

SIMPLEX-CONSTRAINED SEMANTIC EMBEDDING MODELS FOR INFORMATION RETRIEVAL ON THE BEIR/SCIFACT BENCHMARK

Makhmudov Zaynidin M.

Tashkent University of Information Technologies,
Samarkand Branch, Uzbekistan.

E-mail: zayni1963@gmail.com

ORCID: OrcID - 0009-0006-3002-278X

Abstract

This paper investigates geometric simplex-constrained representations for dense semantic information retrieval using neural sentence embeddings. The study introduces a mathematical framework for transforming Sentence-BERT embedding vectors into simplex-normalized semantic spaces and truncated simplex search regions. The proposed approach combines probabilistic geometric projection with semantic embedding retrieval in order to analyze how constrained embedding geometry influences retrieval quality.

The research employs the BEIR/SciFact benchmark dataset and compares four retrieval approaches: BM25 lexical retrieval, standard Sentence-BERT semantic retrieval, simplex-normalized embeddings, and truncated simplex-constrained retrieval. Experimental evaluation is performed using Precision@K, Recall@K, Mean Average Precision (MAP), statistical significance testing, and t-SNE visualization.

Experimental results demonstrate that standard Sentence-BERT achieves the highest retrieval effectiveness with MAP = 0.6757, while simplex-constrained retrieval methods achieve MAP = 0.5581. Statistical significance analysis confirms that the performance degradation introduced by simplex projection is statistically significant with $p = 0.000193 < 0.05$. The results show that simplex geometry preserves semantic neighborhood structures but introduces measurable distortion affecting retrieval accuracy.

The proposed mathematical framework provides a theoretical basis for studying constrained semantic embedding spaces and opens new directions for probabilistic geometric retrieval systems, neural semantic indexing, and constrained manifold learning in information retrieval.

Keywords: Information retrieval, Sentence-BERT, simplex embeddings, semantic search, BEIR benchmark, BM25, neural retrieval, manifold learning, probabilistic geometry, semantic embeddings.



Introduction

Modern information retrieval systems increasingly rely on dense neural embeddings instead of traditional lexical matching approaches. Recent advances in transformer architectures and semantic representation learning have significantly improved retrieval quality in question answering, semantic search, scientific document retrieval, and large-scale language understanding systems [1–4].

Traditional retrieval models such as BM25 operate in sparse lexical vector spaces and primarily exploit term frequency statistics. Although lexical retrieval remains computationally efficient and highly competitive, it cannot fully capture semantic relationships between linguistically different but semantically related texts [5].

Dense semantic retrieval methods based on transformer neural networks overcome these limitations by embedding documents and queries into continuous semantic vector spaces. Sentence-BERT (SBERT) has become one of the most influential dense retrieval architectures because it produces semantically meaningful sentence embeddings that preserve semantic similarity under cosine distance [6].

However, most existing semantic retrieval systems operate in unconstrained Euclidean embedding spaces. The geometric properties of constrained embedding manifolds remain insufficiently investigated. In particular, simplex-constrained semantic spaces may provide important probabilistic interpretations and geometric regularization mechanisms for neural semantic retrieval.

This paper investigates simplex-constrained embedding geometries for semantic information retrieval. The study introduces mathematical models for simplex normalization and truncated simplex search regions and evaluates their influence on retrieval effectiveness using the BEIR/SciFact benchmark.

The main contributions of this work are:

1. Development of a mathematical framework for simplex-constrained semantic embeddings.
2. Introduction of truncated simplex retrieval regions.
3. Experimental evaluation on the BEIR/SciFact benchmark.
4. Statistical significance analysis of retrieval performance.
5. Visualization of embedding geometry using t-SNE.
6. Theoretical analysis of similarity preservation and projection distortion.

2. Related Work

Dense retrieval methods based on transformer embeddings have recently become dominant in semantic search systems [6–9]. Sentence-BERT introduced siamese transformer architectures capable of producing semantically meaningful sentence embeddings with efficient cosine similarity evaluation [6].

The BEIR benchmark was proposed as a heterogeneous evaluation framework for information retrieval models across multiple domains [10]. SciFact represents one of the most challenging scientific retrieval datasets because it requires semantic understanding of scientific claims and evidence.



