# E-commerce Product Classification

## Problem Statement

**Background Information:** E-commerce has become a prominent industry in recent years, with a wide variety of products available for online purchase. As the number of products continues to grow, it becomes crucial for e-commerce platforms to effectively classify and categorize these products. Accurate product classification enables better search functionality, personalized recommendations, and improved user experience.

One of the key challenges in product classification is the manual effort required to assign accurate categories to each product. This process is time-consuming and prone to errors. To overcome this challenge, machine learning techniques can be employed to automate the product classification task. By training a classifier on a labeled dataset, it becomes possible to predict the category of a product based on its description.

**Objective:** The objective of this project is to develop a machine learning model that can accurately classify e-commerce products into different categories based on their descriptions. The model will be trained on a labeled dataset containing product descriptions and their corresponding categories.

**Dataset Information:** The dataset used for this project is provided in a CSV file format. It consists of two columns: "description" and "label". The "description" column contains the text descriptions of the products, while the "label" column contains the corresponding category assigned to each product. The categories include "Electronics," "Household," "Books," and "Clothing & Accessories." The dataset has been preprocessed to remove missing values and duplicate observations.

**Approach:** The project follows a step-by-step approach to achieve the objective:

1. Data Loading and Preprocessing: The dataset is loaded from the CSV file and processed to handle missing values and duplicate observations. The "label" column is manually encoded to numeric values for further analysis.
2. Exploratory Data Analysis: The dataset is analyzed to gain insights into the distribution of product descriptions across different categories. Histograms are generated to visualize the distribution of the number of characters, words, and average word length in the descriptions for each category.
3. Feature-Target Split: The dataset is divided into features (product descriptions) and the target variable (categories).
4. Train-Test Split: The data is split into training and testing sets. The training set will be used to train the machine learning model, while the testing set will be used to evaluate its performance.
5. Text Preprocessing: The text descriptions undergo several preprocessing steps, including converting to lowercase, removing whitespaces, removing punctuations, removing HTML tags, removing emojis, removing other Unicode characters, converting acronyms, converting contractions, removing stopwords, and performing spell checking.
6. TF-IDF Vectorization: The preprocessed text descriptions are converted into numerical feature vectors using the TF-IDF (Term Frequency-Inverse Document Frequency) technique. This step transforms the text data into a format suitable for machine learning algorithms.
7. Model Selection and Training: Several machine learning classifiers, including Logistic Regression, K-Nearest Neighbors, Decision Tree, Support Vector Machine (SVM), Random Forest, Stochastic Gradient Descent, Ridge Classifier, XGBoost, and AdaBoost, are employed for training and classification. The classifiers are evaluated based on their accuracy scores.
8. Model Evaluation: The trained models are evaluated using performance metrics such as accuracy score. The model with the highest accuracy is selected as the final model for product classification.

**Conclusion:** By developing an accurate product classification model, e-commerce platforms can automate the categorization process and improve user experience. The model trained on the provided dataset will enable the classification of e-commerce products into different categories based on their descriptions. This project showcases

the application of machine learning techniques in the e-commerce domain and highlights the potential benefits of automated product classification.

# Framework

1. **Import Required Libraries:** Begin by importing the necessary libraries for data manipulation, visualization, and machine learning. Commonly used libraries include pandas, numpy, matplotlib, sklearn, and nltk.
2. **Load the Dataset:** Use the appropriate function to load the dataset from the CSV file into a pandas DataFrame. Verify the successful loading of the dataset by inspecting a few rows using the head() or sample() function.
3. **Data Preprocessing:** Perform data preprocessing steps to handle missing values and duplicate observations. If any missing values are present, consider using techniques such as filling with a default value or dropping the rows. Remove duplicate observations using the drop_duplicates() function.
4. **Exploratory Data Analysis (EDA):** Conduct exploratory data analysis to gain insights into the dataset. Generate descriptive statistics and visualize the distribution of product descriptions across different categories using histograms or bar plots. Calculate and plot the average number of characters, words, and average word length in the descriptions for each category.
5. **Feature-Target Split:** Divide the dataset into features (product descriptions) and the target variable (categories). Assign the product descriptions to a variable, such as X, and the corresponding categories to another variable, such as y.
6. **Train-Test Split:** Split the data into training and testing sets using the train_test_split() function from sklearn. Specify the proportion of the dataset to allocate for testing (e.g., 80% for training, 20% for testing). This step ensures the model's performance is evaluated on unseen data.
7. **Text Preprocessing:** Preprocess the text descriptions to convert them into a format suitable for machine learning algorithms. Implement a series of text preprocessing steps, such as converting to lowercase, removing whitespaces, removing punctuations, removing HTML tags, removing emojis, removing other Unicode characters, converting acronyms, converting contractions, removing stopwords, and performing spell checking. You can utilize libraries like nltk or regex for these preprocessing tasks.
8. **TF-IDF Vectorization**: Transform the preprocessed text descriptions into numerical feature vectors using the TF-IDF technique. Apply the TfidfVectorizer class from sklearn to convert the text data into a numerical representation. This

