

NLP Explainability Techniques

Maria Ribalta i Albado

mribalta001@ikasle.ehu.eus

Abstract

In (Danilevsky et al., 2020), five explainability techniques were defined: Feature Importance, Example Driven, Surrogate Model, Provenance-Based and Declarative Induction.

In this project, we review in a theoretical and experimental way each technique. We play with at least one library or algorithm that suits each technique and evaluate how good they are theoretically (advantages and disadvantages of the approach, concerns about the methodology, etc) as well as assess how useful the libraries we used have been (in terms of speed, scalability, etc).

Conclusions show that even though some of the techniques seem more solid than others, there is still a lot of room for improvement in the XAI field and there are still lots of open questions that can be addressed.

1 Introduction

In (Danilevsky et al., 2020), the authors reviewed the state of XAI focused on the NLP field. They defined five explainability techniques in which they classified the explanations that could be given from a model: Feature Importance, Example Driven, Surrogate Model, Provenance-Based and Declarative Induction.

A year later, (Xie et al., 2020) also made a great report describing the state of XAI in general, providing great and deep insights regarding the techniques used and their reliability.

In this project, we have tested, at least, one library or algorithm that suits each technique. We evaluate how good the technique is theoretically (advantages and disadvantages of the approach, concerns about the methodology, etc) as well as assess how good the libraries used have been (in terms of speed, scalability, etc).

The work is available and reproducible in https://github.com/maria-ribalta/nlp_explainability_techniques.

2 NLP Explainability Techniques

In this section, we describe the five explainability techniques in a theoretical margin and the concerns that the literature has regarding them.

Later on, in section 3 we experiment with specific libraries and derive our own conclusions.

2.1 Feature Importance

The Feature Importance technique derives an explanation by investigating the importance scores of different features used to output the final prediction (e.g. latent features learned by neural networks, attention scores, lexical features, etc).

Theoretically, this technique is interesting since it takes features that are intrinsically from the model to show how it reacts or behaves on them. Nonetheless, an interpretable feature does not imply explainability in the prediction, in other words, the fact that we can plot a feature does not mean that is the explanation of a given classification.

As mentioned in (Xie et al., 2020) especially for attention scores, there are several concerns in the literature about its reliability regarding using it for understanding models' predictions.

2.2 Example-Driven

Example-driven approaches explain the prediction of an input instance by identifying and presenting characteristics of the given sample or similar ones.

This technique differs from the previous one in the sense that the elements used for the explanations are not part of the model but rather metrics or features extracted from the input samples themselves.

2.3 Surrogate Model

In Surrogate Models, predictions are explained by learning a second, usually more explainable model, as a proxy: this is explaining the decisions of the

first model through a second. Surrogate model-based approaches are model-agnostic and can be used to achieve either local or global explanations.

In the literature is mentioned that learned surrogate models and the original models may have completely different mechanisms to make predictions, which leads to concerns about the fidelity of surrogate model-based approaches.

2.4 Provenance-Based

Provenance-Based explanations are provided by illustrating some or all of the prediction derivation process. This process is intuitive and effective and the final prediction is the result of a series of reasoning steps.

In the original survey, the authors propose ([Amini et al., 2019](#)). This project consists of the creation of a dataset called [MathQA](#) that includes the steps to solve a mathematical problem in a code-like syntax. The motivation behind it was to provide explainability for [Google's AQuA](#) dataset (a dataset with mathematical problems and four options to choose from).

This technique is understandable but it also implies having a step of the derivational process that can be explained in a human-readable way. The resources to achieve these explanations can be expensive in time, and development and/or require human assessment.

2.5 Declarative Induction

The Declarative Induction technique works with human-readable representations, such as rules, trees, and programs are induced as explanations.

Some good examples of those are: decision trees (which generate a series of rules that step by step derive the data onto branches) or any other type of rule-based systems.

This technique is the most simple and effective in terms of explainability. Since the model is already self-explanatory we don't need to put extra effort into understanding predictions or behaviours.

