# Economic News Articles

---

## Problem Statement

**Background:** In the field of economics, staying up-to-date with relevant news articles is crucial for making informed decisions. However, with the increasing volume of news articles available online, it becomes challenging to manually sort and identify articles that are relevant to specific economic topics. Therefore, developing an automated system to classify the relevance of economics news articles can greatly assist researchers, analysts, and decision-makers in efficiently filtering and accessing the most pertinent information.

**Objective:** The objective of this project is to build a machine learning model that can classify the relevance of economics news articles as either "relevant" or "not relevant" based on their textual content. By accurately predicting the relevance, the model can help in streamlining the process of article categorization and information retrieval.

**Dataset:** The dataset used in this project is stored in a CSV file named 'data.csv'. It contains a collection of economics news articles along with their relevance labels. The dataset includes the following columns:

- text: The textual content of the news articles.
- relevance: The relevance label indicating whether the article is "relevant" or "not relevant" to economics.

**Approach:** The code begins by importing the necessary libraries and dependencies for data processing, visualization, and machine learning. It then loads the dataset from the CSV file and displays its shape, providing an overview of the data.

Next, it performs data cleaning and preprocessing steps to prepare the text data for analysis. The clean function utilizes various techniques such as removing named entities,

converting text to lowercase, removing punctuation and digits, removing stopwords, and lemmatization. These steps help in transforming the raw text into a more standardized and manageable format.

After preprocessing the text, the code uses the TF-IDF (Term Frequency-Inverse Document Frequency) vectorization technique to convert the textual data into numerical features. This transformation allows the machine learning algorithms to process the text data effectively. The TF-IDF vectorizer calculates the importance of each word in the articles based on their frequency and rarity across all the articles.

The code then splits the data into training and testing sets using the train_test_split function. It applies four different classifiers: Gaussian Naive Bayes, Multinomial Naive Bayes, Logistic Regression, and Support Vector Machines (SVM) to train and evaluate the models. For each classifier, it calculates and displays the training and testing accuracy scores, as well as the classification report showing precision, recall, and F1 scores for each class.

Additionally, the code generates confusion matrices to visualize the performance of the classifiers in terms of true positives, true negatives, false positives, and false negatives. It also plots the Receiver Operating Characteristic (ROC) curve and calculates the Area Under the Curve (AUC) to evaluate the model's overall performance.

Finally, the code demonstrates an ensemble approach by using a Voting Classifier, which combines the predictions of multiple classifiers to make a final decision. It shows the training and testing accuracy scores for the ensemble model.

At the end of the code, the predictions for the relevance of all the articles in the dataset are saved to a CSV file named "us-economic-news-relevance.csv".

Conclusion: This project aims to provide a solution for classifying the relevance of economics news articles. By utilizing machine learning algorithms and natural language processing techniques, it offers a way to automate the categorization process and efficiently filter relevant articles from a large collection. The code demonstrates the training and evaluation of multiple classifiers, enabling the selection of the best-performing model for the task at hand.

# Framework

1) **Importing Libraries and Dependencies**
- Import the necessary libraries and dependencies required for data processing, visualization, and machine learning tasks. This may include libraries like Pandas, NumPy, scikit-learn, Matplotlib, and NLTK.

2) **Loading the Dataset**
- Read the dataset from the provided CSV file using a suitable function, such as pd.read_csv().
- Display the shape of the dataset to get an overview of the data.

3) **Data Cleaning and Preprocessing**
- Create a function, such as clean(text), to perform data cleaning and preprocessing on the text data.
- Implement various techniques in the clean() function, such as removing named entities, converting text to lowercase, removing punctuation and digits, removing stopwords, and lemmatization, using libraries like NLTK or spaCy.
- Apply the clean() function to the 'text' column of the dataset to preprocess the text data.

4) **Text Vectorization using TF-IDF**
- Use the TF-IDF (Term Frequency-Inverse Document Frequency) vectorization technique to convert the preprocessed textual data into numerical features.
- Initialize a TF-IDF vectorizer object, such as TfidfVectorizer() from scikit-learn.
- Fit the vectorizer on the preprocessed text data using the fit_transform() function.
- Obtain the feature matrix by transforming the preprocessed text data using the transform() function.

5) **Splitting the Data into Training and Testing Sets**
- Split the dataset into training and testing sets using the train_test_split() function from scikit-learn.
- Specify the feature matrix (TF-IDF representation of text) as the X variable and the 'relevance' column as the y variable.
- Set an appropriate test size and random state for reproducibility.

