

Exploring the linear separability of syntactic and semantic information in BERT embeddings

Qingxia Guo, Saiya Karamali, Lindsay Skinner, and Gladys Wang

University of Washington

{qq07, karamali, skinnel, qinyanw}@uw.edu

Abstract

Syntax and semantics are at the foundation every language, yet distinctions between the two are not readily agreed upon. We seek to explore how representations of these two information sets manifest in BERT embeddings. Specifically we investigate the degree of the linear separability of syntactic and semantic information in BERT embeddings, as well as quantify how important the linear component corresponding to one information set is to solving a classification task on the other information set. We use Iterative Nullspace Projection to decompose word-level BERT embeddings into syntactic, non-syntactic, semantic and non-semantic components to be used in syntactic and semantic classification tasks. Our results show that...

1 Introduction

The boundary between semantics and syntax is a hotly debated topic in linguistics, but do large language models make such a distinction? If they do make such a distinction, do language model embeddings present this information in a way that is easily separated and recognized by humans? The objective of this project is to explore BERT’s (Devlin et al., 2019) reliance on certain syntactic information when handling a semantic task, and vice versa. Specifically, we seek to quantify the importance of linearly-separable syntactic or semantic information when performing semantic or syntactic classification, respectively.

To achieve our goal, we construct a linear probing system for a task and then employ Iterative Nullspace Projection (INLP from here on) (Ravfogel et al., 2020) to generate a new embedding devoid of information learned from the probing task. We then measure the performance of this new embedding on downstream syntactic and semantic classification tasks. The design of our probing procedure follows (Elazar et al., 2020), which employs

INLP to investigate whether BERT uses part-of-speech (POS) information when solving language modeling (LM) tasks. INLP has been used for a variety of tasks ((Ravfogel et al., 2020), (Elazar et al., 2020), (Gonen et al., 2020), etc.) though this is the first case we know of in which INLP has been used to investigate the linear separability of syntactic and semantic information.

A novel method for removing information from an embedding, INLP iteratively trains linear models on a specific classification task, and projects the input on intersection of the nullspaces of those linear models. Our objective is that, by applying the INLP procedure to a syntactic task, we are able to separate the representation into a syntactic space and a non-syntactic space. We then compare the performance of a model that seeks to classify semantic labels using the original BERT embeddings with an otherwise identical model trained on embeddings projected onto the non-syntactic space, in order to see if BERT is using syntactic information when performing the semantic task. Conversely, we can also first probe a semantic task, thus defining a semantic and non-semantic space, and then investigate the performance of embeddings projected onto those spaces when performing a syntactic classification task. Once we derive the semantic space and syntactic space from the experiment, we further investigate on the separability of the two spaces by comparing the embedding projections.

To evaluate the separability of syntax and semantics, we use two tasks: one task for the probing system and INLP procedure, and one task for evaluating performance on embeddings before and after INLP. We choose Combinatory Categorical Grammar (CCG from here on) tagging (Hockenmaier and Steedman, 2007) as the syntactic task and semantic category labeling (Bonial et al., 2014) as the semantic task.

The remainder of the paper proceeds as follows:

Section 2 explores previous work related to our experiment. Section 3 provides a description of the probing and evaluation tasks and gives an overview of the experiment pipeline. Section 4 reviews our experiments and affiliated results. Section 5 discusses the implications of those results. Finally, section 6 gives an overview of the entire process and outlines possible next steps.

2 Related Work

The separation and overlap between syntax and semantics has been of interest to linguists for many years. More recently, with the growing popularity of large language models, computational linguists have begun to explore how large language models deal with the boundaries of these information sets in word and sentence embeddings.

Huang et al., 2021 use paraphrase pairs and new target syntax to train a semantic encoder, syntactic encoder and decoder to learn separate representations of the semantic and syntactic information contained in BART embeddings, in order to create semantically equivalent paraphrases with the new syntactic structure. Alongside the encoders they also train an adversarial syntax discriminator to try and predict the source syntax from the semantic embeddings, thus encouraging the disentanglement of the semantic and syntactic information by training the semantic embedder to remove as much syntactic information as possible. Their results show that one can achieve some removal of syntactic information from semantic embeddings, so disentanglement of some information is possible. Though they do not achieve perfect separation of the two information sets. Other non-linear approaches to syntactic-semantic information disentanglement have been carried out in (Chen et al., 2019)

Unlike the aforementioned studies, we seek to explore the linear separability of syntactic and semantic information in large language model embeddings at the word level. To accomplish this task we apply the Iterative Nullspace Projection method to syntactic (CCG) and semantic labeling tasks in order to define the syntactic and semantic components of BERT embeddings that will be used in our downstream classification tasks.

