# Kaggle survery questions clustering

## Problem Statement

**Background:** Kaggle is a platform that hosts data science competitions and provides a community for data scientists to collaborate and learn from each other. As part of this platform, Kaggle conducts surveys to gather insights into the data science community, including their preferences, tools, and opinions. These surveys generate a wealth of data that can be analyzed to understand trends and patterns within the data science field.

**Dataset Information:** The dataset used in this project consists of Kaggle survey questions. Each row represents a question asked in a survey, and the columns contain the question text. The dataset aims to capture the various topics and themes covered in the Kaggle surveys.

**Objective:** The objective of this project is to cluster the Kaggle survey questions based on their similarity in terms of content. Clustering the questions allows us to identify groups of questions that are related or cover similar topics. This can provide valuable insights into the main themes and areas of interest within the Kaggle survey data.

**Approach:** To cluster the survey questions, we utilize a technique called Latent Dirichlet Allocation (LDA), which is a topic modeling algorithm. LDA assumes that each document (in this case, a survey question) is a mixture of topics, and each topic is a probability distribution over words. By applying LDA to the survey questions, we can uncover the underlying topics and group similar questions together.

**Key Steps:**

- Import the necessary libraries for data processing, visualization, and topic modeling.
- Read the Kaggle survey questions dataset.
- Prepare the corpus by tokenizing the questions and removing stop words and punctuation.
- Create a dictionary of the tokens and save it for future reference.

- Convert the corpus into a bag-of-words representation and save it.
- Apply the TF-IDF (Term Frequency-Inverse Document Frequency) transformation to the corpus.
- Build an LDA model using the transformed corpus.
- Obtain the topic distributions for each question in the corpus.
- Visualize the topics and their associated words using the pyLDAvis library.
- Perform correlation analysis on the topic distributions to identify similarities between topics.
- Generate a cluster map to visualize the clusters of questions based on topic similarity.

**Dataset Limitations:**

- The dataset consists only of survey questions and does not include other survey details or responses.
- The quality and accuracy of the clustering results depend on the quality of the survey questions and the underlying assumptions of the LDA algorithm.
- The number of topics is predefined and may not accurately represent the true underlying topics in the dataset.

**Potential Applications:**

- The clustering of survey questions can help in organizing and categorizing large amounts of survey data.
- It can provide insights into the main themes and topics covered in the Kaggle surveys.
- The identified clusters can be used to improve survey design and question formulation.
- Researchers and data scientists can use the clustering results to explore specific topics of interest within the survey data.

Overall, this project aims to leverage the power of topic modeling and clustering techniques to analyze and understand the Kaggle survey questions, enabling data-driven insights and improved survey design in the data science community.

# Framework

Import the required libraries: Begin by importing the necessary libraries for data processing, visualization, and topic modeling. Some of the libraries used in the provided code include pandas, numpy, plotly, gensim, nltk, seaborn, matplotlib, and pyLDAvis.

**Read the dataset:** Load the Kaggle survey questions dataset using the pd.read_csv() function. This step ensures that the data is available for further analysis and processing.

Prepare the corpus: Extract the survey questions from the dataset and create a corpus. Iterate through the columns of the dataset to collect the question texts and store them in a list. This step is necessary for tokenization and subsequent topic modeling.

**Tokenization and stop word removal:** Tokenize the survey questions by splitting them into individual words. Convert the text to lowercase and remove stop words, punctuation, and special characters using the NLTK library's stopwords module and the string module. The resulting tokenized questions are stored as a list of lists.

**Create a dictionary:** Build a dictionary of the tokenized questions using the Gensim library's corpora.Dictionary() function. The dictionary maps each unique word to a unique numerical ID. Save the dictionary for future reference using the dictionary.save() method.

**Convert corpus to bag-of-words representation:** Convert the tokenized questions into a bag-of-words representation using the dictionary created in the previous step. This step transforms each question into a list of tuples, where each tuple contains the word ID and its frequency in the question.

**Apply TF-IDF transformation:** Apply the TF-IDF (Term Frequency-Inverse Document Frequency) transformation to the bag-of-words corpus using the Gensim library's models.TfidfModel() function. This step assigns a weight to each word based on its frequency in the corpus.

**Build an LDA model:** Build an LDA (Latent Dirichlet Allocation) model using the transformed corpus, dictionary, and the number of desired topics. Use the Gensim library's models.LdaModel() function, passing the corpus, dictionary, and the number of topics as parameters.