However, these kind of models are quite limited; for instance: rule-based systems are not scalable and adding and reviewing rules to maintain a certain accuracy and avoid problems of concept drift is a never-ending task.

3 Experiments

3.1 Models

In this project we have used:

- A `google-bert/bert-base-uncased` model which we have fine-tuned for sentiment classification. This has been used in the experiments of the following techniques: Feature Importance, Example Driven and Surrogate Model.
- For the Provenance-based approach, we required a generative model so we used `meta-llama/Llama-2-7b-chat-hf`.
- We have used a Random Forest classifier from `sklearn` in the experiments of Surrogate Model and Declarative Induction.

3.2 Data

The data used to fine-tune the Bert model is the `poem_sentiment` dataset ([Sheng and Uthus, 2020](#)) from Hugging Face. It consists of verses of different classified from 0 to 3 according to the sentiment: 0 for negative, 1 for positive, 2 for neutral and 3 for mixed feelings.

Since the dataset was very small, the original test and train datasets were merged. For our experiments and examples, we have used the song [All Too Well \(10 Minute Version\) \(Taylor's Version\)](#) ([From the Vault](#)), which we have manually labelled with the training's dataset criteria.

Specifically in this report, we will show how each technique behaves with the verse number 93: "*From when your Brooklyn broke my skin and bones*", which is labelled with the tag 0 (negative sentiment). Other verses and their explanation results are available in the public code.

3.3 Procedure

We have chosen a library and/or algorithm for each explainability technique and checked how they worked and how good were the explanations.

To evaluate it, we have performed **human evaluation** and focused on the following aspects:

- Addressing the advantages and disadvantages of each technique (in a model and library agnostic focus).
- Addressing strong and weak points of the specific libraries or method selected (implementation details, speed of the algorithm, limitations, etc). Here we focus on the model and library used.

In general terms, we will evaluate if the explanation allows us to see if the model predictions are aligned with a human interpretation.

3.4 Feature Importance

Attention scores are easy to plot and are one of the most common features to show when explaining certain models.

The BertViz library (Vig, 2019) allows to automatically plot the attention scores and heads of any Transformer architecture. We have plotted it using our fine-tuned Bert model. The results can be seen in figure 1.

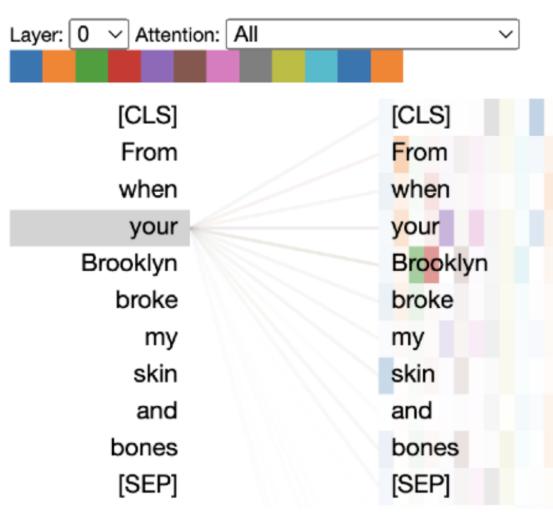


Figure 1: BertViz Attention plot

The BertViz library is smooth and easy to use. The results are fast to compute and the interactivity allows the user to play around.

Despite this, the attention plots are not fully interpretable and very few conclusions can be derived from the charts obtained (especially from the head's plot, available in Appendix A).

As a human, it is hard to check how the plots provided by the library are interpretable or related to the output of the network. Given that this technique focuses on the model's features, the selection of bad features can lead to this type of results.

A limitation that must also be mentioned is that the usage of this library is limited to transformer architectures.