INLP, introduced in (Ravfogel et al., 2020), is a method to define a linear guarding function that masks all the linear information in a word embedding that may be used for a downstream classification

task. In the original paper the authors use this method to remove gender bias from BERT embeddings of biographical descriptions and then measure how easy it is to determine an individual’s gender from the guarded embedding by using various downstream classification methods. Beyond using the INLP method to guard protected attributes, the authors hypothesize several additional use cases for this procedure, including information disentanglement.

The authors of (Elazar et al., 2020) use INLP for exactly this task. They use INLP in order to separate and guard certain linguistic information sets from BERT embeddings in order to better understand what information is being used by large language models, and not just what is encoded. The main premise behind this paper is that if a particular property is used to solve a task, then the removal of that property should negatively influence the model’s ability to solve that task. Conversely, if the removal of a property has little influence on the model’s ability to perform a task then we know that property is not a significant contributing factor in the model’s ability to perform that task. Specifically, (Elazar et al., 2020) seeks to quantify the importance of the information sets used for part-of-speech tagging, syntactic dependency labeling, named entity recognition and syntactic constituency boundaries on BERT’s ability to perform the language modeling task.

We take a similar approach to (Elazar et al., 2020) by separating the information sets used for CCG tagging and semantic labeling from word-level BERT embeddings. However, as we are interested in the linear separability of these two information sets, we will test how the removal of these information sets impacts the embeddings’ performance on semantic labeling and CCG tagging, respectively, rather than language modeling.

3 Methods

We construct two separate probing tasks to isolate the syntactic and semantic information in word-level BERT embeddings. The embeddings are separated into syntactic and non-syntactic, and semantic and non-semantic components via INLP which is described in section 3.1. These embedding components are then combined to form new embeddings, which are evaluated on the same tasks that were used for probing.

3.1 The Iterative Null-Space Projection method

The INLP method first introduced in (Ravfogel et al., 2020), is used to create a guarding function that masks all the linear information contained in a set of vectors, X , that can be used map each vector to $c \in C$, where C is the set of all categories. This is accomplished by training a linear classifier, a matrix W , that is applied to each $x \in X$ in order to predict the correct category c with the greatest possible accuracy. In other words, Wx defines a distribution over the set of categories C and we assign x to the class $c \in C$ which is allotted the greatest probability by Wx . Note that the classifier’s accuracy must be greater than that achieved by guessing the majority category, otherwise x contains no linear information relevant for the categorization task and thus no guarding function is needed. Once W is determined, for any $x \in X$ we can remove the information that W uses to predict c by projecting x onto the null-space of W , $N(W) = \{x | Wx = 0\}$. Call this projection function P_1 and let $\hat{x} = P_1(x)$. This removes all of the linear information in x that W used to predict the category c .

However, this process does not necessarily remove all of the linear information in x that could be used to predict c . For example, x may contain redundant information and W may have only used one set of this information for its prediction. In this case, the redundant information would still be present in \hat{x} . Thus, we must repeat the above process, defining a new linear classifier \hat{W} that uses \hat{x} to predict c . If \hat{W} is still able to predict c with a greater than majority class guess accuracy, then we know that \hat{x} contained linear information about c . As above, we project \hat{x} onto the null-space of \hat{W} via the projection function P_2 and define a new $\hat{x} = P_2(P_1(x))$.

We iteratively apply this process until no linear information remains in \hat{x} , i.e. a linear classifier is unable to predict the correct category c with any probability greater than that achieved by guessing the majority class. The final $\hat{x} = P_n(P_{n-1}(\dots P_1(x)))$ contains no linear information about the categories in C and we call $P(x) = P_n(P_{n-1}(\dots P_1(x)))$ the guarding function.

We will pair the INLP method with the probing tasks described in sections 3.3 and 3.4 in order to create two guarding functions that will enable us to isolate the linear components of BERT embeddings

that contain syntax-specific and semantics-specific information.

3.2 Data

We use the English Parallel Meaning Bank v4.0 (Abzianidze and Bos, 2017) to test the linear separability of the semantic and syntactic information in word-level BERT embeddings. This dataset is consist of gold standard and silver standard of word-level semantic tags. The gold standard contains 5,438 sentences with annotations that are manually verified while the silver standard contains 62,739 sentences with autogenerated annotations. All of our experiments will be conducted on gold standard data.

