# Predicting User Purchase Intent: Deep Contextual Sequence Modeling for E-Commerce Platforms

Dr. Sadik Khan[1]

[1] *Dept. of CSE, IET, Bundelkhand University, Jhansi, India.*

| Article Info | Abstract: |
|---|---|
| | Online retail environments are constantly evolving, generating rich streams of user interaction data in the process. Predicting which users are likely to transition from casual browsing to making an actual purchase has become a critical factor in crafting targeted marketing, dynamic pricing, and personalized shopping experiences. This paper delves into the use of deep contextual sequence modeling to forecast purchase intent on e-commerce platforms. By drawing on state-of-the-art architectures such as Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, Gated Recurrent Units (GRUs), and Transformers—combined with attention mechanisms—our approach captures both short-term user actions and longer-term patterns across multiple sessions. We enhance these models by embedding contextual information at the product level, encoding product metadata and user history more effectively. After an in-depth review of existing literature, we introduce a hybrid system that leverages attention-based components along with product embeddings to achieve robust user purchase intent predictions. Our experiments, conducted on a large public dataset, underscore significant gains in accuracy and F1 scores when compared to simpler methods and other sequence-based baselines. A comparative analysis provides further evidence of the system's strengths, illustrated by tables and charts that detail performance benchmarks. Ultimately, this research demonstrates how integrating contextualized deep sequence models can yield impactful and precise predictions, helping online retailers better address user needs in real time.<br><br>*Keywords:* User Purchase Intent, Deep Learning, Contextual Sequence Modeling, E-Commerce, Recurrent Neural Networks, Transformers, Attention Mechanism, LSTM, Product Embeddings |

# 1. Introduction

E-commerce has reshaped how people explore and purchase products, creating a rapidly expanding environment where each click or search leaves behind data points that shed light on consumer interests. In this setting, accurately predicting when and how users decide to buy products can greatly affect both the user's shopping experience and a company's bottom line. Data-driven strategies that can decipher user intent not only boost conversion rates but also enable real-time personalization and more nuanced sales approaches. Yet, traditional machine learning models—which often rely on static, aggregated features—frequently prove insufficient for capturing the temporal and contextual richness that characterizes contemporary user behavior online.

A key driver in predicting user intent involves the sequence of activities users engage in during their journey on a platform. Standard methods might, for instance, look at the total number of product page views or whether someone added items to their cart. Although these behaviors provide insights, they overlook the dependencies and transitions between events. In contrast, sequence modeling methods can preserve the order and duration of user interactions, highlighting patterns that basic aggregation might miss. Tools such as RNNs and their variants (LSTM and GRU) are particularly suited to handling sequential data because they maintain "memory" of what happened in previous time steps, allowing them to form richer context for current predictions.

However, RNN-based models can encounter difficulties with very long user sessions, where traditional recurrent networks risk losing track of information collected many steps beforehand. To address that, techniques like LSTMs and GRUs have become prevalent in many text and speech applications. More recently, attention mechanisms and Transformer architectures have sparked intense interest, notably because they can learn dependencies across a sequence without needing to process inputs strictly in a time-ordered fashion. This opens possibilities for capturing longer-range patterns and contextual details more efficiently. Within e-commerce, these benefits are especially relevant since user behavior can be erratic, with periods of inactivity interspersed with bursts of high engagement.

In parallel, product-level context also plays an important role. Even if two users browse the same number of items, the categories, brands, or attributes of those items can vastly differ. Embedding these attributes in a latent space can help the model recognize semantic relationships among products. For instance, if users consistently transition between shoes and apparel or show a preference for certain brands, embeddings can reflect that similarity, enhancing predictions about their eventual purchase actions.

Despite these advances, the field still grapples with constraints around real-time data, multi-session user behaviors, and the integration of product metadata. Some studies address user intent by training models on session-level data without incorporating the user's longer-term patterns across multiple sessions. Others may focus solely on the user's demographic features, missing out on the specificity and nuance found in the user's unique browsing history. The research discussed here seeks to overcome these obstacles by introducing a deep contextual sequence modeling framework that can combine multi-session data with product embeddings and attention-driven architectures.

