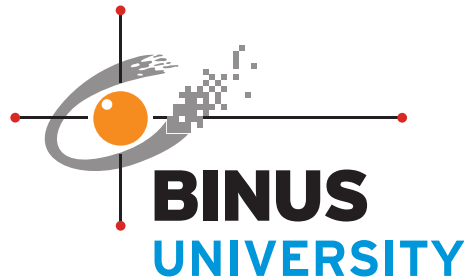


# **Creating Chatbot Using Natural Language Processing**



**Lecturer :  
D4017\_Jude Joseph Lamug Martinez, MCS**

**Arranged by :  
2702368122\_Vammy Johannis Jiang**

**Algorithm and Programming  
Faculty of Computer Science  
Binus International University 2022/2023**

## Table of Contents

<b><i>Chapter 1 – Project Specification .....</i></b>	<b><i>3</i></b>
1.1 Introduction .....	3
1.2 Idea Inspiration .....	3
1.3 How it Works.....	3
<b><i>Chapter 2 – Solution Design.....</i></b>	<b><i>4</i></b>
2.1 Solution Overview .....	4
2.2 Solution Description.....	4
<b><i>Chapter 3 - Implementation and Explanation .....</i></b>	<b><i>5</i></b>
3.1 Modules Used .....	5
3.2 How it works .....	6
<b><i>Chapter 4 - Evidence of Working Program .....</i></b>	<b><i>11</i></b>
<b><i>Chapter 5 - Lesson Learned.....</i></b>	<b><i>12</i></b>
<b><i>Chapter 6 - Source Code and Video Demo .....</i></b>	<b><i>13</i></b>
<b><i>References.....</i></b>	<b><i>14</i></b>

# **Chapter 1 – Project Specification**

## **1.1 Introduction**

This project was made using python programming language. The main purpose of this project is to make a fully functional chatbot which can answer question about Binus university. The chatbot is made by using natural language processing as its model.

## **1.2 Idea Inspiration**

The idea came from my brother when I did not know what I wanted to make for the final project. He proposed that I should ask Chatgpt for project ideas. I proceeded to ask Chatgpt for project ideas and natural language processing caught my eye immediately so I decided to make a chatbot based on natural language processing.

## **1.3 How it Works**

The project uses Natural Language Toolkit (NLTK) for text tokenization and stemming. The chatbot project utilizes a JSON file to define user intents with associated pattern and responses. The project also use PyTorch library for constructing and training a basic feedforward neural network. In which the json file is used as training data for the model after the pre-processing process.

## Chapter 2 – Solution Design

### 2.1 Solution Overview

The program consists of 5 components :

- Chat.py
- Model.py
- Nltk\_utils.py
- Train.py
- Intents.json

### 2.2 Solution Description

1. Nltk\_utils.py  
The utility tools for tokenization, stemming and bag of words part of pre-processing of a natural language processing language.
2. Model.py  
The model of the neural network module of the chatbot.
3. Intents.json  
The training data for our chatbot with consist of tags, patterns, and responses
4. Train.py  
The training module for the chatbot. It consists of ChatClassifier which is used for training and saving the chatbot model. And ChatDataSet where it is a PyTorch Dataset for handling and training data.
5. Chat.py  
The main program used to activate the chatbot consists of the main body of the chatbot itself.

