

HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY



SCHOOL OF ELECTRONICS AND TELECOMMUNICATIONS

NATURE LANGUAGE PROCESSING

Sentence Similarity Learning by Lexical Decomposition and Composition - Q&A Search Application

Instructor: PhD Do Thi Ngoc Diep

Students : Trinh Quoc An - 20224269

Nguyen Van Hieu - 20224312

HANOI, 1-2025

Contents

1	Overview	2
2	Introduction	2
3	Model Overview	3
4	Model Implementation	4
5	Experiment	8
5.1	Training Model	8
5.2	Test Set Evaluation	9
5.3	Q&A Search Application	10
6	Conclusion	10
6.1	Key Achievements	10
6.2	Limitations	11
7	Other Models	11

1 Overview

This project implements a **Two-channel CNN** model to predict the similarity score (ranging from 0 to 1) between two Vietnamese sentences. The main idea is:

- Preprocess and tokenize both sentences.
- For each word in the first sentence, find the best-matching word in the second sentence, and *decompose* each word vector into *similar* and *dissimilar* components (using linear decomposition).
- Combine (*stack*) all the *similar* components into one matrix and all the *dissimilar* components into another matrix, effectively creating two channels for a CNN.
- A two-channel CNN then learns features from both matrices. It outputs a single similarity score after passing through a fully-connected layer and a sigmoid.

We have referred to the method from the article: [1] *Sentence Similarity Learning by Lexical Decomposition and Composition* (Zhiguo Wang, Haitao Mi, Abraham Ittycheriah).

2 Introduction

Sentence similarity is a key measure in both Natural Language Processing (NLP) and Information Retrieval (IR), used to assess how closely two sentences align in meaning. Common applications include:

- **Paraphrase identification:** Determining whether two sentences express the same idea.
- **Question answering:** Ranking candidate answers by relevance.

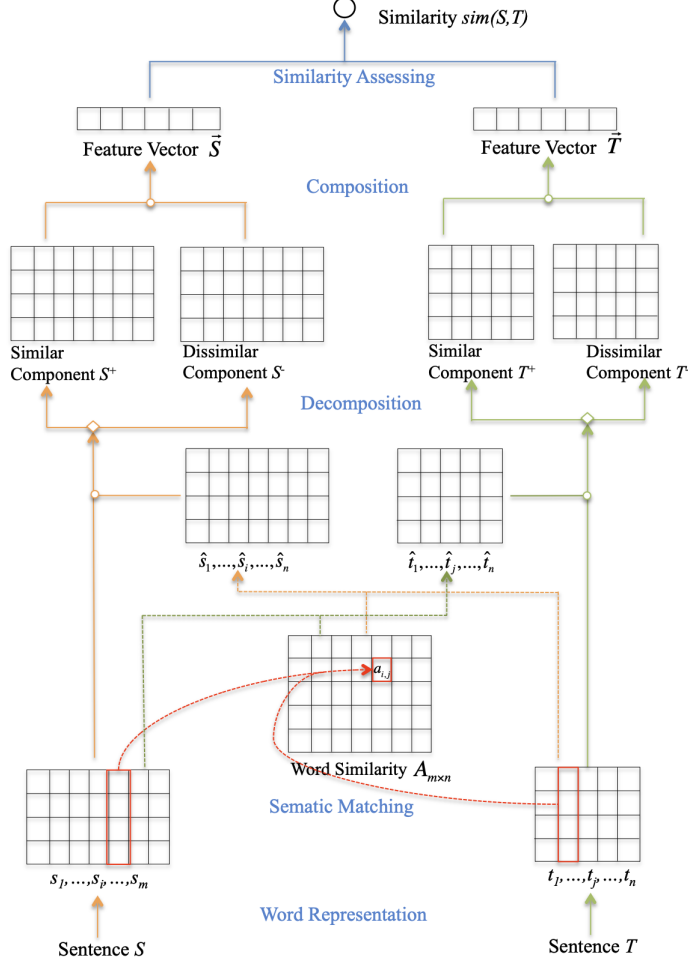
Despite considerable progress, three core challenges remain:

1. **Lexical Gap:** Different words or phrases can carry the same meaning (e.g., *irrelevant* vs. *not related*).
2. **Multi-level Similarity:** Semantic overlap must be evaluated at the word level, phrase level, and syntactic level.
3. **Dissimilarity Cues:** Parts of a sentence that do not match exactly can still provide critical information about nuance or specificity (e.g., *sockeye* vs. *salmon*).