## 6) Training and Evaluating Classifiers

- Initialize and train multiple classifiers using different algorithms available in scikit-learn, such as Gaussian Naive Bayes, Multinomial Naive Bayes, Logistic Regression, and Support Vector Machines (SVM).
- For each classifier, fit the model on the training data and evaluate its performance on the testing data.
- Calculate and display the training and testing accuracy scores for each classifier.
- Generate a classification report showing precision, recall, and F1 scores for each class using the classification_report() function.
- Visualize the performance of the classifiers using confusion matrices to see true positives, true negatives, false positives, and false negatives.
- Plot the Receiver Operating Characteristic (ROC) curve for each classifier and calculate the Area Under the Curve (AUC) to assess the overall model performance.

## 7) Ensemble Approach: Voting Classifier

- Implement an ensemble approach using a Voting Classifier to combine the predictions of multiple classifiers and make a final decision.
- Initialize a Voting Classifier object, such as VotingClassifier() from scikit-learn, and provide the list of trained classifiers as estimators.
- Fit the Voting Classifier on the training data and evaluate its performance on the testing data.
- Calculate and display the training and testing accuracy scores for the ensemble model.
- Saving the Predictions
- Save the predictions for the relevance of all the articles in the dataset to a CSV file, such as "us-economic-news-relevance.csv", using the to_csv() function from Pandas.

## 8) Conclusion

- Summarize the key findings and outcomes of the project, emphasizing the importance of automated article relevance classification in the field of economics.
- By following this detailed outline, someone can effectively write the code for the project, implementing data cleaning, text preprocessing, feature engineering,

model training, evaluation, and ensemble techniques to classify the relevance of economics news articles.

# Code Explanation

The code you provided is an example of a machine learning project that aims to classify the relevance of economics news articles. It utilizes popular libraries like Pandas, scikit-learn, and NLTK to preprocess the text data, extract features using TF-IDF vectorization, train multiple classifiers, and apply an ensemble approach.

Here's a breakdown of the code and its workflow:

1) **Importing Libraries and Dependencies**
   - The first step is to import the necessary libraries and dependencies that will be used throughout the code. These libraries include Pandas, NumPy, scikit-learn, Matplotlib, and NLTK. These powerful tools will help us process and analyze the data efficiently.

2) **Loading the Dataset**
   - After importing the required libraries, the code loads the dataset. In this case, the dataset is stored in a CSV file. The pd.read_csv() function is used to read the CSV file and load it into a Pandas DataFrame. This DataFrame is a tabular structure that allows us to manipulate and analyze the data easily.

3) **Data Cleaning and Preprocessing**
   - Before training our machine learning models, we need to preprocess the text data to remove any noise or unnecessary information. The code defines a function called clean(text) which performs various preprocessing techniques.
   - Within the clean() function, different operations are carried out, such as removing named entities, converting text to lowercase, removing punctuation and digits, removing stopwords, and lemmatization. These operations help in standardizing the text and removing irrelevant information.
   - The clean() function is then applied to the 'text' column of the dataset using Pandas' apply() function. This ensures that the text data in each article is processed and cleaned appropriately.

4) **Text Vectorization using TF-IDF**
   - After cleaning the text data, we need to convert the textual information into numerical features that can be used by machine learning algorithms. The code utilizes the TF-IDF (Term Frequency-Inverse Document Frequency) vectorization technique to achieve this.

- TF-IDF calculates the importance of each word in a document relative to the entire dataset. It assigns higher weights to words that are frequent in a document but rare in other documents.
- To perform TF-IDF vectorization, the code initializes a TF-IDF vectorizer object from scikit-learn. This object is capable of converting the preprocessed text data into a numerical representation.
- The vectorizer is then fitted on the preprocessed text data using the fit_transform() function. This process calculates the TF-IDF scores for each word in the dataset and creates a feature matrix.
- The feature matrix represents each article as a set of numerical features derived from the text. These features will be used as input for our machine learning models.

**5) Splitting the Data into Training and Testing Sets**
- To evaluate the performance of our models, we need to split the dataset into training and testing sets. The code uses the train_test_split() function from scikit-learn to achieve this.
- The feature matrix (TF-IDF representation of text) is assigned to the X variable, and the 'relevance' column, which indicates the relevance of each article, is assigned to the y variable.
- By specifying an appropriate test size and random state, the dataset is divided into two sets: one for training the models and another for testing their performance.

**6) Training and Evaluating Classifiers**
- With the data prepared, the code proceeds to train and evaluate multiple classifiers. These classifiers are popular machine learning algorithms available in scikit-learn, including Gaussian Naive Bayes, Decision Tree, Random Forest, and Support Vector Machine (SVM).
- For each classifier, the code follows a consistent workflow: creating an instance of the classifier, fitting it to the training data, making predictions on the test data, and evaluating the performance using metrics such as accuracy, precision, recall, and F1-score.
- The evaluation results are printed for each classifier, allowing us to compare their performance and select the most suitable one for our task.

