# Paper Subject Prediction

## Problem Statement

**Background:** In the field of research and academia, there is a vast amount of scholarly papers being published every day. Classifying these papers into specific subject categories can be a challenging and time-consuming task. Automating the process of subject prediction can greatly assist researchers, educators, and students in quickly finding relevant papers in their respective domains.

**Objective:** The objective of this project is to build a machine learning model that can predict the subject category of a given scholarly paper based on its summary. The model will be trained on a dataset of papers with their corresponding subject labels, and it will learn to associate specific keywords or patterns in the paper summaries with the appropriate subject categories.

**Dataset:** The dataset used in this project is called "arxiv_data.csv." It contains a collection of scholarly papers from the arXiv repository. The dataset includes two main columns: "summaries" and "terms." The "summaries" column contains the brief summaries of the papers, while the "terms" column contains the subject categories assigned to each paper. Each paper can have multiple subject categories assigned to it, represented as a list of terms.

| summaries | terms |
|------------------------------------------------|----------------------------------|
| This paper presents a novel approach... | [machine learning, computer vision] |
| In this study, we propose a new... | [natural language processing] |
| The objective of this research... | [data mining, pattern recognition] |

In the example above, the first paper is about a novel approach in the field of machine learning and computer vision. The second paper focuses on a new technique in natural language processing. The third paper deals with research in data mining and pattern recognition.

**Approach:** To solve this problem, we will use a machine learning approach. The project code provided implements a multi-label classification model using TensorFlow and Keras. The model is trained on the dataset, where the paper summaries are used as input features, and the subject categories are the target labels. The code performs exploratory data analysis to understand the dataset, handles data preprocessing, splits the dataset into training, validation, and test sets, and finally trains a shallow multi-layer perceptron (MLP) model.

The MLP model consists of multiple dense layers with ReLU activation functions, followed by a final dense layer with a sigmoid activation function. The model is trained using the binary cross-entropy loss function and the Adam optimizer. The training history is plotted to visualize the model's performance.

# Framework

1) **Import Dependencies:** Begin by importing the necessary dependencies for the project, including TensorFlow, Keras, pandas, numpy, and matplotlib.

2) **Load and Explore the Dataset:**
- Use the pd.read_csv() function to load the dataset from the "arxiv_data.csv" file.
- Print the first few rows of the dataset using the head() function.
- Use the len() function to determine the total number of rows in the dataset.
- Check for duplicate titles in the dataset and remove them if any exist.
- Perform exploratory data analysis (EDA) to gain insights into the dataset:
- Identify the number of unique terms and the occurrence of each term.
- Filter out rare terms by removing rows with terms that appear only once.
- Convert the string representation of terms to a list using the literal_eval() function.
- Split the dataset into training, validation, and test sets using the train_test_split() function.

3) **Data Preprocessing:**
- Use TensorFlow's StringLookup to encode the terms as multi-hot vectors.
- Create a vocabulary lookup table using the StringLookup layer.
- Invert the multi-hot encoded labels to obtain the original terms using a helper function.
- Determine the maximum sequence length for the input summaries.
- Define the batch size and padding token.

4) **Create Datasets:**
- Create a function, make_dataset(), to convert the dataframes into TensorFlow datasets.
- Use the lookup() function to binarize the labels.
- Shuffle the training dataset and batch the datasets.
- Apply the make_dataset() function to create training, validation, and test datasets.

5) **Text Vectorization:**
- Create a TextVectorization layer to vectorize the text data.
- Use the adapt() function to fit the vectorizer on the training dataset.
- Map the vectorizer to the datasets using the map() function.

6) **Model Definition and Training:**

- Define the model architecture by creating a function, make_model(), that returns a sequential model.
- Compile the model with binary cross-entropy loss and the Adam optimizer.
- Train the model using the fit() function, passing the training and validation datasets and the number of epochs.
- Plot the training and validation loss and binary accuracy using the plot_result() function.

