## RESEARCH ARTICLE

# Machine Learning Approach to Classification of Online Users by Exploiting Information Seeking Behavior

**MOONA KANWAL**[1,2], **(Member, IEEE), MUZAMMIL AHMAD KHAN**[2], **(Member, IEEE), NAJMA ISMAT**[2], **NAJEED A. KHAN**[1], **AND AFTAB A. KHAN**[3]

[1]Computer Science and Information Technology Department, NED University of Engineering and Technology, Karachi 75270, Pakistan
[2]Department of Computer Engineering, Sir Syed University of Engineering and Technology, Karachi 75300, Pakistan
[3]College of Education, Health and Human Services, Longwood University, Farmville, VA 23901, USA

Corresponding author: Moona Kanwal (mkkhan221@cloud.neduet.edu.pk)

**ABSTRACT** In today's world, technology has engulfed the internet with an excessive amount of unfiltered, spontaneous, and incessant data from multiple sources. Complex algorithms are designed to present information effectively based on user intent. The online experience of users is a combination of various behaviors exhibited to seek information, including searching, sharing, and verifying information. However, this multi-faceted user behavior is yet to be explored comprehensively. This research contributes towards proposing a user intent-machine learning model for classifying users based on their online search, share, and verification behavior, identifying different types of users based on their online engagement, and demonstrating that dynamic online interactions can be classified based on their searching, sharing, and verifying behavior. User feedback on online behavior and practices is gathered through a questionnaire, encompassing participants from diverse gender, occupational, and age demographics. Following the extensive feature engineering, the significant features are presented to K-Mean clustering to identify user intent classes or profiles and their characteristics. A supervised learning Linear Discriminant Analysis Classifier (LDAC) is then trained on data to classify these classes. The proposed framework successfully predicted the user intent class with 80% accuracy. The model is further tested on users' dynamic interaction data gathered through a second user study. The information search, share, and verify activity data is transformed to fit the model and labeled by human raters using the user profiles resulting from clustering. The research achieves an Inter-rater reliability (IRR) of 60%, whereas the model predicted the user with 67% accuracy. This research indicates that a user's purpose in seeking information, their willingness to share information on social media, and their inclination to view information as credible can all contribute to understanding their intentions, identifying behavioral similarities, and can be used to recognize intent through dynamic interactions that can be used in targeted marketing, and search engine optimization.

**INDEX TERMS** User intent, cluster, browsing preference, information sharing, user behavior, search reasons, human behavior.

The associate editor coordinating the review of this manuscript and approving it for publication was Hai Dong.

## I. INTRODUCTION

Since the earliest civilization, humans have been found to have a thirst for knowledge. They gained information and

shared it with others. The execution of information exchange was primarily done through word of mouth. The term epistemology (the theory of knowledge) can be traced from Greek philosophy. For example, Plato (429–347 B.C.E.) wanted to know how we acquire knowledge. Locke (1632-1704) wanted to see the operation of human understanding. Piaget (1896-1980) developed a theory of genetic epistemology or the idea of cognitive development, and Vygotsky thought that we gain knowledge through our social culture. He is well-known for his socio-cultural theory of cognitive development [1]. In today's era, the internet and social media have taken over the social culture, and Vygotsky's followers promote digital access to knowledge [2], creating a universal culture worldwide. Traditional means of information dissemination, like newspapers, magazines, books, television, and radio, are being replaced by readily available and accessible digital sources of information, including web channels, online social networks, podcasts, etc. [3]. While people appreciate the new media era that has diluted the physical barriers in communication and provided them with a platform to share their opinions [3], [4], people are also getting overwhelmed by the information bomb that now explodes on their screens. Users increasingly struggle to focus on one topic and filter the relevant information. The insight into how the user interacts with the internet to retrieve information, share it on social media, or trust information can help filter information effectively to suit user intent. Authors in [5] summarized through survey that an approach towards finding users' online intention follows two steps: (1) analyzing the user perspective and approach towards online platforms and using technology to generate a general profile of user online preferences, and (2) developing learning models to identify users based on online information preferences.

This research takes inspiration from [5] and targets the following research questions:
a) Can user's intentions be classified based on user information-seeking behavior on the internet?
b) Can the user intent classifier predict user intent type using user dynamic behavior?

Many studies and models aim to classify and predict user intention in seeking information online [6]. One of the most widely used behavioral models is the Theory of Planned Behaviour (TPB) [7], which proposes that user intention is influenced by attitude, subjective norm, and perceived behavioral control Another popular model is the Theory of Reasoned Action (TRA) [7], which suggests that user intention is determined by attitude and subjective norms. These behavior models place user intention as an antecedent of behavior, and rightly so, as their behavior reflects user intent. Therefore, to predict user intention, it is important to understand behavior, especially in an online environment where user interactions with the internet represent behavior, and interactions can be captured as different data types and through various means.

Using a machine learning approach, these models have inspired this research to determine user intention through user search behavior, information dissemination, and information trust behavior. This is done in a series of steps. Initially, user information is collected through a survey questionnaire. The behavioral similarities are identified at the preprocessing phase, and the K-Mean Clustering is employed to delineate user intent groups. The clustering process identifies five distinct groups based on user intentions. Subsequently, a Linear Discriminant Analysis (LDA) classifier is trained on the data to classify user intent clusters. For model validation, Inter-Rater Reliability is used on the user's dynamic interactions captured via another user study. The user profile is generalized from the observations of the machine learning model and utilized in annotating the user's dynamic interactions.

## A. USING MACHINE LEARNING TO ANALYZE USER SEARCH BEHAVIOR

One of the trending research topics is to capture user intention in acquiring information and provide ease in bringing that information to the user screen. The research stretches from user navigation, individual preferences, likes and dislikes, search results, and relevance of search queries [8] to semantic analysis of websites for customizing search results to better suit the user intent [9]. The user interaction and preferences on social media are also studied to understand user behavior. In the study [10], the users' physical attributes are used to analyze user interaction in image searching and content to design a search intent system. Another study [11] analyzes clickstream data in predicting the intention of shopping online using deep learning models. In [12], search personalization using search and click history is explored, and [13] applies reinforcement learning to model user intent using information visualization. The study's authors [14] have developed a taxonomy for searching source code in computer programming. Moreover, Conversational Information Seeking (CIS) or providing information needs via conversation with the search engine is another paradigm used to assess user intention in searching. Authors in [15] and [16] explored the search intent and user behavior among others in conversational search.

Research on studying and modeling user intention in online searching also focuses on user engagement and behavior on social media with millions of active users [17]. Online social platforms have also become a dynamic source of information, as was noted by the researchers [18] during Covid'19. This research and many others have shaped search engine results, social platform information feeds, and online marketing [19].

## B. USER'S INTENT TO DISSEMINATE INFORMATION

The need for users to propagate information is also essential, as it can aid in filtering out misinformation and regulating the dissemination of authentic information. There was once a time when information dissemination could only be done through designated publishing bodies like newspapers, official news channels, book publishing houses, etc. With the availability of technology and cheap network access to social media platforms, sharing information has become easier and

accessible to everyone. Even the traditional means of information dissemination are using web-based technologies to compete with social media platforms. It is also true that the information available via authenticated print media and news platforms is still considered more reliable than the social platforms. Social platforms, in some way, act like digital word of mouth. Recent surveys show that social platforms are the prime source of information gathering and exchange [20]. The authors [21] investigated differences in gender behavior and intention in information sharing. According to the authors, the information shared on social platforms varies from status updates and liking a post to advertisement sharing, specially on social media like Facebook or Twitter. Furthermore, the intention to share information differs for men and women based on social ties and commitments. The study [22] discusses factors like extraordinary circumstances or times of crisis (e.g., the Covid19 pandemic), social influence, or user attitudes as causes of an increase in the use of social platforms to access and share information being the only way to connect to rest of the world. Another study [23] investigated Taiwan's Instagram users' social behavior and used big data analytics and k-mean to cluster users and generate user profiles for social media and commerce development.

### C. MACHINE LEARNING IN VERIFYING INFORMATION

The Internet is an information hub where the public can access timely, efficient, and up-to-date information. Information supply and circulation have increased drastically with easy access to social platforms like Facebook, Twitter, YouTube, etc. But unfortunately, not every piece of information is reliable. Information verification has become a challenge for news agencies and end users with the spread of fake information. Research shows that users easily fall for false information if it supports their viewpoint [24].

A recent study showed that certain types of social media users are highly motivated to verify information versus the users whose primary intent on social media is to seek entertainment [25]. Authors in [26] conducted a study exploring the user's purpose in verifying the information and the methods employed. The research discovered various factors like source credibility and headline content to be prima facie in ascertaining the veracity of information. The users' intentions were investigated in [27] who spread fake information on social media intentionally or unintentionally, and an influence graph was proposed to determine the degree of unintentional spreaders of fake news. In [28], the authors compared social media, news media television, and newspaper credibility perceived by users. They concluded through the survey data that young adults perceive news on social media as more credible than others. Another research [29] also used news and social media textual and temporal content along with user responses to predict fake news using the BART model. Similarly, several research studies have targeted fact-checking and credibility assessment of internet and social media content to facilitate the end user to make informed decisions about

the information they seek online. The research ranged from claim detection systems [30], predictive methods [31], creating multimodal data repositories [32], [33], and annotating data [34]. User perception of information verification can be used in tuning automated filters.

This research investigates user online behavior and practices towards online searching, information propagation, and verification and builds user profiles that can contribute to designing a better user experience.

### D. RESEARCH CONTRIBUTION

The research contributions are as follows.
- Explored user behavior and practices employed online to search, disseminate, and verify information. All three aspects provide a complete picture of user intention to acquire knowledge and have not been covered as three dimensions of user intent to seek information.
- Developed a learning model to classify users.
- Tested classifier on user dynamic behavior.