### 7) Ensemble Learning with Voting Classifier

- To further improve the classification performance, the code employs an ensemble learning technique called Voting Classifier.
- The Voting Classifier combines the predictions from multiple individual classifiers and selects the most popular prediction as the final output.
- The code initializes a Voting Classifier object and adds the trained individual classifiers to it. The Voting Classifier then makes predictions on the test data.
- The performance of the Voting Classifier is also evaluated and printed, providing an overall assessment of the ensemble approach.

### 8) Visualizing the Results

- Finally, the code includes a visualization step using Matplotlib. It creates a bar chart to display the performance metrics of each classifier and the ensemble classifier.
- This visualization helps in comparing the different models and identifying the most effective one based on the evaluation metrics.

By understanding the code and its workflow, you can now dive into each function and its purpose, enabling you to modify and extend the code to suit your specific needs. Remember to experiment and iterate, as machine learning is an iterative process that often involves tweaking various components to achieve the best results. Good luck with your project!

# Future Work

1) **Data Augmentation**
- One potential area of improvement is data augmentation. By augmenting the existing dataset with additional synthetic examples, we can increase the diversity and quantity of training data, which can lead to better model performance.
- To implement data augmentation, you can explore techniques such as synonym replacement, backtranslation, or word embedding-based methods. These techniques can generate new examples by modifying or adding new text based on existing samples.

2) **Advanced Text Preprocessing**
- Enhancing the text preprocessing step can further improve the quality of the features extracted from the text data. You can consider incorporating additional techniques such as spell correction, handling negations, or handling word contractions.
- For spell correction, libraries like pySpellChecker or autocorrect can be utilized. Negations can be handled by adding a "not_" prefix to words following negation terms (e.g., "not relevant"). Word contractions can be expanded using contraction mapping dictionaries.

3) **Experimenting with Different Classifiers**
- Although the provided code utilizes popular classifiers, there are numerous other algorithms available in scikit-learn and other libraries that can be explored.
- To improve classification performance, you can experiment with algorithms such as gradient boosting machines (e.g., XGBoost or LightGBM), deep learning models (e.g., recurrent neural networks or transformers), or even ensemble approaches like stacking.
- Simply import the desired classifier, initialize an instance, and follow the same workflow as shown in the provided code to train, evaluate, and compare the performance.

4) **Hyperparameter Tuning**
- The performance of machine learning models heavily depends on the choice of hyperparameters. Fine-tuning the hyperparameters can significantly enhance the model's predictive power.
- Utilize techniques like grid search or random search to systematically explore the hyperparameter space of the selected classifiers.

- Define a grid of hyperparameter values to search over, and let scikit-learn automatically fit and evaluate the models with different combinations of hyperparameters. Select the best-performing combination based on evaluation metrics.

5) **Handling Class Imbalance**

- If the dataset exhibits class imbalance, where one class has significantly more samples than the other, it can impact the model's ability to learn and generalize effectively.
- Apply techniques to address class imbalance, such as oversampling the minority class (e.g., using SMOTE) or undersampling the majority class (e.g., random undersampling).
- Implement the chosen technique using libraries like imbalanced-learn, and evaluate the impact on the model's performance.

6) **Model Deployment and Serving**

- Once you have obtained a satisfactory model, you can deploy it to a production environment, allowing others to interact with it and make predictions on new data.
- Use frameworks like Flask, Django, or FastAPI to create a web application that exposes an API endpoint for accepting new article texts and returning the predicted relevance.
- Ensure the model is properly packaged and versioned, and set up a hosting service or infrastructure to deploy the application.

By following these steps, you can continue to improve the economics news article classification project. Remember to document your experiments, keep track of the changes you make, and evaluate the impact of each improvement. This iterative process will help you refine the model and achieve better performance over time. Good luck with your future work!

# Concept Explanation

**Concept Explanation: Support Vector Machines (SVM)**

Imagine you're at a zoo, and there's a grumpy zookeeper who wants to separate a group of zebras from a group of giraffes. The zookeeper doesn't know much about zebras or giraffes, but luckily, you're a smart AI assistant who can help!

You notice that the zebras have stripes and the giraffes have long necks. So, you come up with a brilliant idea: you draw an imaginary line on the ground that separates the two groups. This line should be the best possible way to separate zebras from giraffes based on their features.

Now, here's where Support Vector Machines (SVM) come into play. In the SVM world, this imaginary line is called a "decision boundary." The goal of SVM is to find the best decision boundary that maximally separates the two groups while keeping them as far apart as possible.

But how does SVM determine the best decision boundary? It's like a game of balancing acts! SVM looks for the widest "safety margin" between the two groups. It wants to draw the decision boundary in such a way that it has the biggest gap between zebras and giraffes. Think of it as a red carpet rolled out between the two groups, creating a buffer zone.