## Chapter 3 - Implementation and Explanation

### 3.1 Modules Used

1. NumPy (import numpy as np):
  - NumPy is a fundamental package for scientific computing in Python. It provides support for large, multi-dimensional arrays and matrices, along with mathematical functions to operate on these arrays.
2. Random (import random):
  - The random module provides functions for generating pseudo-random numbers. It is used here to introduce randomness, such as selecting responses randomly in the chatbot.
3. JSON (import json):
  - The json module is used for encoding and decoding JSON (JavaScript Object Notation) data. In this code, it is used to load JSON data from a file (intents.json).
4. PyTorch (import torch, import torch.nn as nn, from torch.utils.data import Dataset, DataLoader):
  - PyTorch is an open-source machine learning library for Python. torch is the main PyTorch module, providing multi-dimensional tensors and mathematical operations. torch.nn contains neural network modules and functionalities.
  - Dataset and DataLoader are classes in PyTorch's torch.utils.data module for handling datasets and creating data loaders for efficient batching.
5. Scikit-Learn (from sklearn.feature\_extraction.text import CountVectorizer, from sklearn.model\_selection import train\_test\_split):
  - Scikit-Learn is a machine learning library for Python.
  - CountVectorizer is used for converting a collection of text documents to a matrix of token counts. train\_test\_split is used for splitting datasets into training and testing sets.
6. NLTK (import nltk):
  - NLTK (Natural Language Toolkit) is a powerful library for working with human language data. nltk.download('punkt') downloads the Punkt tokenizer models required for tokenization.
  - PorterStemmer is a stemming algorithm used to reduce words to their root or base form.

## 3.2 How it works

First before we know how the chatbot works, we will need to know what natural language processing is. Natural Language Processing (NLP) is an interdisciplinary subfield of computer science and linguistics that focuses on giving computers the ability to support and manipulate human language. NLP combines computational linguistics, machine learning, and deep learning models to process human language in the form of text or voice data and understand the intent and sentiment behind the language. NLP is used for a wide variety of language-related tasks, including answering questions, classifying text in various ways, and conversing with users. Some key applications of NLP include:

- Sentiment analysis: Classifying the emotional intent of text
- Text classification: Categorizing text data based on various criteria
- Entity recognition: Identifying and extracting specific entities from the text
- Text generation: Creating human-like text using machine learning models
- Speech recognition: Parsing spoken language into words and turning sound into text

NLP is essential for processing and analyzing text and speech data, and it can work through differences in dialects, slang, and grammatical irregularities typical in day-to-day conversations. Companies use NLP for various automated tasks, such as processing, analyzing, and archiving large documents, analyzing customer feedback or call center recordings, running chatbots for automated customer service, and answering who-what-when-where question.

In the context of chatbots, Natural Language Processing (NLP) plays a crucial role in enabling machines to understand and respond to text or voice data, capturing the intent and sentiment behind the language. NLP works in chatbots by using a combination of techniques and algorithms to analyse and process human language.

Here are some key components of how NLP works in chatbots:

- Text processing: NLP involves techniques such as tokenization, stemming, lemmatization, and sentiment analysis to process and understand the input text. Tokenization is the process of breaking down text into individual units, such as words or phrases, while stemming and lemmatization are methods of reducing words down to their roots to group related words together.

- Intent recognition: NLP algorithms are used to identify the underlying intent of the user's input, such as making a purchase, asking a question, or seeking assistance.
- Contextual understanding: NLP models can capture contextual information from the user's input, such as the tone, style, and context of the conversation. This allows the chatbot to generate more human-like and relevant responses based on the user's context and previous interactions.
- Machine learning: NLP often involves the use of machine learning algorithms to improve the chatbot's performance over time.

We will start with the first process of natural language processing which is text processing with functions defined in `nltkutils.py`

```
1  import numpy as np
2  import nltk
3  nltk.download('punkt')
4  from nltk.stem.porter import PorterStemmer
5  stemmer = PorterStemmer()
6
7  def tokenize(sentence):
8      return nltk.word_tokenize(sentence)
9
10
11  def stem(word):
12      return stemmer.stem(word.lower())
13
14
15  def bag_of_words(tokenized_sentence, words):
16      # stem each word
17      sentence_words = [stem(word) for word in tokenized_sentence]
18      # initialize bag with 0 for each word
19      bag = np.zeros(len(words), dtype=np.float32)
20      for idx, w in enumerate(words):
21          if w in sentence_words:
22              bag[idx] = 1
23
24      return bag
```

As we can see from the code above, the code first imports libraries such as `numpy` and `Natural Language Toolkit (NLTK)`. After that, 3 main functions which consist of `tokenize`, `stem`, and `bag of words`. The `tokenize` function utilizes the `NLTK` library to tokenize a given sentence into a list of words or tokens. The `stem` from `NLTK` library is used to perform stemming on a given word. Stemming involves reducing a word to its base or root form by removing prefixes or suffixes. The `bag of words` function creates a bag of words representation for a tokenized sentence based on 0 and 1s.