In response to these challenges, the work summarized here proposes a novel model that decomposes each word based on its semantic matching in the other sentence, then composes both similar and dissimilar components using a two-channel CNN.

3 Model Overview

We have based on the model in the article to build a similar but smaller model, which may not be as accurate as the original model. Methods proposed in [?] demonstrate strong performance in neural dialogue systems and sentence embeddings.



(a) Figure 1: Model overview from article: Sentence Similarity Learning by Lexical Decomposition and Composition.

Model Details

This model addresses three core challenges in sentence similarity: bridging the lexical gap, capturing multi-level similarity (from words to syntax), and leveraging both similar and dissimilar parts.

Given two sentences S and T , the goal is to compute a similarity score $\text{sim}(S, T)$ via the following steps:

1. Word Representation:

- Convert each word in S and T into a vector s_i or t_j using cc.vi.300.vec.
- Words with similar meanings are thus closer in embedding space.

2. Semantic Matching:

- For each word s_i in S , compute a *semantic matching vector* \hat{s}_i by comparing it to all words in T .
- Similarly, each t_j has a matching vector \hat{t}_j based on S .

3. Decomposition:

- Decompose each word s_i (and t_j) into a *similar component* (s_i^+) and a *dissimilar component* (s_i^-) according to how it aligns with \hat{s}_i .
- This step ensures that the model does not ignore any mismatched (dissimilar) information.

4. Composition:

- Assemble the similar components (s_i^+) into matrix S^+ and the dissimilar components (s_i^-) into matrix S^- .
- Likewise, construct T^+ and T^- for sentence T .
- A composition function (often a two-channel CNN) extracts features from these matrices to produce two feature vectors \tilde{S} and \tilde{T} .

5. Similarity Assessing:

- Finally, concatenate \tilde{S} and \tilde{T} and apply a scoring function (e.g., a sigmoid) to predict $\text{sim}(S, T)$ in the range $[0, 1]$.

By simultaneously capturing both similar and dissimilar elements, this approach outperforms methods that only focus on matched portions, providing richer semantic insights.

4 Model Implementation

The `get_vector(word)` function

```
def get_vector(word):  
    if word in model_emb.key_to_index:  
        return model_emb[word]  
    else:  
        return np.zeros(embedding_dim, dtype=np.float32)
```

- **Purpose:** Retrieves the embedding vector of a given word from `model_emb`. If the word is out-of-vocabulary (not found in `key_to_index`), it returns a zero vector.
- **Significance:** Ensures every token has a fixed-size embedding. Out-of-vocabulary words default to zeros, preventing errors in later computations.

Vietnamese Preprocessing and Tokenization

```
def preprocess_vi(sentence):  
    s = sentence.lower().strip()  
    s = ViTokenizer.tokenize(s)  
    tokens = s.split()  
    return tokens
```

- **Steps:**

1. Convert the sentence to lowercase and strip whitespace.
2. Use PyVi (ViTokenizer) to tokenize Vietnamese text, producing tokens with underscores
3. Split into a list of token strings.

Utility Functions: cosine_sim, find_best_match, linear_decompose

cosine_sim(a, b)

```
def cosine_sim(a, b):  
    norm_a = np.linalg.norm(a)  
    norm_b = np.linalg.norm(b)  
    if norm_a < 1e-9 or norm_b < 1e-9:  
        return 0.0  
    return float((a @ b) / (norm_a * norm_b))
```

- **Purpose:** Computes the cosine similarity between two vectors **a** and **b**.
- **Handling edge cases:** If either vector has near-zero length, it returns 0.0 to avoid division by zero.

find_best_match(vec_s, list_vec_t)

```
def find_best_match(vec_s, list_vec_t):  
    best_sim = -1.0  
    best_vec = np.zeros_like(vec_s)  
    for vec_t in list_vec_t:  
        sim = cosine_sim(vec_s, vec_t)  
        if sim > best_sim:  
            best_sim = sim  
            best_vec = vec_t  
    return best_vec
```