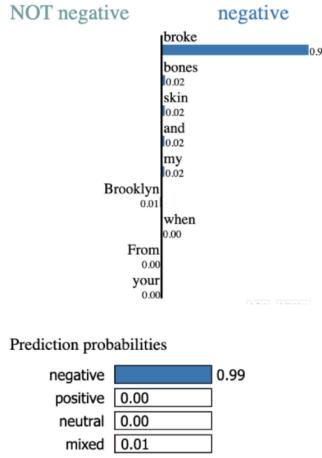
3.5 Surrogate Model

For this technique, we have used the [LIME library](#). Lime (Ribeiro et al., 2016) is the abbreviation for Local Interpretable Model-agnostic Explanations and it is an XAI library that supports surrogate model explanations for classifiers.

We provide the code for using LIME's algorithm with a neural network, which was not available online by the time this project was done.

Text with highlighted words

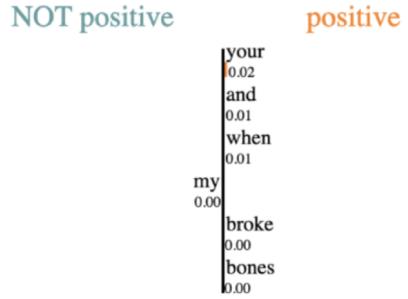
From when your Brooklyn **broke** my skin and bones



(a) LIME's explanation of fine-tuned Bert

Text with highlighted words

From **when** **your** Brooklyn broke my skin and bones



(b) LIME's explanation of undertrained random forest

Figure 2: Comparison of LIME's explainability with two different models

We checked the explanations given by the library comparing two models: our fine-tuned Bert model, which provided a correct classification of the verse (label 0: negative, with 0.99 probability); and an undertrained Random Forest, which provided an incorrect classification of the verse (label 2: neutral, with 0.75 probability).

Explanations can be checked in Figure 2. In 2a we can check how the verb “broke” had a big relevance towards the classification of the verse. In 2b we can see the effects of the undertraining since the words highlighted do not seem to have any relevance regarding the classification of the verse. In fact, they are mostly stop-words, which proves that the model is not working properly, assigning relevance to words that have little to no meaning for classification and at a lexical level.

LIME is restricted to classification algorithms so this limits many other models: regression models, generative models, etc. The explainer takes long to process the explanations for one single sample (3 minutes of execution just one sample) when using our neural model.

3.6 Example Driven

LRP (Bach et al., 2015) stands for layerwise relevance propagation. This method assigns a score to the classification layer and propagates back to each layer until it reaches the input. This way, we get a score that explains the relevance of an input given the predicted class.¹

To test it, we have used the library `interpreting_nlp`. We have checked the explanations of our fine-tuned bert. Results can be checked in Figure 3.

```
Relevance of word embeddings:
From when your Brooklyn broke my skin and bones

Relevance of positional embeddings:
From when your Brooklyn broke my skin and bones

Relevance of type embeddings:
From when your Brooklyn broke my skin and bones

Relevance of combined embeddings:
From when your Brooklyn broke my skin and bones
```

```
Attr Forward Pass Output:
[[ 4.2277    -1.9183239   -1.6355193  -0.58723706]]
LRP Scores:
From: -0.005948931621620199
when: -0.00156334931615186945
your: -0.0096980615862566
Brooklyn: -0.01835773829471589
broke: -0.0627836242039361
my: -0.017896950387639794
skin: -0.02998509989776326
and: -0.00589396115725497
bones: -0.0378684778766487
```

Figure 3: LRP score results

It can be seen that the verb *broke* has the highest LRP score, contributing to the classification of the first class, which is the label assigned to “negative”. The relevance of the verb is highlighted in the plots provided by the library.

¹Address to Appendix B for mathematical details

Example-Driven explanations are one of the most clear and understandable of all those we have tested. They can also have a wide variety of implementations and contexts.

The algorithm has proved to detect the verb form “broke” as crucial for the verse to be labelled as negative. As humans, given the sentence, it is logical to think that the verb contributes to a bad connotation in the verse. In terms of library efficiency, the algorithm has run fast.

