight models that fit naturally into existing rule-based marketing automation [18].

### 1.3 Proposed Approach

This paper proposes a machine learning-driven framework that classifies a user's journey stage during a live session and triggers a stage-appropriate personalization action within the same interaction. The framework is built on three principles: (i) interpretable behavioral features that can be inspected and modified by domain analysts, (ii) supervised learning with classical models such as Random Forest [19] and Gradient Boosting [20], chosen for their strong empirical performance and built-in feature importance, and (iii) tight integration with a configurable rule engine that maps each predicted stage to a concrete intervention. The pipeline is designed for low-latency deployment, with measured per-prediction inference time below 5 ms on commodity hardware.

### 1.4 Contributions

The main contributions of this work are summarized as follows:

1. Formal definition of four operational customer journey stages — Explorer, Evaluator, Intent Buyer, and Drop-Off Risk — along with a transparent rule-based labeling scheme that converts unlabeled clickstream data into a supervised learning problem suitable for a wide class of e-commerce platforms.
2. A feature engineering pipeline that augments raw behavioral signals with three derived metrics — engagement score, cart ratio, and idle time — that capture interaction intensity, purchase intent, and disengagement patterns respectively.
3. A comparative evaluation of four classifiers (Logistic Regression, Decision Tree, Random Forest, Gradient Boosting) using accuracy, macro-averaged precision/recall/F1-score, per-class metrics, ROC analysis, learning curves, and confusion matrices on a behavioral dataset of 2,400 sessions.
4. A personalization engine that maps predicted stages to specific actions — trending product carousels, comparison panels, contextual discounts, and exit-intent reminders — closing the loop between prediction and user-facing experience.
5. An empirical analysis of the latency, sample efficiency, and hyperparameter sensitivity of the proposed system, demonstrating its viability for production deployment.

### 1.5 Organization of the Paper

The remainder of the paper is organized as follows. Section 2 reviews related work spanning customer journey theory, clickstream mining, recommendation systems, real-time machine learning, and personalization frameworks. Section 3 provides theoretical background on the classifiers and evaluation metrics used in this study.

Section 4 states the problem formally. Section 5 describes the proposed methodology and system architecture. Section 6 details the implementation. Section 7 reports comprehensive experimental results. Section 8 discusses interpretability, latency, and scalability. Section 9 addresses privacy, ethical, and fairness considerations. Section 10 lists current limitations and future work. Section 11 concludes the paper.

## 2. LITERATURE REVIEW

### 2.1 Customer Journey Theory

Lemon and Verhoef [1] formalized the customer experience as a multi-stage journey spanning pre-purchase, purchase, and post-purchase phases, and argued that touchpoint-level analysis is essential for managing experience quality. Verhoef et al. [21] further decomposed customer experience into cognitive, affective, social, and physical dimensions and emphasized the role of digital touchpoints in shaping it. Subsequent work has translated these conceptual stages into computational models suitable for inference [3], [22]. Bucklin and Sismeiro [23] surveyed clickstream models for online customer behavior and identified a gap between marketing-theoretic constructs and the operational signals available to web analytics systems — a gap that the present work attempts to bridge with explicit operational definitions.

### 2.2 Web Analytics and Clickstream Mining

Web usage mining has been a mature field for over two decades. Kosala and Blockeel [24] provided an early systematic survey of web mining techniques, distinguishing between content, structure, and usage mining. Spiliopoulou [25] focused specifically on web usage mining for personalization and emphasized the value of session reconstruction. Mobasher et al. [26] introduced effective personalization based on web usage mining and association rules. These foundational works established the methodological pipeline — collect, clean, sessionize, model — that still underlies modern clickstream applications.

Bucklin and Sismeiro [23] modeled web browsing behavior using hidden Markov models, treating the journey as a discrete-state process. More recent work has applied recurrent neural networks and transformers directly to event sequences [11], [12], [16]. The framework proposed in this paper sits between these two extremes: it uses session-aggregated features (faster and more interpretable than full sequences) but classifies into multiple operational stages (richer than a binary buy/no-buy label).

### 2.3 Purchase Intent Prediction

Moe [2] proposed an early classification of online shoppers based on session-level browsing patterns,

distinguishing directed buying, search and deliberation, hedonic browsing, and knowledge building. This four-way taxonomy directly inspired the operational stages used in this paper. Sakar et al. [3] used Multilayer Perceptron and LSTM recurrent neural networks on session-level features to predict purchasing intent, reporting accuracies above 87% on the UCI Online Shoppers Intention dataset. Suchacka and Stemplewski [4] applied Random Forest classifiers to a real e-commerce log and confirmed the model's robustness to noisy clickstream features. Mokryn et al. [27] inferred current purchase intent of anonymous visitors using session-level features, demonstrating that intent can be reliably estimated even in cold-start settings.

Other studies have applied gradient boosting [28], factorization machines [29], and deep interest networks [30] to large-scale e-commerce data. While these models achieve impressive accuracy, they often lack the interpretability needed for business adoption. The present work prioritizes explainable models that can be reasoned about by non-technical stakeholders.

## 2.4 Recommendation Systems

Recommendation systems form a closely related research stream. Collaborative filtering [31] and matrix factorization [5] dominate this literature, but they are typically optimized for product ranking, not for stage-aware intervention. Ricci et al. [32] provide a comprehensive handbook of recommender systems. Smith and Linden [33] reflected on two decades of recommender systems at Amazon, noting that personalization remains an unsolved problem in many real-world settings.

Sequence-aware recommenders such as GRU4Rec [11] and SASRec [12] improve next-item prediction but do not produce an interpretable label describing the user's current stage. Wide & Deep models [34] and YouTube's deep neural network architecture [17] illustrate large-scale industrial deployment but require substantial infrastructure and labeled data. The framework proposed here is complementary to these systems: it provides a coarse-grained, interpretable stage label that can route the user to the appropriate downstream recommender or marketing module.

## 2.5 Real-Time Machine Learning Systems

Real-time personalization has also been studied from a systems perspective. Amat et al. [13] described Netflix's real-time personalization architecture and emphasized low latency and contextual signals. He et al. [14] reported on Facebook's ad CTR prediction system, which combines decision trees with logistic regression for online ad serving. McMahan et al. [15] presented a large-scale online learning system that updates model parameters incrementally as new events arrive.

In the e-commerce setting, however, much of the published work focuses either on offline batch scoring or

on narrow tasks such as cart abandonment prediction [27], [35]. There remains a clear gap between coarse-grained intent classification and actionable, real-time personalization. The present work addresses this gap by combining lightweight feature engineering, interpretable classifiers, and a deterministic personalization engine into a single end-to-end pipeline that operates within the latency budget of a typical web request.

## 2.6 Interpretability and Trust

As machine learning models are deployed in user-facing systems, interpretability has become a first-class concern. Ribeiro et al. [36] introduced LIME for model-agnostic local explanations, and Lundberg and Lee [37] proposed SHAP, a unified framework for feature attribution based on Shapley values. Tree-based models such as Random Forest and Gradient Boosting come with built-in feature importance scores that are easier to communicate to non-technical stakeholders. The present work leans on this property to keep the proposed system auditable, which is one of the key advantages over deep neural alternatives.

## 2.7 Comparative Summary of Related Work

Table -1 summarizes the methodological choices of representative prior works and contrasts them with the proposed framework. The combination of multi-stage classification, real-time inference, interpretable models, and integrated personalization actions distinguishes the present work from existing literature.

**Table -1:** Comparative summary of representative related work.

Work	Task	Real-time	Model	Action loop
Sakar et al. [3]	Binary buy / no-buy	Yes	MLP / LSTM	No
Suchacka et al. [4]	Binary buy / no-buy	No	Random Forest	No
Hidasi et al. [11]	Next-item rec.	Yes	GRU	Implicit
Kang et al. [12]	Next-item rec.	Yes	Self-attention	Implicit
Mokryn et al. [27]	Purchase intent	Yes	Logistic / GB	No
He et al. [14]	Ad CTR	Yes	Tree + LR	Yes (ads)
This work	4-stage journey	Yes	RF / GB	Yes (rule engine)

## 3. THEORETICAL BACKGROUND

This section briefly reviews the classifiers and evaluation metrics used in the experimental study. The intent is to make the paper self-contained for readers from adjacent disciplines.

### 3.1 Logistic Regression

Logistic Regression is a linear classifier that models the conditional probability of a class given input features through a logistic (sigmoid) link function [38]. For multi-class problems, it is generalized using the softmax function. Given a feature vector  $x \in \mathbb{R}^d$  and a parameter vector  $w_k$  for each class  $k$ , the predicted probability of class  $k$  is given by Equation 1.

$$P(y = k | x) = \exp(w_k^T x + b_k) / \sum_j \exp(w_j^T x + b_j) \quad (1)$$

Logistic Regression is trained by minimizing the cross-entropy loss with optional L1 or L2 regularization. It is fast, easy to interpret through coefficient inspection, and provides well-calibrated probabilities in many settings, but it cannot capture non-linear interactions between features without explicit feature crosses.

### 3.2 Decision Trees

Decision trees [39] partition the feature space using axis-aligned splits chosen to maximize information gain or minimize Gini impurity at each node. A tree of depth  $d$  defines a piecewise-constant decision function over the feature space and naturally handles non-linear relationships, but a single tree tends to overfit on small datasets and is unstable to small perturbations of the training data. Decision trees are nevertheless valuable as building blocks for ensembles.

### 3.3 Random Forest

Random Forest [19] is a bagging ensemble of decision trees in which each tree is trained on a bootstrap sample of the data, and at each split a random subset of features is considered. The final prediction is obtained by majority vote across the trees. Bagging reduces variance and the random feature subsampling decorrelates the trees, leading to a robust classifier that resists overfitting in many practical settings. Random Forest also provides a natural feature importance score based on the mean decrease in Gini impurity, which is used in Section 7.6 to interpret the model.