**Get topic distributions:** Obtain the topic distributions for each question in the corpus using the LDA model. This step creates a double wrapper over the original corpus, transforming it from bag-of-words representation to TF-IDF representation and then applying the LDA model to assign topic probabilities to each question.

**Visualize topics and words:** Use the pyLDAvis library to visualize the topics and their associated words. The pyLDAvis.gensim.prepare() function takes the LDA model, corpus, dictionary, and a multidimensional scaling (MDS) method as input. It generates an interactive visualization that displays the topics, their associated words, and their relative importance.

**Perform correlation analysis:** Conduct a correlation analysis on the topic distributions to identify similarities between topics. In the provided code, the sns.clustermap() function from the seaborn library is used to create a cluster map, visualizing the correlation between topics. The resulting plot helps in understanding the relationships and similarities between the topics.

**Finalize and present the results:** At this point, you can finalize the code by adding necessary comments, improving code readability, and ensuring that all dependencies are met. Present the final results, including the cluster map and any other relevant visualizations or insights derived from the analysis.

By following this outline, you can write the code for the Kaggle survey questions clustering project. It covers the key steps, libraries, and techniques needed to perform topic modeling and clustering on the survey questions dataset.

# Code Explanation

**Importing libraries:** The code begins by importing the necessary libraries required for data processing, visualization, and topic modeling. These libraries provide functions and tools that will be used throughout the code.

**Reading the dataset:** The next step is to read the dataset using the pd.read_csv() function. This function reads the contents of the 'questions.csv' file and stores it in a pandas DataFrame. The DataFrame represents tabular data, similar to a spreadsheet, and allows for easy manipulation and analysis.

**Preparing the corpus:** In this step, the code extracts the survey questions from the dataset and creates a corpus. A corpus is a collection of text documents. The code iterates through the columns of the dataset and collects the question texts, storing them in a list. This list becomes the corpus for further analysis.

**Tokenization and stop word removal:** Tokenization is the process of splitting text into individual words or tokens. In this code, the survey questions are tokenized by splitting them into individual words. The text is then converted to lowercase to ensure consistency. Stop words, such as common words like "the" or "and," are removed from the tokenized text using the NLTK library's stopwords module and the string module, which contains punctuation marks.

**Creating a dictionary:** The code builds a dictionary of the tokenized questions using the Gensim library's corpora.Dictionary() function. This dictionary maps each unique word to a unique numerical ID. It serves as a lookup table for the words in the corpus and is used for further processing.

**Converting corpus to bag-of-words representation:** In this step, the tokenized questions are converted into a bag-of-words representation using the dictionary created in the previous step. A bag-of-words representation represents text as a collection of unique words and their frequencies in the document. The resulting representation is a list of tuples, where each tuple contains the word ID and its frequency in the question.

**Applying TF-IDF transformation:** The code applies the TF-IDF (Term Frequency-Inverse Document Frequency) transformation to the bag-of-words corpus using the Gensim library's models.TfidfModel() function. TF-IDF assigns a weight to each word based on its frequency in the corpus. This step helps to highlight important words that

are unique to specific questions and downplay common words that appear in many questions.

**Building an LDA model:** LDA (Latent Dirichlet Allocation) is a topic modeling technique that discovers topics in a collection of documents. The code builds an LDA model using the transformed corpus, the dictionary, and the desired number of topics. The Gensim library's models.LdaModel() function is used for this purpose. LDA assigns topic probabilities to each question in the corpus.

**Getting topic distributions:** The code obtains the topic distributions for each question in the corpus using the LDA model. This step creates a double wrapper over the original corpus, transforming it from the bag-of-words representation to the TF-IDF representation and then applying the LDA model to assign topic probabilities to each question.

**Visualizing topics and words:** The code uses the pyLDAvis library to visualize the topics and their associated words. The pyLDAvis.gensim.prepare() function takes the LDA model, the corpus with topic probabilities, the dictionary, and a multidimensional scaling (MDS) method as input. It generates an interactive visualization that displays the topics, their associated words, and their relative importance. This visualization helps in understanding the main themes and patterns in the survey questions.

**Clustering the topics:** Lastly, the code performs clustering on the topics using the correlation matrix of the topic distributions. The sns.clustermap() function from the seaborn library is used to create a hierarchical cluster map that visually groups similar topics together based on their correlation. This step provides insights into the relationships between different topics.

The overall workflow of the code involves reading the dataset, preprocessing the text by tokenizing and removing stop words, creating a dictionary and a bag-of-words representation, applying TF-IDF transformation, building an LDA model, visualizing the topics, and performing clustering on the topics.

