# Bridging the Gap: A Comparative Analysis of Sparse Autoencoders and Latent Dirichlet Allocation for Document Modeling

Ran D. Zilca

*Abstract*—This paper presents a comprehensive comparison between Sparse Autoencoders (SAEs) and Latent Dirichlet Allocation (LDA) in the context of document modeling. We derive the mathematical foundations of both methods, highlighting their similarities and differences. Our analysis reveals that despite their distinct origins, SAEs and LDA share fundamental similarities in their objective of finding sparse, low-dimensional representations of documents. Experimental results on the 20 Newsgroups dataset demonstrate the comparable performance of SAEs and LDA in topic extraction and document representation tasks. We show that SAEs can be configured to mimic key properties of LDA, producing interpretable "topics" and sparse document representations. This work bridges the gap between neural network-based and probabilistic approaches to document modeling, offering insights for future hybrid models.

## I. Introduction

Document modeling is a fundamental task in natural language processing, with applications ranging from information retrieval to text classification. Two prominent approaches to this task are Sparse Autoencoders (SAEs) and Latent Dirichlet Allocation (LDA). While these methods originate from different paradigms - SAEs from neural networks and LDA from probabilistic modeling - they share the common goal of discovering latent structures in document collections.

SAEs are neural networks designed to learn compact, sparse representations of input data. When applied to document modeling, they can extract features that correspond to underlying themes or topics in the text. LDA, on the other hand, is a generative probabilistic model that explicitly represents documents as mixtures of latent topics.

Despite their different foundations, we hypothesize that SAEs and LDA, when properly configured, can perform remarkably similar functions in document modeling. This paper aims to:

1) Provide a detailed mathematical comparison of SAEs and LDA.
2) Demonstrate how SAEs can be designed to mimic key properties of LDA.
3) Empirically compare the performance of SAEs and LDA on common document modeling tasks.

Our analysis offers insights into the connections between neural and probabilistic approaches to document modeling, potentially paving the way for new hybrid models that combine the strengths of both paradigms.

## II. Background

### A. Sparse Autoencoders

Sparse Autoencoders (SAEs) are a type of artificial neural network designed to learn efficient data representations [3]. An SAE consists of an encoder that maps the input to a hidden representation, and a decoder that attempts to reconstruct the input from this representation. The "sparse" aspect comes from constraints or regularization that encourage the hidden representations to be sparse, i.e., to have many zero or near-zero values.

Formally, for an input vector $x \in \mathbb{R}^V$, where $V$ is the vocabulary size, the encoder produces a hidden representation $z \in \mathbb{R}^K$:

$$z = \sigma(W_1 x + b_1) \tag{1}$$

where $W_1$ is the weight matrix, $b_1$ is the bias vector, and $\sigma$ is an activation function, typically ReLU or sigmoid.

The decoder then attempts to reconstruct the input:

$$\hat{x} = \sigma(W_2 z + b_2) \tag{2}$$

The SAE is trained to minimize the reconstruction error while maintaining sparsity in the hidden layer.

### B. Latent Dirichlet Allocation

Latent Dirichlet Allocation (LDA) is a generative probabilistic model for collections of discrete data such as text corpora [2]. LDA represents documents as mixtures of latent topics, where each topic is characterized by a distribution over words.

The generative process for LDA is as follows:

1) For each topic $k$, draw a word distribution $\phi_k \sim$ Dirichlet($\beta$)
2) For each document $d$:
   a) Draw a topic distribution $\theta_d \sim$ Dirichlet($\alpha$)
   b) For each word $w_n$ in the document:
      i) Draw a topic $z_n \sim$ Multinomial($\theta_d$)
      ii) Draw a word $w_n \sim$ Multinomial($\phi_{z_n}$)

The goal is to infer the latent topic structure (the $\phi_k$ and $\theta_d$ distributions) from the observed words in the documents.

## III. Mathematical Derivation

To demonstrate the similarity between SAEs and LDA, we will derive their objective functions and highlight the parallels.

## A. Sparse Autoencoder Objective

The objective function for an SAE typically consists of two terms: the reconstruction error and a sparsity constraint:

$$L_{SAE} = \underbrace{\sum_{i=1}^{V}(x_i - \hat{x}_i)^2}_{\text{Reconstruction Error}} + \underbrace{\lambda \sum_{j=1}^{K}|z_j|}_{\text{Sparsity Constraint}} \tag{3}$$

where $\lambda$ is a hyperparameter controlling the strength of the sparsity constraint.

## B. LDA Objective

The objective of LDA is to maximize the likelihood of the observed words given the model parameters. The log-likelihood for a single document can be written as:

$$\log P(w|\alpha, \beta) = \log \int_{\theta} \left( \prod_{n=1}^{N} \sum_{k=1}^{K} P(w_n|\phi_k)P(z_n = k|\theta) \right) P(\theta|\alpha)d\theta \tag{4}$$

## C. Reformulating SAE in LDA-like Terms

To elucidate the similarity between SAEs and LDA, we can reformulate the SAE objective in probabilistic terms. This reformulation will highlight the parallels between the two models in their approach to document modeling.