Example of each function is shown below

```
# Example Sentence
sentence = "Chatbots are becoming increasingly popular."

# Tokenization
tokenized_sentence = tokenize(sentence)
# Output: ['Chatbots', 'are', 'becoming', 'increasingly', 'popular', '.']

# Stemming
stemmed_sentence = [stem(word) for word in tokenized_sentence]
# Output: ['chatbot', 'are', 'becom', 'increas', 'popular', '.']

# Predefined Set of Words
predefined_words = ['chatbot', 'are', 'become', 'increasingly', 'popular']

# Bag of Words
bow = bag_of_words(tokenized_sentence, predefined_words)
# Output: [1, 1, 1, 1, 1]
```

Next is the model.py component :

```
import torch
import torch.nn as nn

class NeuralNet(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(NeuralNet, self).__init__()
        self.l1 = nn.Linear(input_size, hidden_size)
        self.l2 = nn.Linear(hidden_size, hidden_size)
        self.l3 = nn.Linear(hidden_size, num_classes)
        self.relu = nn.ReLU()

    def forward(self, x):
        out = self.l1(x)
        out = self.relu(out)
        out = self.l2(out)
        out = self.relu(out)
        out = self.l3(out)
        return out
```

The code above showcases a neural network class using PyTorch for a classification task. The network is defined by 3 layer of linear unit (ReLU) activation functions. This neural network is a feed forward neural network used for the training on the chatbot later on. The role of this neural network is in it ability to learn and comprehend complex relationships within input data such as patterns in user queries or sentences.