step calculates the TF-IDF values for each word in the dataset and creates a feature matrix.

9.  **Model Selection and Training:** Choose a set of machine learning classifiers suitable for text classification tasks, such as Logistic Regression, K-Nearest Neighbors, Decision Tree, Support Vector Machine (SVM), Random Forest, Stochastic Gradient Descent, Ridge Classifier, XGBoost, or AdaBoost. Iterate over each classifier, create an instance of the model, and train it using the training set generated in step 6. Evaluate the model's performance using appropriate metrics, such as accuracy score, cross-validation, or precision/recall.

10. **Model Evaluation:** Compare the performance of the trained models using evaluation metrics. Calculate the accuracy score for each model on the testing set. Select the model with the highest accuracy score as the final model for product classification.

11. **Conclusion and Further Improvements:** Summarize the results obtained from the project, highlighting the accuracy and effectiveness of the chosen model for e-commerce product classification. Discuss potential areas of improvement, such as fine-tuning model parameters, exploring different feature engineering techniques, or using more advanced deep learning models for better performance.

12. **Documentation and Presentation:** Document the code by adding comments and explaining the code's functionality, especially for complex sections. Create a detailed report or presentation summarizing the project, including the problem statement, approach, results, and future work.

By following this detailed outline, one can write the code for the E-commerce Product Classification project, effectively implementing data preprocessing, exploratory data analysis, feature engineering, model training, evaluation, and documentation stages.

# Code Explanation

**Section 1:** Importing the Required Libraries In this section, the necessary libraries for data manipulation, visualization, and machine learning are imported. These libraries provide various functions and tools that will be used throughout the code. For example, pandas is used for data manipulation, numpy for numerical operations, matplotlib for data visualization, and sklearn for machine learning algorithms.

**Section 2:** Loading the Dataset This section focuses on loading the dataset from a CSV file into a pandas DataFrame. The dataset contains information about e-commerce products and their respective categories. By using the appropriate function, the dataset is read and stored in the DataFrame. This allows for easy access and manipulation of the data.

**Section 3:** Data Preprocessing Data preprocessing is an essential step in any machine learning project. In this section, the code handles missing values and duplicate observations. Missing values are commonly dealt with by either filling them with a default value or removing the rows entirely. Duplicate observations are removed to ensure that the data is clean and representative of the actual dataset.

**Section 4:** Exploratory Data Analysis (EDA) EDA involves analyzing and visualizing the dataset to gain insights and better understand the data. Descriptive statistics are generated to summarize the dataset, providing information like the mean, median, and standard deviation of the product descriptions. Visualizations such as histograms or bar plots are created to show the distribution of product descriptions across different categories. These visualizations help in identifying patterns and potential relationships in the data.

**Section 5:** Feature-Target Split In this section, the dataset is divided into features and the target variable. The features are the product descriptions, and the target variable is the corresponding category of each product. This split allows the machine learning model to learn from the descriptions and predict the categories accurately.

**Section 6:** Train-Test Split To evaluate the model's performance on unseen data, the dataset is split into training and testing sets. The training set is used to train the model, while the testing set is used to evaluate its performance. This split ensures that the

model's accuracy is measured on data it has never seen before, providing a more reliable evaluation metric.

**Section 7:** Text Preprocessing Text preprocessing involves transforming the raw text descriptions into a format suitable for machine learning algorithms. Several preprocessing steps are performed, such as converting text to lowercase, removing whitespaces and punctuation, and removing stopwords (commonly used words that carry little meaning). These steps help in reducing noise and standardizing the text data, making it easier for the model to learn meaningful patterns.

**Section 8:** TF-IDF Vectorization TF-IDF vectorization converts the preprocessed text descriptions into numerical feature vectors. It calculates the importance of each word in the dataset based on its frequency and rarity across all documents. The resulting feature matrix contains numerical representations of the text data, which can be used as input for machine learning models.