*1) Probabilistic Interpretation of SAE:* Let us consider the following probabilistic interpretation of the SAE:

1) Assume the reconstruction follows a Gaussian distribution:

$$P(x|\hat{x}) = \mathcal{N}(x|\hat{x}, \sigma^2 I) \tag{5}$$

where $x$ is the input document vector, $\hat{x}$ is the reconstructed document vector, and $\sigma^2$ is the variance of the reconstruction error.

2) Assume a Laplace prior on the hidden units for sparsity:

$$P(z) = \prod_{j=1}^{K} \text{Laplace}(z_j|0, b) \tag{6}$$

where $z$ is the hidden representation, $K$ is the number of hidden units, and $b$ is the scale parameter of the Laplace distribution.

*2) Deriving the SAE Objective:* Given these assumptions, we can derive the SAE objective as follows:

1) The negative log-likelihood of the SAE model is:

$$-\log P(x, z) = -\log P(x|\hat{x}) - \log P(z) \tag{7}$$

2) Expanding the Gaussian log-likelihood term:

$$-\log P(x|\hat{x}) = \frac{1}{2\sigma^2} \sum_{i=1}^{V}(x_i - \hat{x}_i)^2 + \frac{V}{2} \log(2\pi\sigma^2) \tag{8}$$

where $V$ is the vocabulary size.

3) Expanding the Laplace prior term:

$$-\log P(z) = \sum_{j=1}^{K} \frac{|z_j|}{b} + K \log(2b) \tag{9}$$

4) Combining these terms and dropping constants, we get:

$$-\log P(x, z) \propto \frac{1}{2\sigma^2} \sum_{i=1}^{V}(x_i - \hat{x}_i)^2 + \frac{1}{b} \sum_{j=1}^{K}|z_j| \tag{10}$$

This formulation closely resembles the original SAE objective function:

$$L_{SAE} = \underbrace{\sum_{i=1}^{V}(x_i - \hat{x}_i)^2}_{\text{Reconstruction Error}} + \underbrace{\lambda \sum_{j=1}^{K}|z_j|}_{\text{Sparsity Constraint}} \tag{11}$$

where $\lambda = \frac{2\sigma^2}{b}$ controls the strength of the sparsity constraint.

*3) Comparison with LDA:* Now, let's compare this reformulated SAE objective with the LDA objective. Recall that the log-likelihood for a single document in LDA is:

$$\log P(w|\alpha, \beta) = \log \int_{\theta} \left( \prod_{n=1}^{N} \sum_{k=1}^{K} P(w_n|\phi_k)P(z_n = k|\theta) \right) P(\theta|\alpha)d\theta \tag{12}$$

We can draw the following parallels:

1) Document Reconstruction: In SAE, this is represented by the reconstruction error term $\sum_{i=1}^{V}(x_i - \hat{x}_i)^2$. In LDA, it's captured by $\prod_{n=1}^{N} \sum_{k=1}^{K} P(w_n|\phi_k)P(z_n = k|\theta)$, which models the probability of generating each word in the document.
2) Latent Representation: In SAE, this is the hidden layer $z$. In LDA, it's the topic distribution $\theta$ for each document.
3) Sparsity: In SAE, sparsity is enforced through the L1 regularization term $\sum_{j=1}^{K}|z_j|$. In LDA, sparsity in the topic distributions is induced by the Dirichlet prior $P(\theta|\alpha)$.
4) Probabilistic Interpretation: While LDA is inherently probabilistic, we've shown that SAE can also be interpreted in probabilistic terms, bridging the gap between these two approaches.

*4) Implications:* This reformulation reveals that SAEs and LDA, despite their different origins, share fundamental similarities in their approach to document modeling:

1) Both methods aim to find a lower-dimensional representation of documents (topics in LDA, hidden representations in SAE).
2) Both incorporate mechanisms to encourage sparsity in these representations.
3) Both involve a trade-off between faithfully reconstructing the original documents and finding meaningful, generalizable patterns.

The key difference lies in how these objectives are optimized: LDA uses probabilistic inference techniques, while SAE uses gradient-based optimization. However, the underlying goals and the structure of the objectives are remarkably similar, suggesting that insights and techniques from one approach could potentially be applied to the other.

## IV. EXPERIMENTAL SETUP

To empirically compare SAEs and LDA, we conducted experiments on the 20 Newsgroups dataset, a popular benchmark for document classification and topic modeling tasks.

### A. Dataset

The 20 Newsgroups dataset consists of approximately 20,000 newsgroup documents, partitioned evenly across 20 different newsgroups. We preprocessed the data by removing headers, footers, and quotes, and represented documents using term frequency vectors.