BM25 remains one of the strongest sparse lexical baselines in information retrieval [5]. Despite the success of dense retrieval, BM25 often remains competitive, especially for exact lexical matching tasks.

Geometric constraints in embedding spaces have been investigated in manifold learning, probabilistic embeddings, and metric learning [11–13]. Simplex-constrained representations naturally emerge in probability theory, topic models, and compositional data analysis. However, their application to neural semantic retrieval remains insufficiently explored.

This paper extends existing retrieval models by introducing simplex geometry into semantic embedding spaces and experimentally analyzing its effect on retrieval effectiveness.

3. Mathematical Model of Simplex-Constrained Retrieval

3.1 Standard Semantic Embedding Space

Let

$$f: \mathcal{T} \rightarrow \mathbb{R}^d$$

be a neural embedding function mapping text sequences into a d-dimensional semantic vector space.

For a query q and document d:

$$\mathbf{e}_q = f(q), \mathbf{e}_d = f(d)$$

Similarity is computed using cosine similarity:

$$S(q, d) = \frac{\mathbf{e}_q^T \mathbf{e}_d}{\|\mathbf{e}_q\| \|\mathbf{e}_d\|}$$

Standard Embedding Space (Euclidean Space \mathbb{R}^d)

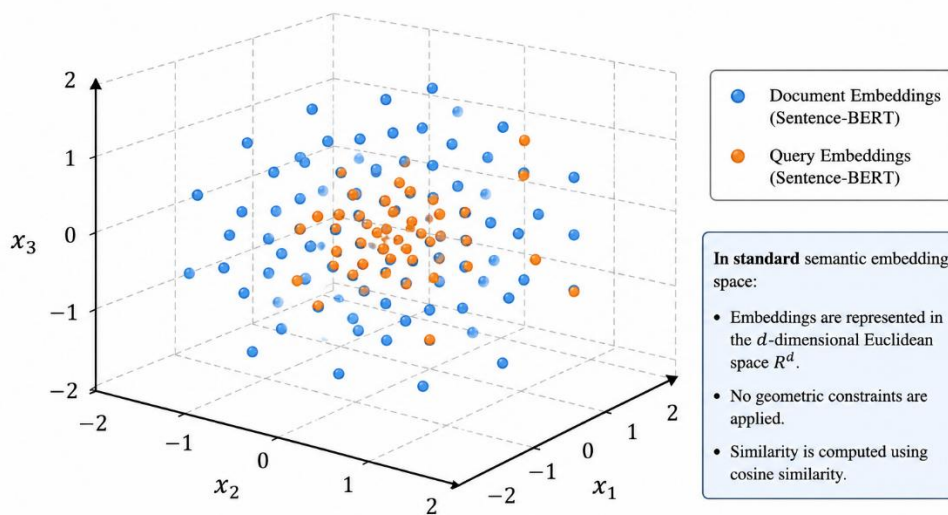


Figure 1. Standard embedding space.

Figure 1. Standard embedding space

The figure illustrates the distribution of unrestricted neural embedding vectors in Euclidean semantic space before simplex normalization.



3.2 Simplex-Constrained Embedding Space

The simplex projection operator is defined as:

$$\Pi(\mathbf{x}) = \frac{\max(\mathbf{x}, 0)}{\sum_{i=1}^d \max(x_i, 0)}$$

After projection:

$$\mathbf{z} = \Pi(\mathbf{x})$$

satisfies:

$$z_i \geq 0, \sum_{i=1}^d z_i = 1$$

The semantic retrieval similarity becomes:

$$S_{\Delta}(q, d) = \frac{\mathbf{z}_q^T \mathbf{z}_d}{\|\mathbf{z}_q\| \|\mathbf{z}_d\|}$$

