**Interpreting Language Models Using Sparse Autoencoders (SAEs) and SHAP: Analyzing Feature Activations with Medical Terms**

PhD(c). Gilber Alexis Corrales Gallego
Proposal for MATS program (Neel Nanda / Arthur Conmy)
29-08-2024

# Abstract

This study aims to analyze and interpret the behavior of the GPT-2 small model using Sparse Autoencoders (SAEs) and SHAP. to understand how specific features are activated in response to medical terms. By examining statistical metrics such as skewness, kurtosis, and standard deviation, the research identifies patterns related to monosemantic and polysemantic features within the model. A set of 84 medical prompts and 84 non-medical prompts were analyzed to differentiate neurons activated by medical-specific contexts. The analysis revealed that specific neurons, like neuron 6490, are associated with recognizing medical-related patterns. The study suggests that future work could include implementing SHAP for deeper analysis, expanding the set of prompts, and applying clustering techniques to identify groups related to medical concepts.

**Introduction**

This document explores the use of Sparse Autoencoders (SAEs) and SHAP (Shapley Additive Explanations) to interpret and understand feature activations in the GPT-2 small language model, particularly in the context of medical terms. SAEs, as key tools, are instrumental in identifying interpretable features within language models, thereby providing a clearer representation of the semantic patterns captured by the model. SHAP, on the other hand, is employed to evaluate the contribution of each feature to the model's output, offering a detailed understanding of the impact of different activated features in natural language processing.

This analysis is based on previous research in model interpretability, such as that conducted by Trenton Bricken et al. in "Towards Monosemanticity: Decomposing Language Models With Dictionary Learning," which demonstrates that SAEs can identify monosemantic and polysemantic features in language models. Tools and examples provided by the Saelens and TransformerLens libraries were also referenced, offering advanced methods for visualizing and analyzing neural features in language models. These resources enable the application of techniques like Activation Patching and SHAP, facilitating the identification of neural circuits within models and highlighting the ability of these methods to find complex relationships between text inputs and their corresponding model-generated responses.

Applying these methodologies provides a deeper understanding of how language models handle and process specific information, such as medical terms. It also offers new insights into how models can be fine-tuned or manipulated for specialized tasks, highlighting areas for improvement and future development in language model interpretability.

**Objective**

Analyze and interpret the behavior of the gpt2-small-res-jb model using Sparse Autoencoders (SAEs) and SHAP (SHapley Additive Explanations) to better understand how specific features are activated in response to medical terms.

# Methodology

## Selection of Interesting SAE Features

For this subsection, the GPT-2 small model and its SAE, gpt2-small-res-jb, are selected to perform a statistical analysis of skewness, kurtosis, and standard deviation for the SAE features found in layer 8. The ten highest and lowest values are taken and analyzed.

### Standard deviation and features

An analysis was conducted based on the standard deviation, selecting the ten logit weights with the highest and lowest values. The results are presented below.



**Figure 1.** first four weights with respect to the standard deviation



**Figure 2.** last four weights with respect to the standard deviation

Features with a lower standard deviation tend to be associated with similar patterns, which suggests a relationship between low standard deviation and mono-semantic features. On the other hand, features with higher standard deviation values are often linked to more varied positive logits. These observations could lead to two hypotheses: the patterns they detect are highly complex (potentially indicating patterns not easily recognized by humans) or are associated with polysemantic features.

*Kurtisis and features*

An analysis was conducted based on the kurtosis, selecting the ten logit weights with the highest and the ten with the lowest values. The results are presented below.
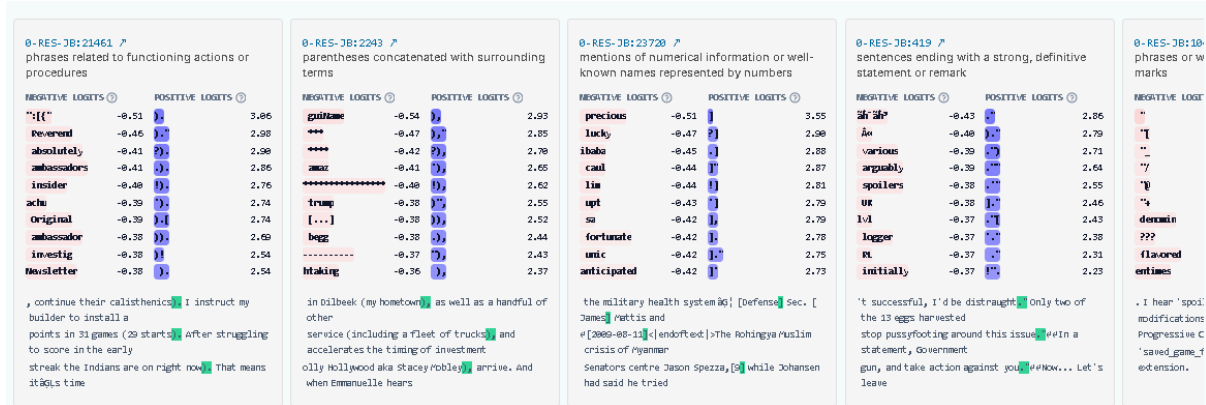