### B. Implementation Details

We implemented the SAE using TensorFlow and the LDA model using scikit-learn. The key configurations are as follows:

---
**Algorithm 1** Sparse Autoencoder Configuration
---
1: Input dimension: 5000 (vocabulary size)
2: Hidden layer dimension: 100
3: Activation function: ReLU
4: Sparsity constraint: L1 regularization
5: Optimizer: Adam
6: Loss function: Binary cross-entropy + L1 regularization
7: Training epochs: 50
8: Batch size: 256

---

---
**Algorithm 2** LDA Configuration
---
1: Number of topics: 100
2: Dirichlet parameter for document-topic distribution ($\alpha$): 1/100
3: Dirichlet parameter for topic-word distribution ($\beta$): 1/100
4: Maximum iterations: 50
5: Learning method: Online variational Bayes

---

## V. RESULTS AND DISCUSSION

### A. Topic Coherence

We evaluated the quality of topics extracted by both SAE and LDA using the topic coherence metric. Specifically, we used the CV coherence measure [1] as implemented in gensim's CoherenceModel. The SAE achieved a mean topic coherence score of 0.3867, while LDA achieved a score of 0.4146. This suggests that both LDA and SAE are capable of producing reasonably coherent topics (with LDA having a slight edge), demonstrating that SAEs can be configured to mimic key properties of LDA in document modeling tasks. The close performance also indicates that there might be potential in developing hybrid approaches that combine the strengths of both probabilistic and neural network-based models for topic discovery.

| Model | Classification Accuracy |
|-------|------------------------|
| SAE | 4.48% |
| LDA | 48.86% |

TABLE I
DOCUMENT CLASSIFICATION ACCURACY USING SAE AND LDA REPRESENTATIONS

| Model | Sparsity (% of zero values) |
|-------|------------------------------|
| SAE | 0.9996 |
| LDA | 0.8679 |

TABLE II
SPARSITY OF LEARNED REPRESENTATIONS

*1) Handling Rare Words in Coherence Calculation:* During the computation of topic coherence, we encountered warnings related to division by zero. This issue typically arises when certain words in the extracted topics do not appear in the reference corpus or appear very infrequently. Such rare words can lead to unreliable probability estimates in the coherence calculation.

To address this issue, we implemented the following strategies:

1) Preprocessing: We ensured that very rare words were removed from the corpus during the preprocessing stage. Specifically, we used a minimum document frequency threshold when creating our vocabulary.
2) Robust Coherence Calculation: In our coherence computation function, we added safeguards to handle cases where word probabilities are zero or near-zero. This involved adding a small smoothing factor to prevent division by zero and taking the logarithm of probabilities only when they are positive.
3) Topic Filtering: For topics containing words that do not appear in the corpus, we excluded these topics from the coherence calculation or replaced the problematic words with the next most probable words in the topic.

These measures helped to ensure more robust and meaningful coherence scores, although it's important to note that the presence of such rare words in the extracted topics may indicate a need for further refinement of the topic models or preprocessing steps.

### B. Document Classification

We used the learned representations from both models as features for a document classification task, using a Naive Bayes classifier. The results are presented in Table I.

The results indicate that XXX [CHECK THE CODE - IS THIS CORRECT??).

### C. Sparsity Analysis

We analyzed the sparsity of the learned representations for both models. For SAE, we considered a value to be zero if its absolute value was below 1e-5. For LDA, we considered a topic proportion to be zero if it was below 0.01.

The sparsity analysis reveals a significant difference in how SAE and LDA represent documents. The extreme sparsity of SAE explains its poor classification performance and suggests

that modifications to its architecture or training process might be necessary to make it more competitive with LDA for document modeling tasks. The results also highlight the importance of carefully balancing sparsity and informativeness in learned representations for both topic modeling and downstream tasks like classification.

## VI. Conclusion

In this paper, we have presented a comprehensive comparison between Sparse Autoencoders (SAEs) and Latent Dirichlet Allocation (LDA) for document modeling. The mathematical derivation revealed fundamental similarities in their objectives, despite their different origins. However, the experimental results on the 20 Newsgroups dataset demonstrated differences in sparsity, as reflected in degraded performance of SAE compared with LDA on document classification problem. The findings suggest that insights and techniques from one approach could potentially be applied to the other, opening up new avenues for hybrid models that combine the strengths of both neural and probabilistic approaches to document modeling.

Future work could explore:

1) Developing hybrid models that integrate SAE and LDA approaches, possibly using more traditional probabilistic approach for expandability of more robust neural models.
2) Investigating the impact of different sparsity constraints on SAE performance in document modeling tasks
3) Extending this comparison to other document modeling techniques and datasets

## References

[1] M. Röder, A. Both, and A. Hinneburg, "Exploring the space of topic coherence measures," in *Proceedings of the eighth ACM international conference on Web search and data mining*, 2015, pp. 399–408.
[2] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, vol. 3, no. Jan, pp. 993–1022, 2003.
[3] A. Ng, "Sparse autoencoder," *CS294A Lecture notes*, vol. 72, no. 2011, pp. 1–19, 2011.