# Future Work

Enhance Text Preprocessing: Improve the text preprocessing step by considering additional techniques such as stemming or lemmatization to further normalize the text. This can help in reducing the vocabulary size and improving topic modeling accuracy.

**Explore Different Topic Modeling Algorithms:** Experiment with other topic modeling algorithms, such as Non-negative Matrix Factorization (NMF) or Hierarchical Dirichlet Process (HDP), to compare their performance with LDA. This exploration can provide insights into the suitability of different algorithms for the survey questions dataset.

**Optimize Number of Topics:** Conduct an analysis to determine the optimal number of topics in the dataset. This can be done by evaluating different numbers of topics and selecting the one that maximizes coherence or minimizes topic overlap. It is important to strike a balance between having enough topics to capture meaningful themes and avoiding excessive fragmentation.

**Evaluate Topic Coherence:** Measure the coherence of the generated topics to assess their interpretability and quality. Coherence metrics, such as c_v, u_mass, or C_npmi, can be used to evaluate the semantic similarity between words within each topic. This analysis helps in refining the topic modeling process and selecting the most meaningful topics.

**Explore Topic Evolution Over Time:** If the survey questions dataset contains a temporal aspect, analyze the evolution of topics over time. Split the dataset into different time periods and apply topic modeling separately to identify how the topics change or evolve. This analysis can provide valuable insights into the changing trends or concerns within the survey data.

**Leverage Topic Labels and Metadata:** If available, incorporate topic labels or metadata associated with the survey questions into the analysis. This additional information can provide context and assist in better understanding the topics. It can also enable more targeted clustering and analysis based on specific attributes.

**Perform Cluster Analysis:** Extend the clustering analysis beyond topics and explore clustering the survey questions themselves based on their content similarity. This can be achieved using techniques such as K-means clustering or hierarchical clustering. The

resulting clusters can reveal distinct question patterns or help in identifying subgroups within the dataset.

**Interactive Visualization and Exploration:** Develop an interactive dashboard or visualization tool to allow users to explore and interact with the clustering results. This can include features such as topic drill-down, keyword search, or filtering based on metadata. The goal is to provide a user-friendly interface for exploring the survey questions and gaining insights.

**Step-by-Step Guide to Implement Future Work:**

**Enhance Text Preprocessing:**

- Research and implement additional text preprocessing techniques like stemming or lemmatization.
- Modify the code to incorporate the chosen technique(s) for normalizing the text.

**Explore Different Topic Modeling Algorithms:**

- Research alternative topic modeling algorithms such as NMF or HDP.
- Replace the LDA model with the chosen algorithm and evaluate its performance.

**Optimize Number of Topics:**

- Implement a loop to iterate over a range of topic numbers.
- Use coherence metrics or other evaluation measures to assess the quality of topics for different topic numbers.
- Select the optimal number of topics based on the evaluation results.

**Evaluate Topic Coherence:**

- Implement coherence metrics calculation using libraries like gensim.models.CoherenceModel.
- Apply the coherence metrics to evaluate the quality of topics generated by the model.

**Explore Topic Evolution Over Time:**

- Analyze the temporal aspect of the survey questions dataset.
- Split the dataset into different time periods.

- Apply topic modeling separately to each time period and compare the topics.

**Leverage Topic Labels and Metadata:**

- Incorporate topic labels or metadata associated with survey questions into the code.
- Use the additional information to enhance the topic modeling and clustering analysis.

**Perform Cluster Analysis:**

- Research clustering algorithms such as K-means or hierarchical clustering.
- Apply the chosen clustering algorithm to the survey questions dataset.
- Analyze the resulting clusters and interpret the patterns or subgroups.

**Interactive Visualization and Exploration:**

- Explore visualization libraries like Plotly or Dash.
- Develop an interactive dashboard or visualization tool to present the clustering results.
- Implement features such as topic drill-down, keyword search, or metadata filtering for user interaction.

By following this step-by-step guide and implementing the suggested future work, you can enhance the Kaggle survey questions clustering project and gain deeper insights from the dataset. Remember to continuously evaluate and iterate on the results to refine the analysis.

# Concept Explanation

**Concept Explanation: Topic Modeling with LDA - Unlocking the Hidden Themes of Survey Questions**

Imagine you have a huge stack of papers with survey questions from people all over the world. You want to understand the underlying themes and patterns hidden in these questions, but reading them one by one would take forever! So, what can you do? Enter the magical algorithm called Latent Dirichlet Allocation, or LDA for short.