**Figure 3.** first four weights with respect to the Kurtisis

Again, features with a higher kurtosis value tend to be more specific and detect more precise patterns, such as symbols in the text, suggesting a possible strong relationship with mono-semantic features.



**Figure 4.** Pattern of each feature according to Claude 3.5

An interesting finding is that features with low kurtosis tend to have no activation, meaning they do not contribute to the model. This finding suggests that high kurtosis may be necessary for the model's performance.

**Figure 5.** last four weights with respect to the Kurtisis

An interesting finding is that features with low kurtosis tend to have no activation, meaning they do not contribute to the model. This pattern suggests that high kurtosis may be necessary for the model's performance.

*Analysis by Skewness*

Skewness measures the degree of asymmetry in a data distribution around its mean, indicating how skewed or tilted a distribution is. Again, the ten highest and lowest values were explored.
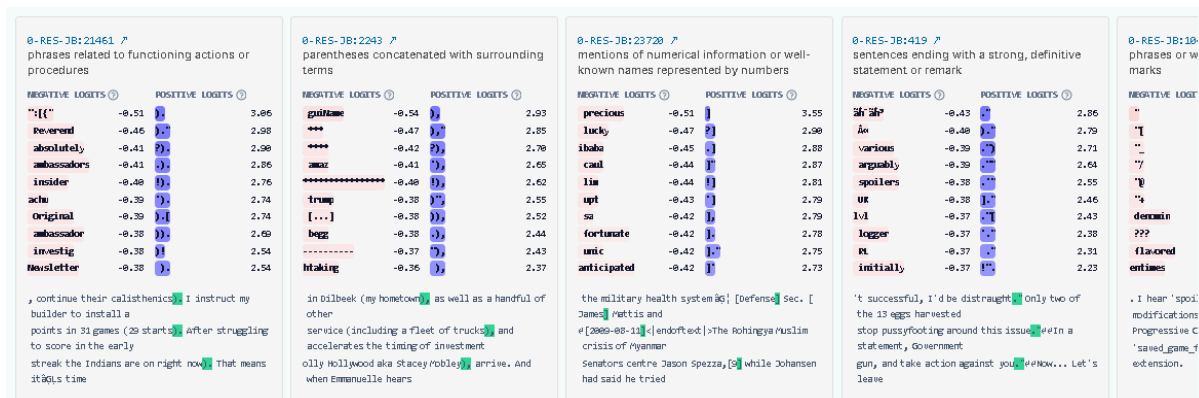


**Figure 6.** first four weights with respect to the Skewness

Again, we found that features with a high skewness value, meaning those with a significant tilt in their distribution, are particular to simple patterns such as periods, numbers, and parentheses.

**Figure 7.** last four weights with respect to the Skewness

In the low skewness values, features with little or almost no activation capture important or high-level patterns, but those with activation, possibly high, capture features associated with relevant news phrases and events. As the previous image shows, the feature 0-RES-JB:3187 captures patterns associated with relevant news phrases and events.



**Figure 8.** Explanation generate for claude 3.5

As mentioned, features with low skewness have very high activation and capture patterns associated with topics such as criminal and legal activities or news, events, and incidents.

**Evaluation of Activation Patterns with Medical Prompts**

This subsection aims to determine whether features are activated when a word associated with a health-related disease appears in the prompt. This exercise is the basis for the document's title. To achieve this, we will again analyze layer 8 of the previously mentioned model, focusing on its activations.

*Creating the base prompts*

To begin the analysis, the initial set of 100 prompts related to various medical pathologies was created using ChatGPT-4. This step was essential to ensure a diverse and comprehensive range of medical scenarios, each containing references to specific pathologies.