## II. RESEARCH METHODOLOGY

This research aims to investigate online activities and preferences that users employ in searching, sharing, and verifying information and develop a user intent model that can summarize the general characteristics of users concerning online behavior and practices. The proposed methodology is presented in Figure 1a. The process is divided into two phases. The first phase includes developing a Machine Learning (ML) model to classify users based on their online search, sharing, and verifying behavior. The model is an extension of the prior work discussed in [35]. Phase II involves testing the ML model on users' dynamic interactions to acquire information.

The machine learning pipeline for Phase 1, illustrated in Figure 1b, is developed after considering related work with a particular focus on the categorical nature of the data. The literature review also revealed that the existing models have yet to cover user search and social and verify information aspects as complete user information-seeking behavior. Other works include different data types, like textual and vector, and use state-of-the-art models like BERT [36] or BART [29] deep net models. However, due to privacy, text extraction from social media and visual data are hard to obtain from participants. This concern has also inspired this research to use user feedback on internet experience as they perceive it and test it to validate model prediction accuracy on user dynamic activities, keeping the features simple and data easy to acquire.

The investigation uses the prior work [35] of collecting user feedback on their online activities, preferences, and use of social networks for information dissemination. The data is pre-processed and encoded. Three new attributes, Search Openness, Online Extravert, and Information Conscientious, are introduced. These attributes are computed based on existing data and represent scores for searching, sharing, and verifying behavior. K-mean Clustering is used on the new
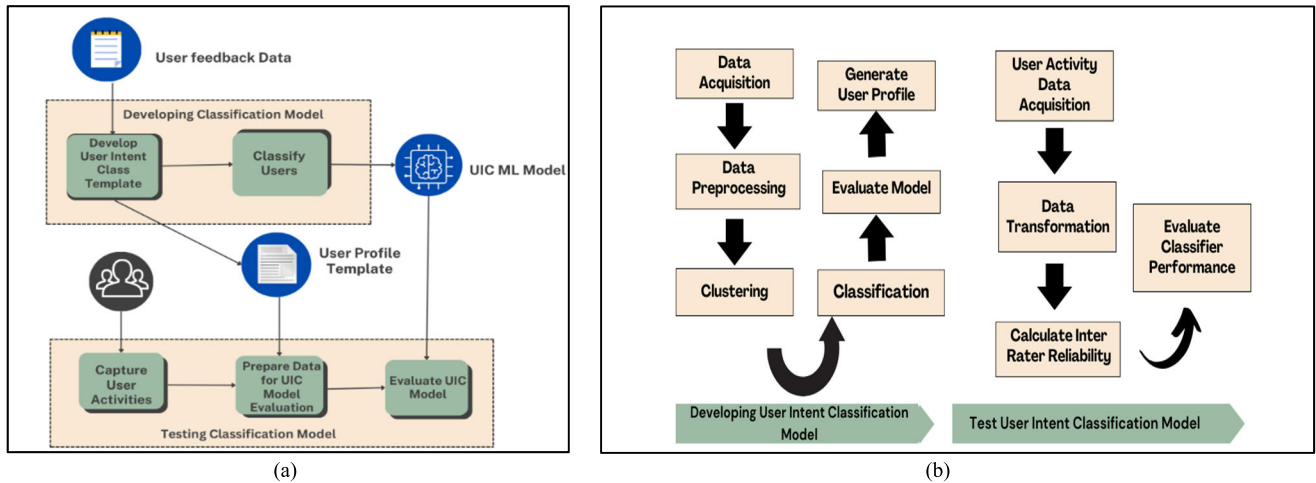
**FIGURE 1.** (a) Proposed research methodology. (b) Proposed process flow.

attributes to group similar data. Labels are assigned to clusters based on the user characteristics. The study then uses these clusters from previous research and adds to the original data, and then different machine learning classifiers are used to classify users. Different evaluation metrics are used to understand models' performance, whereas SHAP [37] class-wise scores validation is used for the model efficacy and feature contribution in class.

Testing the ML model in phase II involves capturing user interactions and transforming them into model features. User attributes are then analyzed to understand user intent as learned by the machine. The user profile created in prior work is used here to label phase II data. Finally, the model is used to predict users, the results of which are validated using Inter-Rater Reliability (with two human raters).

PyCaret,[1] Python pandas, and scikit-learn[2] are used for feature engineering, data modeling, and classification. Pearson Chi-square Test is run using the Python SciPy[3] stats module. Visualization is done using PyCaret and Origin-Pro2023.[4]

### A. DEVELOPING USER INTENT ML MODEL
The first phase of this research involves developing a machine-learning model that can classify users based on their online behavior and answer RQ1.

The following sections explain the machine learning pipeline shown in Figure 1b.

### 1) DATA ACQUISITION
Data acquisition is the most critical, systematic, yet tedious part of research. Data can be collected through various means, whether manual, electronic, or experiments, and may involve

other living beings or machines. For developing the user intent model in research phase I, data specification is not fulfilled by the existing datasets available from prior research; hence, the qualitative survey-based data collection process is governed.

The user study (US1) is based on qualitative questions and case studies in the form of a questionnaire and represents user feedback about online behavior. It is divided into four parts: (1) users' demographic data, (2) searching behavior and action, (3) user intent toward information dissemination, and (4) user perception and actions on verification of information. The questions are carefully crafted and approved by the domain experts before being shared with the participants. The participants are from varied backgrounds and are invited using social platforms, community groups, and universities.

Three hundred initial responses are received against the threshold of 200. The final participants are mainly from South Asia and North America, with mostly Urdu or Hindi as their native language. Both genders are considered, with ages ranging from 18 to above 60 years. Six occupations are mainly considered: science and technology, finance, insurance and commerce, education, and Unemployment. Any other occupations are regarded as Others. Table 1 presents the initial categories and features. The intention head in the table represents the user behavioral areas explored through US1. Data collection is discussed in detail in [35].

### 2) DATA PREPROCESSING
Data obtained from User Study (US1), which has both ordinal and nominal categorical features [38], requires pre-processing to clean, transform, and encode it into a format suitable for machine learning models [39].

As discussed in previous work [35], the data is structured and categorical with textual data. To process the data easily, all textual data in sentence form are transformed into short words. The category's skewness is decreased by combining skewed options with other options, eliminating options,

---

[1] https://pycaret.readthedocs.io
[2] https://scikit-learn.org
[3] https://scipy.org
[4] https://www.originlab.com

**TABLE 1.** Initial variables and categories.

| Intention | Main variables | Categories |
|---|---|---|
| Demographic | Gender | Male, Female |
| | Age | Adult (18 – 40 yrs.), Older Adults (> 40 yes) |
| | Occupation | Science and Technology, Finance and Insurance, Education, Unemployed/Homemaker, Others. |
| Preferences | Browsing Preference | News Channels-Internationals, News Channel-Local, Talk Shows, Financial Websites, E-commerce Websites, Educational and Technology Websites, Entertainment Websites, Nothing Particular |
| | Information medium | Internet Search, YouTube, Newspaper, TV and Radio Channels, Social Media, Word of Mouth |
| | Language | English, Urdu/Hindi, Punjabi, Arabic, Chinese, French, Others |
| | Search Engine | Google, Bing, Yahoo, Social Media, Others |
| Frequency | Search Frequency | Daily, Weekly, Monthly |
| Reason | Search Reasons | Any Info on a Particular Topic, More Information on a Particular Topic, Online Activity (shopping, downloading, gaming), Specific Website |
| Interest | Share Info on Online Social Network | Educational Information, Disaster, Emergency, and Other Breaking News; Entertainment Content; Finance Information; General News, Gossip, Conspiracy, and History Information, Interviews and Videos, Job Postings, Pray and Wellness Mantras; Medicinal Remedies. Religious Content; Science and Technology Info, None |
| Preference | Information Share Medium | Social Media, Blogs, Private Social Groups, Private Messages, Don't Share |
| Frequency | Information Verification Frequency | Always, Frequent, Sometimes |
| Preference | Medium for Verifying Information | Ask an Expert, Ask on Forums (e.g., Reddit, Quora), Ask a Friend, Lookup on Information Repositories, Lookup Scholarly Articles, Search Fact-Checking Websites, Search using Search Engines, Undisclosed. |

**TABLE 1.** *(Continued.)* Initial variables and categories.

| Criteria | Information Credibility Criteria | The Article has High Likes, Date Provided, The Article is Well Written, The Author Profile is Provided, The Author is Well Known, The Publisher is Well Known, References are Provided, Agrees with the Content, Content is Considered Credible. |
|---|---|---|
| | Trusted Online Source Criteria | Complete Information Provided, Content has Clarity, Domain Extension, Recommended by Friend, About and Contact Details Provided, Website has Good UX, Good Content, High Likes or Reviews, Dated Recent, In Top Search results, None. |
| Attitude | Case1 | Confirm from Friend, Forward to Friend, Ignore Post, Post on Social Media, Verify. |
| | Case 2 | Discard, Investigate, Share. |

or creating range bins; e.g., age ranges created in the questionnaire are further combined to create two bins in years, adults (18-40) and old adults (> 40).

Figure A-I in the supplementary file shows the distribution of skewed categories in sample features data. The skewness is considerably removed after merging some categories while eliminating features that show no variance.

Figure A-II in the supplementary file shows a distribution graph of sample features after preprocessing. Furthermore, the categorical data is converted to numeric format using hot encoding [40] for machine learning processing. Hot encoding converts categorical data into a binary vector where each category value is assigned a separate column. All entries with a value are marked one, and the rest are marked zero. This is repeated for all the values of categories. This encoding scheme is selected amongst other efficient encoding schemes like Hash encoding, Factorization [40], Target encoding [41], etc., as it doesn't increase data volume and doesn't pose data loss. Hot encoding does create sparse data, but it becomes a problem for a large number of columns [41], which is not the case here. The data columns are further merged to reduce skewness.