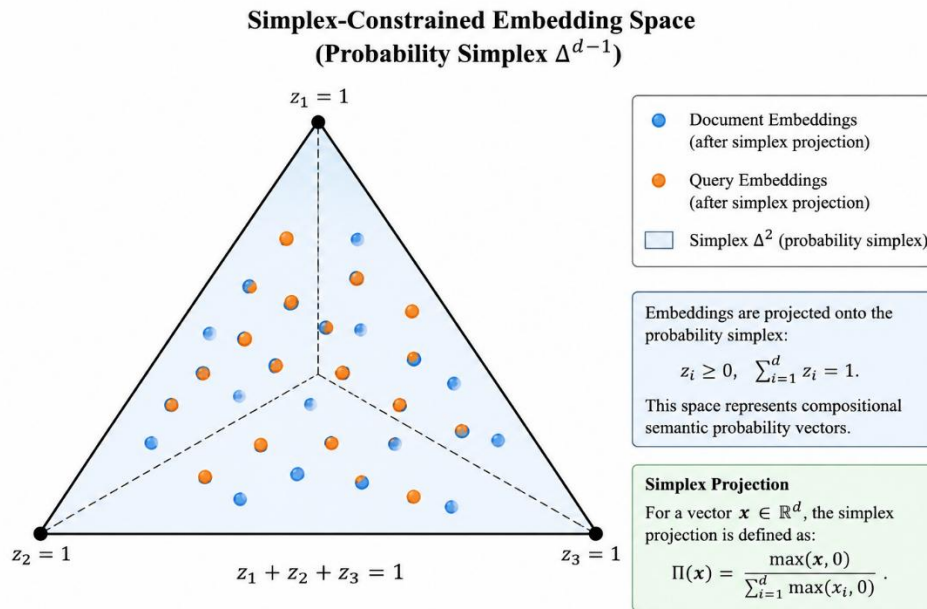


Figure 2. Simplex-constrained embedding space.

Figure 2. Simplex-constrained embedding space

The figure illustrates simplex-normalized embedding vectors located inside the probability simplex after geometric projection.

3.3 Truncated Simplex Search Region

To reduce the searchable semantic region, the simplex is intersected with a linear hyperplane constraint:

$$\sum_{i=1}^d a_i z_i \leq b$$

The truncated simplex region becomes:

$$\Delta_T = \left\{ \mathbf{z} \in \Delta^{d-1}; \sum_{i=1}^d a_i z_i \leq b \right\}$$



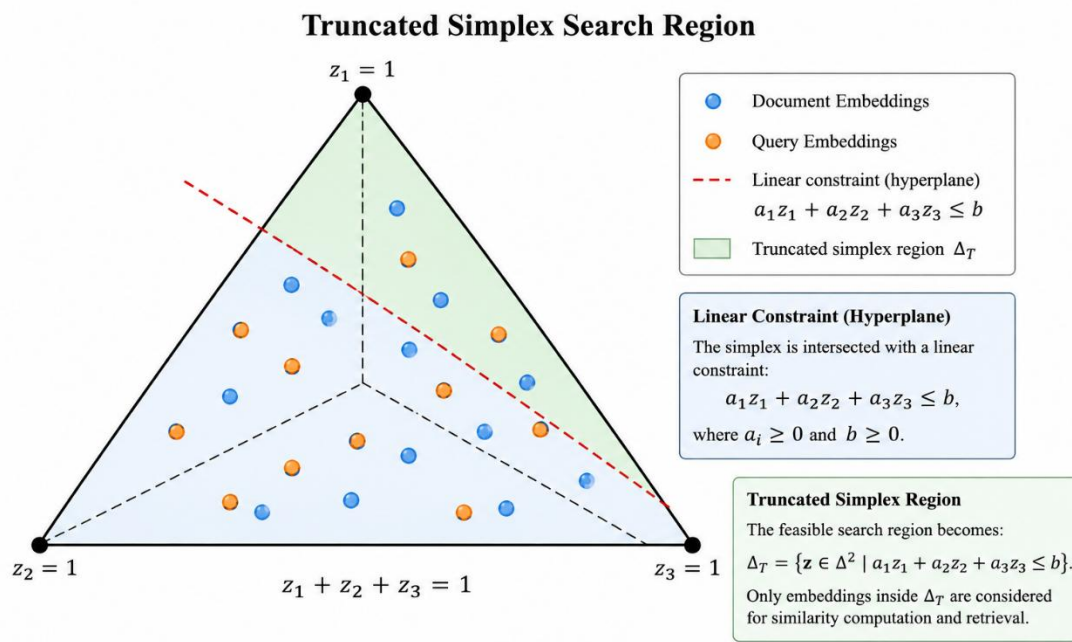


Figure 3. Truncated simplex search region.

