# Extract PDF Text and Generate Embeddings Using Hugging Face



**Estimated time needed:** 15 minutes

## What you will learn

You will learn how to read PDF files in Node.js using the pdf-parse library to extract text content from the PDF. During the extraction process you will clean the text by removing newline characters and merging the text into a single continuous paragraph.

You will interact with external APIs, specifically Hugging Face's machine learning API, to convert extracted text into embeddings. This process involves sending text to the Hugging Face inference client, which then returns an array of embeddings generated by a pre-trained model, `sentence-transformers/all-MiniLM-L6-v2`.

## Learning objectives

After completing this lab, you will be able to:

- Utilize the pdf-parse library to read a PDF file and extract its text content.

- Use the Hugging Face API to convert the extracted text into vector embeddings.

- Establish and utilize an inference client from Hugging Face to perform feature extraction.

- Implement asynchronous JavaScript using async/await to ensure appropriate ordering of the I/O operations such as reading PDF file.

## Prerequisites

- Intermediate competency with JavaScript Node.js

# Important notice about this lab environment

Skills Network Cloud IDE (based on Theia and Docker) is an open-source Integrated Development Environment (IDE) that provides an environment for hands-on labs in course and project-related labs.

Please be aware that sessions for this lab environment are not persistent. Every time you connect to this lab, a new environment is created for you.
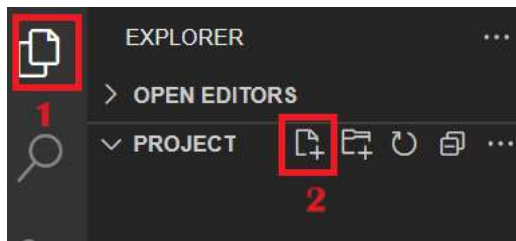
**You will lose data if you exit the environment without saving to GitHub or another external source.**

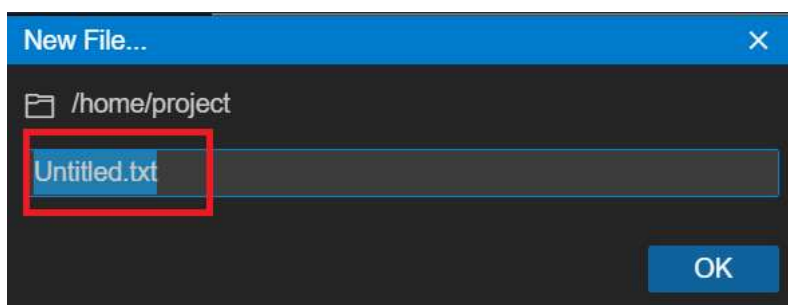Plan to complete these labs in a single session to avoid losing your data.

# Extract text from PDF and convert to embeddings

### Task 1: Create the JavaScript File

- Select **Explorer** on the left side of the terminal window, as shown at number 1 in the following screenshot. Then, within the **Project** folder, select **New File** as shown at number 2 in the following screenshot.



4. A pop-up box should display with the default file name **Untitled.txt**. Change the default name to `extract_pdf.js`.

## Task 2: Install necessary packages

1. You need to install `pdf-parse` and `@huggingface/inference`.

- fs for file system operations (built-in with Node.js).
- pdf-parse to extract text from PDF files.
- @huggingface/inference to interact with the Hugging Face API for machine learning tasks.

```
npm install pdf-parse @huggingface/inference
```

2. Execute the following command to get PDF test file from which you will extract the text. This PDF contains a **Food Menu**.

```
wget https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/6z4EmYqfgIxyYLGukHedvA/foodMenu.pdf
```

## Task 3: Import and Require

Import the necessary modules at the top of your file and initialize the Hugging Face inference client with your API key you obtained from the **Set Up Hugging Face API** lab.

```
const fs = require("fs");
const pdf = require("pdf-parse");
const { HfInference } = require("@huggingface/inference");
const hf = new HfInference("your_huggingface_api_key");
```

Note: The free version of the Hugging Face API key only allows a limited number of accesses in a given time period, although the documentation doesn't make that limit clear. It does seem like the limit is on an hourly basis, so if you get an error after frequent accesses to their API, you may need to wait up to an hour until you can continue your work.