- **Purpose:** Finds the word vector **t_j** in **list_vec_t** that has the highest cosine similarity with the vector **vec_s**.
- **Significance:** Determines which word in sentence **T** best matches a given word in sentence **S**.

```
linear_decompose(s_i, s_i_hat)

def linear_decompose(s_i, s_i_hat):
    alpha = cosine_sim(s_i, s_i_hat)
    s_plus = alpha * s_i_hat
    s_minus = (1 - alpha) * s_i
    return s_plus, s_minus
```

- **Purpose:** Splits the vector `s_i` into two components:
 - `s_plus`: the similar component weighted by the cosine similarity `alpha`.
 - `s_minus`: the remaining portion $(1 - \alpha) * s_i$.
- **Significance:** Captures how much of `s_i` is "covered" or matched by `s_i_hat`.

The Dataset and `collate_fn`

Similarity Dataset

```
class SimilarityDataset(Dataset):
    def __init__(self, df):
        super().__init__()
        self.df = df.reset_index(drop=True)

    def __len__(self):
        return len(self.df)

    def __getitem__(self, idx):
        row = self.df.iloc[idx]
        sent1 = str(row["sentence1"])
        sent2 = str(row["sentence2"])
        label = float(row["similarity"]) # range [0..1]
        return sent1, sent2, label
```

- **Purpose:** Stores the data from a DataFrame containing three columns: `sentence1`, `sentence2`, and `similarity`.
- **Significance:** Conforms to PyTorch's `Dataset` interface for easy batching.

The `collate_fn(batch)` function

```
def collate_fn(batch):
    max_len = 20
    ... // you can see more detail in source code
    return input_tensor, labels_tensor
```

- **Purpose:** This custom `collate_fn` prepares a batch of input tensors and labels for a two-channel CNN. It tokenizes each sentence, retrieves word embeddings, decomposes each embedding into "similar" and "dissimilar" components, truncates/pads them to a fixed length, and finally stacks them into a $(B, 2, \text{max_len}, \text{emb_dim})$ tensor along with the corresponding similarity labels.

- **Significance:** Organizes the data into two channels (similar/dissimilar) suitable for the CNN input.

Two-Channel CNN Model

```
class TwoChannelCNN(nn.Module):
    def __init__(self, emb_dim, num_filters=64, kernel_size=3):
        super().__init__()
        self.conv = nn.Conv2d(
            in_channels=2,
            out_channels=num_filters,
            kernel_size=(kernel_size, kernel_size),
            padding=(1,1)
        )
        self.pool = nn.AdaptiveMaxPool2d((1,1))
        self.fc = nn.Linear(num_filters, 1)
        self.sigmoid = nn.Sigmoid()

    def forward(self, x):
        feat = self.conv(x)          # => (B, num_filters, ?, ?)
        feat = nn.functional.relu(feat)
        feat = self.pool(feat)       # => (B, num_filters, 1, 1)
        feat = feat.squeeze(-1).squeeze(-1) # => (B, num_filters)
        out = self.fc(feat)          # => (B, 1)
        out = self.sigmoid(out)      # => (B, 1), [0..1]
        return out
```

- **Purpose:** The TwoChannelCNN class applies a 2D convolution over the two input channels (similar and dissimilar), pools the feature maps, and produces a single score (0–1) representing the predicted similarity between two sentences.

Main Function

```
# Main function to train and save the model
if __name__ == "__main__":
    df_all = pd.read_csv("data.csv", sep=";", decimal=",")
    ...
    def eval_loss(dloader):
        ...
        return np.mean(losses)

    num_epochs = 10
    for epoch in range(num_epochs):
        ...
        print(f"Epoch {epoch+1}/{num_epochs}, "
              f"TrainLoss={np.mean(train_losses):.4f}, "
              f"ValLoss={val_l:.4f}")
    torch.save(model.state_dict(), "model_weights_vi.pt")
```


- **Purpose:** The code reads a CSV file (data.csv) containing pairs of sentences and their similarity scores.

It shuffles the data, splits it into an 80% training set and 20% validation set, and creates SimilarityDataset objects for both.

Two PyTorch DataLoaders (train_loader and val_loader) are built, using a custom collate_fn that processes the data into two channels (similar/dissimilar).