The original dataset does not include CCG Tags, but Abzianidze and Bos, 2017 utilized CCG parser to produce CCG tags. We follow a similar procedure to apply a CCG parser (Yoshikawa et al., 2017) to develop word-level CCG tags. Once we obtain both CCG tags and semantics tags for the dataset, we can perform syntactic and semantics probing task as desired.

3.3 Syntactic probing task

The syntactic probing task involves training a linear classifier on the final layer BERT embeddings in order to predict the CCG tag associated with each word. We will use this classifier in the INLP algorithm in order to create a guarding function for the information that is necessary to complete the CCG labeling task. For a given embedding, x , the projection that results from applying this guarding function, P_{syn} , to the embedding will represent the non-syntactic information contained in the embedding and will from now on be referred to as the “non-syntactic component” of the embedding, $v_{nosyn} = P_{syn}x$. We can then determine the “syntactic component” of the embedding by taking the difference of the embedding vector with the non-syntactic component, $v_{syn} = x - v_{nosyn}$.

3.4 Semantic probing task

Similar to the above, the semantic probing task involves training a linear classifier on the final layer BERT embeddings in order to predict the semantic tag (described in the data section) associated with each word. This classifier is used in the INLP algorithm in order to create a guarding function, P_{sem} , for the information necessary to complete the Semantic tag labeling task. As described in the Syntactic probing task section, we shall use the

resulting guarding function to decompose the original embedding into a “non-semantic component”, $v_{nosem} = P_{sem}x$, and a “semantic component”, $v_{sem} = x - v_{nosem}$.

3.5 Evaluation tasks

Our goal is to determine which information sets captured in the BERT embeddings are relevant for our evaluation tasks. We thus use the components derived from the probing tasks to create new embeddings that isolate specific types of information. These embeddings are then evaluated on the syntactic and semantic tasks that were used for probing, and their performance is compared to that of the original embeddings. We also compare the performance of each model trained on one of these embeddings with another trained on new embeddings that are created by randomly removing the same number of dimensions from the original embeddings as are removed by the INLP guarding function. In doing so we can test the extent to which the loss of the particular information set of interest is responsible for the drop in performance, as opposed to a general loss of information.

The new embeddings to be tested include the syntactic component, the non-syntactic component, the semantic component and the non-semantic component derived from the probing tasks. Additionally, we can create an embedding that contains syntactic information and removes semantic information by linearly projecting the syntactic component onto the non-semantic component, $v_{syn-sem} = P_{sem}v_{syn}$. Using a similar process, we can create an embedding that contains semantic information and removes the syntactic information present, $v_{sem-syn} = P_{syn}v_{sem}$. Finally, we can create an embedding that contains the semantic information captured by the syntactic component, by linearly projecting the syntactic component onto the semantic component, $v_{syn*sem} = \langle v_{syn}, v_{sem} \rangle u_{sem}$ where u_{sem} is the unit vector for v_{sem} . And similarly, we can create an embedding that contains the syntactic information captured by the semantic component, $v_{sem*syn} = \langle v_{sem}, v_{syn} \rangle u_{syn}$.

We will assess each of these embedding types, the original BERT embeddings and the embeddings created by randomly removing directional information on the CCG and Semantic labeling tasks that were used in the probes.

4 Results

We first probe our tasks on the original embeddings, and linear classifiers can successfully predict both CCG tags and semantics tags (around 85% testing accuracy) as resulted in table 2. Then we apply INLP method to the linear classifiers to derive guarding matrices P_{syn} and P_{sem} which will project the embeddings to the nullspace of syntactic space and semantic space respectively. With some linear tranformation with the original embeddings and the guarding matrices, we are able to extract the embeddings described in table 1.

expression	description
v_{orig}	Original BERT embeddings
$v_{nosem} = P_{sem}v_{orig}$	Gained after INLP with semantics task
$v_{nosyn} = P_{syn}v_{orig}$	Gained after INLP with syntactic task
$v_{sem} = v_{orig} - v_{nosem}$	Semantics representation
$v_{syn} = v_{orig} - v_{nosyn}$	Syntactic representation
$v_{syn-sem} = P_{sem}v_{syn}$	Contains syntactic features only
$v_{sem-syn} = P_{syn}v_{sem}$	Contains semantic features only
$v_{syn*sem} = \langle v_{syn}, v_{sem} \rangle \cdot \frac{v_{syn}}{\ v_{sem}\ }$	Syntactic representations projected on semantic space
$v_{sem*syn} = \langle v_{syn}, v_{sem} \rangle \cdot \frac{v_{syn}}{\ v_{syn}\ }$	Semantics representations projected on semantic space

Table 1: Description of Embeddings

To further compare the impact of INLP, we conduct experiments with random directions removed. For each embeddings, we record the number of removed directions using $Null(P)$ where P is the projection matrix. For embeddings not derived from matrix multiplication, we obtain the number from $Null(M)$ where M is the embedding matrix with size (768,instance number).¹ Then we create embeddings with the same number of directions randomly removed and probe the tasks. All the testing accuracy of our final experiments are in table 2.