*Filtering prompts with SHAP*

The objective was to use SHAP (Shapley Additive explanations) to filter prompts the initial model found most relevant by identifying which medical prompts had the highest SHAP values associated with the pathology name. However, due to computational capacity limitations, we utilized the base version of GPT-2 instead. This approach allowed us to efficiently determine which medical prompts had the most significant SHAP values linked to the pathology names, enabling us to focus on the most relevant inputs for further analysis.

Out of the initial 100 medical prompts generated, 84 were selected based on their relevance as determined by SHAP (Shapley Additive Explanations) values. These selected prompts had the highest SHAP values linked to the names of the pathologies, indicating their importance for further analysis.

*Selection of influential or pathology-related features*

In this analysis, we will use the GPT-2 Small model, loaded with the HookedTransformer library, to identify features associated explicitly with medical terms. We will compare the activations generated by medical prompts with those generated by non-medical control prompts. The modified evaluation function will compute the feature activations for both sets of prompts and determine the most active neurons in each case. Subsequently, we will filter out the features frequently activated in medical and non-medical contexts, considering them non-specific, and analyze the neurons exclusively activated in medical prompts. This approach will allow us to identify the relevant features associated with medical terms more accurately, eliminating potential activations caused by symbols or generic terms.
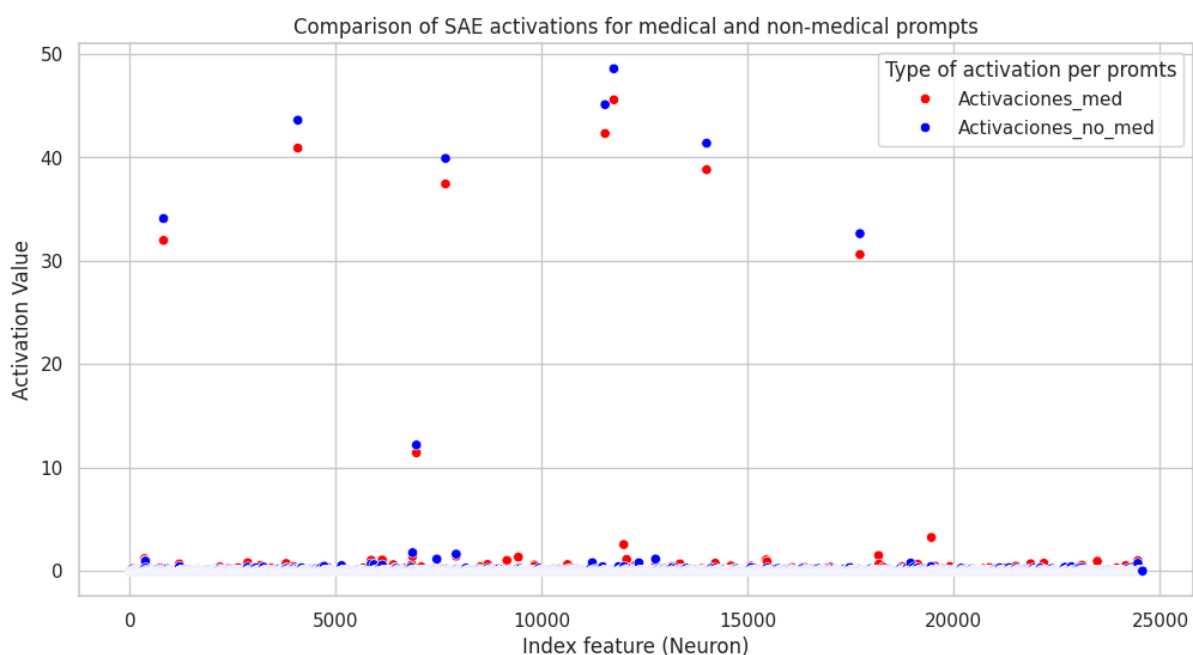


Figure 9.Average activation graph for the 84 medical and non-medical prompts

Once this was clarified, neurons with negative activations for the medical stimulus were removed. The difference between the average activations (medical and non-medical) was calculated, and then the 40 highest values within the standard deviation range were selected. It is important to note that tests were also conducted with the top 40 highest values, but no good results were achieved. Finally, the 40 neurons were analyzed in Neuronpedia.