Table 2 lists some of the features after transformation (for readability, the naming convention defined above is not used). Scaling and normalization are not required for binary data, but for transformation to numeric data for dimension reduction and scoring, normalization is performed on the dataset. The initial dataset has 30 main variables. After preprocessing the data, it is reduced to 25, which is expanded to 50 variables after encoding. The participants' responses are also reduced to 255 after preprocessing.

Feature selection is an essential step in data analysis as it eliminates redundant or insignificant features that may cause overfitting of machine learning models. The outcome also

**TABLE 2.** Example categories after feature engineering.

| Main variables | Categories |
|---|---|
| Gender | Male, Female |
| Age | Adult (18 – 40 yrs), Older Adults (> 40 yes) |
| Occupation | Science and Technology, Finance and Insurance, Education, Unemployed/Homemaker, Others. |
| Browsing Preference | News Channels, Financial/E-commerce, Educational/Technology, and Entertainment websites |
| Information medium | Internet Search, YouTube, Newspaper-TV-Radio, People, Social Media |
| Share Info on Online Social Network | Specific Discipline (History, Science, Technology, Education, Finance), News (All types), Entertainment Content, News, Body, and Soul |
| Share Medium | Social Media/Blogs, Private /Social Groups, Message |
| Verification Medium | Ask an Expert/Friend, Search the Internet/Social Media. |
| Info credibility criteria | Article Quality/Popularity, Author Profile Provided, the Publisher is Credible. |

**TABLE 3.** P adjusted values for categories.

| Main Variables | $\alpha_c$ | $\alpha_t$ |
|---|---|---|
| Gender | 0.025 | 0.15 |
| Age | 0.025 | 0.15 |
| Occupation | 0.01 | 0.06 |
| Information Medium Used | 0.0125 | 0.025 |
| Browsing Preference | 0.0125 | 0.025 |
| Share Information on OSN | 0.0125 | 0.025 |
| Information sharing medium | 0.025 | 0.15 |
| Information verification medium | 0.025 | 0.15 |
| Information verification Criteria | 0.016 | 0.033 |

reduces the data dimensions, which may help improve the accuracy of the machine learning model [42]. The filtering method [43] is chosen for initial analysis, and univariant statistical analysis, like the Hypothesis Test for Independence, is used to help distinguish dependent and independent variables [38].

The chi-square test for independence ($\chi 2$) is used in hypothesis testing to find the significant correlations between variables. The chi-square test ($\chi 2$) is performed on the one hot encoded pre-processed data. The p-values are derived for all the variables paired with each other. The chi-square test ($\chi 2$) is calculated as

$$\sum \left( x^2 = \frac{(O - E_{ij})^2}{E_{ij}} \right) \quad (1)$$

where O = Observed (the actual count of cases in each cell of the data), E = Expected value (calculated below), $\chi 2$= The cell Chi-square value, $\sum \chi 2$ = Formula instruction to sum all the cell Chi-square values $\chi 2$ i,j = i,j is the correct notation to represent all the cells, from the $i^{th}$ row, $j^{th}$ column.

$$E_{ij} = \frac{M_i \times M_j}{n} \quad (2)$$

where E is Expected value of $i^{th}$ variable with $j^{th}$ variable (each cell), $M_i$ is $i^{th}$ Row marginal value, $M_j$ is $j^{th}$ Column marginal value, n is total sample size.

The standard $\alpha$ of 0.05 required adjustment since we have multiple classes under each variable [44]. Therefore, the

p-value is adjusted using the Bonferroni-Adjusted method defined as.

Bonferroni-Adj p-value

$$= \frac{\text{target alpha level}}{n - \text{rank number in pair (by degree of significance)} + 1} \quad (3)$$

where: Target alpha level is the overall alpha level (usually .05), n is number of tests.

Bonferroni Adjustment divides the nominal p-value by the number of tests performed simultaneously or, in other words, by the number of classes within the variable. This ensures that the main variable's overall significance doesn't exceed the nominal p-value [45]. Thus, the significance of features is determined by p-value $\leq \alpha_c$, whereas $\alpha_c < $ p-value $ < \alpha_t$ is considered medium significance. Then, feature selection is performed based on the adjusted p-values and correlations obtained from the analysis.

Table 3 outlines the adjusted $\alpha_c$ and $\alpha_t$. Figure A in the Supplementary file shows the adjusted p values and feature correlations. Almost all the features indicate some level of correlation in the feature set. The features related to cases 1 and 2 are excluded from the final analysis and are used separately to compute the scores, which will be discussed next.

### 3) COMPUTING FEATURE SCORES

Various statistical, mathematical, and machine learning models have been used to explore data patterns and reduce large sets of categorical features into numerical data types, like [43] or [44], which compute feature scores based on weighted item scores, ensuring data integrity. A similar technique is used in this research to reflect the data in numeric and reduced form and use learning models to find patterns. Three new

variables are introduced and are used as computed scores. They are *Search Openness (SO), Online Extravert (OE)*, and *Information Conscientious (IC)* [35]. SO summarizes search behavior features, whereas share behavioral features (including sharing attitude) are computed as OE. IC summarizes information verification mediums, criteria, and verification attitude in cases 1 and 2. The demographic features and the new calculated scores are used in user behavioral pattern recognition.

The scores for SO, SE, and IC are calculated as the sum of the weighted mean of features:

$$Sc_V = \sum_{j=1}^{j=k} WM_{ij} \qquad (4)$$

where Sc is score, V is SO, SE and IC, k is total no. of categories.

The weighted mean (WM) is defined as:

$$WM_{ij} = R_{ij*}M_j \qquad (5)$$

where R is Participant response on $j^{th}$ category, $WM_{ij}$ is Weighted Mean at the $i^{th}$ response and $j^{th}$ category.

The mean of the binary category is defined as:

$$M_j = \frac{Fc_j}{S_n} \qquad (6)$$

where Fc is frequency of $j^{th}$ category.

S is n number of samples, M is Mean of $j^{th}$ category.

The new score and demographic features are presented for cluster analysis.

#### 4) IDENTIFYING USER GROUPS BY EMPLOYING CLUSTERING TECHNIQUE

Machine Learning offers clustering techniques for unsupervised data, i.e., data with unknown targets [46]. Clustering techniques discover similar patterns in data and group them based on it. It is hence selected to examine similarities in participants' data that formulate clusters and determine the labels for reasonable clusters. After computing feature scores, the new data, along with original demographic features, are first normalized to remove biasing, and then outliers are removed. Normalization uses the Z-score method [55] and the Isolation Forest or I-Forest technique to remove outliers. Normalized data (data reduced further to 201 rows) consisting of demographic data, SO, OE, and IC scores, mentioned in the previous section, is presented to K-Mean clustering. This is discussed in detail in [35].

Different clustering techniques are tried and evaluated on average silhouette score [47] and Elbow curve. K-Mean, an unsupervised technique, has been selected. K-Mean finds centroids of K clusters closest to the samples [48]. The number of clusters K here is arbitrary, and the K-Elbow curve is used to identify optimal clusters k= 5. Silhouette scores are used to validate the cluster separation. Figure 2a indicates that the average silhouette score of 0.3 shows the fair separation of clusters. The analysis also shows no negative score in clusters,
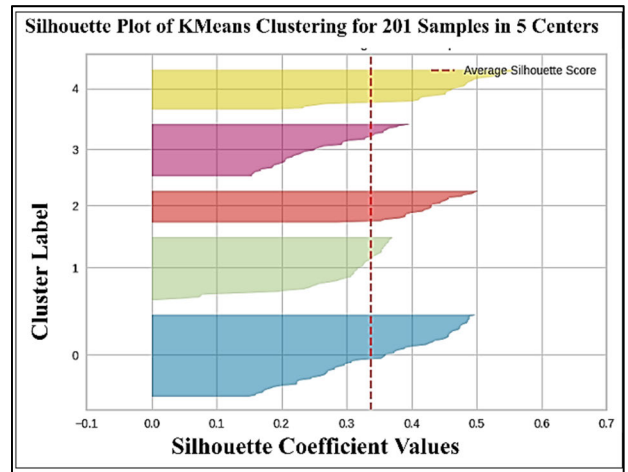


**FIGURE 2.** Silhouette plot for clusters [35].

indicating that all the data members are well-placed in clusters. The other techniques include Hierarchal, Spectral, and Birch; Table 1 in the Supplementary file provides silhouette scores for different techniques using different cluster counts.

The resultant cluster feature obtained from K-Mean is inserted into the original binary data. Then, the scores calculated from case 1 and case 2 are inserted as features separately, and the SE, SO, and IC scores are removed. Finally, dimension reduction is performed, and data is presented to the classifier. Since the clusters represent the user profile based on user intent, the cluster feature is named User Intent Class (UIC).

#### 5) BUILDING USER INTENT CLASSIFIER

Classification is a supervised machine learning method that splits data based on the target or class provided. It then uses that knowledge to attach a label to the unseen data. Several techniques fall under classification, such as classifying data into two classes, called binary classification, or more than two classes, called multiclass classification [49]. Classification has been used along with unsupervised learning to obtain groups based on similarity and identify what these groups represent, e.g., musical instrument classification [50], cancer classification [51], or predicting credit approval [52]. The second process in the methodology's first phase, as shown in Figure 1a, is to use classification on User Intent Class (UIC) to predict user groups and label them based on SO, OE, and IC. Using Kernel PCA, the new data, with the UIC as the target, is reduced to Principal Components (PC).

Principal Component Analysis (PCA) [53] is the most common and effective multivariate data analysis and feature reduction technique. Since PCA doesn't support categorical data, normalization is performed first [54]. Kernel PCA [55] is a variant that uses the Radial Basis Function (RBF) to project non-linear features to lower-dimension space. This study uses Kernel PCA, resulting in seven principal components.