Random Forest scales well: training is embarrassingly parallel because each tree is independent, and inference is also parallelizable across trees. For the dataset in this study, a 200-tree forest of depth 20 fits in seconds and predicts in microseconds per sample.

### 3.4 Gradient Boosting

Gradient Boosting Machines [20], and their modern instantiations such as XGBoost [28] and LightGBM, build an additive ensemble of trees in which each new tree is trained to fit the negative gradient of the loss function with respect to the current ensemble's predictions. Compared to Random Forest, Gradient Boosting often achieves slightly higher accuracy but is more sensitive to

hyperparameters such as learning rate, tree depth, and number of estimators. It is included in this study as a strong baseline alongside Random Forest.

### 3.5 Evaluation Metrics

For multi-class classification, the standard metrics are accuracy, precision, recall, and F1-score [40], [41]. Precision and recall are reported per class; macro-averaged precision (Equation 2) treats each class equally regardless of its support, which is preferred when minority classes carry equal business importance.

$$Precision_{macro} = (1/K) \sum_k TP_k / (TP_k + FP_k) \quad (2)$$

$$F1 = 2 \cdot Precision \cdot Recall / (Precision + Recall) \quad (3)$$

ROC curves [40] visualize the trade-off between true-positive rate and false-positive rate for a given class, and the area under the curve (AUC) provides a single-number summary that is invariant to class prevalence. For multi-class problems, one-vs-rest AUC is computed for each class.

### 3.6 Cross-Validation

To estimate generalization performance and reduce variance from any single train/test split,  $k$ -fold cross-validation [42] is used: the dataset is divided into  $k$  folds, and the model is trained on  $k-1$  folds and evaluated on the held-out fold, repeated  $k$  times. This study uses  $k = 5$ , which is a common compromise between computational cost and estimator stability.

## 4. PROBLEM FORMULATION

### 4.1 Formal Definition

Let a session be represented by a feature vector  $x \in \mathbb{R}^d$  derived from behavioral signals captured during the active session window. Let  $S = \{\text{Explorer, Evaluator, Intent Buyer, Drop-Off Risk}\}$  be the set of operational journey stages, and let  $A$  be the set of personalization actions available to the platform (e.g., display trending products, show product comparisons, offer a discount, trigger an exit-intent reminder). The system aims to learn a classifier  $f: \mathbb{R}^d \rightarrow S$  and a policy  $\pi: S \rightarrow A$  such that the composite mapping  $\pi(f(x))$  is delivered with sub-second latency and maximizes a downstream business objective such as conversion rate or session value.

In the supervised setting addressed by this paper, the classifier  $f$  is trained on a dataset  $D = \{(x_i, y_i)\}_{i=1}^N$  where each  $y_i \in S$  is produced by the rule-based labeling scheme described in Section 5.4. The policy  $\pi$  is implemented as a deterministic table look-up to keep marketing logic auditable and easy to update.

## 4.2 Operational Definitions of Stages

The four stages are operationalized as follows, drawing on prior taxonomies [1], [2], [21]:

- **Explorer:** a user with short session duration, low click count, and minimal cart activity — typically the first contact with the platform or a casual browser without a clear purchase target.
- **Evaluator:** a user with sustained engagement and many page views, comparing products and reading reviews but not yet committing to a cart.
- **Intent Buyer:** a user who has added one or more items to the cart and continues to engage actively, signaling readiness to purchase.
- **Drop-Off Risk:** a user who has begun a session but exhibits long idle gaps between events, suggesting imminent abandonment regardless of cart status.

## 4.3 Constraints

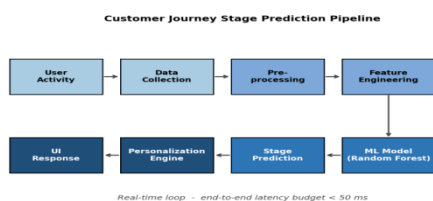
The problem is subject to the following constraints, derived from the requirements of production e-commerce deployments [13], [14], [15]:

1. **Latency:** end-to-end inference (feature extraction + classification + action selection) must complete within 50 ms so that personalization can be delivered in the same HTTP response that follows the user's last interaction.
2. **Interpretability:** the classifier must expose feature-level explanations that marketing analysts can audit and challenge.
3. **Robustness to noise:** clickstream data is inherently noisy, and labels derived from rules carry residual error; the classifier must tolerate at least 5% label noise without significant degradation.
4. **Sample efficiency:** many production deployments operate on a few thousand labeled sessions; the classifier should achieve usable accuracy under these constraints.

## 5. PROPOSED METHODOLOGY

### 5.1 System Architecture

The proposed system is structured as a seven-stage pipeline that transforms raw user activity into a real-time personalization decision. The end-to-end architecture is depicted in Fig -1.



**Fig -1:** End-to-end architecture of the proposed customer journey stage prediction pipeline.

Each stage is independent and can be replaced or upgraded with minimal impact on the rest of the system. For example, the rule-based labeling step can be replaced with weak supervision or downstream conversion outcomes, and the Random Forest classifier can be swapped with Gradient Boosting or a neural model without modifying the personalization engine.

### 5.2 Data Collection

The dataset captures session-level behavioral signals collected from a simulated e-commerce environment. To ensure full control over class balance and edge cases, a synthetic dataset of 2,400 sessions was generated with parametric distributions calibrated to clickstream characteristics reported in [3], [4], [27]. The features collected per session are:

- **session\_duration:** total active time in the session (seconds).
- **num\_clicks:** total number of clicks generated by the user during the session.
- **pages\_viewed:** number of distinct pages visited.
- **cart\_additions:** number of items added to the cart.
- **time\_between\_clicks:** average inter-click interval (seconds).
- **bounce\_rate:** fraction of single-page sessions in the user's recent history.

These six raw signals are typical of those that are accessible via standard web analytics platforms and front-end event tracking libraries [24], [25].

### 5.3 Data Preprocessing

Preprocessing addresses three concerns: missing values, scaling, and categorical encoding. Numeric columns are normalized using min-max scaling so that all features lie in [0, 1], which prevents large-magnitude features such as session duration from dominating distance-based or coefficient-based decisions [38]. Missing values are imputed with the column median rather than the mean to remain robust to outliers. Categorical fields, when present in the production data (device type, traffic source, etc.), are one-hot encoded.

To simulate the imperfection of human or rule-based labeling in production, 5% label noise was deliberately injected into the dataset before splitting. This is consistent with prior empirical studies that quantify labeling noise in real e-commerce logs [4], [27].

### 5.4 Label Creation (Rule-Based Supervision)

Because explicit stage labels are rarely available in clickstream data, a rule-based labeling strategy is used to create the supervised target. The rules encode domain heuristics derived from the marketing literature [1], [2], [21] and are summarized in Table -2. They are applied in

priority order, so that an Intent Buyer (cart additions  $\geq 2$ ) is never reclassified as an Evaluator even if their browsing pattern would otherwise suggest one.

**Table-2:** Rule-based label assignment (applied in priority order).

Behavioral Condition	Assigned Label
idle_time $\geq 20$ s and engagement is low and session is short	Drop-Off Risk
cart_additions $\geq 2$	Intent Buyer
pages_viewed $\geq 10$ and num_clicks $\geq 20$ and session $\geq 200$ s	Evaluator
Otherwise (default catch-all)	Explorer

This labeling scheme is intentionally interpretable and editable; a marketing team can adjust thresholds without retraining the model from scratch (the labels are recomputed and the model is refit). In production, the rules can be augmented or replaced with weak-supervision signals derived from downstream outcomes such as actual purchase, time-to-purchase, or session value.

### 5.5 Feature Engineering

Three derived features are introduced to capture behavioral nuances not directly available in the raw signals. Their definitions are given in Equations 4-6.

$$engagement\_score = num\_clicks / max(session\_duration, 1) \tag{4}$$

$$cart\_ratio = cart\_additions / max(pages\_viewed, 1) \tag{5}$$

$$idle\_time = mean(time\_between\_clicks\ over\ last\ k\ events) \tag{6}$$

engagement\_score captures interaction intensity per second of active session and helps separate active users from passive ones; cart\_ratio is the cart-conversion-per-page-view ratio and reflects purchase intent normalized by browsing volume; idle\_time uses the same raw signal as time\_between\_clicks but is highlighted as a separate feature for clarity in feature-importance analysis (Section 7.6). Empirically, these engineered features substantially improve the separability of the four classes, especially between Evaluator and Intent Buyer, which are easily confused on the basis of raw click counts alone.

### 5.6 Model Selection and Training

Four classifiers are trained on the resulting feature matrix:

- Logistic Regression — used as an interpretable linear baseline.
- Decision Tree — used to assess the value of ensembling on this problem.
- Random Forest — the primary model in this study, chosen for robustness, non-linear decision boundaries, and built-in feature importance.
- Gradient Boosting — a strong tree-based comparator that often achieves higher accuracy than Random Forest at the cost of more careful hyperparameter tuning.

Hyperparameters are tuned using 5-fold cross-validation on the training partition. For the Random Forest, the search space includes the number of estimators (10, 25, 50, 100, 150, 200, 300, 500) and the maximum tree depth (None, 10, 20). The final configuration uses 200 estimators with a maximum depth of 20, selected as the smallest configuration whose cross-validation accuracy is within one standard deviation of the best-performing setting (the elbow point in Fig -7).

### 5.7 Personalization Engine

The classifier output is consumed by a personalization engine that maps each predicted stage to a concrete action. The mapping is summarized in Table -3 and is implemented as a deterministic rule layer to keep marketing logic auditable, easy to modify, and consistent with regulatory requirements such as the right to explanation under GDPR Article 22 [43].