**Section 9:** Model Selection and Training This section involves selecting a suitable machine learning model for text classification and training it on the training set. Different classifiers, such as Logistic Regression, Decision Tree, or Random Forest, can be used. Each classifier is instantiated, trained on the training set, and then evaluated using appropriate metrics, such as accuracy or cross-validation. This step helps in identifying the model that performs the best for the given task.

**Section 10:** Model Evaluation After training and evaluating multiple models, this section focuses on comparing their performances. The accuracy score is calculated for each model on the testing set. The model with the highest accuracy score is considered the best-performing model for product classification. This evaluation step helps in determining the effectiveness of the trained models and selecting the most suitable one.

**Section 11:** Predicting Categories Finally, the selected model is used to predict the categories for new, unseen product descriptions. The preprocessed and vectorized text data is fed into the trained model, and it outputs the predicted category for each product. These predictions can be used for various purposes, such as organizing products on an e-commerce website or automating categorization tasks.

By following this code's workflow, one can effectively preprocess the data, explore and visualize it, train and evaluate different models, and ultimately make predictions on new product descriptions.

# Future Work

**1. Expand and Enhance the Dataset** To improve the accuracy and generalizability of the model, it's important to have a diverse and comprehensive dataset. Collect additional product descriptions from various sources or consider incorporating user-generated content, such as customer reviews or comments. Ensure that the dataset covers a wide range of product categories to provide a more robust training and testing environment.

**2. Advanced Text Preprocessing** Consider implementing more advanced text preprocessing techniques to further improve the quality of the text data. This could involve techniques such as stemming (reducing words to their base form) or lemmatization (reducing words to their dictionary form) to handle different word variations. Additionally, explore techniques for handling misspelled words or abbreviations commonly found in product descriptions.

**3. Feature Engineering** Incorporate additional features that could enhance the model's performance. For example, consider extracting numerical features from the product descriptions, such as the length of the description or the presence of specific keywords related to certain categories. These additional features can provide valuable information to the model and potentially improve its predictive accuracy.

**4. Experiment with Different Machine** Learning Models Explore a wider range of machine learning models beyond the ones used in the initial code. Consider models specifically designed for text classification tasks, such as Support Vector Machines (SVM), Naive Bayes, or deep learning models like Recurrent Neural Networks (RNNs) or Transformer-based models like BERT. Compare the performances of these models using appropriate evaluation metrics and select the one that achieves the best results.

**5. Hyperparameter** Tuning Optimize the selected machine learning model by tuning its hyperparameters. Hyperparameters are adjustable settings that control the learning process of the model. Use techniques like grid search or random search to systematically explore different combinations of hyperparameters and identify the optimal configuration. This step can significantly improve the model's performance.

**6. Implement Model Deployment** Once the best-performing model is identified, deploy it into a production environment for real-world usage. This could involve

integrating the model into an e-commerce platform or creating a standalone application for product categorization. Consider using frameworks like Flask or Django to build a web-based application that accepts product descriptions as input and returns predicted categories as output.

**7. Continuous Model Improvement** Regularly update and retrain the model to adapt to changing trends and new product categories. As new data becomes available, periodically retrain the model to incorporate the latest information and improve its performance over time. Implement a process for collecting user feedback and iteratively improving the model based on user suggestions and corrections.

**Step-by-Step Implementation Guide:**

1. Collect a diverse dataset of product descriptions and their corresponding categories.
2. Preprocess the text data by converting it to lowercase, removing stopwords, punctuation, and other irrelevant characters.
3. Split the dataset into features (product descriptions) and the target variable (categories).
4. Split the data into training and testing sets for model evaluation.
5. Vectorize the preprocessed text data using techniques like TF-IDF or word embeddings.
6. Select and train various machine learning models on the training set.
7. Evaluate the performance of each model using appropriate metrics on the testing set.
8. Choose the best-performing model based on evaluation results.
9. Fine-tune the hyperparameters of the selected model using techniques like grid search or random search.
10. Deploy the selected model into a production environment, such as an e-commerce platform or a standalone application.
11. Regularly update and retrain the model with new data to improve its accuracy and adaptability.

By following this step-by-step guide, you can enhance the project by expanding the dataset, improving preprocessing techniques, experimenting with different models, tuning hyperparameters, deploying the model, and continuously improving its performance.

# Concept Explanation

Let's dive into the concept of the algorithm used in this project in a friendly and funny manner. Imagine you have a magical assistant named AIce Cream, who loves categorizing products just as much as it loves ice cream flavors! Let's call the algorithm "AIce Cream Classifier" for now.