7) **Model Evaluation and Inference:**

- Evaluate the model's performance on the test dataset using the evaluate() function.
- Create a model for inference by combining the vectorizer and the trained model.
- Create a small inference dataset from the test dataset.
- Perform inference on the samples in the inference dataset:
- Print the input summary and the original labels.
- Obtain the predicted probabilities for each label.
- Sort the probabilities and retrieve the top three predicted labels.
- Print the predicted labels for comparison.

By following this outline, someone can write the code for the "Paper Subject Prediction" project based on the provided code. They can customize and expand upon the code to further improve the model's performance and explore additional evaluation metrics and techniques.

# Code Explanation

1) **Import Dependencies:** The code begins by importing necessary dependencies. These are libraries that provide pre-built functions and tools to help us accomplish our tasks. In this case, we are using TensorFlow and Keras for building and training our machine learning model, as well as other libraries like pandas, numpy, and matplotlib for data manipulation and visualization.

2) **Load and Explore the Dataset:** This part of the code focuses on loading and exploring the dataset. The dataset is stored in a CSV file named "arxiv_data.csv". We use the pd.read_csv() function to read the file and create a pandas dataframe, which is like a table that holds our data. We then print the first few rows of the dataframe to get a glimpse of the data.

3) **Perform Exploratory Data Analysis (EDA):** After loading the dataset, we want to understand it better through exploratory data analysis. We start by checking the number of rows in the dataset using the len() function. Next, we identify and remove any duplicate titles in the dataset. We then analyze the occurrence of terms and filter out rare terms that appear only once. Lastly, we convert the string representation of terms into a list using the literal_eval() function.

4) **Data Preprocessing:** Before training our model, we need to preprocess the data. This involves transforming the text and labels into a suitable format for machine learning. We start by setting up some variables like the maximum sequence length and batch size. These values help ensure that our data is processed efficiently. We also define a padding token, which helps standardize the length of input sequences.

5) **Create Datasets:** To train our model, we need to create training, validation, and test datasets. We split the filtered dataset into these three sets using the train_test_split() function. Then, we use TensorFlow's StringLookup to encode the terms as multi-hot vectors, which is a way to represent categorical data. This helps our model understand the different subjects or labels associated with each paper.

6) **Text Vectorization:** In this step, we transform the text data (the summaries of the papers) into a numerical format that can be processed by our machine learning model. We use the TextVectorization layer provided by Keras to convert the text into vectors. This helps the model understand the patterns and relationships within the text.

7) **Model Definition and Training:** Now it's time to define our machine learning model and train it using the training dataset. We define a shallow multi-layer perceptron (MLP) model with several dense layers. These layers are like building blocks that process and transform the data. We compile the model with a loss function and an optimizer, which help guide the learning process. Then, we train the model using the training dataset and evaluate its performance on the validation dataset. We track the loss and binary accuracy during training and visualize them using plots.

8) **Model Evaluation and Inference:** After training our model, we evaluate its performance on the test dataset to see how well it generalizes to new, unseen data. We print the categorical accuracy score, which indicates the model's accuracy in predicting the correct subjects or labels. We also create a model for inference, which combines the text vectorizer and the trained model. This allows us to make predictions on new samples. We create a small dataset for inference, print the input summaries, original labels, and predicted labels for comparison.

By understanding each step and function in the provided code, we can see how the project progresses from loading and preprocessing the data to training and evaluating the machine learning model. It's an exciting journey that involves transforming text data into numerical representations, building a model, and making predictions on new samples.

# Future Work

In this project, we have built a machine learning model to predict the subjects or categories of research papers based on their summaries. To further enhance the project and explore additional improvements, we can consider the following future work:

**1. Data Augmentation**: One way to improve the performance of the model is to augment the dataset. Data augmentation involves generating new samples by applying various transformations to the existing data. For text data, this can include techniques like word replacement, synonym substitution, or sentence shuffling. By increasing the diversity and quantity of data, we can improve the model's ability to generalize and make more accurate predictions.