**TABLE 4.** Classification models test metrics results.

| Model | Accuracy | AUC | F1 |
|---|---|---|---|
| LDA | 0.8128 | 0.9246 | 0.8102 |
| Extra Trees Classifier | 0.8264 | 0.9196 | 0.8180 |
| Logistic Regression | 0.8197 | 0.9250 | 0.8169 |
| Ridge Classifier | 0.8192 | 0.0000 | 0.8081 |

PC and target are fed to different classifiers, and the best model, based on the F1 score, Confusion Matrix, and AUC, is selected. Finally, the SHapely Additive explanation (SHAP) is calculated to observe the model strength concerning features. Details of these steps are provided in the sections below.

### a: CLASSIFICATION MODELS

The classifiers tested on the principal components are listed in Table 4. This section briefly explains the models and the experiment setup.

*Linear Discriminant Analysis (LDA)* [49] is a linear classifier that fits the class's Gaussian densities to data, applying Bayes' rule. The LDA model in this study used the Gini impurity to measure classification probability.

*Logistic Regression* [56] is a variant of Linear Regression that uses the logistic function sigmoid to estimate probabilities.

*Extra Trees Classifier* [56] uses decision trees on several random data sub-samples and accumulates all the results to generate a final prediction. This improves the prediction accuracy while handling overfitting.

*Ridge Classifier* [49] is based on linear regression, where the model is penalized for preventing overfitting and explaining the multicollinearity of features.

### b: EXPERIMENTAL SETUP

*Stratified K Fold cross-validation* maximizes the classifier performance while leveraging a lower sample count. K fold divides the data into k subsets and holds out one of the sub-sample sets as test data while training on all other subsets combined. The model is evaluated, and the step is repeated for all subsets. The final performance measure is the result of the mean of all the models' results. This way, all the samples get trained or tested for the model throughout the process. Stratified K Fold uses the same method, ensuring that each subset represents all classes. This cross-validation technique is used as the data was low in the count. Different subsets are tried for this study, and the best result is obtained using k=4.

*The training/Test Data split* is set to a 70/30 ratio. Table 4 provides the mean values of all the fold results.

*Scikit-learn Random grid search* is used to tune the hyperparameters of each model. Grid search [57] is an exhaustive search that aims to improve model performance by tuning its hyperparameters. It evaluates the model by using a combination of different parametric values.

Evaluation metrics used in this setup included *Accuracy, F1 score, Confusion Matrix,* and *Area Under the Curve* [55] to evaluate the model's efficacy. SHapley Additive exPlanation (SHAP) is a feature selection metric that evaluates classification models based on their features' contribution. Kernel Shap is a variant of Shap that incorporates classic Shap and Local Interpretable Model-agnostic Explanation (LIME), which assigns value to a feature on its contribution to a single prediction [37]. Shap evaluates the model and kpc contribution under each UIC in model efficiency. Figure 3a to 3c illustrates the kpc Shap value against each cluster.

All the classifiers provided good results; however, LDA is selected based on evaluation results on the class level and mean of all class results. Figures 4a to 4d present LDA results. The LDA model is then presented with a small set of data set aside to be used as prediction test data unseen by a model with the F1-score =0.81.

### 6) GENERATING USER PROFILE

As part of previous work [35], the User Intent Classes (UIC) generated by the K-Mean clustering technique are further explored for significant user characteristics in each cluster. The average participation of categories in clusters is computed for this purpose. The features $f_i$ where i=1- n total features, the participation of $f_i$ in $UIC_j$ where j= 0-4, is calculated as

$$UIC_j \, feature = Avg \, f_i \geq \sum Avg \, f_j \qquad (7)$$

i.e., the feature's average participation is greater or equal to the feature's total average participation in all UIC. Figures B-I to IV in the Supplementary file illustrate the average feature participation in clusters.

UIC0 is named Focused, which includes adults, primarily males, with science and technology as the occupation sector and focuses on searching and sharing information but verifying selected content. NetVenturer or UIC1 has a higher frequency of adult females in education, finance, and insurance occupations who use the internet for serious and casual activity but have a higher need to verify data. UIC2 or Aware includes males who use the internet for targeted content, update information, and use OSN for relaxation and staying connected to the world. The Committed label is assigned to UIC3, which is a late-age adult male who uses the internet and social media to explore; however, they believe in investigating content before accepting or sharing. UIC4 consists of more female adults who are either unemployed or homemakers and use the internet and social media for casual surfing, socializing, and occasionally verifying information. UIC4 is labeled as Casual Surfer.

The user characteristics resulting from [35] are generalized to obtain a template to annotate test data. Table 4 presents the user intent template. The User profile was shown to domain experts who verified the transition process and generic translation, after which it was used as ground truth or golden rule for labeling test user data.

**TABLE 5.** User profile template developed from UIC characteristics extracted from clustering [35].

| Demograph | | | Search Openness (SO) | | Online Extravert (OE) | | Information Conscientiousness (IC) | | Cases | | Label | Short label |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Gender | Age | Occupation sector | Info Medium (IM) | Browsing Preference (GP) | Share Content Type (SCT) | Share Medium (SM) | IC Medium (ICM) | Information Credibility Criteria (ICC) | Verify Content (VC) | Share Any Content (SAC) | | |
| Male | Adult | Sci. & Tech. | General Social (GS/S) | Navigational Informational Transactional (NIT/NT) | Particular Casual News Wellbeing (PCN/PCNW/CN) | Private | General (G/0) | Content | yes | no | **Focused** | F |
| Female | Adult Late Age Adult | Education | General Social Traditional (GST) | Navigational Informational Transactional (NIT/NT) | Particular Casual Wellbeing (PCW) CW/C/W/PCW/PC /PW) | all | General Specialist (GS) | Content Publisher (CP) | no | no | **NetVenturer** | NV |
| Male | Late Age Adult | Finance, Insurance | General Traditional (GT/G/T) | Informational Transactional (IT/I/T) | Particular News Wellbeing (PNW) | Private | General Specialist (GS) | Content Author Publisher (CAP/AP) | no/yes | yes | **Aware** | A |
| Male | Late Age Adult | Others | General Social (GS/S) | Navigational Informational (NI) | Particular Casual News Wellbeing (PCN/PCNW/CN) | all | Specialist (S) | Content Author (CA/A) | yes | no | **Committed** | Co |
| Female | Adult | Unemployed | GS (GS/S) | Navigational (N) | Particular News (P/N/0/PN) | all | General Specialist (GS) | Publisher (P) | no | yes | **Casual Surfer** | CS |

## B. TESTING USER INTENT CLASS MODEL

The UIC classifier is trained and tested on data derived from a user feedback study on online information searching, sharing, and verification. The model is further tested on user dynamic behavior comprising actions and activities collected from users' search history and online social network activities.

### 1) USER DYNAMIC INTERACTIONS DATA ACQUISITION

A second user study (US2) is conducted to obtain online users' search, share, and verification activity logs. Participation is voluntary, and the invitation is sent to university students, professional groups, and alumni circles.

The number of participants is capped at 30. The participants are provided an orientation session and are informed about the nature of the study, its duration, and data privacy. Finally, 20 volunteers are accepted based on the demographic distribution used in the User Intent Class (UIC) Model, their commitment to volunteer till the completion of the study, and their informed consent to use their data. The study includes five tasks: browsing, searching, and social media activities. The tasks are open-ended as well as targeted. For example, in one task, participants are asked to search for information related to work using any online medium. In another task, participants are requested to surf on a social platform of their choice and list content type and any action, such as liking or sharing, they take on that content. The final participants are trained on the tasks and Google data service. Due to privacy concerns, the participants are asked to review their activity data before submitting it, whereas the social media data are recorded on the form provided by the participants. The activity data is recorded on the user's Google account and submitted via JotForm. Besides performing the given task, participants also provided a few months of activity log.
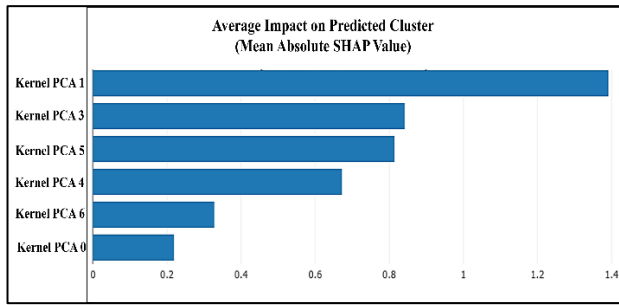
### 2) DATA TRANSFORMATION