A TwoChannelCNN model is instantiated on the chosen device (CPU or GPU), and a mean squared error loss (MSELoss) is used.

The training loop runs for num_epochs, each time performing a forward pass, computing the loss, backpropagating (loss.backward()), and updating weights with Adam.

After each epoch, the validation loss is computed via the eval_loss function. Finally, the model's trained weights are saved to model_weights_vi.pt.

Prediction Function

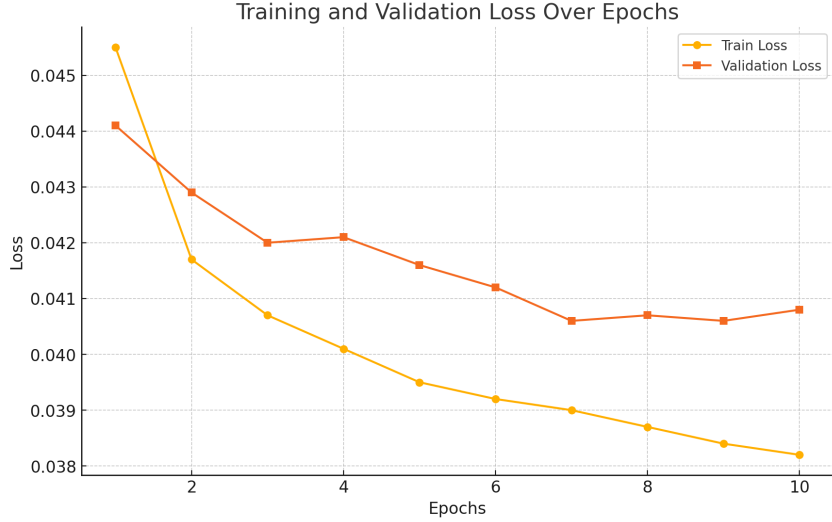
```
def predict_similarity(model, sentence1, sentence2, max_len=20):
    ...
```

- **Purpose:** This function predict_similarity calculates the similarity score between two sentences by tokenizing and embedding the sentences, matching and decomposing embeddings into "similar" and "dissimilar" components, preparing a two-channel input for the model, and passing it through a trained model in evaluation mode to generate the final similarity score.

5 Experiment

5.1 Training Model

The training process demonstrated steady convergence over 10 epochs, with both training loss (TrainLoss) and validation loss (ValLoss) decreasing consistently. The small gap between TrainLoss and ValLoss indicates effective generalization, with no signs of overfitting or underfitting. Validation loss stabilized around 0.0406–0.0412 after epoch 6, suggesting the model has reached optimal performance on the validation set. Minor fluctuations observed are within acceptable bounds due to the stochastic nature of training. The results suggest the model is well-trained and suitable for deployment or further experimentation.



(a) Figure 2: Training result.

5.2 Test Set Evaluation

We trained and evaluated several iterations of our sentence similarity model using different training data configurations and varying the number of training epochs. The primary goal of this evaluation was to select the model that generalizes best to unseen data, as measured by the F1-score on a held-out test set. A balanced F1-score is preferred in our application because we aim to minimize both false positives (incorrectly predicting high similarity) and false negatives (incorrectly predicting low similarity).

Initially, we trained a baseline model (`model_weights_vi`) using [describe initial training data] and 10 epochs. Subsequent models (`model_weights_vi_1` through `model_weights_vi_7`) explored variations in training data size, data augmentation techniques (if any), and the number of training epochs. For example, `model_weights_vi_1` used [describe data for this model] and [number] epochs, while `model_weights_vi_3` incorporated [describe changes, e.g., data augmentation] and trained for [number] epochs. [Optionally describe other key variations in data or training setup].

Table 1 summarizes the performance of each model on the test set. The model `model_weights_vi_5`, trained with [describe the specific data and epochs for this model], achieved the highest F1-Measure of 0.7724, indicating the best balance between Precision and Recall. While `model_weights_vi_3` achieved a remarkably high Recall (0.9713), its significantly lower Precision (0.6122) suggests a tendency to overpredict similarity, which is undesirable for our use case. [You can add a sentence or two explaining why high recall but low precision is not desirable for your *specific* application].