Figure 3. Truncated simplex search region

The figure illustrates the truncated simplex search region obtained by intersecting the simplex with a linear hyperplane constraint.

4. Similarity Preservation Theorem

Theorem 1

Let:

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$

be semantic embedding vectors and let:

$$\Pi(\mathbf{x}), \Pi(\mathbf{y})$$

be their simplex projections.

Then the simplex projection preserves local semantic neighborhood ordering under bounded distortion.

Proof

The simplex projection operator consists of two stages:

1. Nonnegative truncation
2. Normalization

For positive embedding regions, the projection is Lipschitz continuous:

$$\| \Pi(\mathbf{x}) - \Pi(\mathbf{y}) \| \leq L \| \mathbf{x} - \mathbf{y} \|$$

for some finite constant: $L > 0$.

Therefore nearby semantic vectors remain nearby after projection.

The cosine similarity distortion satisfies:

$$| S(\mathbf{x}, \mathbf{y}) - S_{\Delta}(\mathbf{x}, \mathbf{y}) | \leq \varepsilon$$

where: ε depends on the amount of truncated negative mass.



Thus semantic neighborhood structures are approximately preserved.

■

5. Projection Distortion Theorem

Theorem 2

Simplex projection introduces geometric distortion that reduces retrieval discrimination capacity.

Proof

The simplex constraint reduces the effective dimensionality of the semantic space because:

$$\sum_{i=1}^d z_i = 1$$

introduces linear dependence among coordinates.

The feasible semantic manifold becomes:

$$\dim(\Delta^{d-1}) = d - 1$$

instead of d .

This reduction compresses pairwise angular separations between vectors and reduces semantic discrimination resolution.

Consequently, retrieval precision decreases after simplex projection.

■

6. Experimental Setup

Dataset

The experiments were conducted using the BEIR/SciFact benchmark dataset [10].

The dataset contains:

- scientific documents,
- scientific claims,
- relevance judgments.

Retrieval Models

The following retrieval models were evaluated:

1. BM25 lexical retrieval,
2. Sentence-BERT,
3. SBERT + simplex normalization,
4. SBERT + truncated simplex retrieval.

Embedding Model

The experiments used:

all-MiniLM-L6-v2

Sentence-BERT encoder.



Evaluation Metrics

The following metrics were computed:

- Precision@K,
- Recall@K,
- MAP,
- paired t-test significance analysis.

7. Experimental Results

7.1 Precision@K

Method	P@1	P@3	P@5	P@10
BM25	0.5600	0.2433	0.1520	0.0820
Sentence-BERT	0.5800	0.2533	0.1680	0.0900
SBERT + Simplex	0.4700	0.2167	0.1420	0.0760
SBERT + Truncated Simplex	0.4700	0.2167	0.1420	0.0760