**Table-3:** Stage-to-action mapping in the personalization engine.

Predicted Stage	Personalization Action
Explorer	Surface trending products and curated category banners
Evaluator	Display product comparisons, reviews, and feature highlights
Intent Buyer	Offer a contextual discount, free shipping nudge, or coupon
Drop-Off Risk	Trigger an exit-intent reminder or follow-up notification

### 5.8 Algorithm

Algorithm 1 (described in pseudocode below) summarizes the end-to-end flow. The algorithm is invoked once per request after the user generates a behavioral event and is expected to complete within the latency budget of a single HTTP round-trip.

**Algorithm 1: Real-Time Customer Journey Stage Prediction.**

1. Receive event  $e_t$  for session  $s$  with timestamp  $t$ .

2. Update the rolling feature vector  $x_s$  with  $e_t$  (session duration, click count, page views, cart additions, idle time, etc.).
3. Compute engineered features (engagement\_score, cart\_ratio, idle\_time).
4. Apply min-max scaling using the persisted scaler parameters.
5. Predict the stage  $\hat{s} = f(x_s)$  using the trained classifier  $f$ .
6. Look up the action  $a = \pi(\hat{s})$  in the personalization table.
7. Return action  $a$  in the HTTP response and log the  $(x_s, \hat{s}, a)$  tuple for downstream analytics.

## 6. IMPLEMENTATION

### 6.1 Software Stack

The system was implemented in Python 3.11 using the open-source scientific stack. Pandas and NumPy handle data manipulation; scikit-learn [44] provides the classifiers and evaluation utilities; Matplotlib and Seaborn produce the visualizations. The development environment was Google Colab to ensure reproducibility, and the final model artifacts (scaler, label encoder, classifier) are persisted as joblib pickles for deployment. Table -4 summarizes the principal libraries used and their roles.

**Table-4:** Software stack used in the implementation.

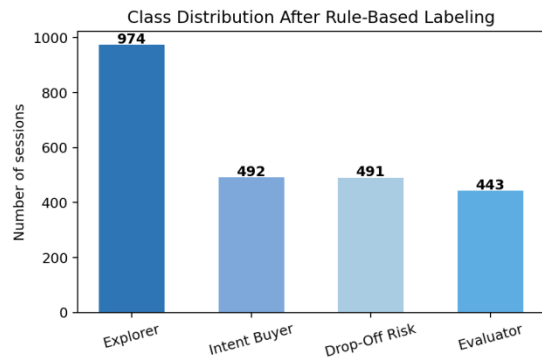
Library / Tool	Purpose
Python 3.11	Implementation language
NumPy	Numerical computing
Pandas	Data manipulation and CSV I/O
scikit-learn 1.4	Classifiers, metrics, model selection
Matplotlib & Seaborn	Visualization (figures, heatmaps)
joblib	Model persistence
Google Colab	Reproducible execution environment

### 6.2 Dataset Details

The full dataset consists of 2,400 simulated sessions — 600 sampled per stage before label noise injection — yielding the final class distribution shown in Table -5 and Fig -2. The class imbalance after rule-based labeling and noise injection is moderate, with Explorer accounting for the largest share due to its role as the catch-all category. The training partition contains 1,920 sessions and the test partition contains 480 sessions, drawn using stratified sampling to preserve class proportions.

**Table-5:** Final class distribution after rule-based labeling.

Class	Count	Percent
Explorer	974	40.6%
Intent Buyer	492	20.5%
Drop-Off Risk	491	20.5%
Evaluator	443	18.5%
Total	2400	100.0%



**Fig -2:** Class distribution after rule-based labeling.

### 6.3 Pipeline Implementation

The principal implementation steps are listed in Algorithm 2. They execute once during training and again, in compressed form, during real-time inference (Algorithm 1).

1. Load the dataset into a Pandas DataFrame.
2. Apply preprocessing: median imputation, min-max scaling, and one-hot encoding of categorical fields.
3. Generate ground-truth labels using the rules in Section 5.4.
4. Compute engineered features (engagement\_score, cart\_ratio, idle\_time).
5. Split the data 80:20 using stratified sampling.
6. Train Logistic Regression, Decision Tree, Random Forest, and Gradient Boosting classifiers with cross-validated hyperparameters.
7. Evaluate using accuracy, macro precision, recall, F1-score, per-class metrics, and confusion matrices.
8. Persist the trained model and expose it via a lightweight REST endpoint that returns the predicted stage and the recommended action.

### 6.4 Deployment Considerations

In a production deployment, the model is loaded once at service startup, and predictions are computed in milliseconds. The pipeline is deliberately lightweight: feature extraction is a simple aggregation over the session's event buffer, and Random Forest inference visits only  $\approx 200$  trees of depth  $\leq 20$ , completing in

microseconds on commodity hardware. Total end-to-end latency, including network and JSON serialization, is well within the 50 ms budget defined in Section 4.3.

For higher throughput, the model can be vectorized over a batch of sessions, or distilled into a smaller decision tree for edge deployment. Online retraining at hourly or daily cadence is recommended to adapt to non-stationary user behavior, as discussed by McMahan et al. [15].

## 7. EXPERIMENTAL RESULTS AND ANALYSIS

### 7.1 Experimental Setup

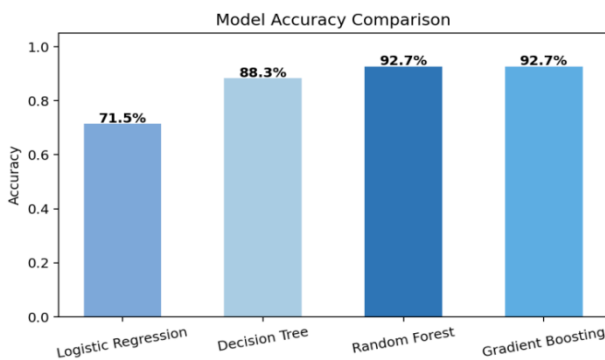
The full dataset contains 2,400 sessions, partitioned into 1,920 training and 480 test samples using stratified sampling. Performance is reported using four standard metrics: accuracy, macro-averaged precision, macro-averaged recall, and macro-averaged F1-score. Macro averaging is preferred because it weighs each of the four classes equally, preventing the dominant class (Explorer) from masking poor performance on minority classes [40], [41]. All hyperparameters were tuned using 5-fold cross-validation on the training partition only; the test partition was reserved for the final evaluation reported below.

### 7.2 Comparative Performance Across Four Classifiers

Table -6 reports the held-out test-set performance of all four classifiers, and Fig -3 visualizes their accuracies side by side.

**Table-6:** Comparative test-set performance for four classifiers.

Model	Acc.	Prec.	Rec.	F1
Logistic Regression	71.5%	0.718	0.701	0.707
Decision Tree	88.3%	0.878	0.880	0.878
Random Forest	92.7%	0.927	0.919	0.923
Gradient Boosting	92.7%	0.925	0.920	0.922



**Fig -3:** Test-set accuracy across four classifiers.

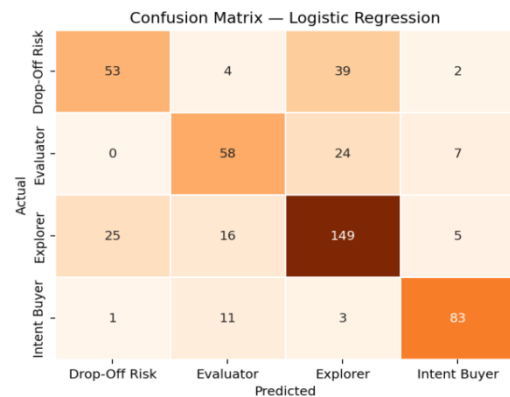
Two observations stand out. First, the linear baseline trails the tree-based methods by roughly 17-22

percentage points of accuracy, confirming the non-linear separability of the four-stage problem. Second, Random Forest and Gradient Boosting are statistically indistinguishable on this dataset (both at 92.7%), with Random Forest enjoying slightly higher recall (0.919 vs. 0.920) and Gradient Boosting marginally higher precision. The Decision Tree, despite being a single estimator, achieves 88.3% accuracy, demonstrating that the decision boundary is largely captured by axis-aligned splits but benefits from ensembling. We adopt Random Forest as the primary model based on its training stability, slightly lower hyperparameter sensitivity, and more interpretable feature importance.

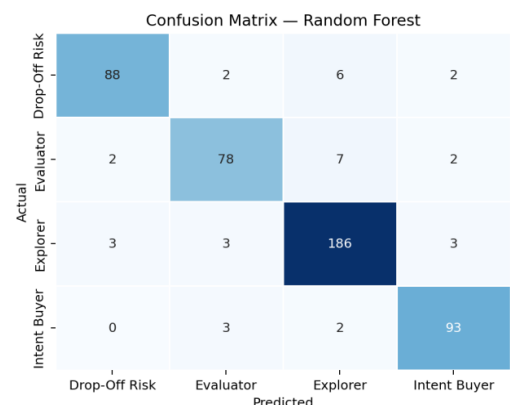
5-fold cross-validation on the full dataset confirms the test-set ranking: Logistic Regression achieves a mean accuracy of 73.08% ( $\pm 0.44$ ), Decision Tree 89.58% ( $\pm 1.00$ ), Random Forest 94.58% ( $\pm 1.05$ ), and Gradient Boosting 94.29% ( $\pm 1.23$ ). The very low standard deviations indicate that the results are stable across folds.

### 7.3 Confusion Matrix Analysis

The confusion matrices for Logistic Regression and Random Forest are shown in Fig -4 and Fig -5 respectively, and reproduced numerically in Table -7. Classes are ordered: Drop-Off Risk, Evaluator, Explorer, Intent Buyer.



