

ANALYSING THE LINEARITY OF LINGUISTIC RELATIONS IN LANGUAGE MODEL EMBEDDING SPACES

Anonymous authors

Paper under double-blind review

ABSTRACT

We propose a framework to analyse how strongly different linguistic relations are linearly encoded in language model embedding spaces. We formalise linear encoding via a constrained linear approximation over related and unrelated word pairs and apply this to an extended BATS dataset covering inflectional, derivational, lexicographic, and encyclopedic relations in GloVe, RoBERTa, and ModernBERT. Our experiments show near-perfect linear encodings for inflectional and derivational relations, but substantially higher errors for lexicographic and encyclopedic relations, especially for one-to-many and many-to-many associations. We also find that RoBERTa and ModernBERT generally encode relations more linearly than GloVe. These results indicate that our framework can reveal which relational structures are most linearly accessible in embeddings, offering a compact tool for probing and comparing relational geometry across models.

1 INTRODUCTION

Large language models (LLMs) and other deep learning-based natural language processing models function by transforming the input text into high-dimensional numerical vectors called embeddings, in which meaning is represented in a distributed way. What such an embedding vector represents depends on its numerical values and its position relative to other embedding vectors in the embedding space. This distributed, high-dimensional coding makes language processing models powerful but also opaque, because it is hard to see what information—especially about linguistic relationships—is encoded where, and how it influences model behaviour.

Probing methods are a widely used way to study what kinds of information are encoded in a deep learning based language processing model’s internal representations by testing what information can be recovered from these representations (Conneau et al., 2018; Nedumpozhimana & Kelleher, 2024). This work proposes a novel framework to analyse the latent representations of language processing models, and goes beyond the standard probing methods in two key ways. First, rather than simply testing whether particular information is present in an embedding, our framework can be used to understand how it is encoded in the embedding space, and more specifically, whether it is represented in a linear or non-linear form. The theoretical inspiration for our approach is the linear representation hypothesis (Park et al., 2024), which proposes that, for at least some linguistic properties, models organise their internal space linearly. Such linearly encoded linguistic properties may be more accessible to the model’s downstream computation as compared to other information in the model’s representations, and so may disproportionately influence the model’s behaviour. Thus, by identifying which information is encoded linearly, we can begin to explain which information strongly drives model behaviour and how we might safely intervene on this behaviour. Second, while much existing work has applied probing to individual concepts or token-level properties, we focus on linguistic relations (such as syntactic or semantic relations). Adopting a relation-based rather than concept-based perspective is both novel and advantageous because a relational view asks how models represent the links between elements in text, which drive many downstream behaviours.

The framework is defined for arbitrary linguistic relations (word-to-word, word-to-sentence, and sentence-to-sentence), and can handle varying relational complexity (one-to-one, one-to-many, and many-to-many).

2 LINEARLY ENCODED RELATIONS

We formalise the concept of the linearity of a relation by defining that any relation r is linearly encoded in the embedding space if there exist two linear operators that map representations of a pair to the same embedding vector if and only if that pair is related. In linear algebra, a linear relation is one where related pairs (e_1, e_2) in a module M over a ring R satisfy a linear equation $f_1 e_1 + f_2 e_2 = 0$, where f_1 and f_2 are two elements in the ring R (Lang, 2002). In our case, we consider the embedding space as a Module of all d -dimensional vectors (\mathbb{R}^d) over the ring of all $d \times d$ square matrices ($\mathbb{R}^{d \times d}$). Note that the space of square matrices over matrix addition and matrix multiplication is a ring, and therefore, d -dimensional vectors over $d \times d$ matrices are a module.

Suppose r be the target linguistic relation and t_1 and t_2 are related linguistic expressions (i.e., $(t_1, t_2) \in r$). Let E be the embedding mapping that maps any linguistic expression to a d -dimensional embedding vector in the embedding space (\mathbb{R}^d) and let e_1 and e_2 be the two d -dimensional embedding vectors of linguistic expressions t_1 and t_2 represented as column matrices. Then, if the relation r is linearly encoded in the embedding space, then there exist two $d \times d$ square matrices L_r and R_r that correspond to the relation r that maps both e_1 and e_2 to the same vector. To align this definition with the standard definition of a linear relation, we can multiply the R_r operator matrix by -1 , so that it will obey the linear equation $L_r e_1 + R_r e_2 = 0$. For more notational simplicity, we can concatenate e_1 and e_2 to create a single $2d$ -dimensional vector e_{12} , and column wise concatenate L_r and R_r to create a single $d \times 2d$ matrix M_r . Then we can formally define that if a relation r is linearly encoded in the embedding space defined by the embedding mapping E , then there exists an M , such that:

$$M_r e_{12} = 0 \iff (t_1, t_2) \in r \quad (1)$$

3 LINEAR APPROXIMATION

Since many linguistic relations will not satisfy the exact linear encoding condition in Equation 1, we next define a linear approximation that quantifies how closely a relation can be represented linearly.

