

# NLP Essentials: Removing Stopwords and Performing Text Normalization using NLTK and spaCy in Python

[INTERMEDIATE](#)[NLP](#)[PYTHON](#)[TECHNIQUE](#)[TEXT](#)[UNSTRUCTURED DATA](#)

## Overview

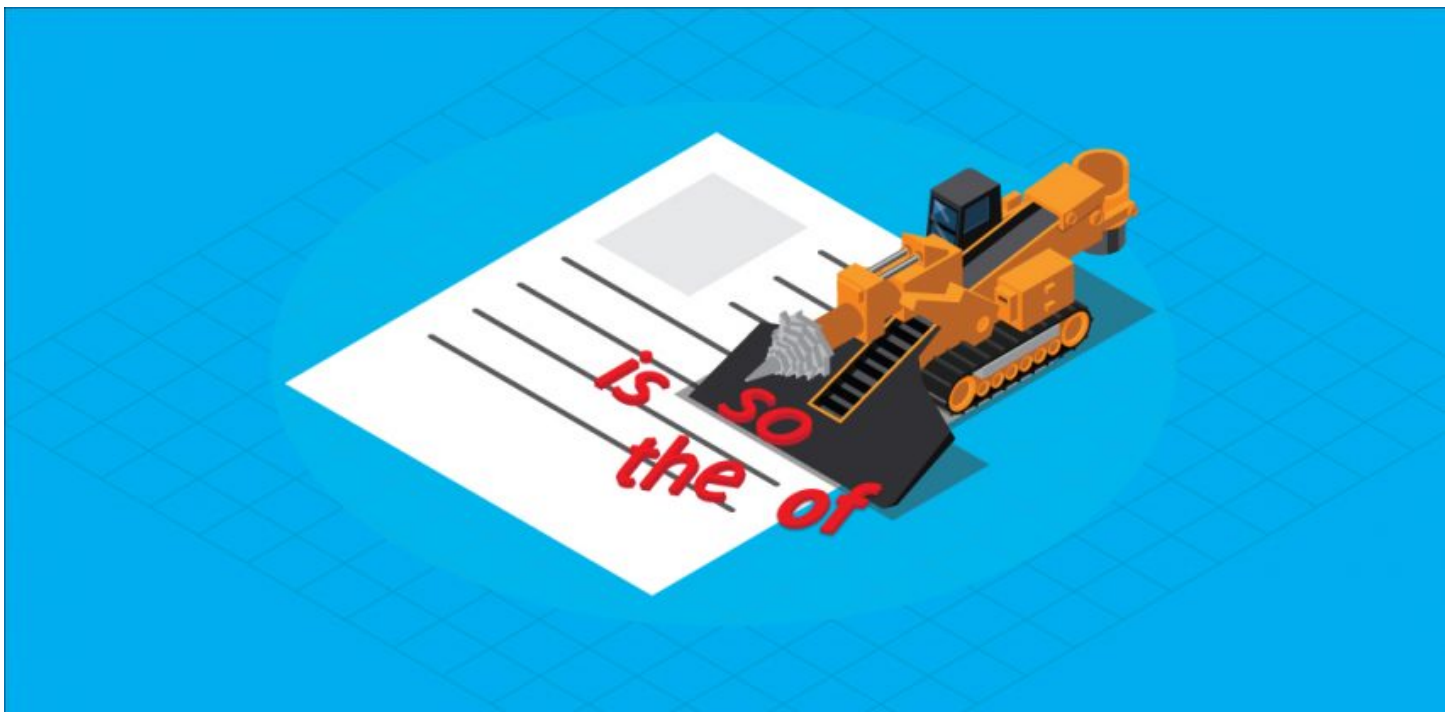
- Learn how to remove stopwords and perform text normalization in Python – an essential Natural Language Processing (NLP) read
- We will explore the different methods to remove stopwords as well as talk about text normalization techniques like stemming and lemmatization
- Put your theory into practice by performing stopwords removal and text normalization in Python using the popular NLTK, spaCy and Gensim libraries

## Introduction

Don't you love how wonderfully diverse [Natural Language Processing\\_\(NLP\)](#) is? Things we never imagined possible before are now just a few lines of code away. It's delightful!

But working with text data brings its own box of challenges. Machines have an almighty struggle dealing with raw text. We need to perform certain steps, called preprocessing, before we can work with text data using NLP techniques.