But wait, there's a catch! SVM also has some strict rules. It only cares about a few animals, the "support vectors," which are the ones closest to the decision boundary. These animals have the responsibility to ensure that the decision boundary stays in place. They're like the bodyguards of our red carpet, making sure that no zebras sneak into the giraffe area or vice versa.

Now, let's bring this back to our project. In the code you provided, SVM is used to classify economics news articles as relevant or irrelevant. Just like the zebras and giraffes, our SVM algorithm wants to find a decision boundary that best separates the relevant and irrelevant articles based on their features.

The features could be things like the presence of specific words or patterns in the articles. SVM learns from the provided dataset, trying to find the optimal decision boundary that maximizes the separation between the two classes.

Once the SVM has learned this decision boundary, we can use it to classify new, unseen articles. We feed an article into the SVM, and it will predict whether the article is relevant or irrelevant based on its position relative to the decision boundary.

So, in a nutshell, SVM helps us draw a line in the sand (or grass, in this case!) between relevant and irrelevant articles. It finds the widest gap between the two classes, assigns bodyguards (support vectors), and then uses this decision boundary to make predictions.

I hope this fun explanation helped you understand the concept of Support Vector Machines (SVM) a bit better! Remember, SVM is just one of the many algorithms in the vast world of machine learning, but it's a pretty cool one!

# Exercise Questions

**1. Question: What is the purpose of using TF-IDF vectorization in the text classification task? How does it help improve the performance of the SVM model?**

**Answer:** TF-IDF (Term Frequency-Inverse Document Frequency) vectorization is used to convert text documents into numerical vectors. It assigns weights to words based on their importance in a document and across the entire dataset. In the context of the SVM model, TF-IDF helps improve performance by capturing the relative importance of words. Words that appear frequently in a specific document but less frequently in the overall dataset are considered more informative. This way, the SVM model can better distinguish between relevant and irrelevant articles based on the significant words associated with each class.

**2. Question: Can you explain the concept of the "kernel" in SVM? How does it affect the SVM's ability to separate data?**

**Answer:** In SVM, the kernel is a function that transforms the input data into a higher-dimensional space, where it becomes easier to find a linear decision boundary. The kernel trick allows us to perform classification in this transformed space without explicitly calculating the coordinates of the data points. The choice of kernel affects the SVM's ability to separate data. For example, a linear kernel assumes that the data can be separated by a straight line, while a non-linear kernel like the Radial Basis Function (RBF) allows for more complex decision boundaries. By choosing an appropriate kernel, we can increase the SVM's flexibility in capturing complex relationships in the data.

**3. Question: What is the significance of the hyperparameters C and gamma in SVM? How do they impact the model's performance and generalization?**

**Answer:** The hyperparameters C and gamma play a crucial role in SVM. The parameter C controls the trade-off between maximizing the margin and minimizing the classification errors. A smaller C value leads to a larger margin but allows more misclassifications, while a larger C value reduces the margin but enforces stricter classification. On the other hand, gamma defines the influence of each training example. A smaller gamma value considers a broader influence, resulting in smoother decision boundaries, while a larger gamma value makes the SVM focus more on nearby points, potentially leading to more complex decision boundaries. The choice of C and gamma

should be carefully tuned to balance between overfitting and underfitting, ultimately impacting the model's performance and generalization on unseen data.

## 4. Question: What are the advantages and limitations of using SVM for text classification tasks? Are there any scenarios where SVM may not be the best choice?

**Answer:** SVM offers several advantages for text classification tasks. It can handle high-dimensional data efficiently, works well in cases of limited training samples, and can handle non-linear relationships through kernel functions. Additionally, SVM provides good generalization performance and works well with balanced datasets. However, SVM also has limitations. It may not perform optimally when dealing with extremely large datasets due to computational constraints. Furthermore, SVM is not inherently probabilistic, meaning it doesn't provide direct probabilities for predictions. In scenarios where interpretability or probabilistic outputs are essential, other models like logistic regression or Naive Bayes may be more suitable choices.

## 5. Question: How can we evaluate the performance of the SVM model in the text classification task? What evaluation metrics can we use, and what do they signify?

**Answer:** To evaluate the performance of the SVM model in text classification, we can use various evaluation metrics. Some commonly used ones include accuracy, precision, recall, and F1-score. Accuracy measures the overall correctness of the model's predictions. Precision indicates the proportion of correctly classified relevant articles among all the articles predicted as relevant, highlighting the model's ability to avoid false positives. Recall (also known as sensitivity) measures the proportion of correctly classified relevant articles out of all the actual relevant articles, focusing on avoiding false negatives. F1-score is the harmonic mean of precision and recall, providing a balanced measure that considers both false positives and false negatives. These evaluation metrics help us understand the SVM model's performance in terms of classification accuracy, false positives, false negatives, and overall effectiveness in identifying relevant articles.