LDA is like a magician that can automatically group similar questions together and reveal the hidden topics without you having to read every single question. It's like having a superpower!

**Here's how LDA works:** Let's say you have a collection of survey questions about different topics like sports, technology, food, and movies. LDA takes all these questions and tries to figure out the themes that might be present in the entire collection. It does this by assuming that each question is a combination of different topics and each topic is a mixture of words.

To make sense of this, let's take an example. Suppose we have a bunch of questions, and LDA starts its magic. It might find a topic that seems to be about sports. This topic could have words like "game," "player," "score," and "team." Another topic could be about technology with words like "computer," "software," "programming," and "internet." LDA assigns a probability to each word, indicating how likely it is to be associated with a particular topic.

Now, let's say we give LDA a new survey question. It will analyze the words in that question and assign probabilities to different topics. For example, it might assign a high probability to the sports topic and a lower probability to the technology topic. Based on these probabilities, LDA helps us understand the dominant theme of the question.

But how does LDA actually do this magic trick? It uses a statistical approach and iteratively adjusts its probabilities until it finds the best combination of topics and words that explain the survey questions. It's like LDA is constantly learning and improving its understanding of the topics.

Once LDA has finished its magic, it presents us with the discovered topics and their associated words. We can then explore these topics and see which questions belong to each topic. This allows us to uncover the hidden patterns and understand what people are most curious about.

Now, isn't that amazing? LDA helps us unlock the hidden themes of survey questions, making our analysis faster and more efficient. It's like having a personal assistant who reads and categorizes thousands of questions for us!

So, the next time you have a pile of survey questions to analyze, remember the magic of LDA. Let it work its wonders and reveal the secret topics lurking within the words. Happy topic modeling!

# Exercise Questions

**What is the purpose of using the LDA algorithm in this project? How does it help in analyzing survey questions?**

**Answer:** The LDA algorithm is used in this project to perform topic modeling on survey questions. It helps in analyzing the questions by automatically discovering hidden topics or themes within the question dataset. LDA identifies the underlying patterns in the questions and groups them based on similarity, making it easier to understand the main themes discussed in the survey.

**Why is it important to preprocess the survey questions before applying the LDA algorithm? What preprocessing steps are performed in this project?**

**Answer:** Preprocessing the survey questions is important to improve the quality of the analysis. In this project, several preprocessing steps are performed before applying LDA. These steps include removing common words and punctuation, converting all text to lowercase, and tokenizing the questions into individual words. Preprocessing helps in removing noise and irrelevant information, allowing LDA to focus on the meaningful patterns in the questions.

**What is the role of the tf-idf (Term Frequency-Inverse Document Frequency) model in this project? How does it contribute to the topic modeling process?**

**Answer:** The tf-idf model is used in this project to transform the raw text data into numerical representations that capture the importance of words within each question and across the entire dataset. By calculating the tf-idf scores, which consider both the frequency of a word within a question and its rarity across all questions, the model assigns higher weights to words that are more informative and distinctive. This helps LDA in identifying the key words and topics that are significant in the survey questions.

**What are the steps involved in performing topic modeling using LDA in this project? Explain each step briefly.**

**Answer:** The steps involved in performing topic modeling using LDA in this project are as follows:

- Preprocessing: The survey questions are preprocessed by removing common words, punctuation, and converting text to lowercase.

- Creating a Dictionary: A dictionary is created to map each unique word in the corpus to a unique numerical ID.
- Creating a Corpus: The questions are transformed into a bag-of-words representation, where each question is represented as a list of (word ID, word count) pairs.
- Applying tf-idf: The tf-idf model is applied to the corpus to assign weights to each word based on its importance.
- Training LDA: The LDA model is trained on the tf-idf weighted corpus, which discovers the underlying topics in the questions.
- Analyzing Topics: The learned topics and their associated words are extracted from the trained LDA model, providing insights into the main themes of the survey questions.

**How can the correlation heatmap be useful in analyzing the topics generated by LDA? What information does it provide?**

**Answer:** The correlation heatmap generated in this project shows the pairwise correlation between the topics discovered by LDA. It provides insights into how closely related the topics are to each other based on the similarity of their word distributions. The heatmap helps in identifying clusters of related topics and understanding the relationships between them. It can be useful in exploring the coherence and distinctiveness of the topics, which can guide further analysis and interpretation of the survey question themes.