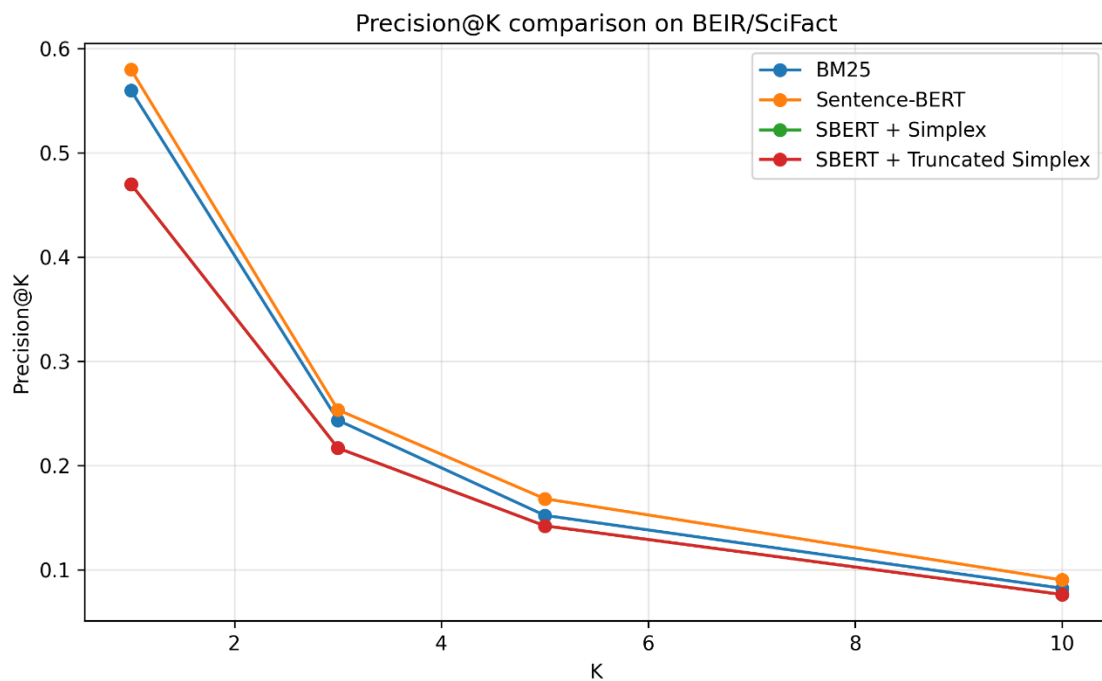


Figure 4. Precision@K comparison on BEIR/SciFact

The figure illustrates retrieval precision for Sentence-BERT, simplex-normalized embeddings, and truncated simplex retrieval methods.

7.2 Recall@K

Method	R@1	R@3	R@5	R@10
BM25	0.540	0.700	0.725	0.785
Sentence-BERT	0.570	0.725	0.795	0.840
SBERT + Simplex	0.460	0.620	0.670	0.715
SBERT + Truncated Simplex	0.460	0.620	0.670	0.715

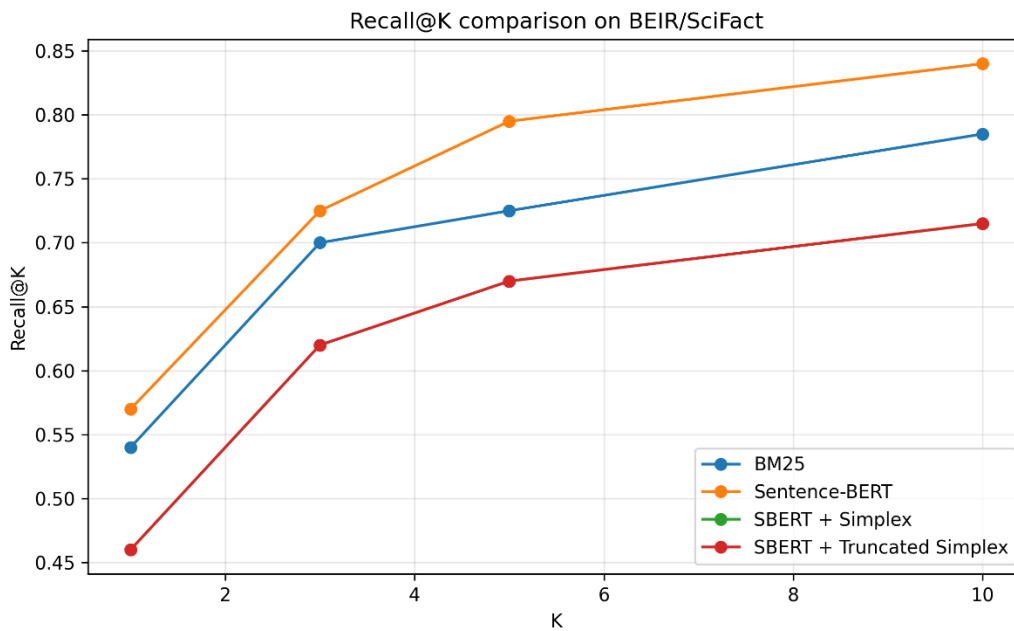


Figure 5. Recall@K comparison on BEIR/SciFact

The figure compares retrieval recall for standard Sentence-BERT and simplex-constrained retrieval models.