The improvements observed with `model_weights_vi_5` suggest that [explain the likely reason for the improvement - e.g., the increased data size, the data augmentation technique, or the optimal number of epochs]. This highlights the importance of carefully tuning the training process and data selection to achieve optimal performance for sentence similarity tasks.

Model	Precision	Recall	F1-Measure
model_weights_vi	0.8135	0.2836	0.4206
model_weights_vi_1	0.7036	0.7582	0.7299
model_weights_vi_2	0.7317	0.6526	0.6899
model_weights_vi_3	0.6122	0.9713	0.7510
model_weights_vi_4	0.6845	0.8729	0.7673
model_weights_vi_5	0.6882	0.8800	0.7724
model_weights_vi_6	0.7903	0.4824	0.5991
model_weights_vi_7	0.7218	0.7611	0.7409

Table 1: Test set evaluation results (threshold = 0.7).

5.3 Q&A Search Application

Once we have the model, we start testing on the Q&A dataset to find the questions that have the highest similarity to the user’s question. Our model will list the questions in the Q&A dataset that have a high similarity value (greater than 0.7) and display them on the screen along with the answer.

```

Question: hang động được hình thành như thế nào

Top most similar questions with answers:
1. hang động sông băng được hình thành như thế nào? - Similarity: 0.75
Answers:
- Một hang động sông băng chìm một phần trên sông băng Perito Moreno.
- Mặt tiền băng cao khoảng 60 m
- Các khối băng trong hang động sông băng Titlis
- Hang động sông băng là hang động được hình thành bên trong lớp băng của sông băng.
- Hang động sông băng thường được gọi là hang động băng, nhưng thuật ngữ này được sử dụng chính xác để mô tả các hang động nền đá chứa băng quanh năm.

```

(a) Figure 4: Q&A search application.

6 Conclusion

This project successfully implemented a model to measure the semantic similarity between Vietnamese sentences by leveraging a novel approach inspired by lexical decomposition and composition methods. The use of a two-channel convolutional neural network (CNN) enabled the model to effectively capture both the similar and dissimilar semantic components of sentence pairs.

6.1 Key Achievements

We incorporated the FastText pre-trained word embeddings (cc.vi.300.vec) to represent Vietnamese words in a high-dimensional vector space, facilitating the effective semantic comparison of words. The linear decomposition method allowed us to split each word vector into its similar and dissimilar components based on their semantic alignment with the other sentence. This step addressed the lexical gap challenge by preserving both matched and unmatched information. A two-channel CNN architecture was designed and implemented to extract features from the decomposed components, enabling a richer representation of sentence pairs. We trained and validated the model on a Vietnamese dataset of 5000 sentence pairs, achieving promising results, although limited by the size of the dataset. The model was applied to a practical Q&A search application,

which demonstrated its potential in real-world scenarios by ranking questions based on similarity.

6.2 Limitations

The performance of the model is limited by the size and diversity of the training dataset. A larger and more varied corpus would likely improve the generalization and accuracy of the model. The similarity scores produced by the model are not fully optimal, suggesting room for improvement in hyperparameter tuning, architecture design, and data augmentation. Handling out-of-vocabulary words with zero vectors, while practical, may have reduced the quality of embeddings in certain cases, especially for rare or domain-specific words.

In conclusion, this project demonstrates the feasibility and effectiveness of using lexical decomposition and two-channel CNNs for sentence similarity tasks in Vietnamese. While the results are promising, further improvements in data size, model complexity, and domain adaptation are needed to fully realize the potential of this approach in practical applications.

7 Other Models

In addition to the method we just implemented above, here we will mention 2 more methods.

Model 1: Context-Based Semantic Similarity

This method leverages the context around sentences to predict semantic similarity. The approach combines both context-aware representations and unsupervised semantic learning. [2]

Core Components:

1. Contextualized Semantic Representation:
 - A sentence is understood not only by its content but also by how it interacts with its surrounding context.
 - The semantic representation of a sentence is built by predicting the likelihood of its surrounding context.
2. Two Learning Approaches:
 - Discriminative Model:
 - Predicts how natural or logical a sentence appears when paired with its left and right contexts.
 - Uses conditional probabilities $P(C|S)$, where C is the context and S is the sentence.
 - Generative Model:
 - Predicts the likelihood of generating words in a sentence based on its surrounding context.
 - Uses conditional probabilities $P(S|C)$, where C is the context and S is the sentence.
3. Measuring Semantic Similarity:
 - Each sentence is represented as a vector in a contextualized semantic space.

