# Master Thesis

## Explorative study of explainable artificial intelligence techniques for sentiment analysis applied for English language

by Rohan Kumar Rathore

Advisor: Dr. Anton Kolonin

Scientific workshop "Big Data Analytics"
Novosibirsk State University, Russia

March 16, 2021

# Introduction

- Artificial intelligence : Artificial agents achieving goals smartly
- Machine learning : Algorithmic models responsible for smartness
- Explainable artificial intelligence : Techniques to explain the models
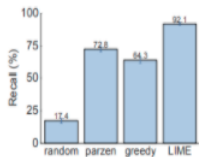
# Outline
Explainable artificial intelligence (XAI) techniques for sentiment analysis

- Model development
  - Sentiment analysis model on IMDB movie reviews dataset
- Technique I
  - Local interpretable model-agnostic explanations (LIME): Explaining with surrogate models
- Technique II
  - Layer-wise relevance propagation (LRP): Explaining with propagated weights relevance scores of the network
- Technique III
  - Artificial neural network decision tree algorithm (Ruleex ANN-DT): Explaining by extraction of decision trees from artificial neural networks
- Performance Analysis
  - Simulatability test: A model is simulatable when a person can predict its behavior on new inputs
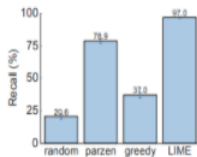
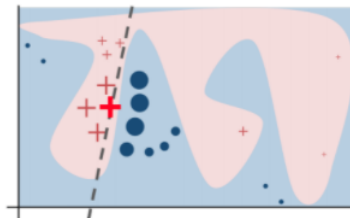$$\xi(x) = \underset{g \in G}{\text{argmin}} \ \mathcal{L}(f, g, \pi_x) + \Omega(g)$$



(a) Sparse LR

(b) Decision Tree

**Recall on truly important features**



**Toy example to present intuition for LIME.**

Input text: *"This movie was beyond disappointment. Well acted story that means nothing. The plot is ridiculous and even what story there is goes absolutely nowhere. It truly isn't worth a nickel, buffalo or otherwise..pun intended!"*
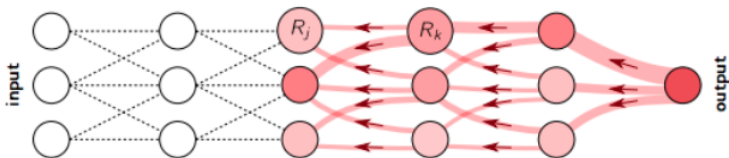
POS WORD CONTRIBUTE:

worth Well truly

NEG WORD CONTRIBUTE:

ridiculous disappointment nothing even plot acted means

**Fig. 10.2.** Illustration of the LRP procedure. Each neuron redistributes to the lower layer as much as it has received from the higher layer.

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$$

**Basic Rule**

$$R_j = \sum_k \frac{a_j \cdot \rho(w_{jk})}{\epsilon + \sum_{0,j} a_j \cdot \rho(w_{jk})} R_k,$$

**Epsilon Rule**

$$R_j = \sum_k \frac{a_j \cdot (w_{jk} + \gamma w_{jk}^+)}{\sum_{0,j} a_j \cdot (w_{jk} + \gamma w_{jk}^+)} R_k$$

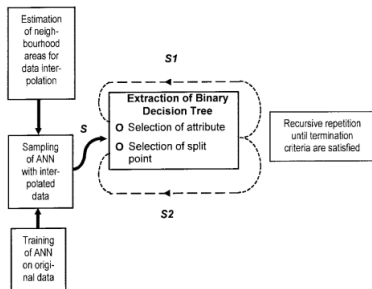**Gamma Rule**

Input text: *"This movie was beyond disappointment. Well acted story that means nothing. The plot is ridiculous and even what story there is goes absolutely nowhere. It truly isn't worth a nickel, buffalo or otherwise..pun intended!"*