Some relations can be approximately encoded linearly; when the embeddings contain noise, this can often be corrected with slight modifications. Whereas some relations can only be linearly encoded by excluding extreme instances. Even for relations that are not exactly linearly encodable, it can still be informative to quantify the degree of linearity they exhibit.

One can see that if we are not considering unrelated pairs of expressions, then a trivial solution (i.e., $M_r = 0$) exists for any relation. To avoid such a trivial solution, we also need to consider unrelated pairs along with related pairs. In the case of related pairs, based on the condition 1, $M_r e_{12}$ should be a 0 vector, and therefore its Euclidean norm will be 0. Whereas, in the case of unrelated pairs (\bar{t}_1, \bar{t}_2) , $M_r e_{\bar{12}}$ should not be a 0 vector ($e_{\bar{12}}$ is the concatenated embedding of \bar{t}_1 and \bar{t}_2), and therefore its Euclidean norm will be strictly greater than 0. In that case, by scaling M_r we can make sure that the Euclidean norm will be greater than or equal to 1. Therefore, we rewrite the condition 1 as:

$$\|M_r e_{12}\|^2 = 0, \forall (t_1, t_2) \in r \text{ and } \|M_r e_{\bar{12}}\|^2 \geq 1, \forall (\bar{t}_1, \bar{t}_2) \notin r \quad (2)$$

Based on this condition, we define the linear approximation of a linguistic relation r as the matrix \tilde{M}_r such that,

$$\tilde{M}_r = \underset{(t_1, t_2) \in r}{\operatorname{argmin}} \left\{ \sum \|M e_{12}\|^2 \right\} \text{ such that, } \|M e_{\bar{12}}\|^2 \geq 1, \forall (\bar{t}_1, \bar{t}_2) \notin r \quad (3)$$

Practically, it is not possible to consider all related pairs and all unrelated pairs, and therefore, we select p related pairs and n unrelated pairs. We can concatenate the embeddings of p related pairs to create a $2d \times p$ matrix P_r , and n unrelated pairs to create a $2d \times n$ matrix N_r . Now, we can restate the optimisation problem as,

$$\tilde{M}_r = \underset{\{tr((MP_r)^T(MP_r))\}}{\operatorname{argmin}} \left\{ \|Mv\|^2 \geq 1, \forall v \in \operatorname{columns}(N_r) \right\} \quad (4)$$

This optimisation has a quadratic objective function with $2d^2$ variables and n quadratic constraints. To simplify this optimisation, we apply singular value decomposition on P_r , such that $P_r = U\Sigma V^T$. Let's assume the SVD of M is QSR^T , then we constrain our search space of M such that the right singular matrix of M is the same as the left singular matrix of P_r (i.e., $R = U$). Then we can rewrite the objective function of the above optimisation problem as,

$$\text{tr}((MP_r)^T(MP_r)) = \text{tr}((QSU^T U\Sigma V^T)^T(QSU^T U\Sigma V^T)) \quad (5)$$

$$= \text{tr}((Q\Sigma V^T)^T(Q\Sigma V^T)) = \text{tr}(V\Sigma^T S^T Q^T Q\Sigma V^T) \quad (6)$$

$$= \text{tr}(V\Sigma^T S^T S\Sigma V^T) \quad (7)$$

Here, the Σ will be a $2d \times p$ matrix and S will be a $d \times 2d$ matrix, and therefore the number of non-zero diagonal entries of Σ will be at most $2d$, and that of S will be at most d . Let the diagonal entries of Σ be $\sigma_1, \sigma_2, \dots, \sigma_{2d}$ and the diagonal entries of S be s_1, s_2, \dots, s_d , then $\Sigma^T S^T S\Sigma$ will be a diagonal matrix with entries $\sigma_1^2 s_1^2, \sigma_2^2 s_2^2, \dots, \sigma_d^2 s_d^2$. The trace of a matrix is the sum of its singular values; therefore, the $\text{tr}(V\Sigma^T S^T S\Sigma V^T)$ will be $\sigma_1^2 s_1^2 + \sigma_2^2 s_2^2 + \dots + \sigma_d^2 s_d^2$. Let $\mathbf{x} = (s_1^2, s_2^2, \dots, s_d^2)$ and $\mathbf{c} = (\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2)$. Then we can rewrite equation 7 in terms of \mathbf{x} and \mathbf{c} as,

