# DON BOSCO INSTITUTE OF TECHNOLOGY

Bengaluru, Karnataka - 74

**TEAM COUNT: FOUR**

**Project Title: Real-Time Language Translation Using Neural Machine Translation**

**Team Lead Name: Deepu M S**

**Team lead CAN ID: CAN_33695861**

1. **Name: Deepu MS**

   **CAN ID: CAN_ 33695861**

   **ROLE: FRONT END DEVELOPER**

2. **Name: Divya S**

   **CAN ID: CAN_ 34140735**

   **ROLE: RESEARCHER**

3. **Name: Spandana R**

   **CAN ID: CAN_ 33695563**

   **ROLE: BACK ENDDEVELOPER**

4. **Name: Shama S P**

   **CAN ID: CAN_ 33701582**

   **ROLE:MACHINE LEARNING ENGINEER**

## Advanced Market Segmentation Using Deep Clustering

PHASE 3: Model Training and Evaluation

**3.1 Overview of Model Training and Evaluation:**

For the Advanced Market Segmentation using Deep Clustering project, the key algorithms are:

**Overview**
Neural Machine Translation (NMT) is a deep learning approach to machine translation. This project aims to develop an NMT model for real-time language translation.

**Model Training**

1.Data Collection: Gather a large dataset of paired texts in the source and target languages.

2. Data Preprocessing: Preprocess the data by tokenizing, normalizing, and splitting it into training, validation, and testing sets.

3. Model Architecture: Implement a sequence-to-sequence NMT model with an encoder-decoder structure, using recurrent neural networks (RNNs) or transformers.

4. Training: Train the model using a suitable optimizer and loss function, such as cross-entropy loss.

**Model Evaluation**

1.Metrics: Evaluate the model using metrics such as BLEU score, perplexity, and translation accuracy.

2.Testing: Test the model on a held-out test set to assess its performance.

3. Comparison: Compare the performance of the NMT model with other machine translation approaches.

**Real-Time Translation**

1.Model Deployment: Deploy the trained model in a real-time translation system.

2.Input Processing: Process user input, such as text or speech, and prepare it for translation.

3. Translation: Use the NMT model to translate the input text in real-time.

**Data Collection**

1.Customer Data: Collect customer data, including demographics, behavior, and transaction history.

2. Feature Engineering: Extract relevant features from the data, such as purchase frequency and average order value.

**Model Training**

1.Deep Clustering Model: Implement a deep clustering model, such as a variational autoencoder (VAE) or a deep embedding clustering (DEC) model.

2. Training: Train the model using the customer data and features.

**Model Evaluation**

1.Metrics: Evaluate the model using metrics such as silhouette score, Calinski-Harabasz index, and Davies-Bouldin index.

2. Cluster Analysis: Analyze the clusters obtained from the model to identify patterns and insights.

**Market Segmentation**

1.Segment Identification: Identify distinct market segments based on the clusters obtained from the model.

2.Segment Characterization: Characterize each segment using demographics, behavior, and transaction history.

3. Targeted Marketing: Develop targeted marketing strategies for each segment to improve customer engagement and retention.Model Training Data Collection

Gather a large dataset of paired texts in the source and target languages.

Import pandas as pd

```
# Load dataset
train_ data = pd. read_ csv('train.csv')
test_ data = pd. read_ csv('test.csv')
```

**Data Preprocessing**

Preprocess the data by tokenizing, normalizing, and splitting it into training, validation, and testing sets.

```
from nltk .tokenize import word_ tokenize
from sk learn. model_ selection import train_ test_ split

# Tokenize   and  normalize data
Train _tokenized = train_ data .apply (word_ tokenize)
Test _tokenized = test_ data .apply(word_ tokenize)

# Split data into training, validation, and testing sets

X_ train, X_ val, y_ train, y_ val = train_ test_ split(train_ tokenized, test_ tokenized, test_ size=0.2,
random_ state=42)
```

**Model Architecture**

Implement a sequence-to-sequence NMT model with an encoder-decoder structure, using recurrent neural networks (RNNs) or transformers.

**Model Evaluation**

```
from sklearn.metrics import accuracy_score, f1_score

# Evaluate model on validation set

y_pred = model.predict(X_val)

y_pred_class = np.argmax(y_pred, axis=1)

print('Validation Accuracy:', accuracy_score(y_val, y_pred_class))

print('Validation F1 Score:', f1_score(y_val, y_pred_class, average='macro'))
```

**Real-Time Translation**

```python
from flask import Flask, request, jsonify

app = Flask(__name__)

# Load trained model
model.load_weights('nmt_model.h5')

@app.route('/translate', methods=['POST'])
def translate():
    # Get input text from request
    input_text = request.get_json()['text']

    # Preprocess input text
    input_tokenized = word_tokenize(input_text)

    # Translate input text using NMT model
    output = model.predict(input_tokenized)

    # Return translated text as JSON response
    return jsonify({'translated_text': output})

if __name__ == '__main__':
    app.run(debug=True)
```

**Advanced Market Segmentation using Deep Clustering**
Overview
This project aims to develop a deep clustering approach for advanced market segmentation.

**Data Collection**

Collect customer data, including demographics, behavior, and transaction history.

```python
import pandas as pd

# Load customer data
customer_data = pd.read_csv('customer_data.csv')
```

Model Training Data Preprocessing

Preprocess the data by scaling and normalizing it.

```python
from sklearn.preprocessing import StandardScaler
# Scale and normalize data
scaler = StandardScaler()
customer_data_scaled = scaler.fit_transform(customer_data)
```

## Model Architecture

Implement a deep clustering model, such as a variational autoencoder (VAE) or a deep embedding clustering (DEC) model.

```python
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input, Dense, Dropout

# Define model architecture
input_dim = customer_data_scaled.shape[1]
encoding_dim = 10

input_layer = Input(shape=(input_dim,))
encoder = Dense(encoding_dim, activation='relu')(input_layer)
decoder = Dense(input_dim, activation='sigmoid')(encoder)

autoencoder = Model(input_layer, decoder)
```

## Model Evaluation

```python
from sklearn.metrics import silhouette_score

# Evaluate model using silhouette score
silhouette = silhouette_score(customer_data_scaled, autoencoder.predict(customer_data_scaled))
print('Silhouette Score:', silhouette)
```

## Market Segmentation

```python
from sklearn.cluster import KMeans

# Perform k-means clustering on encoded data
kmeans = KMeans(n_clusters=5)
clusters = kmeans.fit_predict(autoencoder.predict(customer_data_scaled))

# Analyze clusters to identify market segments
print('Cluster Labels:', clusters)

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Embedding

# Define model architecture
model = Sequential()
model.add(Embedding(input_dim=10000, output_dim=128, input_length=max_length))
model.add(LSTM(128, return_sequences=True))
model.add(LSTM(64))
model.add(Dense(64, activation='relu'))
model.add(Dense(output_dim=10000, activation='softmax'))
```