WORD CONTRIBUTE:

ridicul disappoint noth absolut worth well act even plot mean

- Selection of Attribute: Similar to CART algorithm of reducing the

$$V_w = \sum_{k=1}^{2} \frac{n_k}{n} \text{Var}(O_k)$$

  entropy

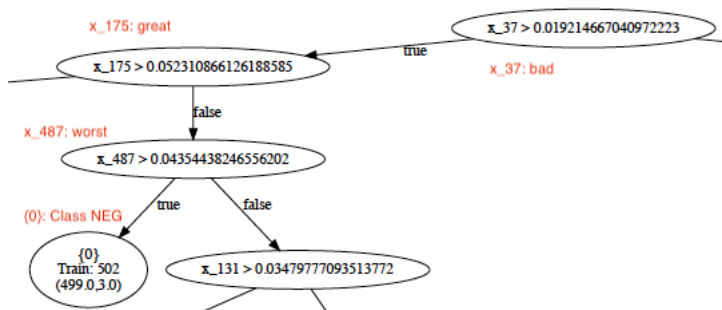- Stopping criteria: Standard deviation or the variance is zero
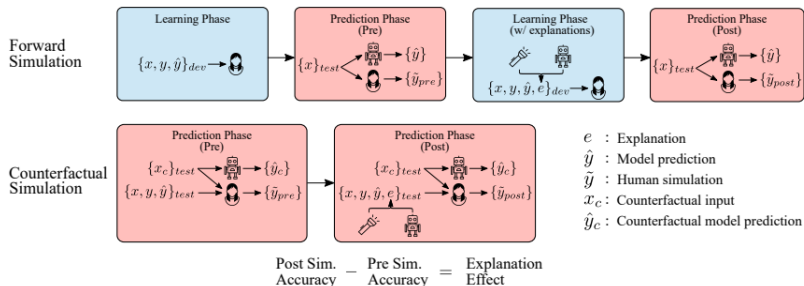
- Statistical pruning technique: chi squared

- Decision Tree representation (showing sub-section here)

# Performance Analysis Method

Simulatability test: Model is simulatable if person can predict its behavior on new inputs

# Results

| XAI Phase | Forward Test | Counterfactual Test | Total |
| --- | --- | --- | --- |
| LIME - Pre | 90.0% | 65.0% | 77.5% |
| LIME - Post | 90.0% | 90.0% | 90.0% |
| **LIME - Change** | **0.0%** | **25.0%** | **12.5.0%** |
| LRP - Post | 90.0% | 65.0% | 77.5% |
| LRP - Pre | 95.0% | 85.0% | 90.0% |
| **LRP - Change** | **5.0%** | **20.0%** | **12.5.0%** |

- A total of 120 data points were collected
- Improvement the accuracy of model prediction capability of the human subject by 12.5%

# Paper Readiness

- International Conference on Data Science and Applications, ICDSA 2021 (Accepted for presentation in conference)
- Improve results (Performance analysis on ANN-DT technique) and apply in other conferences

# References

- Gregor P.J.S., Chris A., Francois S.G., ANN-DT: An Algorithm forExtraction of Decision Trees from Artificial Neural Networks, IEEETransactions on Neural Networks, Vol.10 No.6, 1999
- Montavon G., Binder A., Lapuschkin S., Samek W., Müller K.R.: Layer-Wise Relevance Propagation: An Overview. In: Samek W., Montavon G., Vedaldi A., Hansen L., Müller K.R. (eds.) Explainable AI: Interpreting, Explaining and Visualizing Deep Learning, LNCS, vol. 11700, pp. 193–209. Springer, Cham (2019)
- Ribeiro, M.T., Singh, S., Guestrin, C.: "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16), pp.1135–1144, Association for Computing Machinery, New York (2014)
- GHase, P., Bansal, M.: Evaluating Explainable AI: Which Algorithmic Explanations Help Users Predict Model Behavior?. arXiv:2005.01831 (2020)