Miss out on these steps, and we are in for a botched model. These are essential NLP techniques you need to incorporate in your code, your framework, and your project.



We discussed the first step on how to get started with NLP in [this article](#). Let's take things a little further and take a leap. We will discuss how to remove stopwords and perform text normalization in Python using a few very popular NLP libraries – NLTK, spaCy, Gensim, and TextBlob.

*Are you a beginner in NLP? Or want to get started with machine learning but aren't sure where to begin? We have these two fields comprehensively covered in our end-to-end courses:*

- [Natural Language Processing \(NLP\) Using Python](#)
- [Applied Machine Learning – Beginner to Professional](#)

## Table of Contents

- What are Stopwords?
- Why do we need to Remove Stopwords?
- When should we Remove Stopwords?
- Different Methods to Remove Stopwords
  - Using NLTK
  - Using spaCy
  - Using Gensim
- Introduction to Text Normalization
- What are Stemming and Lemmatization?
- Methods to perform Stemming and Lemmatization
  - Using NLTK
  - Using spaCy
  - Using TextBlob

## What are Stopwords?

Stopwords are the most common words in any natural language. For the purpose of analyzing text data and building NLP models, these stopwords might not add much value to the meaning of the document.

Generally, the most common words used in a text are “the”, “is”, “in”, “for”, “where”, “when”, “to”, “at” etc.

Consider this text string – “There is a pen on the table”. Now, the words “is”, “a”, “on”, and “the” add no meaning to the statement while parsing it. Whereas words like “there”, “book”, and “table” are the keywords and tell us what the statement is all about.



A note here – we need to perform tokenization before removing any stopwords. I encourage you to go through my article below on the different methods to perform tokenization:

- [How to Get Started with NLP – 6 Unique Methods to Perform Tokenization](#)

Here's a basic list of stopwords you might find helpful:

a about after all also always am an and any are at be been being but by came can cant come could did didn't do does doesn't doing don't else for from get give goes going had happen has have having how i if ill i'm in into is isn't it its i've just keep let like made make many may me mean more most much no not now of only or our really say see some something take tell than that the their them then they thing this to try up us use used uses very want was way we what when where which who why will with without wont you your youre

## Why do we Need to Remove Stopwords?

Quite an important question and one you must have in mind.

**Removing stopwords is not a hard and fast rule in NLP. It depends upon the task that we are working on.** For tasks like text classification, where the text is to be classified into different categories, stopwords are removed or excluded from the given text so that more focus can be given to those words which define the meaning of the text.

Just like we saw in the above section, words like **there**, **book**, and **table** add more meaning to the text as compared to the words **is** and **on**.

However, in tasks like [machine translation](#) and [text summarization](#), removing stopwords is not advisable.

Here are a few key benefits of removing stopwords:

- On removing stopwords, dataset size decreases and the time to train the model also decreases
- Removing stopwords can potentially help improve the performance as there are fewer and only meaningful tokens left. Thus, it could increase classification accuracy
- Even search engines like Google remove stopwords for fast and relevant retrieval of data from the database

## When Should we Remove Stopwords?

I've summarized this into two parts: when we can remove stopwords and when we should avoid doing so.

## Remove Stopwords

We can remove stopwords while performing the following tasks:

- Text Classification
  - Spam Filtering
  - Language Classification
  - Genre Classification
- Caption Generation
- Auto-Tag Generation

## Avoid Stopword Removal

- Machine Translation
- Language Modeling
- Text Summarization
- Question-Answering problems

Feel free to add more NLP tasks to this list!

