

NLP Explainability Techniques

Maria Ribalta i Albado - DL4NLP

Motivation

XAI is the discipline that explains why an AI model takes the decision it takes.

In 2020, Danilevsky et al provided a survey summarizing the most relevant XAI for NLP contributions were.

We have taken the explainability techniques on the paper and tried them out as well as proposed some additional techniques that have appeared since then.

Based on..

- A Survey of the State of Explainable AI for Natural Language Processing (Danilevsky et. al, 2020)

Models

- google-bert/bert-base-uncased
- meta-llama/Llama-2-7b-chat-hf
- sklearn's Random Forest

Data

- Fine-tuning** - poem_sentiment
 - verses of poems classified from 0 to 3 according to the sentiment
 - 0: negative, 1: positive, 2: neutral, 3: mixed
- Tests and experiments** - *All Too Well (10 Minute Version) (Taylor's Version) (From the Vault)* (manually labelled with the training's dataset criteria).

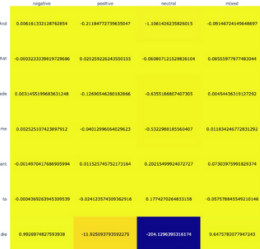
Feature Importance

Layer: 0 Attention: All



Bertviz Attention Visualization System

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 NNs, attention scores, lexical features, etc). **BertViz** is the most famous library to plot attention scores and heads from any BERT model (or similars). We have as well implemented the importance scores of each word of a sentence with **Representation Erasure** scores.



Word Representation Erasure Results

Example Driven

Attr Forward Pass Output:
[[2.465749 -1.7446934 -1.069228 0.37986298]]
LRP Scores:
Just: -0.004472918022042428
between: -0.014479319803017129
us: -0.0030729131287711
;: -0.00022758517016512742
did: -0.0230418528388121
the: -0.0027981546134820574
love: 0.0006156135456628592
affair: -0.001354776939090957
ma: -0.07915791916835851
##in: -0.027547992649463837
you: -0.0026281125802477415
all: 0.031026367439432788
too: -0.0021704092748219767
well: 0.014077428293788757
?: 0.0034378883720426255

Relevance of word embeddings:

Just between us , did the love affair ma ##in you all too well ?

Relevance of positional embeddings:

Just between us , did the love affair ma ##in you all too well ?

Relevance of type embeddings:

Just between us , did the love affair ma ##in you all too well ?

Relevance of combined embeddings:

Just between us , did the love affair ma ##in you all too well ?

LRP results for the sentence: Just between us, did the love affair main you all too well?

Example-Driven approaches explain the prediction of an input instance by identifying and presenting other instances, usually from available labeled data, that are semantically similar to the input instance.

LRP stands for layerwise relevance propagation and is an example-driven technique to showcase the relevance of certain features: positional embeddings, word embeddings, etc. The library **interpret_nlp** allows to extract the LRP scores from any Bert model, in our case, our fine-tuned Bert with the poem_sentiment data.

Surrogate Model

Text with highlighted words

From when your Brooklyn broke my skin and bones

NOT negative negative



Prediction probabilities

negative 0.99

positive 0.00

neutral 0.00

mixed 0.01

LIME's Explainability of a prediction of a Bert model

Text with highlighted words

From when you Brooklyn broke my skin and bones

NOT positive positive



Prediction probabilities

negative 0.19

positive 0.04

neutral 0.75

mixed 0.02

LIME's explainability of a prediction of an undertrained tree 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.

LIME's algorithm consists on generating an input perturbation to evaluate how that affects the prediction and via the surrogate model assesses the explainability.

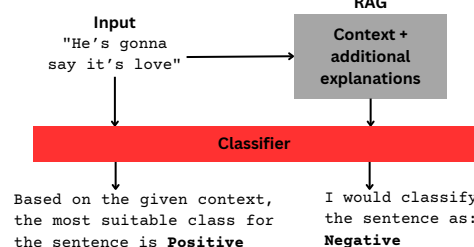
What is a good explanation? Who is the end user of that explanation? What task or objectives are we trying to evaluate? These questions don't permit the standardisation of simple processes for evaluating XAI methods. While there are different types of evaluation (informal examination, comparison to ground truth, human evaluation and others), how to evaluate an XAI approach is still an **open question**. In this project we have applied **human evaluation** to assess each of our approaches, but these cannot be directly compared:

- While **Feature Importance** is direct and easy to implement, recent studies show that, for instance, a feature such as attention is not always interpretable and that attention does not always lead to insight into model prediction.
- Surrogate models** may have completely different mechanisms to make predictions than the original model. This has led in many cases to concerns about the fidelity of this kind of approach.
- Example-Driven** approaches are complex to code and each architecture will require a different implementation, which is costly in many levels.
- The **provenance-based** method is desirable but not always possible (we cannot always build interpretable architectures nor explainable datasets).
- Declarative induction** techniques are simple and easy to understand, yet they are not always scalable (rule-based systems) or complex enough to generalize all tasks (random forests).

XAI is yet a huge field to investigate and we are sure many more techniques will arise with the creation of new models and new architectures,

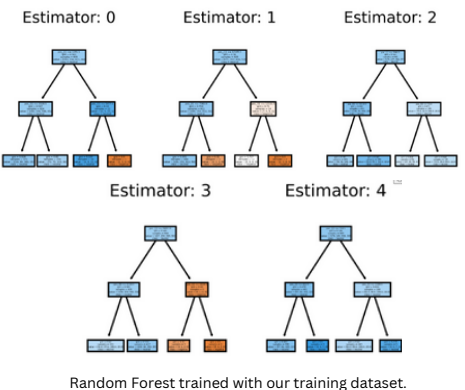
Provenance Based

The **provenance-based** technique is that in which 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. Techniques like **Chain of Thought** or **RAGs** could be considered provenance-based explainability techniques since the steps and additional information the model uses in order to make predictions are human understandable. We generated a toy-RAG system that classifies verses within the context of the song.



Simulation of the process followed by our RAG-toy system to classify the sentence "he's gonna say it's love". In the context of the song, the sentence has a negative connotation.

Declarative Induction



The **declarative Induction** technique is the one that works with human-readable representations, such as rules, trees, and programs are induced as explanations.

Decision trees generate a series of rules that step by step derive the data onto branches. **Random Forests** are the ensemble of different decision trees and are one of the most common and well-known self-explainable models.

Discussion



Find the whole project in Github!