Next is train.py part :

```
1 import numpy as np
2 import random
3 import json
4
5 import torch
6 import torch.nn as nn
7 from torch.utils.data import Dataset, DataLoader
8 from sklearn.feature_extraction.text import CountVectorizer
9 from sklearn.model_selection import train_test_split
10
11 from nltk_utils import tokenize, stem
12 from model import NeuralNet
13
14 # Define a class for the Chat Classifier
15 class ChatClassifier:
16     def __init__(self, intents_file='intents.json', model_file='data.pth'):
17         # Initialize the intents from the JSON file names and placeholders for data
18         self.intents_file = intents_file
19         self.model_file = model_file
20         self.intents = None
21         self.tags = None
22         self.model = None
23         self.vectorizer = None
24         self.x_train = None
25         self.y_train = None
26         self.batch_size = None
27
28     def load_intents(self):
29         # Load intents from a JSON file
30         with open(self.intents_file, 'r') as f:
31             self.intents = json.load(f)
32
33     def process_intents(self):
34         # Process intents data to prepare it for training
35         corpus = []
36         for intent in self.intents['intents']:
37             tag = intent['tag']
38             patterns = intent['patterns']
39             # Extract the corpus with tokenized patterns and associated tags
40             corpus.extend([pattern, tag] for pattern in patterns)
41         # Get unique tags
42         self.tags = list(set(tag for _, tag in corpus))
43
44     # Separate patterns and tags
45     X, y = zip(*corpus)
46     # Tokenize patterns
47     X = [' '.join(tokenize(pattern) for pattern in X)]
48
49     # Use CountVecorizer to convert patterns to a bag-of-words representation
50     self.vectorizer = CountVecorizer(tokenizer.tokenize as a splitter)
51     X = self.vectorizer.fit_transform(X).toarray()
52     # Save the tags for training
53     y = [self.tags.index(tag) for tag in y]
```

```
54
55 def train_model(self, input_size, hidden_size, output_size,
56 num_epochs=1000, batch_size=64, learning_rate=0.01):
57     # Training loop
58     input_size = input_size or self.x_train.shape[1]
59     output_size = output_size or len(self.tags)
60
61     # Initialize the model, loss function, and optimizer
62     model = NeuralNet(input_size, hidden_size, output_size)
63     criterion = nn.CrossEntropyLoss()
64     optimizer = torch.optim.Adam(model.parameters()), lr=learning_rate)
65
66     # Create a PyTorch DataLoader for training data
67     train_dataset = ChatDataset(self.x_train, self.y_train)
68     train_loader = DataLoader(dataset=train_dataset, batch_size=batch_size, shuffle=True, num_workers=0)
69
70     # Training loop
71     for epoch in range(num_epochs):
72         for inputs, labels in train_loader:
73             outputs = model(inputs)
74             loss = criterion(outputs, labels)
75
76             optimizer.zero_grad()
77             loss.backward()
78             optimizer.step()
79
80         # Print loss at regular intervals
81         if epoch % 10 == 0:
82             print(f'Epoch {epoch + 1}/{num_epochs}, Loss: {loss.item():.4f}')
83
84     # Save the trained model
85     self.model = model
86
87 def save_model(self):
88     # Save the trained model's state, input size, hidden size, output size, and tags to a file
89     data = {
90         "model_state": self.model.state_dict(),
91         "input_size": self.model.input_size,
92         "hidden_size": self.model.hidden_size,
93         "output_size": self.model.output_size,
94         "tags": self.tags,
95     }
96     torch.save(data, self.model_file)
97     print(f'Training complete. Model saved to {self.model_file}')
98
99 def train_and_save_model(self):
100     # Load intents, process them, and train the model
101     self.load_intents()
102     self.x_train, self.y_train = self.process_intents()
103     self.batch_size = len(self.x_train) if self.batch_size is None else self.batch_size
104     self.train_model()
105
106 def visualize_training(self, num_epochs=1000):
107     # Visualize the training progress by printing the loss at regular intervals
```

```
110
111 def visualize_training(self, num_epochs=1000):
112     # Visualize the training progress by printing the loss at regular intervals
113     criterion = nn.CrossEntropyLoss()
114     optimizer = torch.optim.Adam(self.model.parameters()), lr=0.01)
115
116     train_dataset = ChatDataset(self.x_train, self.y_train)
117     train_loader = DataLoader(dataset=train_dataset, batch_size=self.batch_size, shuffle=True, num_workers=0)
118
119     for epoch in range(num_epochs):
120         for inputs, labels in train_loader:
121             outputs = self.model(inputs)
122             loss = criterion(outputs, labels)
123
124             optimizer.zero_grad()
125             loss.backward()
126             optimizer.step()
127
128             if epoch % 10 == 0:
129                 print(f'Epoch {epoch + 1}/{num_epochs}, Loss: {loss.item():.4f}')
130
131
132 # Define a custom dataset for PyTorch's DataLoader
133 class ChatDataset(Dataset):
134     def __init__(self, X, y):
135         self.x_samples = X.shape[0]
136         self.x_data = torch.tensor(X, dtype=torch.float32)
137         self.y_data = torch.tensor(y, dtype=torch.long)
138
139     def __getitem__(self, index):
140         return self.x_data[index], self.y_data[index]
141
142     def __len__(self):
143         return self.x_samples
144
145 # Usage
146 classifier = ChatClassifier()
147 classifier.train_and_save_model()
148
149 # Visualize training progress
150 classifier.visualize_training()
```

The train.py code contains the ChatClassifier class that facilitates the training and usage of a neural network-based chatbot model. The class reads intent data from a JSON file, pre-processes it, and trains a feedforward neural network using PyTorch. The training involves converting text patterns into a bag-of-words representation using CountVecorizer and then training a neural network to predict the associated intent tags. The training progress can be visualized by printing the loss at regular intervals. The trained model is saved, allowing for subsequent chatbot interactions. The code also includes a custom dataset class for PyTorch's DataLoader.