- The similarity between two sentences is computed using the cosine similarity of their respective vectors:

$$\text{cosine_similarity}(u, v) = \frac{u \cdot v}{\|u\| \|v\|}$$

- Sentences with closer vectors are more semantically similar.

4. Surrogate Model for Efficiency:

- Directly computing context for all sentences in a dataset is computationally expensive.

- A surrogate model is trained to quickly approximate the similarity between two sentences, bypassing the need for full context computation.

Applications: - This method has been used for tasks such as paraphrase identification and answer retrieval. It performs well in both unsupervised and supervised learning setups, providing rich semantic insights by considering both the matched and unmatched portions of sentences.

Model 2: A Semantic and Word Order-Based Sentence Similarity Method

This method computes sentence similarity by combining two key components: semantic similarity and word order similarity. The approach leverages lexical databases (e.g., WordNet) and mathematical techniques to capture both semantic relationships and structural information in sentences. [3]

Core Components:

1. Semantic Similarity: - Word-Level Similarity:

- Each word in a sentence is mapped to a vector in a lexical database (e.g., WordNet), which organizes words hierarchically.

- The semantic similarity between two words is computed based on:

- The shortest path between their nodes in the hierarchy.

- The depth of the nodes (deeper nodes represent more specific meanings).

- Sentence-Level Similarity:

- A joint word set is created by combining all unique words from both sentences.

- Each sentence is represented as a semantic vector, where each dimension corresponds to the similarity of a word in the joint set with the sentence.

- Cosine similarity between the two semantic vectors is used to measure the semantic similarity of the sentences:

$$\text{Semantic Similarity} = \cos(\theta) = \frac{\vec{S}_1 \cdot \vec{S}_2}{\|\vec{S}_1\| \|\vec{S}_2\|}$$

2. Word Order Similarity:

- Each sentence is represented as a word order vector based on the joint word set.

- If a word exists in the sentence, its position is recorded; otherwise, it is assigned a value of 0.

- Example:

- Sentence 1: "The dog chased the cat."

- Sentence 2: "The cat was chased by the dog."

- Joint word set: "dog," "chased," "cat," "was," "by," "the".

- Word order vector for Sentence 1: [2, 3, 4, 0, 0, 1].

- Word order vector for Sentence 2: [6, 4, 2, 3, 5, 1].
- Word order similarity is calculated as:

$$\text{Word Order Similarity} = 1 - \frac{\|\vec{O}_1 - \vec{O}_2\|}{\|\vec{O}_1 + \vec{O}_2\|}$$

3. Combined Similarity:

- The overall similarity score is calculated as a weighted combination of semantic similarity and word order similarity:

$$\text{Sentence Similarity} = \alpha \cdot \text{Semantic Similarity} + (1 - \alpha) \cdot \text{Word Order Similarity}$$

- Here, α is a weight parameter, typically set to prioritize semantic similarity (e.g., $\alpha = 0.8$).

Advantages: - Handles Synonyms and Paraphrases:

- By incorporating semantic similarity, this method effectively compares sentences with different but semantically related words (e.g., "The dog is fast" vs. "The animal is quick").

- Considers Sentence Structure:

- Word order similarity ensures that the method differentiates between structurally different sentences with similar words (e.g., "The dog bit the man" vs. "The man bit the dog").

- Simplicity:

- The method avoids complex deep learning models, making it computationally efficient and easy to implement.

Applications: This method has been applied to tasks such as conversational agents, information retrieval, and paraphrase detection. Its simplicity and effectiveness make it suitable for systems requiring interpretable and computationally efficient solutions.

References

- [1] Zhiguo Wang and Haitao Mi and Abraham Ittycheriah. "Sentence Similarity Learning by Lexical Decomposition and Composition."
- [2] Xiaofei Sun, Yuxian Meng, Xiang Ao, Fei Wu, Tianwei Zhang, Jiwei Li, and Chun Fan. "Sentence Similarity Based on Contexts"
- [3] Yuhua Li, Zuhair Bandar, David McLean and James O'Shea "A Method for Measuring Sentence Similarity and its Application to Conversational Agents"