$$\text{tr}((MP_r)^T(MP_r)) = \mathbf{c} \cdot \mathbf{x} \quad (8)$$

Similarly, the constraints of the optimisation problem can be rewritten as,

$$\|M\mathbf{v}\|^2 \geq 1 \iff (QSU^T \mathbf{v})^T(QSU^T \mathbf{v}) \geq 1 \quad (9)$$

$$\iff \mathbf{v}^T U S^T Q^T Q S U^T \mathbf{v} \geq 1 \iff \mathbf{v}^T U S^T S U^T \mathbf{v} \geq 1 \quad (10)$$

$$\iff ((\mathbf{v}^T U) \odot (\mathbf{v}^T U)) \mathbf{x} \geq 1 \quad (11)$$

Where \odot is the element-wise multiplication. Now we can rewrite the optimisation in terms of \mathbf{x} ,

$$\tilde{\mathbf{x}} = \text{argmin}\{\mathbf{c} \cdot \mathbf{x}\} \text{ s.t.}, \quad (12)$$

$$\mathbf{A}\mathbf{x} \geq 1, \mathbf{x} \geq 0 \quad (13)$$

Where $\mathbf{A} = (N_r^T U) \odot (N_r^T U)$.

The optimum \mathbf{x} (i.e., $\tilde{\mathbf{x}}$) for a relation r serves two roles: its objective value $\mathbf{c} \cdot \tilde{\mathbf{x}}$, normalised by the number of related pairs, defines the *approximation error* for relation r (see Section 4), and its components determine the linear operator $\tilde{M}_r = \text{diag}(\sqrt{\tilde{\mathbf{x}}})U$ where $\sqrt{\tilde{\mathbf{x}}}$ is the element-wise square root of $\tilde{\mathbf{x}}$. Here we ignore the left singular matrix of \tilde{M}_r (Q) by setting it as the identity matrix. Note that the term Q will cancel out, and therefore the optimisation problem is independent of Q .

4 ARE RELATIONS ENCODED LINEARLY IN REPRESENTATIONAL SPACES?

We conducted a preliminary empirical analysis to check whether the proposed framework can be used to investigate whether linguistic relations are linearly encoded in the representational space of some well-known language processing models. For this experiment, we extended the *BATS* dataset (Gladkova et al., 2016), which contains 40 word-to-word relations, including 10 Inflectional relations, 10 Derivational relations, 10 Lexicographic relations, and 10 Encyclopedic relations. For each of these 40 relations, the *BATS* dataset lists 50 pairs of words for which that relation holds. To create a dataset for our experiments for each relation in *BATS*, we manually created 50 more related pairs, and created an extended *BATS* dataset with 100 related pairs for each relation. Many relations we consider in this experiment are one-to-many or many-to-many, and in such cases, every word is paired individually with every other related word, and hence, the number of data points varies for different relations. Details of the number of related pairs are shown in Table 1.

For each of these relations, we also created a set of unrelated pairs by using the words already present in the dataset. In this process of creating unrelated pairs, we treated the domain (set of all words that come first in the related pairs) and the range (set of all words that come second in the related pairs) as two different categories. By this segregation we avoid assuming that the domain and range of a relation should be the same. We then created the full Cartesian product of domain and range, and treated any pair not in the related set as an unrelated pair. As in the case of related pairs, the number of unrelated pairs also varies from one relation to another relation, and these details are shown in Table 1.

Table 1: Statistics of dataset and linear approximation errors (macro average) of *BATS* relations.

Relations	# Related pairs			# Unrelated pairs			Avg error of approximation		
	Min	Max	Avg	Min	Max	Avg	GloVe	RoBERTa	ModernBERT
Inflectional	100	131	108.6	9900	17030	11765.2	0	0	0
Derivational	101	202	124.7	9999	20200	13909.8	0	0	0
Encyclopedic	104	284	177.2	2751	11539	7544.9	0.4704	0.4685	0.4549
Lexicographic	209	1988	965.1	12814	197900	63923.3	0.9201	0.8139	0.8360

To generate representations of words, we used three models: a non-neural representation model, GloVe (Pennington et al., 2014); a neural representation model, RoBERTa (Liu et al., 2019); and one of the most recent neural representation models, ModernBERT (Warner et al., 2025). While generating the GloVe representation, if the word is not in the vocabulary, we randomly assign a fixed 300-dimensional representation for such words. To generate RoBERTa and ModernBERT representations, we selected the average final layer token embeddings (Note that a word can have multiple tokens) generated by the model from the input word. Then we analysed whether relations can be linearly encoded in these representational spaces by linearly approximating these relations and calculated the error of approximation. The approximation error for a relation is the objective function (Eg. 12) normalised by the number of related pairs; zero error implies perfect linear encoding.