¹This is not necessarily equivalent to the number of direction removed. $Null(M) \geq Null(P)$ for the corresponding projection P , but in practice the numbers are very close

Embedding	Directions Removed	CCG Tagging	Semantics Tagging
v_{orig}	0	84.75%	88.56%
$\text{Rand}(v_{nosem})$	77	84.06%	87.71%
$\text{Rand}(v_{nosyn})$	159	82.26%	87.57%
$\text{Rand}(v_{syn-sem})$	81	83.11%	87.61%
$\text{Rand}(v_{sem-syn})$	161	82.46%	86.04%
$\text{Rand}(v_{syn*sem})$	0	p%	p%
$\text{Rand}(v_{sem*syn})$	0	p%	p%
v_{nosem}	77	27.93%	17.28%
v_{nosyn}	159	23.76%	49.43%
$v_{syn-sem}$	81	34.24%	33.91%
$v_{sem-syn}$	161	26.96%	50.21%
$v_{syn*sem}$	0	p%	p%
$v_{sem*syn}$	0	p%	p%

Table 2: Experiment Result of Different Embeddings

5 Discussion

The INLP method successfully guard information retrieved from probe task. When performing CCG tagging, the model accuracy drops significantly (from 84.06 % to 23.76%), and similarly for semantics tagging. To measure the relevancy of semantics information to syntactic task, we compared the performance of CCG tagging with embeddings v_{nosem} with that of embeddings v_{nosyn} . The both performances have considerable decrease as to the original embeddings. The model is using representation from both space when probe the syntactic task. Meanwhile, the performances of semantics tagging with the two embeddings differ significantly, with v_{nosyn} performs almost three times better than v_{nosem} . The result suggests that the semantics task is less dependent on the syntactic information than the opposite direction, which contradicts the assumption that semantics space will be rel

To investigate the compositionality of semantics space and syntactic space, we also calculated the average cosine similarity scores between embeddings in table 3.

expression	similarity score
v_{orig}, v_{nosem}	0.4183
v_{orig}, v_{nosyn}	0.3609
v_{nosyn}, v_{nosem}	0.2830
$v_{syn-sem}, v_{sem-syn}$	0.2579

Table 3: Similarity of Semantics and Syntactic Representations

6 Conclusion

This is where our conclusion will go. This is where our future work will go.

References

- Lasha Abzianidze and Johan Bos. 2017. Towards universal semantic tagging. In *Proceedings of the 12th International Conference on Computational Semantics (IWCS 2017) – Short Papers*, pages 1–6, Montpellier, France.
- Claire Bonial, Julia Bonn, Kathryn Conger, Jena D. Hwang, and Martha Palmer. 2014. [PropBank: Semantics of new predicate types](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 3013–3019, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Mingda Chen, Qingming Tang, Sam Wiseman, and Kevin Gimpel. 2019. [A multi-task approach for disentangling syntax and semantics in sentence representations](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2453–2464, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Yanai Elazar, Shauli Ravfogel, Alon Jacovi, and Yoav Goldberg. 2020. [Amnesic probing: Behavioral explanation with amnesic counterfactuals](#).

Hila Gonen, Shauli Ravfogel, Yanai Elazar, and Yoav Goldberg. 2020. [It’s not Greek to mBERT: Inducing word-level translations from multilingual BERT](#). In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 45–56, Online. Association for Computational Linguistics.

Julia Hockenmaier and Mark Steedman. 2007. [CCG-bank: A Corpus of CCG Derivations and Dependency Structures Extracted from the Penn Treebank](#). *Computational Linguistics*, 33(3):355–396.

James Y. Huang, Kuan-Hao Huang, and Kai-Wei Chang. 2021. Disentangling semantics and syntax in sentence embeddings with pre-trained language models. In *NAACL*.

Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. [Null it out: Guarding protected attributes by iterative nullspace projection](#).

Masashi Yoshikawa, Hiroshi Noji, and Yuji Matsumoto. 2017. [A* ccg parsing with a supertag and dependency factored model](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 277–287. Association for Computational Linguistics.