We believe that this is the best technique from the ones we have tried since it finds a balance between good and understandable explanations, reliability due to the way it is computed; and good latency and ease of use of the library.

3.7 Provenance Based

Techniques like COT(Wei et al., 2023) (Chain of Thought) or RAG(Lewis et al., 2021) (Retrieval Augmented Generation) could fall into provenance-based explanations.

In COT, the model (LLM) provides step-by-step reasoning to reach an objective or decision.

In RAG systems, we send a query that, before reaching the LLM, is complemented with additional information extracted from an external resource (an updated DB, scraping the web, etc). RAG systems allow better contextualization and more accurate responses.

In this experiment, we have coded a toy-RAG system. We asked an generative LLM (Llama 2) to classify the verses and provide an explanation for it. We have compared how the decision of the model changed if we provided additional information (using the RAG system) or not (sending only the verse).²

Our results showed that only in two cases the prediction of the model was the same with and without additional context, which shows that adding information, (explainability) alters the behaviour of the model.

However, our model had better classification results with the non-contextualized samples than the ones that used the RAG system ($f_1 = 0.66$ and $f_1 = 0.54$ respectively). We attribute it to the simplicity of our system.

One example in which the model benefited from having context is the verse “*From when your Brooklyn broke my skin and bones*”. The LLM labelled it correctly with a 0 (Negative) with the additional

²Address to Appendix C for complete details on the implementation

context provided by the RAG and incorrectly with a 3 (Mixed) when no contextualized information was passed. However, their explanation showed that the LLM had a hallucination behaviour since it referred to the word "my" as the main cause of "Mixed" classification.³

In terms of explainability, even if our experiments didn't come as expected, prompts with additional information allow a better understanding and contextualization of a model's decision since they are understandable by a human. However, we could also argue the fact that generative models work by probability and the additional information even if it changes the given response, it might not affect positively to the correctness of the answer.

3.8 Declarative Induction

Since this project was mainly focused on Deep Learning methods, we have simply provided an example of a Random Forest model to show how declarative induction works. However, no further contributions have been proposed.

Scikit Learn implementation allows the visualization of each of the estimators and their derivation rules, as shown in Figure 4.

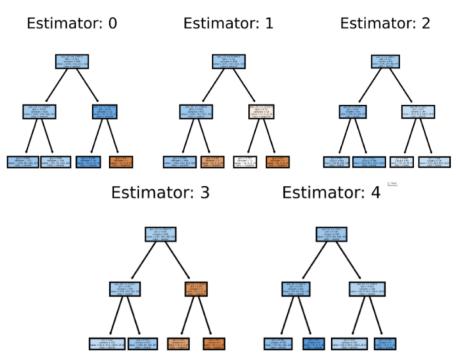


Figure 4: Random Forest

Our implementation of the Random Forest is explanatory but as mentioned, it might not be accurate enough for other tasks. Also, explanations might become incomprehensible if a large number of estimators is used, since it is not possible to check all the rules for each tree.

4 Conclusions and Future Work

We have seen that Feature Importance approaches might involve features that are interpretable but not directly related to the prediction result. Attention

³You can check the prompt and the answer of the LLM in Appendix C

using BertViz, in our case, has not shown good results in terms of explainability.

In Surrogate Models, the explanations shown with LIME's algorithm are logical and seem to make sense. As a human, one can extract good insights from using this library. However, there are concerns regarding the reliability of this approach, since the proxy model might have completely different procedures to do predictions than the original one. Therefore, this approach lacks solidity.

For Example Driven approaches, LRP scores have proved to be useful and interpretable. The scores are given by a mathematical formula that involves knowledge regarding the neuron connections, which makes this approach solid. We believe that this is the best approach of the ones we have tried since it has a good balance in speed, explainability and solidity.

Provenance-based models are easy to understand and provide contextualization and clarity at certain steps of the derivation process. Our RAG toy system has shown that a contextualized answer leads to different decisions than when having no additional context. This approach, however, is not possible or suitable for all tasks, so its scope is limited.