**2. Fine-tuning the Model:** To achieve better performance, we can explore fine-tuning the model architecture and hyperparameters. This involves tweaking the structure of the neural network, adding or removing layers, adjusting the number of neurons, and experimenting with different activation functions. Additionally, we can optimize the learning rate, batch size, and number of training epochs to find the best configuration for our specific task.

**3. Utilizing Pre-trained Language Models:** Pre-trained language models, such as BERT, GPT, or RoBERTa, have shown significant success in natural language processing tasks. We can leverage these models by fine-tuning them on our specific dataset. Pre-trained models have already learned rich contextual representations from large amounts of text data, which can greatly benefit our subject prediction task. Fine-tuning a pre-trained model can improve accuracy and capture more nuanced relationships in the text.

**4. Handling Imbalanced Classes:** In the dataset, we may encounter imbalanced class distributions, where some subject categories have significantly fewer samples compared to others. To address this issue, we can employ techniques such as oversampling the minority classes, undersampling the majority classes, or using class weights during training. These techniques help to ensure that the model learns from all classes effectively, rather than being biased towards the majority class.

**5. Deploying the Model:** After developing and fine-tuning the model, we can deploy it as a production-ready system. This involves creating an API (Application Programming Interface) that can receive input summaries and return the predicted subject categories.

We can use frameworks like Flask or FastAPI to develop a web-based interface where users can interact with the model. Additionally, we may consider containerization using Docker for easier deployment and scalability.

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

To implement the future work for the paper subject prediction project, follow these step-by-step instructions:

1) Data Augmentation: Research and implement data augmentation techniques suitable for text data. Apply these techniques to generate additional samples from the existing dataset.
2) Fine-tuning the Model: Experiment with different model architectures and hyperparameters. Make adjustments to the neural network structure, activation functions, and optimization settings. Train and evaluate the model with these modifications to identify the best configuration.
3) Utilizing Pre-trained Language Models: Choose a pre-trained language model suitable for the subject prediction task. Fine-tune the selected model on the dataset by adapting it to the specific classification task. Train and evaluate the fine-tuned model to assess its performance.
4) Handling Imbalanced Classes: Analyze the class distribution in the dataset and identify any imbalanced classes. Implement techniques such as oversampling, undersampling, or class weighting to address the class imbalance. Train the model using the balanced dataset and evaluate its performance.
5) Deploying the Model: Develop an API using frameworks like Flask or FastAPI to expose the model's functionality. Create an interface where users can input the summaries of research papers and receive predicted subject categories. Consider containerization using Docker for easier deployment and management.

By following these steps, you can enhance the paper subject prediction project by implementing data augmentation, fine-tuning the model, utilizing pre-trained language models, handling imbalanced classes, and deploying the model as a production-ready system.

# Concept Explanation

Alright, buckle up and get ready to dive into the fascinating world of paper subject prediction! Let me break it down for you in a friendly and funny manner.

Imagine you have a massive stack of research papers, each with a title and a summary. Now, your task is to figure out the subject or category of each paper based on its summary. Sounds challenging, right? But don't worry, we've got an algorithm to save the day!

The algorithm we're using is like a super-smart detective who knows how to uncover the hidden secrets of text data. It's called a machine learning algorithm, specifically a deep learning model. Think of it as a Sherlock Holmes of text analysis!

**Here's how our Sherlock Holmes algorithm works:**

1) **Data Detective:** Our algorithm starts by gathering a bunch of research papers with their titles, summaries, and subject labels. It carefully examines each paper, trying to find patterns and connections between the words in the summaries and their corresponding subjects. It's like the algorithm is putting on its detective hat and magnifying glass, looking for clues!

2) **Training Time:** Once our algorithm has studied the papers, it enters a special training phase. During this phase, it learns from the examples in the data and tries to understand the relationships between the words in the summaries and the subjects they belong to. It's like the algorithm is attending a training academy to become the ultimate paper subject predictor!

3) **Neural Network Ninja:** The core of our algorithm is a powerful neural network. This network is made up of interconnected nodes called neurons, and it's trained to recognize and interpret the patterns in the data. It's like our algorithm has become a ninja, using its neural network powers to slice through the complexity of text and extract meaningful information.