This paper is structured as follows. First, we present an in-depth literature review, tracing the evolution from older machine learning paradigms toward the cutting-edge of deep sequence modeling. Next, we detail our proposed framework, from the data preprocessing stage to the architecture of the neural network modules. A schematic diagram illustrates how we stack together product embeddings, sequence encoders, and an attention mechanism. We then turn to the experimental results, showcasing performance comparisons against both traditional and modern baselines and providing insight into why certain architectures outperform others. Finally, we reflect on the findings and address remaining challenges, including cold-start scenarios, data sparsity, and how external factors like shipping costs and product stock influence user decisions. By connecting these insights with real-world e-commerce needs, this work aims to bridge the gap between advanced research ideas and practical, scalable implementations.

## 2. Literature Review

Research on predicting user purchase intent within e-commerce has evolved in tandem with general trends in data analytics and machine learning. Early studies often concentrated on heuristic-driven or rule-based approaches, analyzing how frequently users clicked on products or how much time was spent on product pages. These methods were intuitive yet limited. As the scale of e-commerce expanded, more advanced statistical and machine learning algorithms, such as logistic regression, decision trees, and random forests, began to be applied. While these provided meaningful improvements, they typically collapsed temporal information into static features—for example, summarizing a session by the total number of pages visited—thus discarding the order and timing of events that could be crucial for more accurate intent prediction.

When deep learning started transforming the fields of computer vision and natural language processing, researchers began applying neural architectures to e-commerce data. RNNs, with their ability to process time-series information, naturally became a go-to choice. By reading sequences of user interactions, RNN-based models can predict a future outcome (such as a purchase) based on what has happened in the session so far. Since vanilla RNNs struggle with gradient issues when sequences become long, LSTM networks presented a more robust alternative by introducing gating mechanisms that manage how information is stored, forgotten, and retrieved over time [1]. This approach has been successfully utilized in many situations where user actions are interdependent, like music playlist generation or video consumption.

Simultaneously, GRUs emerged as a leaner choice. They combine certain gates in LSTM cells, reducing parameter counts and frequently training faster while retaining most of LSTM's ability to track long-range patterns [2]. In various studies, LSTMs and GRUs are found to be almost on par in predictive performance, although specifics of the dataset and problem can tip the balance in favor of one over the other. Beyond RNNs, the revolution brought by Transformers and self-attention provided a crucial alternative [3]. Attention-based models sidestep the limitations of recurrent structures by allowing the model to weigh the relevance of each interaction in the entire sequence when predicting. This architecture has been

extremely successful in language models such as BERT and GPT, and it holds equal promise for e-commerce, where certain interactions might be more telling than others regarding user intent.

Another key component in the literature deals with embedding techniques. In the context of e-commerce, items can be represented as vectors that capture something akin to "semantic" similarity. For instance, items from the same category, or those frequently bought together, might end up close in the embedding space. These embeddings can be learned using techniques analogous to Word2Vec and GloVe, which were initially created for text data [4]. By augmenting the raw session data with embedded product representations, researchers have shown improved model accuracy, as the learned embeddings can highlight non-obvious relationships between products that typical ID-based or categorical features might ignore.

Recent work has taken advantage of multi-session data, looking at user actions across multiple visits, rather than focusing on just a single session. This multi-session perspective can be valuable, especially if a user rarely completes a purchase on the first visit. In these cases, aggregated memory across sessions and attention-based weighting of important sessions can yield far more accurate predictions. Hybrid models, combining LSTM/GRU or Transformers with attention, have been proposed to track short-term dynamics alongside a user's broader purchase history [5]. Some also combine multiple objectives—predicting not only a binary purchase outcome but also the likelihood of cart abandonment, or the potential lifetime value of a user—allowing the model to learn shared representations that bolster performance on each task [6].

Despite the clear achievements, a few issues remain. Many deep learning models can be regarded as black boxes, which can be challenging for companies needing interpretability. For instance, if the model mislabels a session, it may not be straightforward to trace which user actions or product signals caused the error. Additionally, the problem of sparse data persists, especially for new products or new users (the so-called cold-start problem). Transfer learning and data augmentation approaches have been tested, but consistent and scalable solutions that can be readily adopted in real-world production systems remain somewhat elusive. Even so, the body of evidence in the literature underscores that advanced sequence-based frameworks, especially those leveraging attention, are among the most promising strategies for modeling e-commerce user behavior in all its complexity.