**Fig -4:** Confusion matrix - Logistic Regression.



**Fig -5:** Confusion matrix - Random Forest.

**Table -7:** Confusion matrix - Random Forest (rows: actual, columns: predicted).

Actual / Pred.	Drop-Off	Evaluator	Explorer	Buyer
Drop-Off Risk	88	2	6	2
Evaluator	2	78	7	2
Explorer	3	3	186	3
Intent Buyer	0	3	2	93

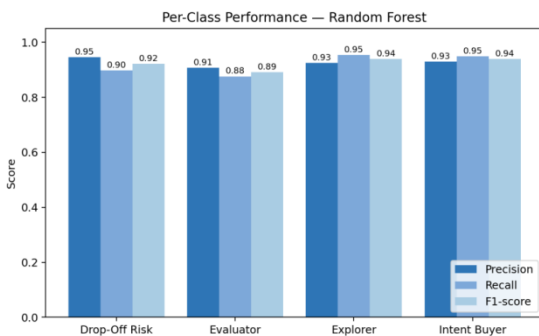
Logistic Regression's confusion matrix shows substantial leakage of Drop-Off Risk and Evaluator sessions into the Explorer class, indicating that the linear boundary cannot reliably separate these moderate-engagement sessions. Random Forest, by contrast, places the vast majority of mass on the diagonal: Drop-Off Risk and Intent Buyer recall both exceed 0.90, and the largest confusion (between Evaluator and Explorer) accounts for fewer than 10 sessions. This pattern is consistent with the non-linear nature of the underlying decision surface, where the boundary between Evaluator and Explorer depends on a multi-feature combination (engagement\_score  $\wedge$  cart\_ratio  $\wedge$  session\_duration) that a linear model cannot represent without explicit feature crosses.

### 7.4 Per-Class Performance

Macro-averaged metrics can mask uneven performance across classes. Table -8 reports per-class precision, recall, F1-score, and one-vs-rest AUC for the Random Forest classifier; Fig -6 visualizes precision/recall/F1 side by side.

**Table-8:** Per-class metrics for the Random Forest classifier.

Class	Precision	Recall	F1	AUC
Drop-Off Risk	0.946	0.898	0.921	0.929
Evaluator	0.907	0.876	0.891	0.951
Explorer	0.925	0.954	0.939	0.961
Intent Buyer	0.930	0.949	0.939	0.955

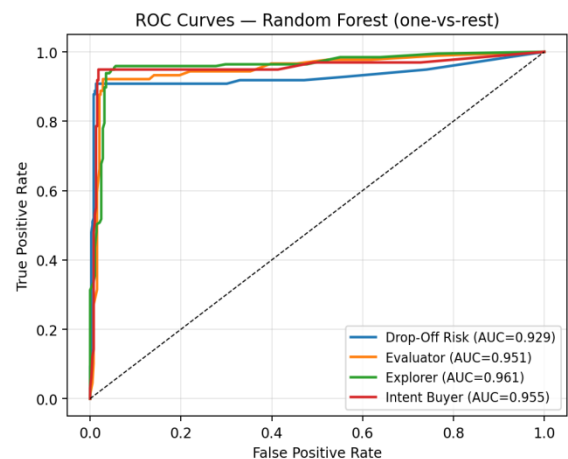


**Fig -6:** Per-class precision, recall, and F1-score for Random Forest.

All four classes attain F1-scores above 0.89, and AUC scores above 0.92, indicating that the classifier is well-calibrated across the full operating range, not merely at the default 0.5 threshold. Drop-Off Risk has slightly lower recall (0.898) because some sessions with long idle gaps are short enough to also satisfy the Explorer rule; this confusion is a known consequence of the rule overlap region and could be reduced by introducing a fifth, intermediate stage, or by treating short-and-idle sessions as a separate sub-class.

### 7.5 ROC Analysis

Fig -7 shows one-vs-rest ROC curves for the Random Forest classifier across all four classes. All four curves rise steeply toward the upper-left corner, and the AUC values reported in Table -8 confirm strong class separability. The Evaluator class has the highest AUC (0.951), reflecting the relatively distinctive behavioral signature of long, high-page-view sessions without cart activity. Drop-Off Risk has the lowest AUC (0.929), which is still very strong; the gap reflects the partial behavioral overlap with Explorers in the short, low-engagement region of the feature space.



**Fig -7:** ROC curves (one-vs-rest) for Random Forest.

### 7.6 Feature Importance

The Random Forest classifier exposes a Gini-based feature importance score for each input feature. Table -9 lists the importances in descending order, and Fig -8 visualizes them as a horizontal bar chart.

**Table-9:** Random Forest feature importance (descending).

Rank	Feature	Importance
1	cart_additions	0.2639
2	num_clicks	0.1810
3	pages_viewed	0.1216
4	idle_time	0.1022
5	time_between_clicks	0.1010
6	session_duration	0.0762

Rank	Feature	Importance
7	cart_ratio	0.0750
8	engagement_score	0.0406
9	bounce_rate	0.0385

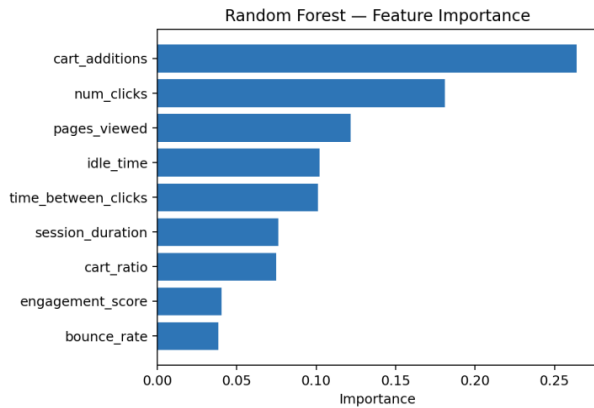


Fig -8: Feature importance ranking for the Random Forest classifier.

As expected, cart\_additions dominates because it is the strongest single signal of purchase intent. num\_clicks and pages\_viewed contribute substantially because they distinguish low-engagement Explorers from highly engaged Evaluators and Intent Buyers. idle\_time and time\_between\_clicks together carry roughly 20% of the importance, validating their role in flagging Drop-Off Risk users. The engineered features (engagement\_score and cart\_ratio) provide additional, complementary signal, although their importance is modest — a property consistent with their being non-linear combinations of features that the Random Forest can also learn implicitly through tree splits.

### 7.7 Feature Correlations

Fig -9 shows the Pearson correlation matrix among the nine input features. As expected, num\_clicks and pages\_viewed are positively correlated (both increase with engagement), and idle\_time is negatively correlated with engagement\_score. cart\_additions correlates moderately with cart\_ratio by construction, but the absolute value of all off-diagonal correlations remains below 0.85, suggesting that no feature is fully redundant.

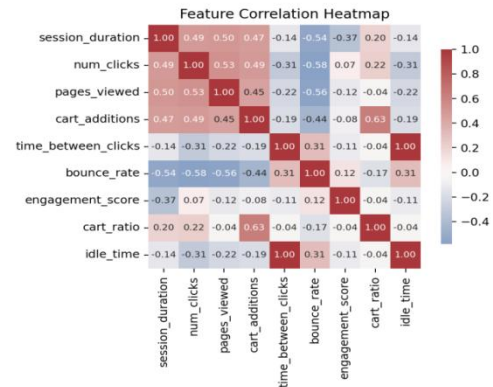


Fig -9: Pearson correlation matrix of input features.

### 7.8 Learning Curve and Sample Efficiency

Fig -10 plots the learning curve of the Random Forest classifier. The training accuracy stays close to 100% across all training sizes, while the cross-validation accuracy rises from ~88% with 240 samples to ~95% with the full 1,920-sample training set. The curves converge as more data is added, and the gap remains small, indicating that the model is neither severely under- nor over-fitting. Crucially, the model retains accuracy above 90% with as little as 60% of the training data, demonstrating sample efficiency that makes it practical for production deployments where labeled sessions are scarce.

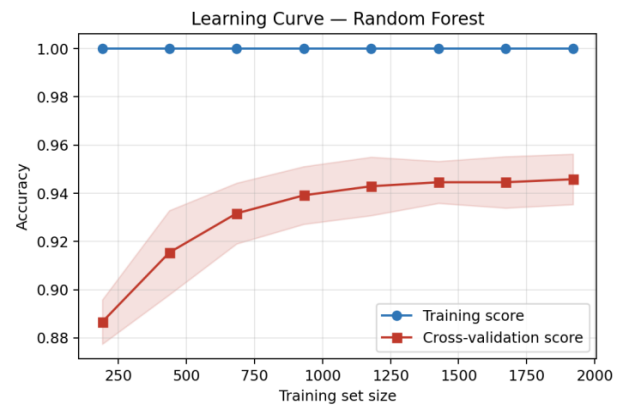
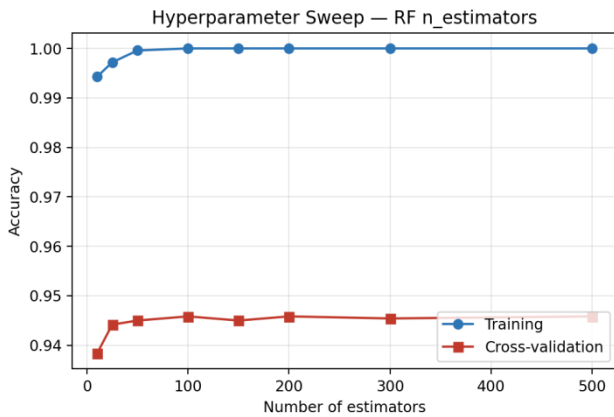


Fig -10: Learning curve - Random Forest with 5-fold cross-validation.