## Different Methods to Remove Stopwords

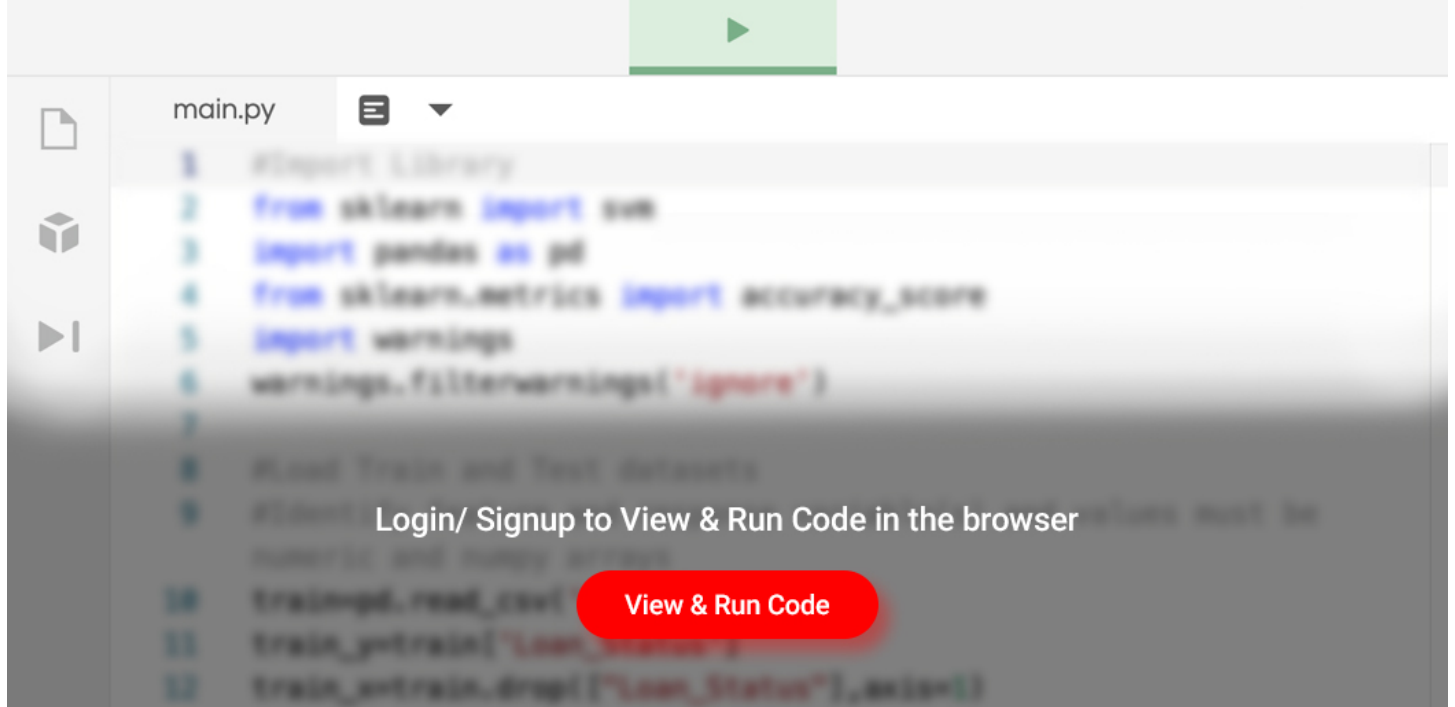
### 1. Stopword Removal using NLTK

NLTK, or the Natural Language Toolkit, is a treasure trove of a library for text preprocessing. It's one of my favorite Python libraries. **NLTK has a list of stopwords stored in 16 different languages.**

You can use the below code to see the list of stopwords in NLTK:

```
import nltk from nltk.corpus import stopwords set(stopwords.words('english'))
```

Now, to remove stopwords using NLTK, you can use the following code block. This is a LIVE coding window so you can play around with the code and see the results without leaving the article!



Here is the list we obtained after tokenization:

He determined to drop his litigation with the monastery, and relinquish his claims to the wood-cutting and fishery rihgts at once. He was the more ready to do this becuase the rights had become much less valuable, and he had indeed the vaguest idea where the wood and river in question were.

And the list after removing stopwords:

He determined drop litigation monastery, relinquish claims wood-cutting fishery rihgts. He ready becuase rights become much less valuable, indeed vaguest idea wood river question.

Notice that the size of the text has almost reduced to half! Can you visualize the sheer usefulness of removing stopwords?

## 2. Stopword Removal using spaCy

spaCy is one of the most versatile and widely used libraries in NLP. We can quickly and efficiently remove stopwords from the given text using SpaCy. It has a list of its own stopwords that can be imported as **STOP\_WORDS** from the **spacy.lang.en.stop\_words** class.