4) **Prediction Magic:** Once our algorithm completes its training, it becomes a prediction wizard! When it encounters a new research paper with a summary, it analyzes the words in the summary and feeds them into its neural network. The network then performs a series of calculations and magically produces a prediction, telling us which subject(s) the paper belongs to. It's like our algorithm is waving a prediction wand and revealing the secrets hidden within the text!

For example, let's say we have a paper with the summary: "This study investigates the effects of caffeine on sleep quality." Our algorithm, with its superpowers, might predict that the subject of this paper is "Sleep Medicine" or "Pharmacology." It's like our algorithm is whispering, "Aha! This paper is all about caffeine and sleep!"

And there you have it! Our paper subject prediction algorithm is like a detective, a trainee, a ninja, and a wizard all in one. It uses the power of machine learning and neural networks to unravel the mysteries of text data and predict the subjects of research papers.

Remember, behind all the complexity and technical jargon, there's always a touch of magic and a dash of humor in the world of algorithms!

# Exercise Questions

1) **Question: What is the purpose of performing exploratory data analysis (EDA) in this project? Provide some examples of the EDA performed.**

**Answer:** The purpose of EDA in this project is to gain insights and understand the characteristics of the dataset before building the model. It helps us identify any data quality issues, understand the distribution of data, and make informed decisions for preprocessing and modeling. Some examples of EDA performed in this project include:

- Checking for duplicate titles and removing them to ensure data quality.
- Analyzing the occurrence of terms to identify rare terms and filter them out.
- Exploring the vocabulary size to understand the complexity of the text data.

2) **Question: What is the role of StringLookup and TextVectorization in the project? How do they contribute to the model training?**

**Answer:** StringLookup and TextVectorization are important components in preprocessing the text data for the model.

- StringLookup is used to convert the categorical labels (paper subjects) into a numerical representation suitable for the model. It creates a mapping between the categorical labels and their corresponding numerical values.
- TextVectorization converts the textual summaries into a numerical representation that can be fed into the model. It tokenizes the text, converts it into sequences of numbers, and applies techniques like n-grams and TF-IDF to capture the importance of words in the summaries. This helps the model understand the text data and make predictions based on it.

3) **Question: Explain the architecture of the shallow MLP (Multi-Layer Perceptron) model used in this project.**

**Answer:** The shallow MLP model consists of three densely connected layers:

- The first layer has 512 neurons and uses the ReLU activation function. It learns complex patterns in the input data.
- The second layer has 256 neurons and also uses the ReLU activation function. It further extracts relevant features from the previous layer's output.
- The final layer has a number of neurons equal to the vocabulary size of the subjects. It uses the sigmoid activation function to produce multi-hot encoded

predictions for each subject. Each output neuron represents a subject, and its activation value indicates the presence or absence of that subject in the paper.

4) **Question: How is the model trained and evaluated in this project? What metrics are used to assess the model's performance?**

**Answer:** The model is trained using the training dataset and evaluated using the validation dataset. During training, the model optimizes the binary cross-entropy loss function using the Adam optimizer. It aims to minimize the difference between the predicted multi-hot encoded labels and the true labels. The model's performance is assessed using two metrics:

- Binary Cross-Entropy Loss: It measures the dissimilarity between the predicted and true labels, giving an indication of how well the model is fitting the data.
- Binary Accuracy: It calculates the percentage of correctly predicted labels, indicating the overall accuracy of the model's predictions.

5) **Question: How is the trained model used for inference? Explain the process of predicting the subject of a new research paper.**

**Answer:** The trained model is used for inference by creating a new model that combines the text vectorization layer and the shallow MLP model. To predict the subject of a new research paper:

- The paper's summary is passed through the text vectorization layer to convert it into a numerical representation.
- The resulting numerical representation is then fed into the shallow MLP model, which produces a multi-hot encoded prediction.
- The predicted probabilities for each subject are ranked, and the top predicted subjects are selected as the final predictions.
- This way, the model utilizes its learned patterns and associations to make predictions on unseen research papers.