So, what does this AIce Cream Classifier do? Well, it's like having a super-powered brain that can read and understand product descriptions to automatically assign them to the right category. It's like having a friend who can instantly tell you if a product is a phone, a shirt, or a fancy kitchen gadget just by reading its description. Impressive, right?

To make it work, AIce Cream Classifier needs to learn from a bunch of examples. It's like training a little ice cream monster to recognize different flavors. You show it a variety of ice cream flavors—vanilla, chocolate, strawberry, you name it—and tell the monster which flavor each one is. Similarly, for our algorithm, we feed it a dataset of product descriptions along with their corresponding categories. We let the AIce Cream Classifier learn from these examples so that it can start making predictions on its own.

Once the AIce Cream Classifier has learned from these examples, it uses its newfound knowledge to categorize new products. You can give it a product description like "A sleek, touch-screen device with a high-resolution camera" and watch it do its magic. The algorithm analyzes the words and phrases in the description, just like how you'd analyze the ingredients and texture of ice cream, and then confidently predicts the most suitable category for the product—let's say it predicts "Electronics" in this case.

But how does AIce Cream Classifier make these predictions? It uses its secret ingredient called "machine learning." Just as you might experiment with different ice cream recipes, the algorithm experiments with different mathematical models and equations to find the best way to understand the product descriptions.

It breaks down the words in the descriptions into tiny pieces, just like you'd break down ingredients into flavors and textures. It looks for patterns and connections between these pieces to figure out which words are important for each category. For example, it might discover that words like "screen," "camera," and "processor" are often associated with the "Electronics" category, just like how "chocolate chips," "fudge swirls," and "creamy texture" are associated with a chocolate ice cream flavor.

Once the AIce Cream Classifier has learned these patterns, it can use them to make predictions on new, unseen product descriptions. It's like giving it a mystery ice cream sample and asking it to identify the flavor based on what it has learned so far. And just like a good ice cream taster, the AIce Cream Classifier gets better and better with more training and practice.

But hey, it's not perfect! Sometimes, just like when you get a brain freeze, the AIce Cream Classifier might make a mistake. It might misclassify a product or get confused by tricky descriptions. That's why it's important to keep improving and updating the algorithm over time. Just like you'd keep trying new ice cream flavors and adjusting your taste preferences, the AIce Cream Classifier can learn from its mistakes and become even more accurate.

So, with the AIce Cream Classifier, you have a friendly and magical assistant who can categorize products with the power of machine learning. It's like having an ice cream-loving superhero helping you sort through all those products and make sense of their descriptions.

I hope this concept explanation brings a smile to your face and helps you understand the algorithm in a fun and friendly way. If you have any more questions, feel free to ask!

# Exercise Questions

**Question 1: Explain the concept of feature extraction and how it is used in the AIce Cream Classifier algorithm.**

**Answer:** Feature extraction is a process where relevant information or characteristics, known as features, are extracted from raw data to be used for analysis or prediction. In the context of the AIce Cream Classifier algorithm, feature extraction involves identifying important words or phrases in product descriptions that can help determine the category of the product.

The AIce Cream Classifier algorithm analyzes the dataset of product descriptions and identifies patterns and connections between words and categories. It looks for specific keywords or phrases that are frequently associated with certain categories. These identified words or phrases act as the features or indicators for each category. For example, words like "camera" and "processor" might be important features for the "Electronics" category.

During the prediction phase, the algorithm examines new product descriptions and checks for the presence of these important features. If the description contains words or phrases that match the identified features, the algorithm assigns a higher probability to the corresponding category. By extracting relevant features, the AIce Cream Classifier algorithm can make accurate predictions based on the presence or absence of these features in product descriptions.

**Question 2: What are some potential challenges or limitations of the AIce Cream Classifier algorithm?**

**Answer:** While the AIce Cream Classifier algorithm is a powerful tool for categorizing products based on their descriptions, it does have certain limitations. Some potential challenges include:

1. Ambiguous descriptions: Product descriptions can sometimes be vague or open to interpretation, making it difficult to determine the correct category. For example, a description like "versatile gadget" might not provide enough context to accurately classify the product.
2. Out-of-vocabulary words: The algorithm relies on the presence of specific words or phrases as features. If a product description contains uncommon or new words

that were not present in the training data, the algorithm may struggle to assign an appropriate category.