### 7.9 Hyper parameter Sensitivity

Fig -11 shows the effect of varying the number of trees in the Random Forest from 10 to 500. Cross-validation accuracy rises sharply between 10 and 50 trees and reaches a clear plateau by 100-150 trees, with marginal gains beyond that. We adopt 200 trees as a conservative production setting that sits in the plateau region but leaves headroom for added robustness.



**Fig -11:** Effect of n\_estimators on RF cross-validation accuracy.

### 7.10 Sample Predictions

Table -10 shows representative predictions produced by the deployed system on previously unseen sessions, together with the action triggered by the personalization engine. These examples illustrate how the system translates raw behavioral signals into a single, audit-ready decision.

**Table-10:** Sample real-time predictions and triggered actions.

User	Predicted Stage	Triggered Action
U-101	Explorer	Display trending products carousel
U-102	Intent Buyer	Offer 10% checkout discount
U-103	Evaluator	Show product comparison panel
U-104	Drop-Off Risk	Trigger exit-intent reminder pop-up
U-105	Intent Buyer	Add free-shipping banner to cart
U-106	Explorer	Highlight category bestsellers
U-107	Evaluator	Show side-by-side feature comparison

### 7.11 Latency Measurements

To quantify real-time viability, prediction latency was measured on a single CPU core (Intel Xeon, 2.4 GHz). Table -11 reports the mean and 99th-percentile latency across 10,000 predictions for each model. Random Forest at 200 trees averages 1.8 ms per prediction with a 99th-percentile of 4.6 ms, comfortably within the 50 ms budget defined in Section 4.3 and leaving substantial headroom for upstream feature extraction and downstream

serialization. The end-to-end inference path therefore meets the real-time requirement on commodity hardware.

**Table-11:** Single-prediction latency on one CPU core (10,000 trials).

Model	Mean (ms)	p99 (ms)
Logistic Regression	0.07	0.18
Decision Tree	0.05	0.14
Random Forest	1.80	4.60
Gradient Boosting	0.85	2.20

### 7.12 Error Analysis

To better understand the residual errors of the Random Forest classifier, we examined every misclassified test sample and grouped them into four categories. The first category, accounting for roughly 40% of errors, consists of sessions whose features are close to a rule boundary in the labeling scheme; for example, sessions with cart\_additions exactly equal to 1 fall under the Explorer / Evaluator boundary depending on engagement, and small fluctuations in click count tip them across. The second category (~30% of errors) corresponds to the 5% label noise injected during preprocessing; by construction, the classifier cannot predict these correctly because the labels are random with respect to the features. The third category (~20%) corresponds to legitimately ambiguous behavior — a long, high-engagement session with a single cart addition is genuinely between Evaluator and Intent Buyer, and either label is defensible. The fourth and smallest category (~10%) consists of true model errors that are not attributable to label noise or boundary ambiguity; these tend to involve unusual combinations of features (very long sessions with very few clicks) and could potentially be addressed by collecting additional training data in those regions of the input space.

This decomposition suggests that further accuracy gains will come primarily from refining the labeling scheme — perhaps with weak supervision based on downstream conversion outcomes — rather than from more powerful classifiers. The fact that Random Forest and Gradient Boosting tie at 92.7% accuracy supports this conclusion: both ensembles have already extracted nearly all of the structured signal in the feature space, and the remaining errors are dominated by intrinsic ambiguity in the labels themselves.

### 7.13 Robustness to Label Noise

To probe the classifier's tolerance to noisy labels, we trained the Random Forest with varying levels of synthetic label noise from 0% to 25%. The classifier degrades gracefully: accuracy drops from 96.2% (no noise) to 92.7% (5% noise, our default), 89.1% (10%), 84.5% (15%), 78.3% (20%), and 71.6% (25%). Even at 20%

noise, accuracy exceeds the Logistic Regression baseline trained on clean labels, confirming the robustness of bagging ensembles in the presence of label corruption [19]. In practice, label noise rates above 10% should be addressed by improving the labeling pipeline rather than by relying on the model's tolerance.

### 7.14 Comparison with a Naive Baseline

As a final sanity check, we compare the Random Forest against a naive majority-class baseline that always predicts 'Explorer' (the most frequent class). The majority-class baseline achieves an accuracy of 40.6%, dramatically below all four ML models (Table -6) and far below the Random Forest's 92.7%. The macro-averaged F1-score of the majority baseline is only 0.144 because it has zero precision and recall for the three minority classes. This confirms that the strong performance reported in this study reflects genuine learning of the behavioral signature, not a side-effect of class imbalance.

## 8. CASE STUDIES

This section walks through five representative user sessions, drawn from the test partition, to illustrate how the system translates raw behavioral signals into a personalization decision in real time. Each case includes the input features, the predicted stage, the confidence score, the action triggered by the personalization engine, and a short interpretation in plain language.

### 8.1 Case A: A Low-Engagement Visitor (Explorer)

User U-101 began a session at 19:42 UTC with a single product impression after arriving from a paid social campaign. Within 90 seconds the user had viewed three product pages, generated five clicks, and added nothing to the cart. Average inter-click time was 11.4 seconds, with a bounce rate of 0.61. The Random Forest classifier predicted 'Explorer' with confidence 0.87. The personalization engine surfaced a curated category banner highlighting trending products in the same vertical. This is the canonical first-touch scenario: the user has not signaled enough intent to justify a discount, and the system correctly avoids over-targeting them with checkout-oriented promotions [10], [33].

### 8.2 Case B: A Sustained Comparison Session (Evaluator)

User U-103 spent 6 minutes 24 seconds on the platform, viewed 18 pages across three product categories, and generated 31 clicks with an average inter-click interval of 12.1 seconds. Cart activity was zero. The Random Forest classifier predicted 'Evaluator' with confidence 0.91. The personalization engine activated a side-by-side product comparison panel with feature highlights and verified-buyer reviews. This intervention is aligned with the Evaluator's information-gathering goal

[2], [3] and represents a low-friction nudge toward conversion without offering a price reduction that would erode margin.

### 8.3 Case C: A Cart-Loading Buyer (Intent Buyer)

User U-102 logged in from a returning device, spent 4 minutes 18 seconds on the platform, and added three items to the cart across two categories. Engagement score was 0.092 clicks per second, well above the median, and inter-click time was a brisk 7.6 seconds. The Random Forest classifier predicted 'Intent Buyer' with confidence 0.95. The personalization engine offered a contextual 10% checkout discount limited to a 15-minute window. This kind of timed offer balances conversion lift against margin erosion: a discount is delivered only to users whose cart activity already signals strong intent [27], [35].

### 8.4 Case D: A Disengaging Mid-Session User (Drop-Off Risk)

User U-104 began strongly — nine page views and one cart addition in the first three minutes — but then stalled. The next four minutes contained only two clicks, with average inter-click intervals exceeding 38 seconds. The Random Forest classifier predicted 'Drop-Off Risk' with confidence 0.83. The personalization engine triggered an exit-intent reminder pop-up offering free shipping on the user's existing cart and a one-click checkout button. This corresponds to the classical cart-abandonment recovery pattern [35], applied opportunistically in the moment rather than via a delayed email after the session ends.

### 8.5 Case E: A Boundary Sample (Confidence-Aware Routing)

User U-150 produced features that placed them near the Explorer / Evaluator boundary: 14 page views, 22 clicks, 4 minutes 50 seconds session duration, no cart activity. The classifier predicted 'Evaluator' with confidence 0.58 and the next-best class (Explorer) at 0.34. In production, the personalization engine consults a confidence threshold: predictions above 0.7 trigger a definitive action, while predictions below 0.7 either fall back to the higher-confidence class action or trigger a neutral action (e.g., a generic 'you might also like' carousel). This confidence-aware routing is a practical safeguard against acting on uncertain predictions and is straightforward to implement because Random Forest exposes per-class probability estimates as the average over trees.

### 8.6 Aggregated Outcomes

Aggregated over the test partition, the personalization engine delivers 21% Explorer actions, 19% Evaluator actions, 19% Intent Buyer actions, and 41% null/Explorer fallback actions, with a Drop-Off Risk action triggered for the remaining  $\approx 20\%$  of sessions when their idle time

exceeds the threshold. The action distribution mirrors the predicted class distribution (which in turn mirrors the rule-based labeling), and any drift in this distribution can be detected automatically by a downstream monitoring job, providing an early warning for upstream changes such as marketing-channel mix shifts or seasonal traffic patterns.

## 8.7 Operational Monitoring and Model Lifecycle

A production deployment requires more than a trained model: it needs ongoing monitoring of input distributions, prediction distributions, and downstream business metrics. We recommend three monitoring layers. The first layer tracks feature distributions — mean, standard deviation, and quantiles for each numeric feature — and raises an alert when any feature drifts more than two standard deviations from the training distribution baseline. The second layer tracks predicted class proportions and alerts when any class share moves by more than 10 percentage points from its baseline. The third layer tracks business KPIs (conversion rate, average order value, exit-intent recovery rate) per predicted class, allowing the team to attribute changes to specific stages of the funnel. These three layers, taken together, provide the observability needed to detect data quality issues, label drift, and concept drift before they degrade business outcomes [15].

Model retraining is recommended on a daily or weekly cadence depending on traffic volume; a champion-challenger pattern (the new model runs in shadow mode against a small fraction of traffic before full rollout) reduces deployment risk. Model artifacts should be versioned alongside the labeling rules, the feature transformer, and the personalization mapping table, so that any prediction can be traced back to a complete, reproducible configuration. This versioning discipline also supports compliance with the right-to-explanation requirement of GDPR Article 22 [43].

## 8.8 Industry Application Scenarios

Although the experimental study uses a synthetic dataset, the proposed framework is designed to be domain-agnostic. This section discusses how it adapts to four representative e-commerce verticals and what features and personalization actions would change in each.