Finally it is chat.py :

```
1 import random
2 import json
3
4 import torch
5
6 from model import NeuralNet
7 from nltk_utils import bag_of_words, tokenize
8
9 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
10
11 with open('intents.json', 'r') as json_data:
12     intents = json.load(json_data)
13
14 FILE = "data.pth"
15 data = torch.load(FILE)
16
17 input_size = data["input_size"]
18 hidden_size = data["hidden_size"]
19 output_size = data["output_size"]
20 all_words = data["all_words"]
21 tags = data["tags"]
22 model_state = data["model_state"]
23
24 model = NeuralNet(input_size, hidden_size, output_size).to(device)
25 model.load_state_dict(model_state)
26 model.eval()
27
28 bot_name = "RBinus"
29 print("Let's chat! (type 'quit' to exit)")
30 while True:
31     sentence = input("You: ")
32     if sentence == "quit":
33         break
34
35     sentence = tokenize(sentence)
36     X = bag_of_words(sentence, all_words)
37     X = X.reshape(1, X.shape[0])
38     X = torch.from_numpy(X).to(device)
39
40     output = model(X)
41     _, predicted = torch.max(output, dim=1)
42     tag = tags[predicted.item()]
43
44     probs = torch.softmax(output, dim=1)
45     prob = probs[0][predicted.item()]
46     if prob.item() > 0.75:
47         for intent in intents['intents']:
48             if tag == intent["tag"]:
49                 print(f"{bot_name}: {random.choice(intent['responses'])}")
50     else:
51         print(f"{bot_name}: I do not understand...")
52
```

Chat.py employs a trained neural network from train.py for a chatbot application. The model is loaded from a saved state, and user input is tokenized and converted into a bag-of-words representation. The neural network predicts the intent of the input, and if the prediction surpasses a confidence threshold, the chatbot provides a relevant response randomly selected from predefined responses associated with the predicted intent. If the confidence falls below the threshold, the chatbot print out I do not understand.

The overall process to activating the chatbot is first, users will need to run the train.py until the neural network finishes training. Then, users can run chat.py to run the chatbot and ask it questions about Binus International.

First, when we run `train.py`

Then, we run `chat.py` and asks it questions :

11

## Chapter 5 - Lesson Learned

In working on this project, I have learned a lot of new things in python :

- Several new libraries (numpy, nltk, PyTorch) as well as how to utilize them to make a very interesting programming.
- Natural language processing was something previously unknown to myself, so in this project I have learned what it is, how does it work, and the processes behind NLP in which in the current world has useful tools in the form of chatgpt and other AI chatbot which utilizes NLP to produce a fluid more human like response.

Overall, I have feel that I have grown as a programmer once again. Moreover, after this project I personally feel more interested to explore new area of coding in which I previously have not been interested to explore. I find the urge to keep improving what I made and hope to make more other projects and interesting tools to help me boost my own programming skill on my programming journey that has just begun a few months ago.

## Chapter 6 - Source Code and Video Demo

GitHub Link : <https://github.com/TanaRuin/Algo-Prog-Final-Project/tree/main>

Video Demo Link :

<https://drive.google.com/drive/folders/1JEFBTFTJYUD24jLqv7CkMdZZ1XMEJg0m?usp=sharing>

## References

*patrickloeber/pytorch-chatbot: Simple chatbot implementation with PyTorch.* (n.d.). GitHub.

<https://github.com/patrickloeber/pytorch-chatbot>

codebasics. End-to-End NLP Project | Build a Chatbot in Dialogflow | NLP Tutorial | S3 E2. *YouTube*. Published online June 23, 2023. Accessed January 7, 2024.

<https://www.youtube.com/watch?v=2e5pQqBvGco&t=9916s>

sentdex. Creating a Chatbot with Deep Learning, Python, and TensorFlow p.1. *YouTube*. Published online November 24, 2017. Accessed January 7, 2024.

[https://www.youtube.com/watch?v=dvOnYLDg8\\_Y&list=PLQVvvaa0QuDdc2k5dwtDTyT9aCja0on8j](https://www.youtube.com/watch?v=dvOnYLDg8_Y&list=PLQVvvaa0QuDdc2k5dwtDTyT9aCja0on8j)

Landbot. Natural Language Processing Chatbot | Quick Overview. *YouTube*. Published online July 13, 2022. Accessed January 7, 2024. [https://www.youtube.com/watch?v=XT923\\_4GY9k](https://www.youtube.com/watch?v=XT923_4GY9k)