Figure 5a provides a mapping scheme of US2 features to UIC model features. The figure illustrates the UIC model features, artifacts extracted, and actions defined to map model features and data extracted from artifacts. The Search behavior features are translated from search activities extracted from the user search, Chrome, and YouTube history, along with participants' input from tasks 1 and 2. Three to four months of history logs records for each participant are reviewed and extracted for transformation. Similarly, sharing behavior features are mapped using Social media activities extracted from Google device activity and participants' own recording of actions on performing relevant tasks (tasks 2,3 and 4). Information verification features, on the other hand, are derived from data extracted from appropriate tasks (Tasks 3, 4, and 5), along with search and browse activities extracted from history. Verification frequency is computed from tasks as well as activity history.
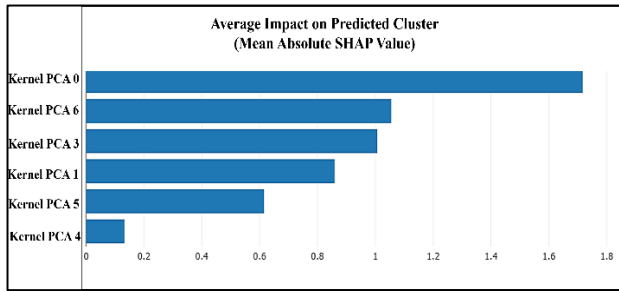
The features extracted from US2 are translated manually using Website category lookup, Website ranking, Search Engine Optimization services, and participants' own recording of the activities while performing tasks in US2. The website category for URLs extracted from the history files is determined using Website Categorization API,[5] which uses IAB Content Taxonomy.[6] The taxonomy offers a 3-tier category system and is used by various research, e.g., [58] to
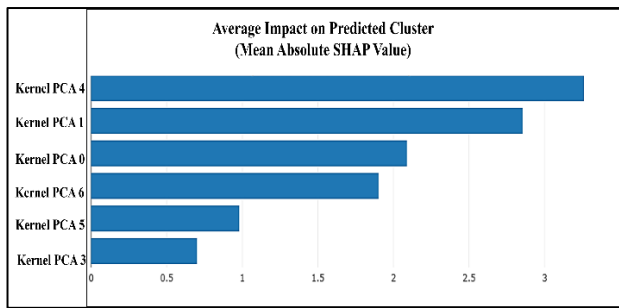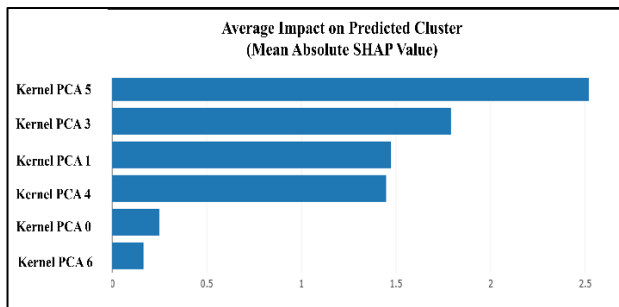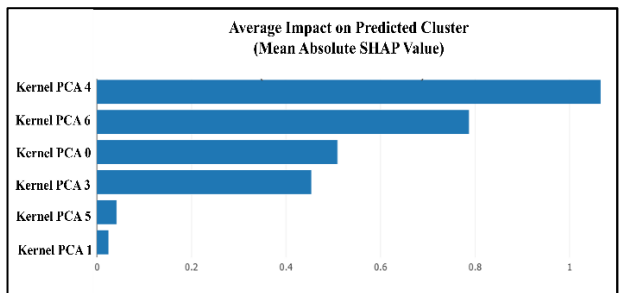
[5] https://website-categorization.whoisxmlapi.com/api
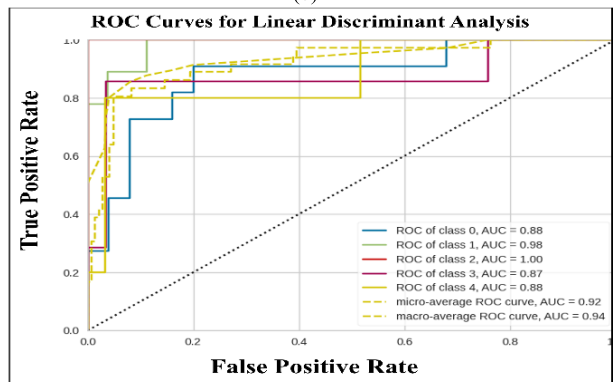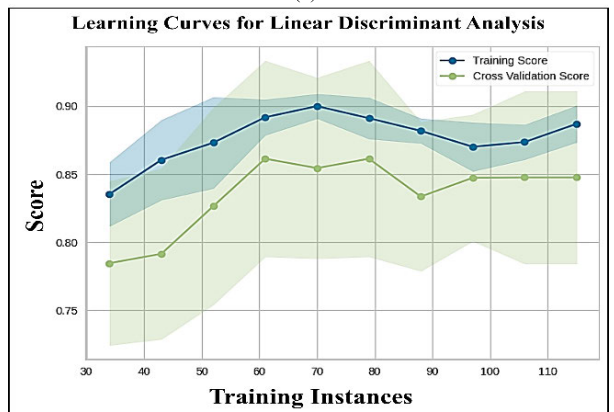[6] https://iabtechlab.com/standards/content-taxonomy/
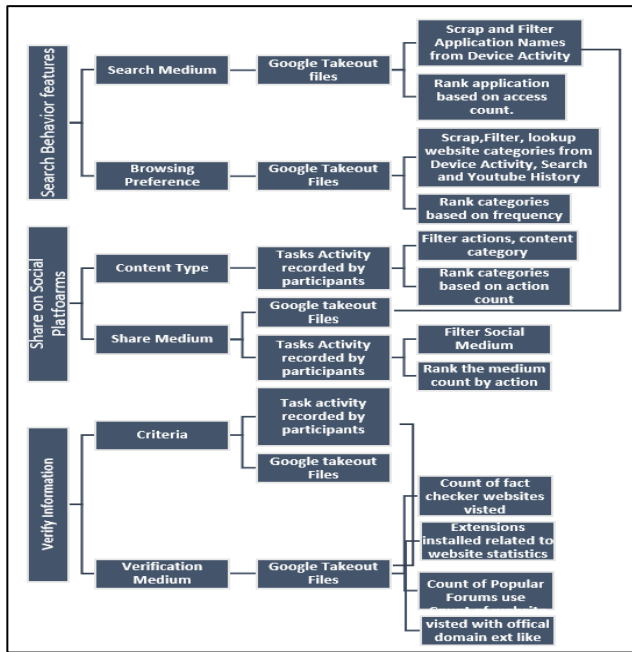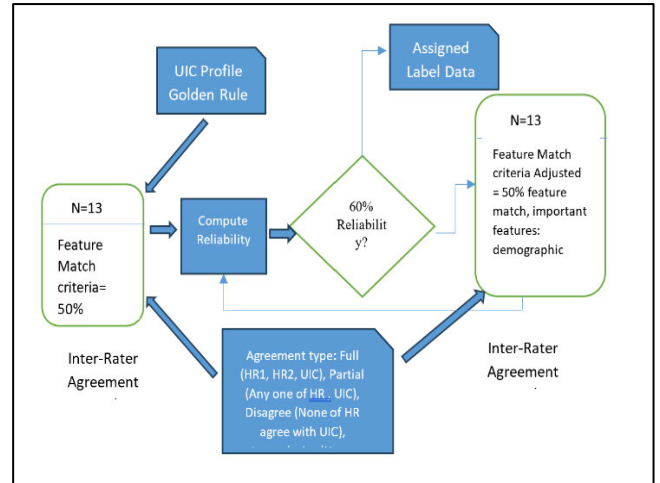
(a)



(b)



(c)



(d)



(e)

**FIGURE 3.** (a) SHAP score for each KPC in UIC0. (b)SHAP score for each KPC in UIC1. (c) SHAP score for each KPC in UIC2. (d) SHAP score for each KPC in UIC3. (e) SHAP score for each KPC in UIC4.
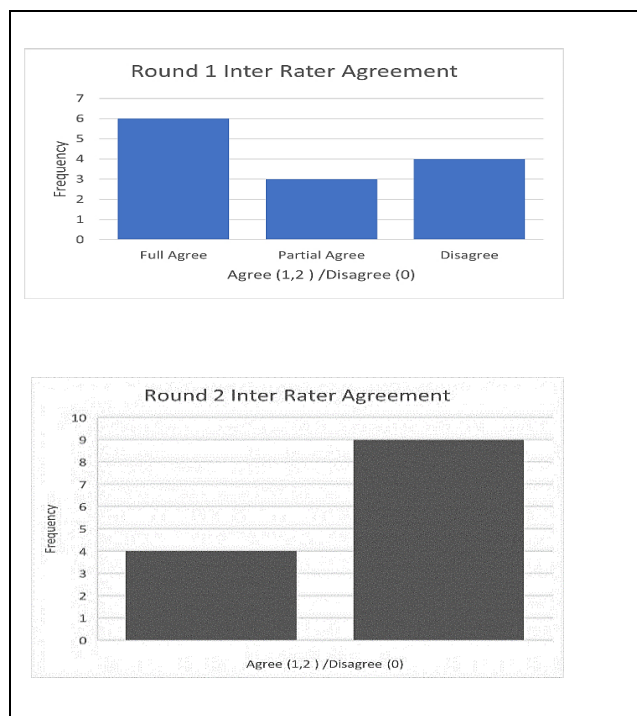


(a)



(b)



(c)



(d)

**FIGURE 4.** (a) Confusion matrix for LDA classifier. (b)Class report for LDA classifier. (c) ROC curves for LDA. (d) LDA learning curves for training and test data.

(a)



(b)



(c)

**FIGURE 5.** (a) User study (US2) data transformation scheme for mapping UIC model features. (b) US2 Annotation And Inter-Rater Reliability Process. (c) Data annotation agreement for rounds 1 and 2.

provide website category lookup service. The three levels of categories are combined into 2 level systems reflecting the categories defined by UIC model features, namely browsing preferences and Sharing on OSN content type. Ahrefs[7]

---

[7]https://ahrefs.com/

and Alexa Web Ranking Services [59] determine the ranking of visited websites. The rank is determined by unique visitors' frequency and visit frequency [60] and provides credibility assessment to some level. Domain details and URL extensions like org, gov, and edu are also used to verify the source veracity; e.g., [61] have been used as one

of the features to check the news source credibility. The fact-checking website also provides services to check website credibility [62].

For verifying and sharing any content cases, task 2 and 3 in US2 are considered explicitly concerning user actions. Share, forward, and post are regarded as a post or share case, whereas browsing or read comments are considered as verify or investigate the case; the rest are regarded as ignore or discarded cases. These transformations map attitude features in Table 2 of the UIC model.