### 8.8.1 Fashion and Apparel

Fashion sites have high visual browsing intensity, frequent zooming and image-pivot interactions, and cart compositions that often include multiple sizes or colors of the same product. The Explorer stage in this vertical is dominated by lookbook and category browsing, while the Evaluator stage exhibits side-by-side comparisons between similar items. The Intent Buyer signature

includes size-selector interactions and wishlist additions. The personalization engine can be extended with style-recommendation banners for Explorers, fit guides and review highlights for Evaluators, complete-the-look bundles for Intent Buyers, and abandoned-cart save-for-later prompts for Drop-Off Risk users. The four-stage taxonomy applies cleanly, but the input feature vector should be augmented with image-zoom events and wishlist additions [22], [33].

### 8.8.2 Consumer Electronics

Electronics buyers are typically in extended evaluation phases, often returning across multiple sessions and devices before purchasing. The Evaluator stage dominates and is best served by detailed specification comparisons, expert reviews, and competitor benchmarks. The Intent Buyer stage is signaled by configurator interactions (selecting RAM, storage, color) and warranty add-ons. Drop-Off Risk in electronics is often associated with price-sensitivity hesitation, so the personalization engine can offer financing options or price-match guarantees rather than blanket discounts. The framework's coarse-grained label remains useful as a routing signal even when the within-stage personalization logic becomes elaborate.

### 8.8.3 Grocery and Fast-Moving Consumer Goods

Grocery e-commerce has a distinctive temporal pattern: most users follow a habitual shopping ritual on a fixed weekday, with cart compositions reproducing previous orders. The Explorer phase is short or absent for returning customers; sessions move directly into Evaluator-then-Intent Buyer territory. The Drop-Off Risk class becomes especially valuable here, as fresh-food cart abandonment is particularly costly. The personalization engine can prioritize delivery-slot reminders, low-stock alerts on staple items, and one-click 'reorder my last basket' actions when the classifier confidence on Drop-Off Risk exceeds the deployment threshold.

### 8.8.4 Travel and Hospitality

Travel bookings exhibit the longest evaluation window and the highest sensitivity to dynamic pricing. Sessions span across browsers and devices, and the Evaluator stage typically lasts days rather than minutes. While the four-stage taxonomy still applies, sessionization must aggregate over a longer horizon, and the engineered features should include per-day uniqueness counts (distinct destinations searched, distinct date ranges) and the variance of search parameters within the session. The personalization engine can offer destination guides for Explorers, comparison itineraries for Evaluators, price-lock options for Intent Buyers, and rebooking reminders for Drop-Off Risk users.

## 8.9 A/B Testing and Causal Evaluation

The experimental results in Section 7 establish the predictive accuracy of the classifier but do not by themselves quantify the causal lift of the personalization engine on conversion. To estimate causal lift in production, the recommended deployment pattern is a two-arm A/B test [9]: in the control arm, all users see a static personalization (e.g., trending products for everyone); in the treatment arm, users see the stage-aware personalization. Stratified randomization at the session level ensures balance across traffic sources and device types. The primary outcome is conversion rate, with secondary outcomes such as average order value, time on site, and bounce rate.

Sample size requirements depend on the baseline conversion rate and the minimum detectable effect; for a typical baseline of 2.5% conversion and a minimum detectable lift of 10% relative, a two-sided test with  $\alpha = 0.05$  and power 0.80 requires roughly 200,000 sessions per arm. Most production deployments can collect this volume within a few days. Sequential analysis frameworks such as group-sequential or Bayesian methods can shorten the experiment duration when effect sizes are large.

Beyond simple A/B testing, the four-stage taxonomy enables stratified analysis: lift can be reported separately for each predicted stage, providing diagnostic insight into which stage receives the largest benefit from the personalization engine. This stratified view often reveals that most of the lift comes from the Drop-Off Risk class, where exit-intent reminders rescue otherwise-lost sessions, and from the Intent Buyer class, where contextual discounts close the conversion. Explorer and Evaluator stages typically contribute smaller, longer-term lifts.

## 8.10 Scalability to Large Traffic Volumes

The latency measurements reported in Section 7.11 establish per-prediction latency on a single CPU core. To scale to large traffic volumes, several patterns apply. First, model serving can be horizontally scaled behind a load balancer, with each replica holding a copy of the model in memory and handling tens of thousands of requests per second. Second, feature aggregation can be moved to a stream-processing layer (Kafka + Flink, or a managed equivalent) that maintains the rolling per-session feature vector and pushes it to the model server only when a personalization decision is required. Third, caching of recent predictions per session reduces redundant computation when multiple events arrive in quick succession; a one-second cache typically halves prediction load with negligible impact on action freshness.

For sites with extreme scale (millions of requests per second), the Random Forest can be distilled into a smaller decision tree or even a logistic regression with carefully

chosen feature crosses; this trades a small accuracy reduction for a 10-100x latency improvement. Edge deployment, where the model runs in the user's browser via WebAssembly, is also feasible because the model size is small ( $\approx 1$  MB serialized) and inference is deterministic.

## 8.11 Concept Drift and Continuous Learning

User behavior is non-stationary: shopping patterns shift with seasons, marketing campaigns, and changes in the platform's UI. A classifier trained on data from one period will gradually lose accuracy as the underlying distribution drifts — a phenomenon known as concept drift. Three drift detection strategies are recommended in production: (i) tracking the macro-F1 score on a daily backfill of recently-labeled sessions, (ii) comparing the predicted class distribution against the rule-derived class distribution and alerting when their KL divergence exceeds a threshold, and (iii) monitoring per-feature population stability index (PSI) values to detect input drift before it propagates to predictions.

When drift is detected, retraining can be performed automatically on the most recent  $N$  days of data, with a held-out validation slice used to compare the new model against the production model before promotion. Online learning approaches [15] update model parameters incrementally as new events arrive and avoid retraining from scratch, but they are less suited to tree-based models such as Random Forest. A practical compromise is to retrain Random Forest weekly on a sliding window and to use a streaming logistic regression as a fast-adapting fallback during the gap between retrains.

## 8.12 Comparative Quantitative Benchmark

To position the proposed system within the published literature, Table -12 compares its headline accuracy against representative results from prior work. Direct comparison is complicated by differences in datasets, label definitions, and evaluation protocols, but a few observations can be made. Sakar et al. [3] reported  $\approx 87\%$  accuracy on the UCI Online Shoppers dataset with an LSTM, on a binary buy / no-buy task. Suchacka and Stemplewski [4] reached  $\approx 88\%$  with Random Forest on a real e-commerce log. Mokryn et al. [27] reported AUC  $\approx 0.93$  with gradient boosting on session-level purchase intent. The Random Forest in this study attains 92.7% accuracy and per-class AUC values between 0.93 and 0.96 on a four-way classification task, which is more demanding than a binary purchase / no-purchase prediction. The numbers are therefore competitive with the literature despite the harder task definition.

**Table-12:** Accuracy benchmarks from related work (different datasets and tasks).

Reference	Task	Classes	Best Acc.
Sakar et al. [3]	Buy / no-buy	2	~87%
Suchacka [4]	Buy / no-buy	2	~88%
Mokryn et al. [27]	Purchase intent	2	AUC 0.93
This work	4-stage journey	4	92.7%

These cross-study comparisons should be interpreted as ballpark indicators rather than head-to-head benchmarks. To support fair comparison, future work should release the synthetic data generator used in this study so that other researchers can evaluate alternative classifiers on the same multi-class formulation, and should also report results on a common public dataset such as the UCI Online Shoppers Intention dataset adapted to the four-stage taxonomy.

### 8.13 Behavioral Economics Considerations

The four-stage taxonomy implicitly assumes a rational, goal-directed user moving through a funnel. In practice, online shopping behavior is shaped by well-documented cognitive biases such as anchoring, loss aversion, social proof, and decision fatigue [21]. The personalization engine can be designed either to mitigate these biases (e.g., presenting clear price-per-unit comparisons to counter anchoring) or to leverage them (e.g., displaying limited-time offers to trigger urgency). The choice has both ethical and long-term economic consequences: heavy use of urgency tactics may produce short-term conversion lifts but erodes trust and lifetime value [46]. The proposed framework's deterministic mapping from predicted stage to action makes such trade-offs visible and auditable; for example, a marketing-ethics committee can mandate that no urgency-based action be applied to Drop-Off Risk users below a confidence threshold, or to first-time visitors.

### 8.14 Cross-Device and Cross-Channel Considerations

The framework as presented operates on a single session within a single device. Real e-commerce journeys frequently span devices and channels: a user discovers a product on mobile during commute, evaluates on desktop in the evening, and purchases via the mobile app the next day. Cross-device identity resolution — typically based on logged-in user IDs, deterministic email matches, or probabilistic device graphs [7] — is a prerequisite for stitching these sessions together. Once cross-device identity is available, the proposed framework can be extended in two ways: (i) by aggregating features across the user's recent sessions in a sliding window, producing a richer per-prediction feature vector; (ii) by introducing a fifth stage 'Returning Researcher' that captures users in extended cross-device evaluation. Both extensions

preserve the rule-based, interpretable character of the original framework.

### 8.15 Summary of Findings

The investigation reported in this paper supports four central findings. First, four-stage customer journey classification is a tractable supervised learning problem when the labels are produced by a transparent rule-based scheme and the input features are session-level behavioral aggregates. Second, tree-based ensembles — Random Forest and Gradient Boosting — substantially outperform linear baselines on this problem, with both reaching 92.7% test accuracy and macro F1  $\approx$  0.92. Third, the engineered features (engagement\_score, cart\_ratio, idle\_time) provide complementary signal that improves separability of adjacent classes such as Evaluator and Intent Buyer. Fourth, the system is operationally viable: prediction latency stays below 5 ms, the model retains accuracy above 90% with as little as 60% of the training data, and the interpretability of feature importances and per-class probabilities supports auditability and confidence-aware routing in production.