# spacy

Here's how you can remove stopwords using spaCy in Python:

```
1  from spacy.lang.en import English
2
3  # Load English tokenizer, tagger, parser, NER and word vectors
4  nlp = English()
5
6  text = """He determined to drop his litigation with the monastery, and relinquish his claims to the wood-cutting and
7  fishery rihgts at once. He was the more ready to do this becuase the rights had become much less valuable, and he had
8  indeed the vaguest idea where the wood and river in question were."""
9
10 # "nlp" Object is used to create documents with linguistic annotations.
11 my_doc = nlp(text)
12
13 # Create list of word tokens
14 token_list = []
15 for token in my_doc:
16     token_list.append(token.text)
17
18 from spacy.lang.en.stop_words import STOP_WORDS
19
20 # Create list of word tokens after removing stopwords
21 filtered_sentence = []
22
23 for word in token_list:
24     lexeme = nlp.vocab[word]
25     if lexeme.is_stop == False:
26         filtered_sentence.append(word)
27 print(token_list)
28 print(filtered_sentence)
```

[view raw](#)

spacy.py hosted with ♥ by GitHub

This is the list we obtained after tokenization:

```
He determined to drop his litigation with the monastery and relinquish his claims to the wood-cutting and \n
fishery rihgts at once. He was the more ready to do this becuase the rights had become much less valuable,
and he had \n indeed the vaguest idea where the wood and river in question were.
```

And the list after removing stopwords:

determined drop litigation monastery, relinquish claims wood-cutting \n fishery rihgts. ready becuae rights become valuable, \n vaguest idea wood river question.

An important point to note – stopword removal doesn't take off the punctuation marks or newline characters. We will need to remove them manually.

Read more about spaCy in this article with the library's co-founders:

- [DataHack Radio #23: Ines Montani and Matthew Honnibal – The Brains behind spaCy](#)

### 3. Stopword Removal using Gensim

Gensim is a pretty handy library to work with on NLP tasks. While pre-processing, gensim provides methods to remove stopwords as well. We can easily import the **remove\_stopwords** method from the class **gensim.parsing.preprocessing**.



Try your hand on Gensim to remove stopwords in the below live coding window:

main.py

```
1 #Import Library
2 from sklearn import svm
3 import pandas as pd
4 from sklearn.metrics import accuracy_score
5 import warnings
6 warnings.filterwarnings("ignore")
7
8 #Load Train and Test datasets
9 #Load
10 traingd.read_csv(
11 train_gtrain["Loan_Status"],axis=1)
12 train_gtrain.drop(["Loan_Status"],axis=1)
```

Login/ Signup to View & Run Code in the browser

View & Run Code

He determined drop litigation monastery, relinquish claims wood-cutting fishery rihgts once. He ready becuae rights valuable, vaguest idea wood river question were.

**While using gensim for removing stopwords, we can directly use it on the raw text.** There's no need to perform tokenization before removing stopwords. This can save us a lot of time.

## Introduction to Text Normalization

In any natural language, words can be written or spoken in more than one form depending on the situation. That's what makes the language such a thrilling part of our lives, right? For example:

- Lisa **ate** the food and washed the dishes.
- They were **eating** noodles at a cafe.
- Don't you want to **eat** before we leave?
- We have just **eaten** our breakfast.
- It also **eats** fruit and vegetables.

In all these sentences, we can see that the word **eat** has been used in multiple forms. For us, it is easy to understand that eating is the activity here. So it doesn't really matter to us whether it is 'ate', 'eat', or 'eaten' – we know what is going on.

Unfortunately, that is not the case with machines. They treat these words differently. Therefore, we need to normalize them to their root word, which is "eat" in our example.

Hence, **text normalization** is a process of transforming a word into a single canonical form. This can be done by two processes, **stemming** and **lemmatization**. Let's understand what they are in detail.

## What are Stemming and Lemmatization?

Stemming and Lemmatization is simply normalization of words, which means reducing a word to its root form.

In most natural languages, a root word can have many variants. For example, the word 'play' can be used as 'playing', 'played', 'plays', etc. You can think of similar examples (and there are plenty).



Stem<sup>ming</sup>  
&  
Lemmat<sup>ization</sup>



# Stemming

Let's first understand stemming:

- Stemming is a text normalization technique that cuts off the end or beginning of a word by taking into account a list of common prefixes or suffixes that could be found in that word
- It is a rudimentary rule-based process of stripping the suffixes ("ing", "ly", "es", "s" etc) from a word

## Lemmatization

Lemmatization, on the other hand, is an organized & step-by-step procedure of obtaining the root form of the word. It makes use of vocabulary (dictionary importance of words) and morphological analysis (word structure and grammar relations).