The framework proposed here aims to contribute to this conversation by melding multi-session user data, product embeddings, and advanced attention-based architectures into a single cohesive system. It addresses the complexities of dynamic user journeys, ensures the model remains mindful of relevant historical actions, and captures crucial context from product metadata. The next section details the methodology, laying out how each part of the pipeline—session segmentation, feature engineering, sequence modeling, and final prediction—fits together.

## 3. Methodology

### 3.1 Overview

The goal of the proposed methodology is to thoroughly integrate user behaviors across time with product-level context, building a deep learning architecture that excels at predicting e-commerce purchase intent. We accomplish this by creating a pipeline that encompasses data preprocessing, feature engineering, sequence encoding, context extraction, and final prediction. The crux of this approach is an attention-based module that sifts through a user's sequence of actions— potentially spanning multiple sessions—and identifies the actions most critical for determining whether they will complete a purchase. Additionally, we incorporate product embeddings that capture nuanced relationships among items, thus helping the model understand how different products interact within a user's browsing history.

Our methodological flow starts with a raw data ingestion stage, followed by session segmentation. We then transform raw interactions into structured sequences, enrich them with product metadata and user-level attributes, and feed the entire package into a neural network. We tried different sequence encoders (LSTM, GRU, and a Transformer variant) to evaluate their strengths in handling time-ordered data. After generating predictions, we compare the results with multiple baselines, including simpler models, to quantify the gains introduced by deep contextual sequence modeling.

### 3.2 System Architecture Diagram

Below is a conceptual representation of the system. It outlines how data moves through preprocessing, enters feature engineering and embedding layers, then flows into the sequence encoding and attention mechanisms before the final

prediction stage. While the diagram is presented here in a textual manner, it can be visualized as a pipeline where raw e-commerce logs get transformed and encoded, and then pass through an attention-augmented deep neural network that outputs the probability of a user making a purchase.

Raw e-commerce interaction data is collected. This data is cleansed, standardized, and segmented into sessions. During feature engineering, we extract user-level attributes (such as user demographics if available), product-level metadata (categories, brand, text descriptions), and session-level statistics. We generate product embeddings that place items in a vector space. The resulting feature set is passed into a sequence encoder, which may be an LSTM, GRU, or Transformer-based module. Then, an attention mechanism highlights the time steps (or user actions) that are most predictive of the final purchase. The output from attention flows into fully connected layers that compute the purchase intent probability.

### 3.3 Data Preprocessing and Session Segmentation

In the initial step, raw logs containing all user interactions—ranging from simple page views to add-to-cart events—are collected. Each record is labeled with a user ID, product ID, timestamp, and a description of the action (for instance, "view," "add to cart," or "purchase"). Given that e-commerce platforms operate continuously, each user may generate multiple sessions over extended periods. We define a session cutoff by setting a threshold for inactivity, so that if a user is inactive for more than a specified time (often 30 minutes), any subsequent interaction is considered the start of a new session.

For each session, we assign a binary purchase label reflecting whether the user checked out items in that session or within a short window afterward. Depending on data availability, user-level contextual information (like location or demographic details) may also be appended to each session. The result is a structured dataset that features one row per session, containing a list of user actions in chronological order. When necessary, we track prior sessions to encode a user's longer-term history, effectively creating sequences that extend beyond a single session.

### 3.4 Feature Engineering and Embedding Generation

After segmenting sessions, we extract various features. Session-level features can include the session duration, the total number of page views, and other aggregated statistics. More granular features can capture transitions between product categories, the frequency of add-to-cart actions, and the presence of discount codes or promotions. User-level features might include membership status, prior spending patterns, or location data, though such information depends on platform-specific data policies.

