

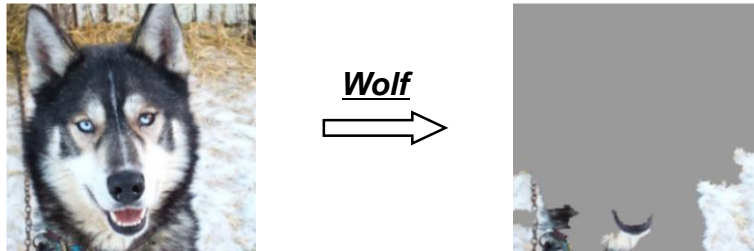
Interpretability of Neural Network Models Used in Data Analysis

Author: Branislav Pecher

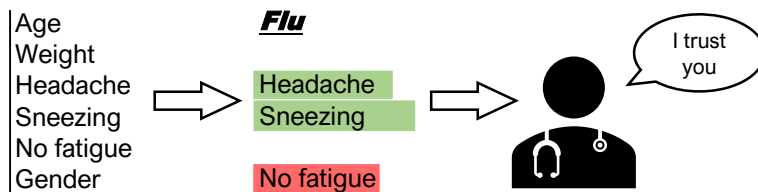
Supervisor: Ing. Jakub Ševcech, PhD.

Why is interpretability important?

- Helps audit and find problems decision process of models



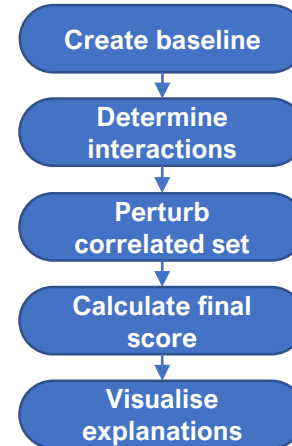
- Enhances trust in the model and therefore its use in practice



- Attribution methods most popular for neural networks
- Shortcomings of current methods:
 - Problems when encountering interactions in data
 - Mostly focused on problems involving image data
 - Explanations sometimes hard to interpret

Proposed method

- Used for explanation of decisions in **text** classification with the use of word embeddings
- Take into consideration interactions when generating explanations – more precise explanations
- Generate insightful explanations – visualization of explanations
- 5 step process



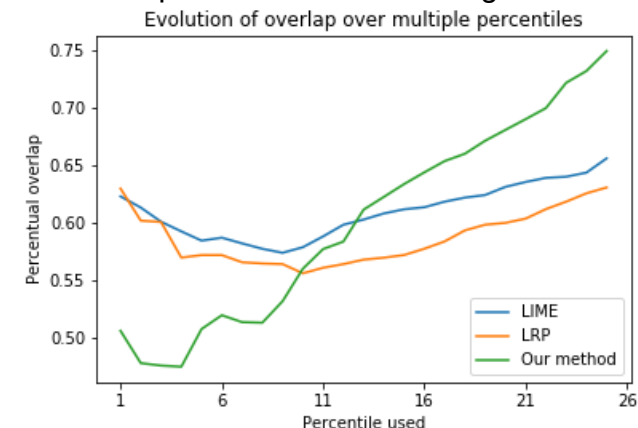
Enjoyable in spite of Leslie Howard's performance. Mr. Howard plays Philip as a flat, uninteresting character. One is supposed to feel sorry for this man; however, I find myself cheering Bette Davis' Mildred. Ms. Davis gives one her finest performances (she received an Academy Award nomination). Thanks to her performance she brings this rather dull movie to life. **Be sure not to miss when Mildred tells Philip exactly how she feels about him.

Enjoyable in spite of Leslie Howard's performance. Mr. Howard plays Philip as a [redacted]. [redacted] One is supposed to feel sorry for this man; however, I find myself cheering Bette Davis' Mildred. Ms. Davis gives one her finest performances (she received an Academy Award nomination). Thanks to her performance she brings this rather [redacted] movie to life. **Be sure not to miss when Mildred tells Philip exactly how she feels about him.

Enjoyable in spite of Leslie Howard's performance. Mr. Howard plays Philip as a flat, uninteresting character. One is supposed to feel sorry for this man; however, I find myself cheering Bette Davis' Mildred. Ms. Davis gives one her finest performances (she received an Academy Award nomination). Thanks to her performance she brings this rather dull movie to life. **Be sure not to miss when Mildred tells Philip exactly how she feels about him.

Experiments and results

- Comparison with other attribution methods – LIME[1] and LRP[2]
- Solution: User experiment to determine ground truth



- Our method is better at finding important features that are otherwise overlooked due to the interactions