



# Data Systems Group

Institute of Computer Science, University of Tartu

XAI-Spring2020 Project Plan  
5<sup>th</sup> March 2020  
Abdul Wahab

**Scope:** A Comprehensive Qualitative and Quantitative Evaluation of Model-Agnostic Local Interpretability Methods.

Phase 1: Benchmarking of Tabular Interpretability methods on 3 BioMedical Tabular Datasets		
Task	Comments	Timeline
1. Consolidating Metrics + Implementing New Metrics	Consolidating all existing and relevant metrics for interpretability + implementing any new metrics that suit the Tabular Dataset type Interpretability.	9-15 <sup>th</sup> March
2. Dataset and Model setup	Gathering the relevant datasets and setting up the model and interpretability methods on these datasets.	16-22 <sup>nd</sup> March
2. CLEAR Benchmarking	Understanding the approach behind this approach method and doing the benchmarking experiments	23-29 <sup>th</sup> March
4. BreakDown Benchmarking	// //	30 <sup>th</sup> March-05 <sup>th</sup> April
5. D-LIME Benchmarking	// //	06 <sup>th</sup> April-12 <sup>th</sup> April
6. Counterfactual Explanations Guided by Prototypes	// //	13 <sup>th</sup> April - 19 <sup>th</sup> April
Phase 2: Benchmarking of Text Interpretability Methods on 3 Datasets		
7. Consolidating Metrics + Implementing New Metrics	Consolidating all existing and relevant metrics for interpretability + implementing any new metrics that suit the Tabular Dataset type Interpretability.	20 <sup>th</sup> April-26 <sup>th</sup> April
8. Dataset and Model setup	Gathering the relevant datasets and setting up the model and interpretability methods on these datasets.	27 <sup>th</sup> April - 03 <sup>rd</sup> May

<b>9. POLAR</b>	Understanding the approach behind this approach method and doing the benchmarking experiments	04 <sup>th</sup> May - 10 <sup>th</sup> May
<b>8. Text Dataset Interpretability Method 2</b>	// //	11 <sup>th</sup> May - 17 <sup>th</sup> May
<b>9. Text Dataset Interpretability Method 3</b>	// //	18 <sup>th</sup> May - 24 <sup>th</sup> May
<b>Phase 3: Benchmarking of Above Methods on Images</b>		
<b>7. Consolidating Metrics + Implementing New Metrics</b>	Consolidating all existing and relevant metrics for interpretability + implementing any new metrics that suit the Tabular Dataset type Interpretability.	15 <sup>th</sup> May - 24 <sup>th</sup> May
<b>8. Dataset and Model setup</b>	Gathering the relevant datasets and setting up the model and interpretability methods on these datasets.	25 <sup>th</sup> May - 31 <sup>st</sup> May
<b>10. Counterfactual Explanations Guided by Prototypes</b>	Understanding the approach behind this approach method and doing the benchmarking experiments	1 <sup>st</sup> June - 6 <sup>th</sup> June
<b>11. Contrastive Explanation (Foil Trees)</b>	// //	7 <sup>th</sup> June - 13 <sup>th</sup> June
<b>12. Convex Density Constraints for Computing Plausible Counterfactual Explanations</b>	// //	14 <sup>th</sup> June - 20 <sup>th</sup> June