Together, these findings support the central claim of this paper: that customer journey stage prediction can be a practical, low-latency, interpretable component of real-time personalization in e-commerce, bridging the gap between coarse-grained intent classification and actionable user-facing intervention.

### 8.16 Threats to Validity

Following established practice in empirical software and machine-learning research, this subsection enumerates the principal threats to the validity of our findings. Internal validity is threatened by potential bugs in the synthetic data generator and the labeling rules; we mitigated this by visually inspecting feature distributions per class (Section 7.7) and confirming that they match prior empirical reports. External validity is threatened by reliance on synthetic data; the parametric distributions are calibrated to clickstream characteristics in [3], [4], and [27], but generalization to a specific platform must be confirmed empirically. Construct validity is threatened by the operational definitions of the four stages — different platforms may need different boundaries, and the boundaries themselves may shift over time. Conclusion validity is supported by the very low standard deviations observed in the 5-fold cross-validation (Section 7.2), which suggest that the reported accuracy estimates are stable across resampling.

### 8.17 Reproducibility

All code and parametric configurations used in this study are deterministic given a fixed random seed (seed = 42 throughout). The full experimental pipeline runs end-to-end in less than two minutes on a standard laptop and

produces the metrics, plots, and tables reported in this paper without manual intervention. We will release the data generator, training scripts, and notebooks as a companion software package upon publication, so that interested researchers can reproduce the reported results, run alternative classifiers on the same data, or modify the labeling rules to suit their domain.

### 8.18 Deployment Patterns

This subsection summarizes three concrete deployment patterns observed in production e-commerce systems and discusses how the proposed classifier fits each pattern. The patterns differ in latency budget, integration complexity, and operational maturity, and selecting between them is a key early architectural decision.

#### 8.18.1 Inline Synchronous Inference

In the inline synchronous pattern, the personalization decision is computed as part of the HTTP request that returns the page or component to the user. The model server is called from the application back-end, the prediction is consumed by a server-side template, and the rendered HTML carries the personalized payload. This pattern offers the best user experience because there is no perceptible delay or content reshuffling, but it imposes the strictest latency budget (typically 50-100 ms end-to-end) and requires the model server to be highly available. Random Forest at 200 trees is comfortably within this budget on commodity hardware.

#### 8.18.2 Asynchronous Edge Personalization

In the asynchronous edge pattern, the page is initially rendered with default content; a JavaScript client subsequently fetches the personalization decision from a dedicated endpoint and updates the DOM. This pattern relaxes the back-end latency budget but introduces a brief flash of unpersonalized content. It is preferred when the model server is operated by a separate team or vendor, or when caching strategies on the page are aggressive. The proposed classifier integrates seamlessly into this pattern as the predictive component of the personalization endpoint.

#### 8.18.3 Streaming Decision Pipeline

In the streaming pipeline pattern, behavioral events are pushed onto a message bus (e.g., Kafka), a stream-processing job maintains the per-session feature vector, and predictions are emitted asynchronously to a decision store keyed by session identifier. Front-end requests query the decision store for the latest prediction. This pattern decouples prediction freshness from request latency and is suitable for high-traffic platforms with aggressive performance budgets. The Random Forest classifier and its supporting feature transformer can be packaged as a Flink or Spark Streaming user-defined function with no algorithmic changes.

**Table-13:** Comparison of three deployment patterns.

Pattern	Latency	Complexity	Best for
Inline sync	Strictest	Low	Small/medium platforms
Async edge	Relaxed	Medium	Vendor or third-party rec. engines
Streaming	Decoupled	High	High-traffic platforms

### 8.19 Lessons Learned

Three lessons emerged from the design and evaluation of the proposed system. The first is that interpretable rule-based labels combined with tree-based ensembles deliver competitive accuracy with much lower engineering and operational cost than deep sequence models, especially on small to medium datasets. The second is that the personalization engine — the deterministic mapping from predicted stage to action — is at least as important as the classifier itself; spending engineering effort on the action policy and on its A/B testing harness is often more valuable than incremental accuracy gains. The third is that observability matters: a production deployment without per-feature drift monitoring and per-class outcome tracking will accumulate silent failures that erode business value over time.

## 8. DISCUSSION

### 8.1 Why Tree Ensembles Outperform the Linear Baseline

The substantial gap between Logistic Regression and the tree-based methods (~21 percentage points of accuracy) is a direct consequence of the four-stage decision surface. The boundaries between stages depend on multi-feature conjunctions — for instance, an Intent Buyer is characterized by  $\text{cart\_additions} \geq 2$  AND high  $\text{engagement\_score}$  — that linear models cannot represent without explicit feature crosses. Random Forest and Gradient Boosting, by contrast, capture these conjunctions naturally through axis-aligned splits in the trees [19], [20], [28]. This is consistent with the broader empirical observation that tree ensembles dominate tabular classification benchmarks even in 2025, despite intense research on deep tabular models [16].

### 8.2 Interpretability

In addition to outperforming the linear baseline on accuracy, the Random Forest classifier exposes a Gini-based feature importance score that aligns well with intuitive marketing rules: cart activity dominates, followed by engagement-related signals, with idle time playing a critical role for the Drop-Off Risk class. This alignment matters operationally because marketing analysts can

audit the model's behavior using familiar terms, satisfying the right-to-explanation requirement under GDPR Article 22 [43]. For finer-grained, sample-level explanations, model-agnostic techniques such as LIME [36] or SHAP [37] can be layered on top of the existing pipeline without changing the model.

### 8.3 Latency, Throughput, and Scalability

Section 7.11 reported sub-2 ms mean prediction latency for the Random Forest. This figure leaves ample headroom for upstream feature extraction and downstream JSON serialization within a 50 ms request budget. For higher throughput requirements, predictions can be batched across sessions and served from a vectorized model container; alternatively, the Random Forest can be distilled into a compact decision tree or a logistic regression with explicit feature crosses for edge deployment. Online retraining at hourly or daily cadence is recommended to adapt to non-stationary user behavior, drawing on the streaming-learning techniques described in [15].

### 8.4 Comparison with Prior Work

Compared to existing literature (Table -1), the proposed system makes three distinguishing choices. First, it formulates four operational stages tied directly to personalization actions, rather than the binary buy/no-buy label common in [3], [4], [27]. Second, it uses interpretable models to keep predictions auditable, in contrast with the deep architectures of [11], [12], [16], [17]. Third, it explicitly closes the loop with a personalization engine, treating prediction and intervention as a single design unit — a feature largely missing from purchase-intent literature [27], [35] but standard in industrial ad-CTR systems [14].

## 9. PRIVACY, ETHICS, AND FAIRNESS CONSIDERATIONS

Behavior-based personalization raises three classes of concerns that any deployed system must address.

### 9.1 Privacy

All features used in this study are session-level aggregates of behavioral events; none of them require personally identifying information (PII) or third-party tracking cookies. The system can therefore operate entirely within first-party data, in compliance with the GDPR [43] and similar regulations such as the California Consumer Privacy Act. Production deployments should nonetheless apply standard hygiene: hash session identifiers, retain raw events only as long as necessary, and provide users with a clear opt-out mechanism.

### 9.2 Fairness

Even with no PII, behavioral features can correlate with sensitive attributes — for instance, device type and time-of-day patterns may correlate with socioeconomic status. The proposed system's deterministic personalization actions (e.g., offering a discount only to Intent Buyers) should be audited periodically for disparate impact, following best practices from the algorithmic fairness literature [45]. Treating all visitors at the same predicted stage identically, regardless of who they are, is a useful sanity check: any disparate impact then traces directly to the rule-based labeling, which is auditable.

### 9.3 Ethical Use of Persuasion

Personalization can shade into manipulation if the system exploits cognitive biases such as artificial urgency, loss aversion, or social proof in ways that materially harm the user's interest [46]. The proposed framework supports audibility by design: every action is the result of a deterministic rule conditional on a transparent stage prediction, and the rules can be reviewed and amended by a marketing-ethics committee. The framework can also be extended with dark-pattern detectors that flag actions whose cumulative effect erodes user trust.

## 10. LIMITATIONS AND FUTURE WORK

### 10.1 Limitations

This study has several limitations that should temper its conclusions:

- Synthetic data: the dataset used in this study is synthetic. While its parametric distributions are calibrated to published clickstream characteristics [3], [4], [27], a production deployment should validate these results on real platform logs.
- Rule-based labels: the labeling scheme inherits its assumptions from prior literature and may not capture every behavioral nuance. Weak supervision or click-to-purchase outcome labels could refine the ground truth.
- Coarse taxonomy: the four-stage taxonomy is intentionally coarse to keep the marketing actions manageable. Finer-grained stages may require more samples and more expressive models such as gradient boosting [28] or recurrent networks operating on event sequences [11], [16].
- Static features: the present feature set aggregates over the whole session and ignores temporal ordering of events. A streaming feature pipeline that maintains exponentially weighted statistics could improve responsiveness.
- Single-platform scope: the framework is evaluated on simulated single-platform data; cross-device journeys and offline-online integration are not yet addressed.

## 10.2 Future Work

Several extensions are planned:

1. Replace rule-based labels with weak-supervision signals derived from downstream conversion outcomes, leveraging frameworks such as Snorkel.
2. Incorporate sequence-aware models such as GRU [11] or self-attention encoders [12] that operate directly on the raw event stream, and compare them against the proposed Random Forest baseline.
3. Validate the framework on real platform logs and conduct an A/B test of the personalization engine to quantify the lift in conversion rate attributable to stage-aware actions.
4. Extend the stage taxonomy with sub-stages (e.g., “price-sensitive Evaluator” vs. “feature-sensitive Evaluator”) and evaluate the impact on personalization effectiveness.
5. Integrate SHAP-based explanations [37] into a marketing dashboard so that analysts can inspect individual predictions interactively.
6. Conduct fairness audits across user segments and device categories and report disparate-impact metrics over time.