Finally, Declarative Induction techniques rely on the fact that a model's architecture or decision path is explainable which is highly desirable but not possible in most cases. Problems of scalability or algorithms being too simple for certain tasks are the main limitation of this technique.

Additionally, there are still many open questions in the field of XAI, since there is no standardized way of evaluating it nor all techniques are directly comparable.

All in all, what is a good explanation? Who is the end user of this explanation? What task or procedure are we trying to explain? These questions are still open and few of the approaches that are used nowadays address them. XAI is, in general, a very immature field yet that has the potential to grow in the following years.

References

- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. MathQA: Towards interpretable math word problem solving with operation-based formalisms. 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

2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. [On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation](#). *PLOS ONE*, 10(7):1–46.

Marina Danilevsky, Kun Qian, Ranit Aharonov, Yannis Katsis, Ban Kawas, and Prithviraj Sen. 2020. [A survey of the state of explainable AI for natural language processing](#). *CoRR*, abs/2010.00711.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. [Retrieval-augmented generation for knowledge-intensive nlp tasks](#).

Marco Túlio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. ["why should I trust you?": Explaining the predictions of any classifier](#). *CoRR*, abs/1602.04938.

Emily Sheng and David Uthus. 2020. [Investigating societal biases in a poetry composition system](#).

Jesse Vig. 2019. [A multiscale visualization of attention in the transformer model](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 37–42, Florence, Italy. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#).

Ning Xie, Gabrielle Ras, Marcel van Gerven, and Derek Doran. 2020. [Explainable deep learning: A field guide for the uninitiated](#). *CoRR*, abs/2004.14545.

A Feature Importance

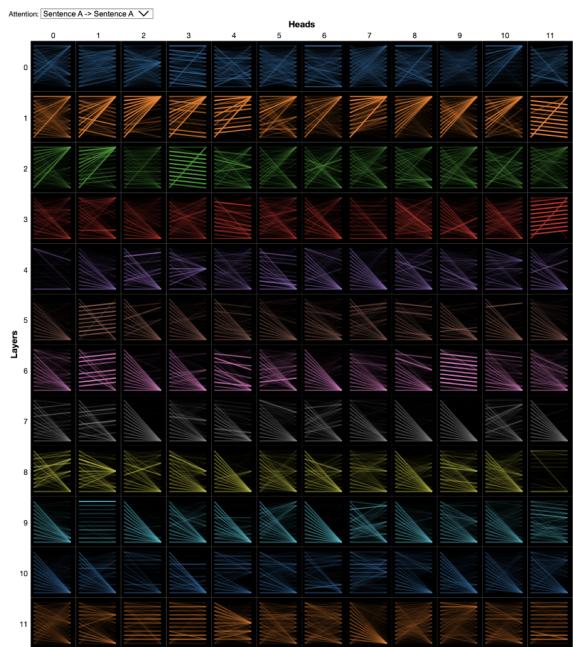


Figure 5: BertViz Attention Heads Plot

B LRP scores

LRP relevance score is defined by:

$$R_i = \sum_j \left(\alpha \cdot \frac{(w_i w_{ij})^+}{\sum_i (w_i w_{ij})^+} - \beta \cdot \frac{(w_i w_{ij})^-}{\sum_i (w_i w_{ij})^-} \right) \cdot R_j \quad (1)$$

- $(x_i)_i$ is the neuron activation at layer l
- $(R_j)_j$ is the relevance score associated to the neurons at layer $l + 1$
- w_{ij} is the weight connectin neuron i to neuron j .
- Where $+$ and $-$ denote the positive and negative parts.
- α and β fulfill that: $\alpha - \beta = 1$

C RAG

The high-level components of a RAG system are a corpus, an input from the user and a similarity measure between the corpus and the user input.

In our case, we defined:

- The input: verses of the test data, lyrics from the song. The same as in the previous experiments.