## Why do we need to Perform Stemming or Lemmatization?

Let's consider the following two sentences:

- **He was driving**
- **He went for a drive**

We can easily state that both the sentences are conveying the same meaning, that is, driving activity in the past. A machine will treat both sentences differently. Thus, to make the text understandable for the machine, we need to perform stemming or lemmatization.

Another benefit of text normalization is that it reduces the number of unique words in the text data. This helps in bringing down the training time of the machine learning model (and don't we all want that?).

## So, which one should we prefer?

**Stemming** algorithm works by cutting the suffix or prefix from the word. **Lemmatization** is a more powerful operation as it takes into consideration the morphological analysis of the word.

Lemmatization returns the lemma, which is the root word of all its inflection forms.

We can say that stemming is a quick and dirty method of chopping off words to its root form while on the other hand, lemmatization is an intelligent operation that uses dictionaries which are created by in-depth linguistic knowledge. **Hence, Lemmatization helps in forming better features.**

## Methods to Perform Text Normalization

# 1. Text Normalization using NLTK

The NLTK library has a lot of amazing methods to perform different steps of data preprocessing. There are methods like *PorterStemmer()* and *WordNetLemmatizer()* to perform stemming and lemmatization, respectively.

Let's see them in action.

## Stemming

```
1  from nltk.corpus import stopwords
2  from nltk.tokenize import word_tokenize
3  from nltk.stem import PorterStemmer
4
5  set(stopwords.words('english'))
6
7  text = """He determined to drop his litigation with the monastery, and relinquish his claims to the wood-cutting and
8  fishery rihgts at once. He was the more ready to do this becuase the rights had become much less valuable, and he had
9  indeed the vaguest idea where the wood and river in question were."""
10
11 stop_words = set(stopwords.words('english'))
12
13 word_tokens = word_tokenize(text)
14
15 filtered_sentence = []
16
17 for w in word_tokens:
18     if w not in stop_words:
19         filtered_sentence.append(w)
20
21 Stem_words = []
22 ps =PorterStemmer()
23 for w in filtered_sentence:
24     rootWord=ps.stem(w)
25     Stem_words.append(rootWord)
26 print(filtered_sentence)
27 print(Stem_words)
```

[view raw](#)

nltkstem.py hosted with ❤ by GitHub

He determined drop litigation monastery, relinquish claims wood-cutting fishery rihgts. He ready becuase rights become much less valuable, indeed vaguest idea wood river question.

He determin drop litig monastri, relinquish claim wood-cut fisheri rihgt. He readi becuas right become much less valuabl, inde vaguest idea wood river question.

We can clearly see the difference here. Now, let's perform lemmatization on the same text.

## Lemmatization

```
1  from nltk.corpus import stopwords
2  from nltk.tokenize import word_tokenize
3  import nltk
4  from nltk.stem import WordNetLemmatizer
5  set(stopwords.words('english'))
6
7  text = """He determined to drop his litigation with the monastery, and relinquish his claims to the wood-cutting and
8  fishery rihgts at once. He was the more ready to do this becuase the rights had become much less valuable, and he had
9  indeed the vaguest idea where the wood and river in question were."""
10
11 stop_words = set(stopwords.words('english'))
12
13 word_tokens = word_tokenize(text)
14
15 filtered_sentence = []
```

```

16 for w in word_tokens:
17     if w not in stop_words:
18         filtered_sentence.append(w)
19 print(filtered_sentence)
20
21
22 lemma_word = []
23 import nltk
24 from nltk.stem import WordNetLemmatizer
25 wordnet_lemmatizer = WordNetLemmatizer()
26 for w in filtered_sentence:
27     word1 = wordnet_lemmatizer.lemmatize(w, pos = "n")
28     word2 = wordnet_lemmatizer.lemmatize(word1, pos = "v")
29     word3 = wordnet_lemmatizer.lemmatize(word2, pos = ("a"))
30     lemma_word.append(word3)
31 print(lemma_word)

```

[view raw](#)

nltklemma.py hosted with ❤ by GitHub

He determined drop litigation monastery, relinquish claims wood-cutting fishery rihgts. He ready becuae rights become much less valuable, indeed vaguest idea wood river question.

He determined drop litigation monastery, relinquish claim wood-cutting fishery rihgts. He ready becuae right become much le valuable, indeed vaguest idea wood river question.