### 3) VALIDATE UIC MODEL PERFORMANCE ON USER INTERACTIONS

The data annotation process is carried out to annotate or assign UIC labels to the US2 data transformed. Inter-Rater Reliability (IRR) is used to validate the annotations. Data Annotation, or data coding or labeling, is a process involving domain experts, specialists, or end users to label the data so it can be used further for automated systems. Data annotation may involve human or sophisticated machine learning algorithms or both [63].

In this research, a hybrid approach, i.e., humans and machines, are used as annotators to label US2 transformed data and validate UIC model prediction. The annotators consist of two humans and the UIC model. Figure 5b outlines the annotation and agreement process. The human raters are technology professionals with graduate degrees. The UIC profile rules are presented to human raters as a template or golden rule, whereas data (N=13) was presented to the UIC model to predict the labels. Inter-Rater Reliability using Fleiss Kappa [64] is computed to reach an agreement between the three annotators or raters. Fliess Kappa ($\kappa$) is used as it supports the nominal or categorical nature of data and is used when raters are fixed and have a fixed number of records.

The Fleiss kappa ($\kappa$) [64] is defined as:

$$\kappa = \frac{P^{'} - P^{'}e}{1 - P^{'}e} \tag{8}$$

where P' Is Observed Agreement And P'E Is Expected Agreement on Random Judgment And Defined By:

$$P^{'} = \frac{1}{pc(c-1)}\left(\sum_{i=1}^{p}\sum_{j=1}^{r} c_{ij}^2 - pc\right) \tag{9}$$

$$P^{'}e = \frac{1}{pc}\sum_{j=1}^{r}\left(\sum_{i=1}^{p} c_{ij}\right)^2 \tag{10}$$

where p represents the 13 participants' records, c represents 5 UIC labels, and r represents three raters. 1-P'e is the degree of agreement reachable above the chance, and P'-P'e is the actual degree of agreement above the chance. $\kappa$ is 1 if a complete agreement is reached and 0 if there is no agreement. The agreeable criterion is agreed upon and varies from case to case. The IRR agreement presented by [65] is followed. The threshold is set to moderate agreement, i.e., 60% or above,

The process is completed in two rounds. The criterion for annotation is set to 50% coverage of rules, whereas all features were given equal weights in the first round. The agreement between human raters and machines is divided into three groups. Full agreement (2) is reached if all three raters agree on the label. Partial agreement (1) is observed if at least one human rater agrees with the machine label. The agreement is inconclusive if human raters agree, but machines disagree. It is considered disagreement (0) if all three raters have different labels or are inconclusive. The reliability threshold is set as 60% or greater. After 1st round, reliability was 40%, below the threshold. In 2nd round, it is decided between the raters that in case of the sample fitting multiple labels, demographic features will be given priority in determining the label. Some of the partial agreements are changed to full agreements. At the end of round 2, reliability is calculated as 60%. Figures D-I and II in the Supplementary file show the intermediate calculations using Fliess Kappa (eq. 8). The annotation agreement for both rounds is shown in Figure 5c. The x-axis in the figure represents the Rater agreement group, and the y-axis represents the agreement count.

Based on the IRR, the model accuracy and F1 score are also computed, and model performance on limited user interaction data is considered acceptable with 67% accuracy and an F1 score of 0.70.

## III. RESULTS AND DISCUSSION

This research aims to identify users' characteristics based on their search preferences, information-sharing intention, and trust in the information. The objective is to develop a framework to identify behavioral patterns and classify users based on online behavior and practices. The user intent machine learning model uses the K-Mean clustering technique to identify behavioral patterns in users, group them into clusters, and use a Linear Discriminant Analysis (LDA) classifier to predict the cluster. The result shows 80% classification accuracy on cross-validation tests and unseen data, indicating that the model can classify users based on their features. Since clusters represent the user intent to search, share, and verify online information, the clusters are termed User Intent Class or UIC. The UIC model is further tested on the user interaction data acquired through a second user study. The user interactions are mapped on the model features and annotated using the UIC profile template generalized from user characteristics. Two human raters annotated the data with an Inter-Rater Reliability (IRR) of 60%. The UIC model predicted the user with 67% accuracy.

The following sections discuss the results in detail.

### A. UIC MODEL PERFORMANCE ON USER FEEDBACK

The clusters generated by K-Mean have an average silhouette coefficient of 0.34. This indicates that the clusters have overlapping membership in some clusters, such as clusters 1 and 3 Figure 2a. Clusters 4, 2, and 0, with silhouette score 0.5, are relatively well separated than clusters 1 and 3. The lower volume of data, outliers interfering in cluster density, and shifting cluster centroid may contribute to a lower silhouette

coefficient. However, the cluster classification shows good results, increasing the performance of the entire model.

Linear Discriminant Analysis (LDA) technique is used to classify UI Clusters. Multiple classification modes are evaluated to determine the best model based on AUC, accuracy, F1, and SHAP values. The models used classic PCA and Kernel PCA for dimension reduction. Extra trees, Ridge, and Logistic Regression also performed well besides LDA, However, LDA performance is better on the individual class level (Figure C I-III in the Supplementary file shows the Precision, Recall, and f1-score for classification models). Similarly, LDA with KPCA showed better Shap values than the other models. Fig 4- a and b indicate every UIC or cluster's f1 score on test data is above 70%, which is satisfactory, though UIC1 or cluster 1 has a higher false negative. However, class true prediction is 100% accurate. Figure 4d shows the learning curve for the LDA model on training and test data. The learning curve shows that the model could improve if presented with more data, but the cross-validation curve indicates the steady performance of the model in the end. The ROC curve in Figure 4c also shows stable convergence of each class. The LDA model also performed well on unseen data with precision, recall, and f1 as 0.79, 0.86, and 0.81, respectively.

The Shap values provide feature significance and participation in the model. The kpc4 (kernelpca4) in Net Venturer Class (UIC1) indicated by Figure 3a, kpc6 (kernelpca6) in Committed Class (UIC3) explained by Figure 3b, and kpc5 and 0 (kernelpca5 and kernelpca0) in Casual Surfer Class (UIC4) highlighted by Figure 3c shows almost no contribution in the model. In contrast, most features show some contribution in the rest of UIC by Figure 3d and Figure 3e. This shows that users from these classes have good feature participation.

## B. UIC MODEL PERFORMANCE ON USER ONLINE INTERACTIONS

The User Intent Class Model is tested and validated on users' browsing and social activities collected through a second user study (US2). The study gathers browsing activities, user social interaction for disseminating information, and actions taken to verify information through various tasks. Users must walk through the tasks presented with scenarios on internet browsers and social media and share the information. The information collected is filtered and mapped to UIC model features. Human Annotators are then consulted to label the data using the UIC user profile template. UIC Model is also used as a machine annotator alongside humans. The Inter-Rater Reliability (IRR) using the Fleiss kappa method is calculated to validate the annotation. The resultant kappa after two rounds of the data annotation process is 60%, which passes the criterion for IRR agreement, thus validating the prediction of the UIC model on user activity data. The UIC model shows an accuracy of 67% and an F1 score of 71% on US2 data, which can be considered acceptable, provided the limited data size with even lower support for each class.

The results also show that the UIC model can predict the user group on the user's online activities associated with information seeking, distributing, and trusting information. The User Profile template is validated through IRR and can be used as the Golden Rule.

## C. COMPARING RELATED WORK

The Internet has impacted users of every walk of life. People use the internet for all purposes, and much research focuses on improving peoples' experience by understanding the need to use the internet. User profiling is one way of grouping users based on their needs, preferences, intentions, and activities; hence, different sets of user profiles that have employed different statistical and machine learning methods and multimodal features have emerged. Reference [38] used linguistic features to classify different types of Twitter users, and [22] used structured data and applied clustering techniques to group Instagram users on commerce activities. This research has targeted user searching behavior, information sharing, and verification intent to generate a generic user profile. Moreover, a user prediction model is also developed to identify user types. The data was derived from user feedback rather than scrapped from the internet as privacy concerns have discontinued many social platforms for sharing data.

## IV. CONCLUSION

This research aims to answer two questions: a) Can the user intention be classified based on their behavior and practices in searching, sharing, and verifying information? b) How can the classification model be validated using user interactions? To achieve this, the research uses a hybrid machine learning approach to successfully clustering similar users according to their online search, sharing, and verification behaviors and facilitating the classification of new users. The unsupervised K-Means Clustering technique groups users with similar features, whereas Linear Discriminant Analysis (LDA) is a supervised classifier to predict user classes. The clustering process resulted in identifying five User Intent Classes (UIC), and LDA classified the test data with approximately 80% accuracy. The analysis resulted in five different user profiles, which are named Committed (those who balance work and leisure), Casual Surfer (search and share for leisure with minimal need to authenticate information), NetVenturer (use the internet for all purposes, but keep eyes open), Focused (Mainly for serious work) and Aware (UpToDate, wise internet usage). These findings can be explored further to tailor online services to use generic user profiles that can be securely used to increase the user experience by providing the information the user generally intends. The research uses user online interaction and activities on the UIC model to test its efficacy. Human raters annotate the interaction data and LDA model, and Inter-rater reliability (IRR) is computed to validate the annotation, User profile, and UIC model prediction. The experts accepted the IRR of 60% and successfully validated the model performance in predicting users on dynamic

data. The research contributes to developing a framework for profiling users based on unique three-dimensional data, i.e., user search intent, information dissemination needs, and content credibility criteria. These aspects of user behavior have been researched separately but not as a component of user intent. Information integrity is an important aspect in this digital era, where information is accessed from all online platforms. The research shows a correlation between search and share, whereas verification adds weight to the user's truth-seeking desire. Insights from this study can help in better design of search filters, Search Engine Optimization (SEO), targeted advertising, and shopping experience based on user profiling.