Product-level features prove particularly beneficial in an e-commerce context. We encode these attributes, such as the product category, subcategory, brand, price range, and textual descriptions, into numeric vectors. One approach is to train product embeddings from scratch by treating product co-occurrences as analogous to word co-occurrences in language models. The objective is for products that are frequently viewed or purchased together to have similar embeddings. Alternatively, if external data is available, pretrained embeddings can serve as a basis. Either way, during training for purchase intent prediction, the embeddings can be fine-tuned to optimize for the final classification goal.

### 3.5 Sequence Encoder and Attention Mechanism

At the core of our framework lies a sequence encoder. We implement and test three variants: LSTM, GRU, and a Transformer-based model. Each user session is viewed as a time series of events, which includes not just the product embeddings but also any relevant contextual features from that specific time step. LSTM and GRU networks handle these events in order, updating hidden states that carry forward the memory of previously observed interactions. They thus learn short- and medium-range dependencies, though LSTMs can often handle slightly longer sequences because of their explicit gating structures.

Transformers approach sequences differently. They deploy self-attention layers that let the model assess relationships among positions in the sequence in parallel, rather than step-by-step. This design speeds up training and can better capture longer-range dependencies because any element in the sequence can relate directly to any other element through learned attention weights [3].

Regardless of the sequence encoder chosen, we place an attention layer on top. Attention weighs the significance of each hidden state in the overall sequence context. In essence, while the sequence encoder captures broad temporal dependencies, the attention module decides which time steps matter most for predicting a purchase outcome. For example, a user's flurry of add-to-cart actions near the end of the session might overshadow product views from the beginning. By integrating attention, we allow the network to focus on the subset of actions that are the strongest signals of purchase intent.

### 3.6 Prediction Layer

Once the attention mechanism condenses the sequence into a weighted context vector, this representation is concatenated with any relevant global context features, such as user-level summaries. We feed the resulting vector into several fully connected layers, each followed by nonlinear activation functions. The final layer uses a sigmoid or softmax activation, making it compatible with binary or multi-class classification, depending on how the problem is framed.

During training, we rely on a loss function such as binary cross-entropy (for purchase vs. non-purchase) or categorical cross-entropy (if multiple outcomes are tracked). We typically use adaptive optimizers like Adam or RMSProp. Hyperparameters—like the number of attention heads in the Transformer, the hidden layer dimensions, and the embedding sizes—are tuned on a separate validation set. After training completes, we select the best model based on validation performance and evaluate it on a holdout test set.

## 4. Results and Analysis

### 4.1 Experimental Setup

To test the proposed framework, we employed a large publicly available dataset of e-commerce interactions spanning various product categories, user demographics, and a wide range of timestamps. After cleaning and filtering, the dataset comprised roughly one million session records from about 100,000 different users, tied to nearly 50,000 unique products. We chose a split of 70% training, 15% validation, and 15% test sets, ensuring that sessions from the same user were restricted to a single set to prevent data leakage.

We built multiple models for comparison: a baseline logistic regression classifier that used only session-level aggregate features, a feedforward neural network trained on the same aggregate features, and three deep sequence models—LSTM with attention, GRU with attention, and a Transformer-based model (which inherently uses attention as well). Product embeddings were shared across these deep models, but only the sequence-based methods utilized them step-by-step during time-ordered processing. Each model was optimized with either a binary cross-entropy loss (when focusing purely on purchase vs. non-purchase) or an extended objective if the classification was multi-class.

### 4.2 Performance Metrics

To capture a comprehensive view of how well the models performed, we evaluated them on accuracy, precision, recall, and F1 score. Because purchase behavior was less common than non-purchase behavior, we gave particular weight to the F1 score, which balances precision and recall. We also considered the Area Under the Receiver Operating Characteristic Curve (AUC), which elaborates on how well each model manages the trade-off between true positives and false positives across various decision thresholds.

### 4.3 Comparison with Baseline Models

In the test set, logistic regression and feedforward neural networks using static session features provided only moderate predictive power. These simpler methods lacked a mechanism to incorporate the time dimension of user actions. When we introduced the LSTM with attention, performance metrics substantially improved, reflecting how sequential analysis and the attention-based focusing of relevant time steps can sharpen purchase intent predictions. GRU with attention yielded results in a similar range, slightly differing in parameters and training time, but generally close to the LSTM's performance.