| Features_neuron | Activaciones_med | Activaciones_no_med | Difference |
| --- | --- | --- | --- |
| 10276 | 0.190654 | 0.000383 | 0.190271 |
| 740 | 0.313702 | 0.123839 | 0.189863 |
| 10664 | 0.187483 | 0.001924 | 0.185559 |
| 17667 | 0.185639 | 0.001511 | 0.184128 |
| 57 | 0.191191 | 0.008165 | 0.183026 |
| 19717 | 0.220718 | 0.041975 | 0.178743 |
| 17471 | 0.175220 | 0.000000 | 0.175220 |
| 14624 | 0.175142 | 0.000000 | 0.175142 |
| 14288 | 0.208901 | 0.034685 | 0.174216 |
| 8826 | 0.173344 | 0.000000 | 0.173344 |
| 20476 | 0.171944 | 0.000000 | 0.171944 |
| 14367 | 0.165671 | 0.000000 | 0.165671 |
| 6585 | 0.163456 | 0.000000 | 0.163456 |
| 19690 | 0.160095 | 0.000000 | 0.160095 |
| 4347 | 0.159153 | 0.000000 | 0.159153 |
| 15719 | 0.156561 | 0.005155 | 0.151406 |
| 4766 | 0.243417 | 0.092564 | 0.150852 |
| 2029 | 0.149327 | 0.000142 | 0.149184 |
| 19548 | 0.161294 | 0.014016 | 0.147278 |
| 12836 | 0.146119 | 0.000000 | 0.146119 |
| 24564 | 0.165707 | 0.019807 | 0.145900 |
| 1667 | 0.170036 | 0.024717 | 0.145319 |
| 6490 | 0.144562 | 0.000000 | 0.144562 |
| 16647 | 0.147324 | 0.003386 | 0.143938 |

| | | | |
|---|---|---|---|
| 21828 | 0.167574 | 0.024935 | 0.142640 |
| 848 | 0.144347 | 0.003005 | 0.141342 |
| 19364 | 0.141108 | 0.001112 | 0.139996 |
| 7344 | 0.149478 | 0.010270 | 0.139209 |
| 13872 | 0.245501 | 0.106667 | 0.138834 |
| 12554 | 0.140330 | 0.003393 | 0.136937 |
| 20955 | 0.136513 | 0.000000 | 0.136513 |
| 15076 | 0.465050 | 0.328780 | 0.136270 |
| 17271 | 0.156657 | 0.021944 | 0.134713 |
| 16687 | 0.135332 | 0.000633 | 0.134699 |
| 14032 | 0.157261 | 0.022669 | 0.134593 |
| 16471 | 0.133981 | 0.000000 | 0.133981 |
| 24232 | 0.142612 | 0.009032 | 0.133580 |
| 23098 | 0.236032 | 0.102520 | 0.133512 |
| 4732 | 0.222698 | 0.089406 | 0.133292 |
| 18896 | 0.225538 | 0.092715 | 0.132823 |

**Tabla 1.** The forty candidate neurons extract patterns related to pathology or medicine



Figure 10. characteristic of the group of forty neurons that recognize patterns related to medicine

Upon analyzing Neuronpedia, neuron 6490 recognized patterns associated with medicine, precisely 'words related to healthcare and medical care,' achieving the search objective. Interestingly, it showed low activation for the proposed prompts.

## Limitations and future developments

If more time were available, it would be interesting to implement SHAP on the GPT-2 small model and analyze both the model outputs and the activations of each autoencoder. However, it should be noted that this would require better computational resources. Testing with a more extensive set of prompts could yield better results and help identify neurons more specialized in specific pathologies.

Future work should focus on analyzing the circuits that activate neuron 6490; better insights might be obtained with more time. Finally, two clustering techniques, such as K-means or t-SNE, could be applied to identify groups related to medical concepts. Manipulating these neurons to increase their activation would help enhance the model.

## References

Bricken, T., Elhage, N., et al. (2023). *Towards Monosemanticity: Decomposing Language Models With Dictionary Learning*. Anthropic. [Link](#).

Nanda, N. (2023). *An Extremely Opinionated Annotated List of My Favourite Mechanistic Interpretability Papers v2*. AI Alignment Forum. [Link](#).

Lundberg, S. M., & Lee, S.-I. (2017). *A Unified Approach to Interpreting Model Predictions*. In Advances in Neural Information Processing Systems (NeurIPS). [Link](#).

**SAELens**: *A Library for Interpreting Language Models Using Sparse Autoencoders*. [Link](#).

**TransformerLens**: *Tools for Analyzing Transformer Models*. [Link](#).