## V. LIMITATION AND FUTURE WORK
The research captured specific demographic data; therefore, more extensive and diverse data is required for further generalization. Diversity was initially aimed for but was restricted by regional reachability and funding restrictions. Payment-based user platforms like Amazon Mechanical Turk can be used to enroll participants of diverse backgrounds. More sophisticated models, like neural nets, can be tried with a more extensive set. The research didn't aim for a real-time model that can also be focused in the future. In addition, further work can be done to capture live user actions and create a mapping framework to test the current model.
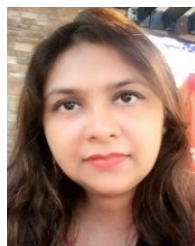
## DECLARATIONS
The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Data collected in this research is part of Ph.D. research and cannot be made available now. However, it can be made available upon reasonable request in the future.

## REFERENCES
[1] L. S. Vygotsky, *Mind in Society: The Development of Higher Psychological Processes*. Cambridge, MA, USA: Harvard Univ. Press, 1978.
[2] K. A. Mills, "Shrek meets vygotsky: Rethinking adolescents' multimodal literacy practices in schools," *J. Adolescent Adult Literacy*, vol. 54, no. 1, pp. 35–45, Sep. 2010.
[3] A. Halevy, C. Canton-Ferrer, H. Ma, U. Ozertem, P. Pantel, M. Saeidi, F. Silvestri, and V. Stoyanov, "Preserving integrity in online social networks," *Commun. ACM*, vol. 65, no. 2, pp. 92–98, Feb. 2022.
[4] X. Feng, X. Wang, and Y. Zhang, "Research on the effect evaluation and the time-series evolution of public culture's Internet communication under the background of new media: Taking the information dissemination of red tourism culture as an example," *J. Comput. Cultural Heritage*, vol. 16, no. 1, pp. 1–15, Mar. 2023.
[5] C. I. Eke, A. A. Norman, L. Shuib, and H. F. Nweke, "A survey of user profiling: State-of-the-art, challenges, and solutions," *IEEE Access*, vol. 7, pp. 144907–144924, 2019.
[6] J. Liu, M. Mitsui, N. J. Belkin, and C. Shah, "Task, information seeking intentions, and user behavior: Toward a multi-level understanding of web search," in *Proc. Conf. Human Inf. Interact. Retr.*, Glasgow, U.K., Mar. 2019, pp. 123–132.
[7] J. L. Hale, B. J. Householder, and K. L. Greene, "The theory of reasoned action," in *The Persuasion Handbook: Developments in Theory and Practice*, vol. 14. Newbury Park, CA, USA: Sage, 2002, pp. 259–286.
[8] J. Shi, P. Hu, K. K. Lai, and G. Chen, "Determinants of users' information dissemination behavior on social networking sites: An elaboration likelihood model perspective," *Internet Res.*, vol. 28, no. 2, pp. 393–418, Apr. 2018.
[9] P. Bedi, S. B. Goyal, A. S. Rajawat, R. N. Shaw, and A. Ghosh, "A framework for personalizing atypical web search sessions with concept-based user profiles using selective machine learning techniques," in *Advanced Computing and Intelligent Technologies* (Lecture Notes in Networks and Systems), vol. 218. Singapore: Springe, 2022.
[10] M. Soleymani, M. Riegler, and P. Halvorsen, "Multimodal analysis of user behavior and browsed content under different image search intents," *Int. J. Multimedia Inf. Retr.*, vol. 7, no. 1, pp. 29–41, Mar. 2018, doi: 10.1007/s13735-018-0150-6.
[11] D. Koehn, S. Lessmann, and M. Schaal, "Predicting online shopping behaviour from clickstream data using deep learning," *Expert Syst. Appl.*, vol. 150, Jul. 2020, Art. no. 113342.
[12] H. Yoganarasimhan, "Search personalization using machine learning," *Manage. Sci.*, vol. 66, no. 3, pp. 1045–1070, Mar. 2020.
[13] T. Ruotsalo, J. Peltonen, M. J. Eugster, D. Głowacka, P. Floreen, P. Myllymaki, G. Jacucci, and S. Kaski, "Interactive intent modeling for exploratory search," *ACM Trans. Inf. Syst.*, vol. 36, no. 4, p. 44, Oct. 2018, doi: 10.1145/3231593.
[14] S. K. Shivakumar, "A survey and taxonomy of intent-based code search," *Int. J. Softw. Innov.*, vol. 9, no. 1, pp. 69–110, Jan. 2021.
[15] P. Ren, Z. Liu, X. Song, H. Tian, Z. Chen, Z. Ren, and M. de Rijke, "Wizard of search engine: Access to information through conversations with search engines," in *Proc. 44th Int. ACM SIGIR Conf. Res. Develop. Inf. Retr.*, Jul. 2021.
[16] A. Salle, S. Malmasi, O. Rokhlenko, and E. Agichtein, "Studying the effectiveness of conversational search refinement through user simulation," in *Advances in Information Retrieval* (Lecture Notes in Computer Science). Cham, Switzerland: Springer, 2021, pp. 587–602, doi: 10.1007/978-3-030-72113-8_39.
[17] S. Dixon. *Statista*. Accessed: Jan. 2024. [Online]. Available: https://www.statista.com/statistics/272014/global-social-networks-ranked-by-number-of-users/
[18] S.-F. Tsao, H. Chen, T. Tisseverasinghe, Y. Yang, L. Li, and Z. A. Butt, "What social media told us in the time of COVID-19: A scoping review," *Lancet Digit. Health*, vol. 3, no. 3, pp. e175–e194, Mar. 2021.
[19] Z. Wu, Y. Liu, Q. Zhang, K. Wu, M. Zhang, and S. Ma, "The influence of image search intents on user behavior and satisfaction," in *Proc. 12th ACM Int. Conf. Web Search Data Mining*, Melbourne, VIC, Australia, Jan. 2019, pp. 645–653.
[20] D. Fraszczak, "Information propagation in online social networks—A simulation case study," in *Proc. 38th Int. Bus. Inf. Manag. Assoc.*, Seville, Spain, 2021, pp. 23–24.
[21] X. Lin and X. Wang, "Examining gender differences in people's information-sharing decisions on social networking sites," *Int. J. Inf. Manage.*, vol. 50, pp. 45–56, Feb. 2020.
[22] A. Mohammed and A. Ferraris, "Factors influencing user participation in social media: Evidence from Twitter usage during COVID-19 pandemic in Saudi Arabia," *Technol. Soc.*, vol. 66, Aug. 2021, Art. no. 101651.
[23] S.-H. Liao, R. Widowati, and C.-J. Cheng, "Investigating Taiwan Instagram users' behaviors for social media and social commerce development," *Entertainment Comput.*, vol. 40, Jan. 2022, Art. no. 100461, doi: 10.1016/j.entcom.2021.100461.
[24] X. Wang, F. Chao, G. Yu, and K. Zhang, "Factors influencing fake news rebuttal acceptance during the COVID-19 pandemic and the moderating effect of cognitive ability," *Comput. Hum. Behav.*, vol. 130, May 2022, Art. no. 107174, doi: 10.1016/j.chb.2021.107174.
[25] R. P. Yu, "How types of Facebook users approach news verification in the mobile media age: Insights from the dual-information-processing model," *Mass Commun. Soc.*, vol. 24, no. 2, pp. 233–258, 2021, doi: 10.1080/15205436.2020.1839104.
[26] S. Edgerly, R. R. Mourão, E. Thorson, and S. M. Tham, "When do audiences verify? How perceptions about message and source influence audience verification of news headlines," *Journalism Mass Commun. Quart.*, vol. 97, no. 1, pp. 52–71, Mar. 2020.
[27] X. Zhou, K. Shu, V. V. Phoha, H. Liu, and R. Zafarani, "This is fake! Shared it by mistake," in *Proc. ACM Web Conf.*, Lyon, France, Apr. 2022.
[28] A. Apejoye, "Comparative study of social media, TV and newspapers news credibility," in *Proc. Int. Conf. Media Commun., Technol. Design*, Apr. 2015.
[29] S. Raza and C. Ding, "Fake news detection based on news content and social contexts: A transformer-based approach," *Int. J. Data Sci. Analytics*, vol. 13, no. 4, pp. 335–362, May 2022.

[30] L. Konstantinovskiy, O. Price, M. Babakar, and A. Zubiaga, "Toward automated factchecking: Developing an annotation schema and benchmark for consistent automated claim detection," *Digit. Threats, Res. Pract.*, vol. 2, no. 2, pp. 1–16, Jun. 2021.

[31] G. S. Cheema, S. Hakimov, and R. Ewerth, "Check_square at CheckThat! 2020: Claim detection in social media via fusion of transformer and syntactic features," 2020, *arXiv:2007.10534*.

[32] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, and H. Liu, "FakeNewsNet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media," *Big Data*, vol. 8, no. 3, pp. 171–188, Jun. 2020.

[33] Y. Li, B. Jiang, K. Shu, and H. Liu, "MM-COVID: A multilingual and multimodal data repository for combating COVID-19 disinformation," 2020, *arXiv:2011.04088*.

[34] M. K. Elhadad, K. F. Li, and F. Gebali, "COVID-19-FAKES: A Twitter (Arabic/English) dataset for detecting misleading information on COVID-19," in *Proc. 12th Int. Conf. Intell. Netw. Collaborative Syst.*, 2020, pp. 256–268.

[35] M. Kanwal, N. A. Khan, and A. A. Khan, "A machine learning approach to user profiling for data annotation of online behavior," *Comput., Mater. Continua*, vol. 78, no. 2, pp. 2419–2440, 2024, doi: 10.32604/cmc.2024.047223.

[36] M. Heidari, J. H. Jones, and O. Uzuner, "Deep contextualized word embedding for text-based online user profiling to detect social bots on Twitter," in *Proc. Int. Conf. Data Mining Workshops (ICDMW)*, Nov. 2020, pp. 480–487.