Here, **v** stands for **verb**, **a** stands for **adjective** and **n** stands for **noun**. The lemmatizer only lemmatizes those words which match the **pos** parameter of the lemmatize method.

Lemmatization is done on the basis of part-of-speech tagging (POS tagging). We'll talk in detail about POS tagging in an upcoming article.

## 2. Text Normalization using spaCy

spaCy, as we saw earlier, is an amazing NLP library. It provides many industry-level methods to perform lemmatization. Unfortunately, spaCy has no module for stemming. To perform lemmatization, check out the below code:

```

1 #make sure to download the english model with "python -m spacy download en"
2
3 import en_core_web_sm
4 nlp = en_core_web_sm.load()
5
6 doc = nlp(u"""He determined to drop his litigation with the monastery, and relinquish his claims to the wood-cutting and
7 fishery rihgts at once. He was the more ready to do this becuae the rights had become much less valuable, and he had
8 indeed the vaguest idea where the wood and river in question were.""")
9
10 lemma_word1 = []
11 for token in doc:
12     lemma_word1.append(token.lemma_)
13 lemma_word1

```

[view raw](#)

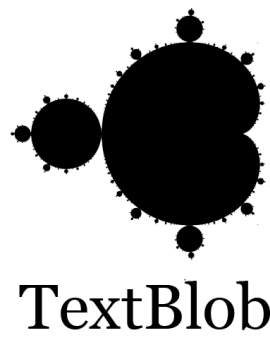
spacylemma.py hosted with ❤ by GitHub

-PRON- determine to drop -PRON- litigation with the monastery, and relinquish -PRON- claim to the wood-cutting and \n fishery rihgts at once. -PRON- be the more ready to do this becuae the right have become much less valuable, and -PRON- have \n indeed the vague idea where the wood and river in question be.

Here *-PRON-* is the notation for pronoun which could easily be removed using regular expressions. **The benefit of spaCy is that we do not have to pass any *pos* parameter to perform lemmatization.**

### 3. Text Normalization using TextBlob

TextBlob is a Python library especially made for preprocessing text data. **It is based on the NLTK library.** We can use TextBlob to perform lemmatization. However, *there's no module for stemming in TextBlob.*



So let's see how to perform lemmatization using TextBlob in Python:

```
1 # from textblob lib import Word method
2 from textblob import Word
3
4 text = """He determined to drop his litigation with the monastery, and relinquish his claims to the wood-cutting and
5 fishery rihgts at once. He was the more ready to do this becuase the rights had become much less valuable, and he had
6 indeed the vaguest idea where the wood and river in question were."""
7
8 lem = []
9 for i in text.split():
10     word1 = Word(i).lemmatize("n")
11     word2 = Word(word1).lemmatize("v")
12     word3 = Word(word2).lemmatize("a")
13     lem.append(Word(word3).lemmatize())
14 print(lem)
```

view raw

textbloblemma.py hosted with ♥ by GitHub

He determine to drop his litigation with the monastery, and relinquish his claim to the wood-cutting and fishery rihgts at once. He wa the more ready to do this becuase the right have become much le valuable, and he have indeed the vague idea where the wood and river in question were.

Just like we saw above in the NLTK section, TextBlob also uses POS tagging to perform lemmatization. You can read more about how to use TextBlob in NLP here:

- [Natural Language Processing for Beginners: Using TextBlob](#)

### End Notes

Stopwords play an important role in problems like sentiment analysis, question answering systems, etc. That's why removing stopwords can potentially affect our model's accuracy drastically.

This, as I mentioned, is part two of my series on 'How to Get Started with NLP'. You can check out [part 1 on tokenization here](#).

And if you're looking for a place where you can finally begin your NLP journey, we have the perfect course for you:

- [Natural Language Processing \(NLP\) Using Python](#)

---

Article Url - <https://www.analyticsvidhya.com/blog/2019/08/how-to-remove-stopwords-text-normalization-nltk-spacy-gensim-python/>



## **Shubham Singh**

A Data Science Enthusiast who loves reading & writing about Data Science and its applications. He has done many projects in this field and his recent work include concepts like Web Scraping, NLP etc. He is a Data Science Content Strategist Intern at Analytics Vidhya. And currently pursuing BTech in Computer Science from DIT University, Dehradun.