The Transformer-based approach outperformed the RNN variants. Its parallelized attention layers likely captured broader interactions within each session more effectively. Even though the difference in metrics compared to LSTM and GRU may

not be massive, it was consistent across multiple runs, suggesting that the Transformer architecture gained a stronger grasp of complex, multi-step user behaviors.

A snapshot of these results can be seen in the following comparative summary:

Model performances on the test set showed that the Transformer-based model, combining multi-head attention with product embeddings, had the highest overall F1 score, followed closely by LSTM and GRU variants. Both logistic regression and the feedforward neural network lagged behind because they could not account for sequence order. This finding highlights the importance of modeling temporal patterns and focusing on pivotal user actions.

### 4.4 Chart Visualization

A clearer view of the results appears in a simple chart that plots each model's F1 score. The baseline models (logistic regression and feedforward) cluster in the lower range, while the sequence models stand out for delivering higher F1 scores. Within the group of sequence models, the Transformer-based method consistently comes out on top, validating the notion that advanced attention mechanisms pay off in these scenarios. This demonstrates that capturing the interactions among time steps in a more flexible, parallelized way can yield extra gains in e-commerce intent prediction tasks.

### 4.5 Error Analysis

To understand where our best-performing model—namely the Transformer-based one—fell short, we conducted an in-depth error analysis on misclassified sessions. We discovered that sessions consisting of only a few events (for example, a user who clicks on a single product and leaves) were particularly tricky for the model, likely because minimal interaction offers limited evidence of future intent. Additionally, new products lacking historical data presented a challenge, indicating that the model heavily relies on both user context and product embeddings to make accurate predictions.

Another noteworthy pattern in the model's errors pertained to users who engaged extensively, added items to their cart, but eventually abandoned the purchase. Such cases often hinge on external factors (shipping fees, discount code unavailability, or personal budget concerns) that the model cannot directly infer from browsing data. These results imply that while advanced sequence models excel at interpreting user behavior, combining them with additional data sources—like price fluctuations or shipping costs—could help reduce overestimation of purchase intent when these external elements are adverse.

## 5. Conclusion

This research analyzed how deep contextual sequence models could enhance the prediction of user purchase intent in e-commerce. By merging sequential user activity data with product embedding layers and attention mechanisms, the presented framework captures both immediate user actions and larger-scale behavior patterns spanning multiple sessions. Through experiments on a substantial public dataset, we demonstrated that sequence-based methods—especially Transformers—outstrip simpler baselines, highlighting the significance of context, time, and product relationships in understanding purchase decisions.

The presented methodology has clear practical implications. Platforms can employ these predictive models to trigger targeted recommendations, personalized promotions, or dynamic pricing at moments when users exhibit high purchase probability. However, challenges remain. Sparse interaction data, newly launched products, and external factors tied to cost or convenience often confound the model's predictions, pointing to the need for additional data streams or specialized architectures. Future work could thus include exploring multi-task learning strategies, real-time model updates, or the integration of external signals such as pricing shifts or competitor promotions. By continuously refining these strategies, e-commerce platforms can move closer to a fully personalized shopping experience that aligns seamlessly with user intent.

## References

1. S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
2. K. Cho et al., "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, Doha, Qatar, 2014, pp. 1724–1734.
3. A. Vaswani et al., "Attention Is All You Need," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Long Beach, CA, 2017, pp. 5998–6008.
4. T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient Estimation of Word Representations in Vector Space," *arXiv preprint arXiv:1301.3781*, 2013.
5. Y. Sun et al., "BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, Beijing, China, 2019, pp. 1441–1450.
6. Dr. Sadik Khan , "Using Web Mining to Trace the Evolution of Language, Ideas, and Cultural Trends Over Time", Int. J. Sci. Inno. Eng. pp. 1-8, 2026-01-01 doi: https://doi.org/10.70849/ijsci03010016114 .
7. G. M. Faris et al., "Multi-Task Neural Networks for E-commerce Customer Lifetime Value and Purchase Intent Prediction," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Anchorage, AK, 2019, pp. 1504–1512.