[37] E. Brusa, L. Cibrario, C. Delprete, and L. G. Di Maggio, "Explainable AI for machine fault diagnosis: Understanding features' contribution in machine learning models for industrial condition monitoring," *Appl. Sci.*, vol. 13, no. 4, p. 2038, Feb. 2023.

[38] A. Agresti, *Categorial Data Analysis*. Hoboken, NJ, USA: Wiley, 2013.

[39] S. A. Alasadi and W. S. Bhaya, "Review of data preprocessing techniques in data mining," *J. Eng. Appl. Sci.*, vol. 12, no. 16, pp. 4102–4107, 2017.

[40] C. Seger, "An investigation of categorical variable encoding techniques in machine learning: Binary versus one-hot and feature hashing," School Elect. Eng. Comput. Sci. (EECS), Uppsala Univ. Library, Sweden, Tech. Rep. URN: urn:nbn:se:kth:diva-237426, DiVA, id: diva2:1259073, 2018.

[41] P. Cerda and G. Varoquaux, "Encoding high-cardinality string categorical variables," *IEEE Trans. Knowl. Data Eng.*, vol. 34, no. 3, pp. 1164–1176, Mar. 2022.

[42] A. Bommert, X. Sun, B. Bischl, J. Rahnenführer, and M. Lang, "Benchmark for filter methods for feature selection in high-dimensional classification data," *Comput. Statist. Data Anal.*, vol. 143, Mar. 2020, Art. no. 106839.

[43] U. M. Khaire and R. Dhanalakshmi, "Stability of feature selection algorithm: A review," *J. King Saud Univ. Comput. Inf. Sci.*, vol. 34, no. 4, pp. 1060–1073, Apr. 2022.

[44] P. S. Wright, "Adjusted P-values for simultaneous inference," *Biometrics*, vol. 48, no. 4, pp. 1005–1013, Dec. 1992, doi: 10.2307/2532694.

[45] D. J. Mundfrom, J. J. Perrett, j. Schaffer, A. Piccone, and M. Roozeboom, "Bonferroni adjustments in tests for regression coefficients," *Multiple Linear Regression Viewpoints*, vol. 32, no. 1, pp. 1–6, 2006.

[46] A. E. Ezugwu, A. M. Ikotun, O. O. Oyelade, L. Abualigah, J. O. Agushaka, C. I. Eke, and A. A. Akinyelu, "A comprehensive survey of clustering algorithms: State-of-the-art machine learning applications, taxonomy, challenges, and future research prospects," *Eng. Appl. Artif. Intell.*, vol. 110, Apr. 2022, Art. no. 104743.

[47] J.-O. Palacio-Niño and F. Berzal, "Evaluation metrics for unsupervised learning algorithms," 2019, *arXiv:1905.05667*.

[48] G. Oyewole and G. Thopil, "Data clustering: Application and trends," *Artif. Intell. Rev.*, vol. 56, pp. 6439–6475, Nov. 2022.

[49] I. H. Sarker, "Machine learning: Algorithms, real-world applications and research directions," *Social Netw. Comput. Sci.*, vol. 2, no. 3, p. 160, May 2021.

[50] A. Pandey, T. R. Nair, and S. B. Thomas, "Combination of K-means clustering and support vector machine for instrument detection," *Social Netw. Comput. Sci.*, vol. 3, no. 2, p. 121, Mar. 2022.

[51] R. Kumar, G. Ganapathy, and J. J. Kang, "A hybrid mod K-means clustering with mod SVM algorithm to enhance the cancer prediction," *Int. J. Internet, Broadcast. Commun.*, vol. 13, no. 2, pp. 231–243, 2021.

[52] C.-H. Weng and C.-K. Huang, "A hybrid machine learning model for credit approval," *Appl. Artif. Intell.*, vol. 35, no. 15, pp. 1439–1465, Dec. 2021.

[53] F. Kherif and A. Latypova, "Principal component analysis," in *Machine Learning*. New York, NY, USA: Academic, 2020, pp. 209–225.

[54] F. L. Gewers, G. R. Ferreira, H. F. D. Arruda, F. N. Silva, C. H. Comin, D. R. Amancio, and L. D. F. Costa, "Principal component analysis: A natural approach to data exploration," *ACM Comput. Surv.*, vol. 54, no. 4, pp. 1–34, 2021.

[55] R. Tan, J. R. Ottewill, and N. F. Thornhill, "Monitoring statistics and tuning of kernel principal component analysis with radial basis function kernels," *IEEE Access*, vol. 8, pp. 198328–198342, 2020.

[56] E. K. Ampomah, Z. Qin, and G. Nyame, "Evaluation of tree-based ensemble machine learning models in predicting stock price direction of movement," *Information*, vol. 11, no. 6, p. 332, Jun. 2020.

[57] G. Hackeling, *Mastering Machine Learning With Scikit-Learn*. Birmingham, U.K.: Packt Publishing, 2017.

[58] V. Ponce-López and C. Spataru, "Social media behaviour analysis in disaster-response messages of floods and heat waves via artificial intelligence," *Comput. Inf. Sci.*, vol. 15, no. 3, pp. 18–36, Jun. 2022.

[59] Y. Akgül, "Accessibility, usability, quality performance, and readability evaluation of university websites of turkey: A comparative study of state and private universities," *Universal Access Inf. Soc.*, vol. 20, no. 1, pp. 157–170, Mar. 2021.

[60] S. Aggarwal, H. Van Oostendorp, Y. R. Reddy, and B. Indurkhya, "Providing web credibility assessment support," in *Proc. Eur. Conf. Cognit. Ergonom.*, Sep. 2014, pp. 1–8.

[61] N. Sitaula, C. K. Mohan, J. Grygiel, X. Zhou, and R. Zafarani, "Credibility-based fake news detection," in *Disinformation, Misinformation, and Fake News in Social Media*. Cham, Switzerland: Springer, 2020, pp. 163–182.

[62] PressBook. (2017). *Pressbook Fact Checker List*. [Online]. Available: https://pressbooks.pub/webliteracy/chapter/fact-checking-sites/

[63] X. Wu, L. Xiao, Y. Sun, J. Zhang, T. Ma, and L. He, "A survey of human-in-the-loop for machine learning," *Future Gener. Comput. Syst.*, vol. 135, pp. 364–381, Oct. 2022.

[64] J. L. Fleiss, "Measuring nominal scale agreement among many raters," *Psychol. Bull.*, vol. 76, no. 5, pp. 378–382, Nov. 1971.

[65] J. R. Landis and G. G. Koch, "The measurement of observer agreement for categorical data," *Biometrics*, vol. 33, no. 1, p. 159, Mar. 1977.

**MOONA KANWAL** (Member, IEEE) received the B.Sc. and M.Sc. degrees in computer engineering from the Sir Syed University of Engineering and Technology, Pakistan, in 2002 and 2007, respectively. She is currently pursuing the Ph.D. degree with the Department of Computer Science, NED University of Engineering and Technology, Pakistan. From 2008 to 2009, she was a Compliance and the SQA Manager with Mixit Technologies, Pakistan. From 2020 to 2022, she was a Professor (part-load) with Seneca College, Canada. Since 2007, she has been an Assistant Professor with the Sir Syed University of Engineering and Technology. She is also a mentor for many projects and startups. She received funding from the Ministry of Science and Technology Pakistan and a scholarship from the Higher Education Commission Pakistan. Her research interests include machine learning, data sciences, human–computer interaction, and human behavior.

**MUZAMMIL AHMAD KHAN** (Member, IEEE) received the B.S. and M.S. degrees in computer engineering from the Sir Syed University of Engineering and Technology, Karachi, Pakistan. He is currently pursuing the Ph.D. degree with the NED University of Engineering and Technology, Karachi. He is an Assistant Professor with the Department of Computer Engineering, Sir Syed University of Engineering and Technology. His research interests include cloud computing, the IoT, cyber security, networks, artificial intelligence, data science, e-commerce, and project management.

**NAJMA ISMAT** received the bachelor's, master's, and Ph.D. degrees, in 1998, 2002, and 2018, respectively. She is currently with the Sir Syed University of Engineering and Technology. She has published several international publications and presented several papers at international and national conferences. Her research interests include mobility, reliability, connectivity, and coverage issues in underwater sensor networks, wireless sensor networks, and the IoT.

**AFTAB A. KHAN** received the master's and Ph.D. degrees in developmental psychology and education from the University of Toronto, Canada. He is currently an Associate Professor of special education with the Department of Education and Counselling, Longwood University. He teaches both graduate and undergraduate courses related to the field of special education. He has over 20 years of experience in the field of special education. He was involved in a large research project on Wisdom Around the World. This includes his experience teaching students at the university level, training professionals at the school level, consulting on autism spectrum disorder, and training parents in the community are a few of the many professional practices he has accomplished.

**NAJEED A. KHAN** received the first M.Sc. degree in mathematics from the University of Karachi, Pakistan, in 1990, the second M.Sc. degree in computer sciences from the NED University of Engineering and Technology, Pakistan, in 1996, and the Ph.D. degree in computer vision from the University of Leeds, U.K., in 2011. He is currently a Professor of artificial intelligence, and an Executive Officer and the Caretaker MoST Chair Professor endowment at the NED University of Engineering and Technology. He has supervised multiple Ph.D. scholars who have completed their degrees. He was awarded multiple research national and international funding grants. He won the HEC National Centre for Artificial Intelligence (NCAI), NED University of Engineering and Technology, as a Co-PI. He is the author of more than 30 publications in JCR and ISI-indexed journals. His research interests include computer vision specific to application areas of medical imaging. He is a Senior Member of the International Association of Computer Science and Information Technology (IACSIT), a member of Asian Council of Science Editors (ACSE), and a Professional Member of the Association of Computing Machinery (ACM).

● ● ●