3. Overfitting: If the algorithm is trained on a dataset that is too small or not representative enough, it may become overly specialized and fail to generalize well to new, unseen data. This phenomenon is known as overfitting, where the algorithm performs well on the training data but poorly on new data.

4. Imbalanced datasets: If the dataset used to train the algorithm is imbalanced, meaning some categories have significantly more examples than others, it may bias the predictions towards the more prevalent categories. This can result in inaccurate predictions for the less represented categories.

5. Contextual understanding: While the AIce Cream Classifier algorithm can identify important keywords, it may struggle to understand the overall context of a product description. It may miss subtle nuances or rely too heavily on specific words without considering the broader context of the description.

Addressing these challenges and limitations requires continuous improvement and refinement of the algorithm, along with regular updates to the training data to ensure better accuracy and robustness.

## Question 3: How would you evaluate the performance of the AIce Cream Classifier algorithm?

**Answer:** Evaluating the performance of the AIce Cream Classifier algorithm involves assessing its accuracy and ability to correctly predict the categories of products. Here are some evaluation measures that can be used:

1. Accuracy: The accuracy of the algorithm can be calculated by comparing its predictions with the known true categories in a test dataset. It is the ratio of correct predictions to the total number of predictions made.

2. Precision and Recall: Precision measures the proportion of correctly predicted positive instances (products correctly classified within a category) out of all instances predicted as positive. Recall, on the other hand, measures the proportion of correctly predicted positive instances out of all actual positive instances in the dataset. Both precision and recall provide a more nuanced view of the algorithm's performance, especially when dealing with imbalanced datasets.

3. F1 Score: The F1 score is a metric that combines precision and recall into a single value. It balances the trade-off between precision and recall and provides a comprehensive evaluation of the algorithm's performance.
4. Confusion Matrix: A confusion matrix visualizes the algorithm's performance by showing the number of correct and incorrect predictions for each category. It helps identify any patterns of misclassification or bias and provides insights into the algorithm's strengths and weaknesses.

These evaluation measures can be used to gauge the performance of the AIce Cream Classifier algorithm and guide further improvements and optimizations.

## Question 4: How can you enhance the AIce Cream Classifier algorithm to improve its accuracy?

**Answer:** There are several ways to enhance the AIce Cream Classifier algorithm to improve its accuracy:

1. Increase Training Data: Providing a larger and more diverse dataset during the training phase can help the algorithm learn a wider range of features and improve its ability to generalize to new data.
2. Feature Engineering: Instead of relying solely on individual words or phrases, feature engineering involves extracting more complex features from the product descriptions. This could include analyzing the syntactic structure, identifying important word combinations, or considering the context of sentences.
3. Word Embeddings: Utilizing word embeddings, such as Word2Vec or GloVe, can capture semantic relationships between words. By representing words in a vector space, the algorithm can better understand the similarities and differences between different words in the product descriptions.
4. Regularization Techniques: Applying regularization techniques like dropout or L1/L2 regularization can prevent overfitting and improve the algorithm's generalization capabilities.
5. Ensemble Methods: Implementing ensemble methods, such as combining multiple classifiers or using a combination of different algorithms, can often lead to improved accuracy by leveraging diverse perspectives and reducing individual biases.

By implementing these enhancements, the AIce Cream Classifier algorithm can become more accurate and reliable in categorizing products based on their descriptions.

**Question 5: Can you propose an alternative approach to the AIce Cream Classifier algorithm?**

**Answer:** Yes, an alternative approach to the AIce Cream Classifier algorithm could involve using a deep learning technique called recurrent neural networks (RNNs). RNNs are particularly useful when dealing with sequential data, such as product descriptions.

In this alternative approach, instead of relying solely on individual words or features, an RNN can take into account the sequential nature of the words in a description. Each word in the description is fed into the RNN one at a time, and the network maintains an internal memory or hidden state, which captures information from previous words.

By processing the entire description through the RNN, it can capture the contextual information and dependencies between words, allowing for a more comprehensive understanding of the description. This can potentially lead to improved accuracy in categorizing products.

However, implementing RNNs requires additional considerations, such as handling vanishing or exploding gradients and selecting appropriate architectures, such as Long Short-Term Memory (LSTM) or Gated Recurrent Units (GRU). Training RNNs also typically requires larger datasets and more computational resources compared to traditional machine learning algorithms.

By exploring alternative approaches like RNNs, we can further enhance the AIce Cream Classifier algorithm and potentially achieve even better performance in categorizing products based on their descriptions.