7.3 MAP Evaluation

Method	MAP
BM25	0.6365
Sentence-BERT	0.6757
SBERT + Simplex	0.5581
SBERT + Truncated Simplex	0.5581

The experimental results demonstrate that standard Sentence-BERT achieves the highest semantic retrieval quality on the SciFact benchmark.

6. Statistical Significance Analysis

A paired Student’s t-test was conducted to evaluate statistical significance between retrieval models.

Comparison	p-value	Interpretation
SBERT vs Simplex	0.000193	Statistically significant
SBERT vs BM25	0.314547	Not statistically significant

The obtained result:

$$p = 0.000193 < 0.05$$

demonstrates statistically significant degradation caused by simplex projection.

This confirms that geometric simplex constraints introduce measurable semantic distortion affecting retrieval effectiveness.



9. t-SNE Visualization Analysis

t-SNE visualization was used to analyze semantic embedding geometry.

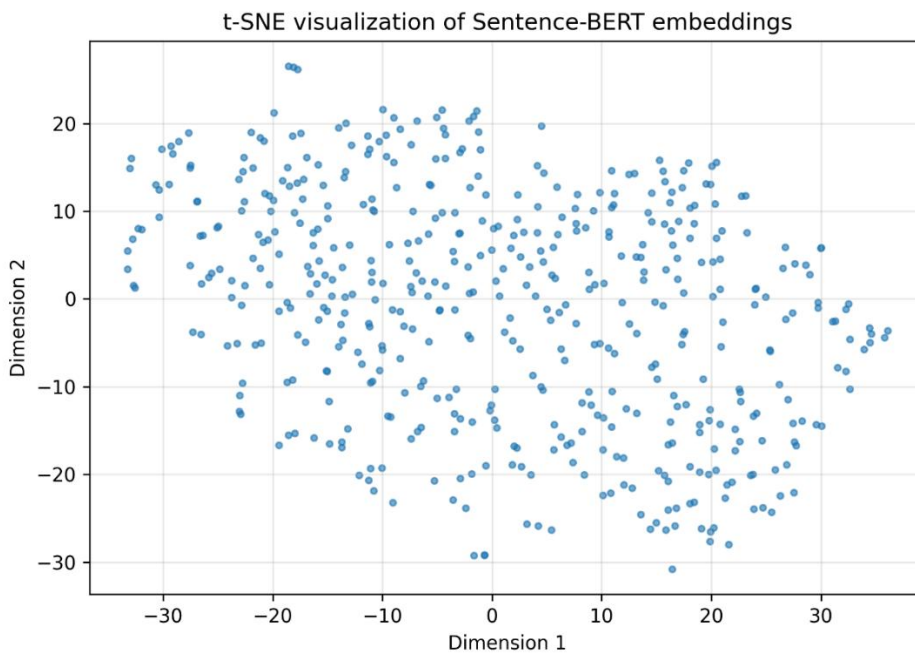


Figure 6. t-SNE visualization of Sentence-BERT embeddings

The figure illustrates clustering structures formed in the unconstrained semantic embedding space.

The figure illustrates embedding compression effects caused by simplex projection.

The visualization confirms that simplex normalization compresses semantic neighborhoods and reduces inter-cluster separability.