## Task 4: Extract text from a PDF file

You need to create to extract text from the PDF file and format it for further processing.

```
const extractTextFromPDF = async (filePath) => {
  try {
    const dataBuffer = fs.readFileSync(filePath);
    const data = await pdf(dataBuffer);
    const text = data.text.replace(/\n/g, " ").replace(/ +/g, " ");
    return text;
  } catch (err) {
    console.error("Error extracting text from PDF:", err);
    throw err;
  }
};
```

Let's review how this funciton works.

- The variable `const extractTextFromPDF` calls the anonomys asynchronous function `async (filePath) => { ... }` which takes a single parameter `filePath`. The string parameter represents the path to the PDF file from which you will extract text.

- A try-catch blocks wraps the entire function to handle any errors that might occur during the execution of the code inside the try block.

- The `dataBuffer` variable stores the results from the Node.js File System, `fs` module's `readFileSync()` method. It reads all content of the file specified by the `filePath` parameter.

- The `data` variable stores the return value from the anonomous asynchronous `pdf()` function from the `pdf-parse` library. It parses the PDF content stored in the `dataBuffer`. The `pdf()` function returns a promise that resolves to an object containing various details about the PDF, including the extracted text as a property called `text`.

- Next, you use the code `const text = data.text.replace(/\n/g, " ").replace(/ +/g, " ");` to format the extracted text.

  - The function `.replace(/\n/g, " ")` replaces all newline characters (\n) with a single space. It converts multi-line text into a single line or paragraph, making it easier to process further.
  - The function `.replace(/ +/g, " ")` replaces multiple consecutive spaces with a single space. It cleans irregular spacing in the text, resulting from the previous replacement or the original PDF formatting.
  - After formatting, `return text` returns the cleaned text.

## Task 5: Convert the extracted text to an embedding

Now, you need to create one more function named `convertTextToEmbedding()` to generate embeddings from the extracted text.

```
const convertTextToEmbedding = async (text) => {
  try {
    const result = await hf.featureExtraction({
      model: "sentence-transformers/all-MiniLM-L6-v2",
      inputs: text,
    });
    // console.log("Embedding Result:", result);
    return result; // Return the embedding array
  } catch (err) {
    console.error("Error converting text to embeddings:", err);
    throw err;
  }
};
```

- The asynchronous `convertTextToEmbedding()` function converts text into embeddings using the Hugging Face API.

- It first tries to call the `featureExtraction()` method with the specified model `sentence-transformers/all-MiniLM-L6-v2` and input text.

- If successful, it returns the resulting embeddings.

- If an error occurs, it logs the error and rethrows it. You can use the optional commented `console.log` line for debugging purposes to log the embedding result.

Now you should create a variable named `filePath` and initialize it with name of the PDF.

```
const filePath = "foodMenu.pdf";
```

## Task 6: Create the main function

Now you need to create the `main()` function to invoke `convertTextToEmbedding()` and `extractTextFromPDF()`.

```
async function main(){
    const text = await extractTextFromPDF(filePath);
    console.log("Extracted Text:", text);
    const embeddings=await convertTextToEmbedding(text);
    console.log(embeddings);
}
```

The main function utilizes typical asynchronous JavaScript operations used for processing PDF text extraction and conversion to embeddings.

- By using async and await, it ensures that each step completes in order, handling the asynchronous nature of file I/O and processing operations gracefully.
- You will find this approach useful in scenarios where operations depend on the completion of previous tasks before proceeding, such as in data preprocessing or document analysis tasks.

Lastly, call the `main()` function in the last line of the `extract_pdf.js` file to kick-off the process.

```
main()
```

## Task 7: Check the output

1. Now open the terminal and perform given command to see the output.

```
node extract_pdf
```

2. You will get the output text extracted from pdf as per given screenshot.

- Output with embeddings



# Conclusion

- The provided code utilizes Node.js and libraries such as `fs` for file system operations, `pdf-parse` for PDF parsing, and the Hugging Face `@huggingface/inference` library for text embeddings.

- The `extractTextFromPDF()` function reads a PDF file synchronously using `fs.readFileSync`, parses it using `pdf-parse`, and cleans up the extracted text by replacing newline characters and multiple spaces.

- The `convertTextToEmbedding()` function utilizes the Hugging Face model `sentence-transformers/all-MiniLM-L6-v2` to convert the extracted text into embeddings, which represent a numerical representation of the text's meaning.

- The `main()` function sequentially calls `extractTextFromPDF()` and `convertTextToEmbedding()`, logging the extracted text and embeddings to the console.

**Author(s)**

Richa Arora