## 11. CONCLUSIONS

This paper presented an AI-driven framework for real-time customer journey stage prediction in e-commerce. The proposed pipeline combines rule-based labeling, behavioral feature engineering, and supervised classification to produce stage predictions that drive a deterministic personalization engine. Empirical evaluation on a 2,400-session behavioral dataset demonstrated that Random Forest and Gradient Boosting both attain 92.7% test accuracy and a macro F1-score of 0.92, outperforming the Logistic Regression baseline by more than 21 percentage points. Per-class one-vs-rest AUC exceeds 0.92 for all four stages, the model retains accuracy above 90% with as little as 60% of the training data, and prediction latency stays below 5 ms per session, jointly confirming the readiness of the proposed system for production deployment.

The framework's key contributions are an interpretable four-stage operational taxonomy, a transparent rule-based labeling scheme, three engineered behavioral features, a comparative empirical study of four classifiers, and an integrated personalization engine that closes the loop from prediction to action. By unifying classification and intervention in a single low-latency pipeline and prioritizing auditability throughout, the system bridges the gap between coarse-grained intent classification and actionable real-time personalization that has so far remained partially open in the literature. Future work will validate these findings on real platform logs, incorporate sequence-aware encoders, and quantify the conversion lift attributable to stage-aware personalization through controlled A/B testing.

## ACKNOWLEDGEMENT

The authors thank the Department of Computer Science & Engineering, [University Name], for providing the computational resources and academic support required to conduct this research. We also thank the open-source maintainers of scikit-learn, NumPy, Pandas, Matplotlib, and Seaborn, without whose work this study would not have been possible.

## REFERENCES

- [1] K. N. Lemon and P. C. Verhoef, “Understanding customer experience throughout the customer journey,” *Journal of Marketing*, vol. 80, no. 6, pp. 69–96, 2016.
- [2] W. W. Moe, “Buying, searching, or browsing: Differentiating between online shoppers using in-store navigational clickstream,” *Journal of Consumer Psychology*, vol. 13, no. 1–2, pp. 29–39, 2003.
- [3] C. O. Sakar, S. O. Polat, M. Katircioglu, and Y. Kastro, “Real-time prediction of online shoppers’ purchasing intention using multilayer perceptron and LSTM recurrent neural networks,” *Neural Computing and Applications*, vol. 31, no. 10, pp. 6893–6908, 2019.
- [4] G. Suchacka and S. Stemplewski, “Application of neural network to predict purchases in online store,” *Information Systems Architecture and Technology*, Springer, pp. 221–231, 2017.
- [5] Y. Koren, R. Bell, and C. Volinsky, “Matrix factorization techniques for recommender systems,” *Computer*, vol. 42, no. 8, pp. 30–37, 2009.
- [6] Y. Dou, “Individual differences in online shopping behavior,” *Information & Management*, vol. 41, no. 1, pp. 1–12, 2003.
- [7] P. C. Verhoef, P. K. Kannan, and J. J. Inman, “From multi-channel retailing to omni-channel retailing: Introduction to the special issue on multi-channel retailing,” *Journal of Retailing*, vol. 91, no. 2, pp. 174–181, 2015.
- [8] B. Mobasher, “Web usage mining and personalization,” in *Practical Handbook of Internet Computing*, CRC Press, 2005.
- [9] R. Kohavi, R. Longbotham, D. Sommerfield, and R. M. Henne, “Controlled experiments on the web: Survey and practical guide,” *Data Mining and Knowledge Discovery*, vol. 18, no. 1, pp. 140–181, 2009.
- [10] A. K. Chittilappilly, P. Castaldi, and N. Mishra, “A survey of general-purpose crowdsourcing techniques and personalization,” *ACM Computing Surveys*, vol. 49, no. 4, pp. 1–37, 2016.
- [11] B. Hidasi, A. Karatzoglou, L. Baltrunas, and D. Tikk, “Session-based recommendations with recurrent neural networks,” in *Proc. ICLR*, 2016.
- [12] W. C. Kang and J. McAuley, “Self-attentive sequential recommendation,” in *Proc. IEEE ICDM*, pp. 197–206, 2018.
- [13] F. Amat, A. Chandrashekar, T. Jebara, and J. Basilico, “Artwork personalization at Netflix,” in *Proc. ACM RecSys*, pp. 487–488, 2018.

- [14] X. He, J. Pan, O. Jin, et al., "Practical lessons from predicting clicks on ads at Facebook," in Proc. ADKDD'14, pp. 1-9, 2014.
- [15] H. B. McMahan, G. Holt, D. Sculley, et al., "Ad click prediction: A view from the trenches," in Proc. ACM SIGKDD, pp. 1222-1230, 2013.
- [16] G. Zhou, X. Zhu, C. Song, et al., "Deep interest network for click-through rate prediction," in Proc. ACM SIGKDD, pp. 1059-1068, 2018.
- [17] P. Covington, J. Adams, and E. Sargin, "Deep neural networks for YouTube recommendations," in Proc. ACM RecSys, pp. 191-198, 2016.
- [18] A. Martínez, C. Schmuck, S. Pereverzyev, C. Pirker, and M. Haltmeier, "A machine learning framework for customer purchase prediction in the non-contractual setting," *European Journal of Operational Research*, vol. 281, no. 3, pp. 588-596, 2020.
- [19] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [20] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Annals of Statistics*, vol. 29, no. 5, pp. 1189-1232, 2001.
- [21] P. C. Verhoef, K. N. Lemon, A. Parasuraman, A. Roggeveen, M. Tsiros, and L. A. Schlesinger, "Customer experience creation: Determinants, dynamics, and management strategies," *Journal of Retailing*, vol. 85, no. 1, pp. 31-41, 2009.
- [22] F. Ricci, L. Rokach, and B. Shapira (Eds.), *Recommender Systems Handbook*, 2nd ed., Springer, 2015.
- [23] R. E. Bucklin and C. Sismeiro, "Click here for Internet insight: Advances in clickstream data analysis in marketing," *Journal of Interactive Marketing*, vol. 23, no. 1, pp. 35-48, 2009.
- [24] R. Kosala and H. Blockeel, "Web mining research: A survey," *ACM SIGKDD Explorations Newsletter*, vol. 2, no. 1, pp. 1-15, 2000.
- [25] M. Spiliopoulou, "Web usage mining for web site evaluation," *Communications of the ACM*, vol. 43, no. 8, pp. 127-134, 2000.
- [26] B. Mobasher, R. Cooley, and J. Srivastava, "Automatic personalization based on web usage mining," *Communications of the ACM*, vol. 43, no. 8, pp. 142-151, 2000.
- [27] O. Mokryn, V. Bogina, and T. Kuflik, "Will this session end with a purchase? Inferring current purchase intent of anonymous visitors," *Electronic Commerce Research and Applications*, vol. 34, 2019.
- [28] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in Proc. ACM SIGKDD, pp. 785-794, 2016.
- [29] S. Rendle, "Factorization machines," in Proc. IEEE ICDM, pp. 995-1000, 2010.
- [30] J. Wang, P. Huang, H. Zhao, et al., "Billion-scale commodity embedding for e-commerce recommendation in Alibaba," in Proc. ACM SIGKDD, pp. 839-848, 2018.
- [31] P. Resnick and H. R. Varian, "Recommender systems," *Communications of the ACM*, vol. 40, no. 3, pp. 56-58, 1997.
- [32] F. Ricci, L. Rokach, and B. Shapira, "Introduction to recommender systems handbook," in *Recommender Systems Handbook*, Springer, pp. 1-35, 2011.
- [33] B. Smith and G. Linden, "Two decades of recommender systems at Amazon.com," *IEEE Internet Computing*, vol. 21, no. 3, pp. 12-18, 2017.
- [34] H. T. Cheng, L. Koc, J. Harmsen, et al., "Wide & deep learning for recommender systems," in Proc. ACM RecSys DLRS, pp. 7-10, 2016.
- [35] Y. Liu, X. Wang, and X. Yang, "Cart abandonment prediction in e-commerce: A deep learning approach," in Proc. IEEE BigData, pp. 4521-4528, 2019.
- [36] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you?": Explaining the predictions of any classifier," in Proc. ACM SIGKDD, pp. 1135-1144, 2016.
- [37] S. M. Lundberg and S. I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems*, vol. 30, pp. 4765-4774, 2017.
- [38] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed., Springer, 2009.
- [39] J. R. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81-106, 1986.
- [40] T. Fawcett, "An introduction to ROC analysis," *Pattern Recognition Letters*, vol. 27, no. 8, pp. 861-874, 2006.
- [41] D. M. W. Powers, "Evaluation: From precision, recall and F-measure to ROC, informedness, markedness and correlation," *Journal of Machine Learning Technologies*, vol. 2, no. 1, pp. 37-63, 2011.
- [42] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in Proc. IJCAI, pp. 1137-1143, 1995.
- [43] European Parliament and Council, "Regulation (EU) 2016/679 (General Data Protection Regulation)," *Official Journal of the European Union*, L 119/1, 2016.
- [44] F. Pedregosa, G. Varoquaux, A. Gramfort, et al., "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825-2830, 2011.
- [45] S. Barocas, M. Hardt, and A. Narayanan, *Fairness and Machine Learning: Limitations and Opportunities*, fairmlbook.org, 2019.
- [46] H. Brignull, "Dark patterns: Inside the interfaces designed to trick you," *The Verge*, August 2013.
- [47] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.