From our empirical analysis, we found that both Inflectional and Derivational relations are linearly encoded in the representational spaces of all three models (with 0 approximation error). However, for Lexicographic and Encyclopedic relations, none of the models has a perfect linear encoding. We also found that, although the average values of errors of linear approximations are comparable for all three models, the RoBERTa and ModernBERT average scores are better than the GloVe average score. This shows that in more recent and powerful language models, relations are encoded more linearly. However, when we compare RoBERTa with ModernBERT, RoBERTa is better on Lexicographic relations, and ModernBERT is better on Encyclopedic relations.

Generally, we found that relations that are one-to-one are more likely to encode linearly in representational space. For example, all inflectional and derivational relations (morphological relations) are one-to-one, and we found near-perfect linear approximations for these relations. However, for relations with one-to-many or many-to-many related pairs, i.e., Lexicographic and Encyclopedic semantic relations, we observed that it is harder to find a linear approximation. For example, for one of the lexical relations, ‘*part-whole*’, which is a many-to-many relation, we got an above 1 average error of approximation for all three representation models (GloVe: 1.3017, RoBERTa: 1.1070, and ModernBERT: 1.1690). However, for the lexical relation with relatively fewer many-to-many related pairs, ‘*antonyms-binary*’, we got lower approximation errors (GloVe: 0.3309, RoBERTa: 0.3246, and ModernBERT: 0.3257). We observed a similar pattern in Encyclopedic relations.

When we further analysed the Encyclopedic relations, we found that relations that are non-deterministic or non-exclusive (one-to-many) are hard to approximate linearly compared to relations with strong, nearly one-to-one associations between entities. For example, in the case of ‘*country-language*’, languages like Malayalam and Hindi have a strong association with India and are therefore more linearly encoded in the embedding space than English, whose association with India is diffuse and non-exclusive. Similarly, for the ‘*thing-colour*’ relation, many related pairs such as ‘banana’ and ‘green’ are context-dependent/non-deterministic because a ‘banana’ can be ‘green’, but it can also be ‘yellow’, and for these pairs, we obtained higher approximation errors.

5 CONCLUSION

In this work, we proposed a framework for analysing the linearity of linguistic relations and applied it to 40 word-to-word relations of varying complexity in GloVe, RoBERTa, and ModernBERT. We found that inflectional and derivational relations admit near-perfect linear encodings, whereas lexicographic and encyclopedic relations—especially one-to-many and many-to-many mappings—yield substantially higher approximation errors, with RoBERTa and ModernBERT generally encoding relations more linearly than GloVe. This shows that our framework can pinpoint which relational structures are most linearly accessible in current language models and provides a practical tool for comparing relational geometry across architectures.

REFERENCES

- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. What you can cram into a single \$!#* vector: Probing sentence embeddings for linguistic properties. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2126–2136, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1198. URL <https://www.aclweb.org/anthology/P18-1198>.
- Anna Gladkova, Aleksandr Drozd, and Satoshi Matsuoka. Analogy-based detection of morphological and semantic relations with word embeddings: what works and what doesn’t. In Jacob Andreas, Eunsol Choi, and Angeliki Lazaridou (eds.), *Proceedings of the NAACL Student Research Workshop*, pp. 8–15, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-2002. URL <https://aclanthology.org/N16-2002/>.
- Serge Lang. *Algebra*. Springer, 3 edition, 2002.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Vasudevan Nedumpozhimana and John D. Kelleher. Topic aware probing: From sentence length prediction to idiom identification how reliant are neural language models on topic? *Natural Language Processing*, pp. 1–29, 2024. doi: 10.1017/nlp.2024.43.
- Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models, 2024. URL <https://arxiv.org/abs/2311.03658>.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pp. 1532–1543, 2014.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Griffin Thomas Adams, Jeremy Howard, and Iacopo Poli. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2526–2547, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.127. URL <https://aclanthology.org/2025.acl-long.127/>.