- The corpus: the corpus works as the “additional information” that helps contextualise and enrich the query before reaching the LLM. In this case, we took the comments of the song from the [Genius webpage](#). Genius provides context and explanations of the lyrics from songs.
- The similarity measure: the Jaccard similarity score.

Then, the process followed for the RAG system:

1. Receive a user input (the lyrics)
2. Perform a similarity measure (we chose the Jaccard similarity as in the article) to match the most suitable document (the context from Genius)
3. Send to the LLM the query with the added information. Our LLM is Llama2 and we call it via the ‘ollama’ proxy.

The Jaccard similarity used for querying the corpus given the input lyric is defined as:

$$J(i, d) = \frac{|i \cap d|}{|i \cup u|}, i \in I, d \in D \quad (2)$$

Where:

- D is the set of documents, in our case the additional information extracted from Genius.
- I is the user input, in our case, the verse from the test.
- The intersection \cap represents the total number of words that appear in both: the verse and the document.
- The union \cup represents the total number of words that appear in the verse or the document.

We take the document d that has the higher Jaccard Similarity given the input i .

The prompt used as input for queries with contextualized information:

You are a classifier that, given a sentence, says if the sentence is negative, positive, neutral or mixed. The sentence is “{song_lyric}”. Additional info is that: {additional_information}”.

Say which class is more suitable and a short explanation

The prompt used as input for queries without contextualized information:

You are a classifier that, given a sentence, says if the sentence is negative, positive, neutral or mixed. The sentence is “{user_input}”. Say which class is more suitable and a short explanation

Our example of verse 93:

```
{'label': 0,
'lyric': 'From when your Brooklyn broke my skin and bones',
'predicted': 0,
'prompt': 'You are a classifier that, given a sentence, says if the sentence is negative, positive, neutral or mixed. The sentence is "From when your Brooklyn broke my skin and bones". Additional info is that: "sticks and stones may break my bones" mean that it's clear "he put her through some sort of verbal abuse". Say which class is more suitable and a short explanation',
'response': '\n'
'Based on the given sentence, I would classify it as negative. '
'The phrase "From when your Brooklyn broke my skin and bones" '
'suggests that someone has been physically hurt or abused, with '
'"Brooklyn" likely being a person who inflicted the harm. The '
'additional context you provided further reinforces this '
'interpretation, as "sticks and stones may break my bones" is '
'often used to convey the idea that verbal abuse can be just as '
'damaging as physical harm.\n'
'\n'
'Therefore, the most suitable class for this sentence would be '
'"Negative."',
'sample_idx': 5}
```

[+ Code](#) [+ Markdown](#)

Figure 6: Input and output of a prompt with additional information provided by a toy-RAG system. The LLM predicts the label correctly.

```
{'label': 0,
'lyric': 'From when your Brooklyn broke my skin and bones',
'predicted': 3,
'prompt': 'You are a classifier that, given a sentence, says if the sentence is negative, positive, neutral or mixed. The sentence is "From when your Brooklyn broke my skin and bones". Say which class is more suitable and a short explanation',
'response': '\n'
'Based on the sentence provided, I would classify it as "Mixed" '
'because it contains both negative and positive elements.\n'
'\n'
'The phrase "your Brooklyn broke my skin and bones" is negative '
'in tone, as it describes physical harm caused by something '
'(Brooklyn) that is presumably a person or entity. The use of the '
'word "broke" implies damage or injury, which has a negative '
'connotation.\n'
'\n'
'However, the sentence also contains positive elements, such as '
'the use of the word "my," which suggests ownership or attachment '
'to the thing being described. Additionally, the phrase "skin and '
'bones" could be interpreted as a metaphor for something fragile '
'or delicate, which has a neutral connotation.\n'
'\n'
'Therefore, based on the mixed nature of the sentence, I would '
'classify it as "Mixed."',
'sample_idx': 5}
```

Figure 7: Input and output of a prompt without additional information. The LLM does not predict the label correctly.