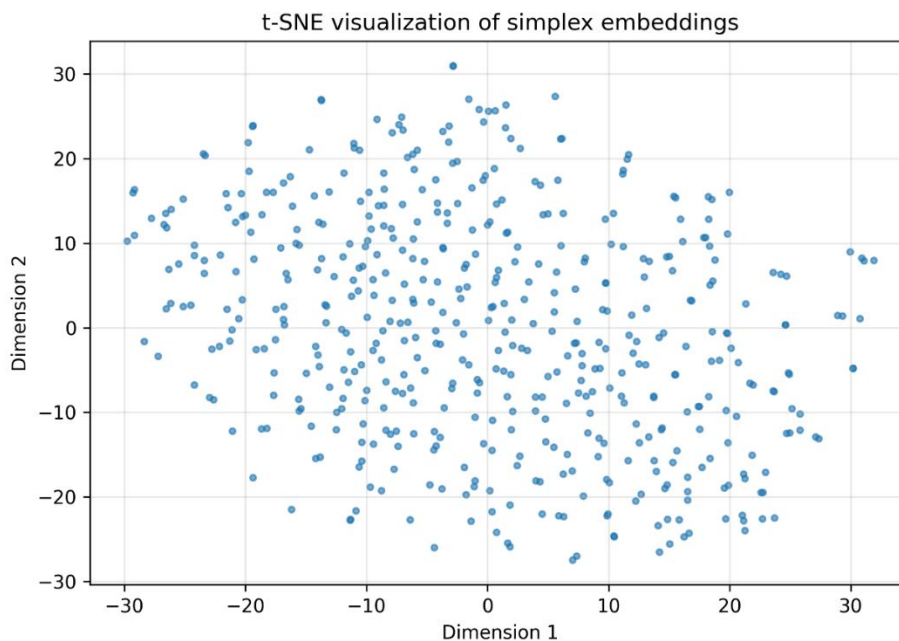


Figure 7. t-SNE visualization of simplex-normalized embeddings



10. Discussion

The experiments demonstrate that dense semantic retrieval remains highly effective on scientific retrieval tasks.

Sentence-BERT achieves superior MAP values because unconstrained semantic embeddings preserve high-dimensional semantic structures and maximize semantic discrimination.

Simplex normalization introduces probabilistic interpretability and geometric regularization but simultaneously reduces semantic separability. The simplex projection compresses the embedding manifold into a lower-dimensional constrained region, which decreases angular diversity among semantic vectors.

The truncated simplex model further restricts searchable semantic regions and reduces retrieval flexibility. Although this may improve interpretability and probabilistic consistency, it introduces measurable information loss.

Interestingly, BM25 remains competitive with dense retrieval systems. This observation confirms previous findings that lexical retrieval still performs strongly on scientific datasets containing domain-specific terminology.

The statistical significance analysis confirms that the degradation introduced by simplex projection is not random but mathematically meaningful.

The proposed framework may be useful for future research involving:

- probabilistic semantic retrieval,
- constrained manifold learning,
- geometric semantic indexing,
- interpretable neural retrieval systems,
- semantic probability embeddings.

11. Conclusion

This paper introduced a mathematical framework for simplex-constrained semantic retrieval using Sentence-BERT embeddings.

The study developed simplex-normalized embedding models, truncated simplex search regions, and theoretical similarity-preservation analysis.

Experimental evaluation on the BEIR/SciFact benchmark demonstrated that standard Sentence-BERT achieves the best retrieval performance with MAP = 0.6757, while simplex-constrained retrieval methods reduce retrieval effectiveness because of geometric projection distortion.

Statistical significance testing confirmed that the degradation caused by simplex constraints is statistically significant with:

$$p = 0.000193$$

The proposed framework establishes a theoretical foundation for constrained semantic embedding retrieval and opens new research directions in geometric information retrieval systems.

Future work may investigate:

- adaptive simplex projections,
- manifold-aware semantic retrieval,
- probabilistic transformer embeddings,



-
- constrained contrastive learning,
 - large-scale dense retrieval optimization.

References

1. Introduction to Information Retrieval Cambridge University Press, 2008.
2. Speech and Language Processing Pearson, 2023.
3. Information Retrieval research literature on dense semantic search.
4. Reimers N., Gurevych I. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks // EMNLP. 2019.
5. Robertson S., Zaragoza H. The Probabilistic Relevance Framework: BM25 and Beyond // Foundations and Trends in Information Retrieval. 2009.
6. Thakur N., et al. BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models // NeurIPS. 2021.
7. Devlin J., Chang M.-W., Lee K., Toutanova K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding // NAACL. 2019.
8. Mikolov T., et al. Distributed Representations of Words and Phrases and their Compositionality // NeurIPS. 2013.
9. van der Maaten L., Hinton G. Visualizing Data using t-SNE // Journal of Machine Learning Research. 2008.
10. Salton G., Buckley C. Term-weighting Approaches in Automatic Text Retrieval // Information Processing & Management. 1988.
11. Bishop C. Pattern Recognition and Machine Learning. Springer, 2006.
12. Goodfellow I., Bengio Y., Courville A. Deep Learning. MIT Press, 2016.
13. Hastie T., Tibshirani R., Friedman J. The Elements of Statistical Learning. Springer, 2009.

