



Interpretability Approaches applied to Predictive Models in Clinical Healthcare

Trends in Bioinformatics
Final Presentation
Tom Martensen, Axel Stebner

Agenda

1. Recap

2. Methods

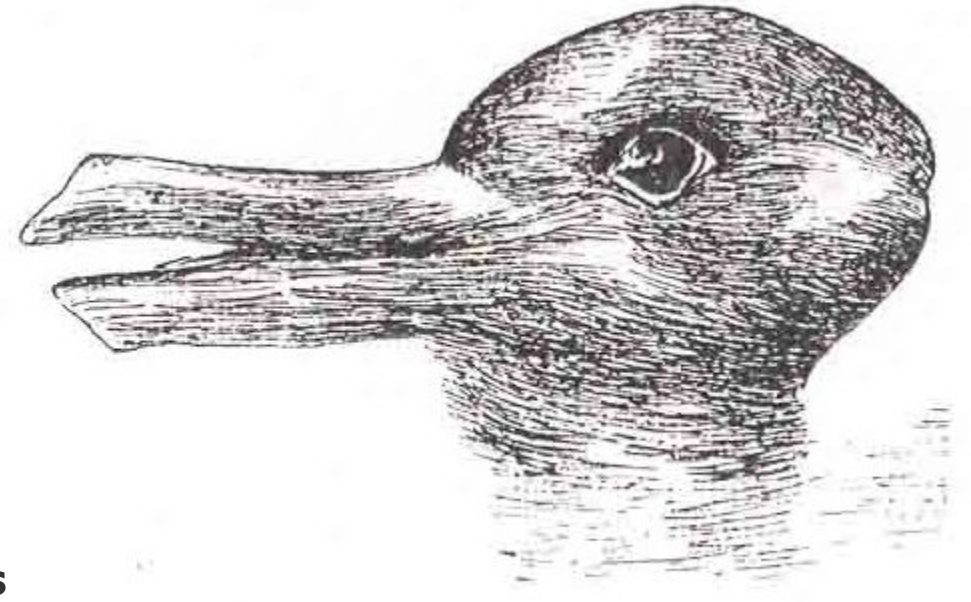
1. Building a Clinical Prediction Model

2. Applying Interpretability Methods in Detail

3. Making Interpretability Available for Domain Experts

3. Results

4. Outlook



Recap: Visions & Objectives

VISION 1

Find and validate medical hypotheses regarding mortality and recovery of AKI

- Train CPM
- Predict patient outcomes
- Gather interpretations
- Derive and evaluate clinical hypotheses

VISION 2

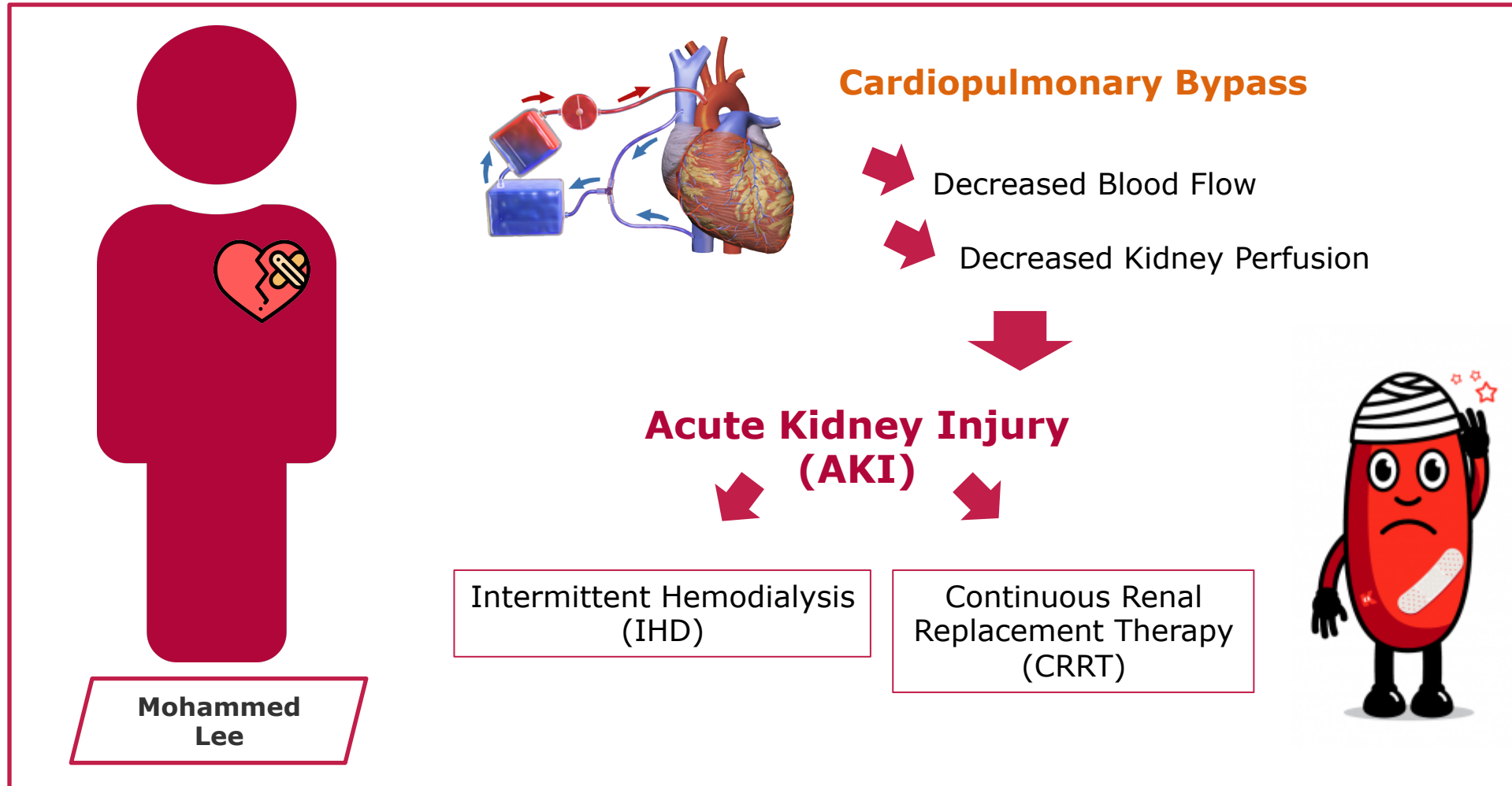
Make interpretations of CPMs available to physicians

- Interpret any CPM
- Make interpretations comparable side-by-side
- Show complexity-faithfulness tradeoff

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 3

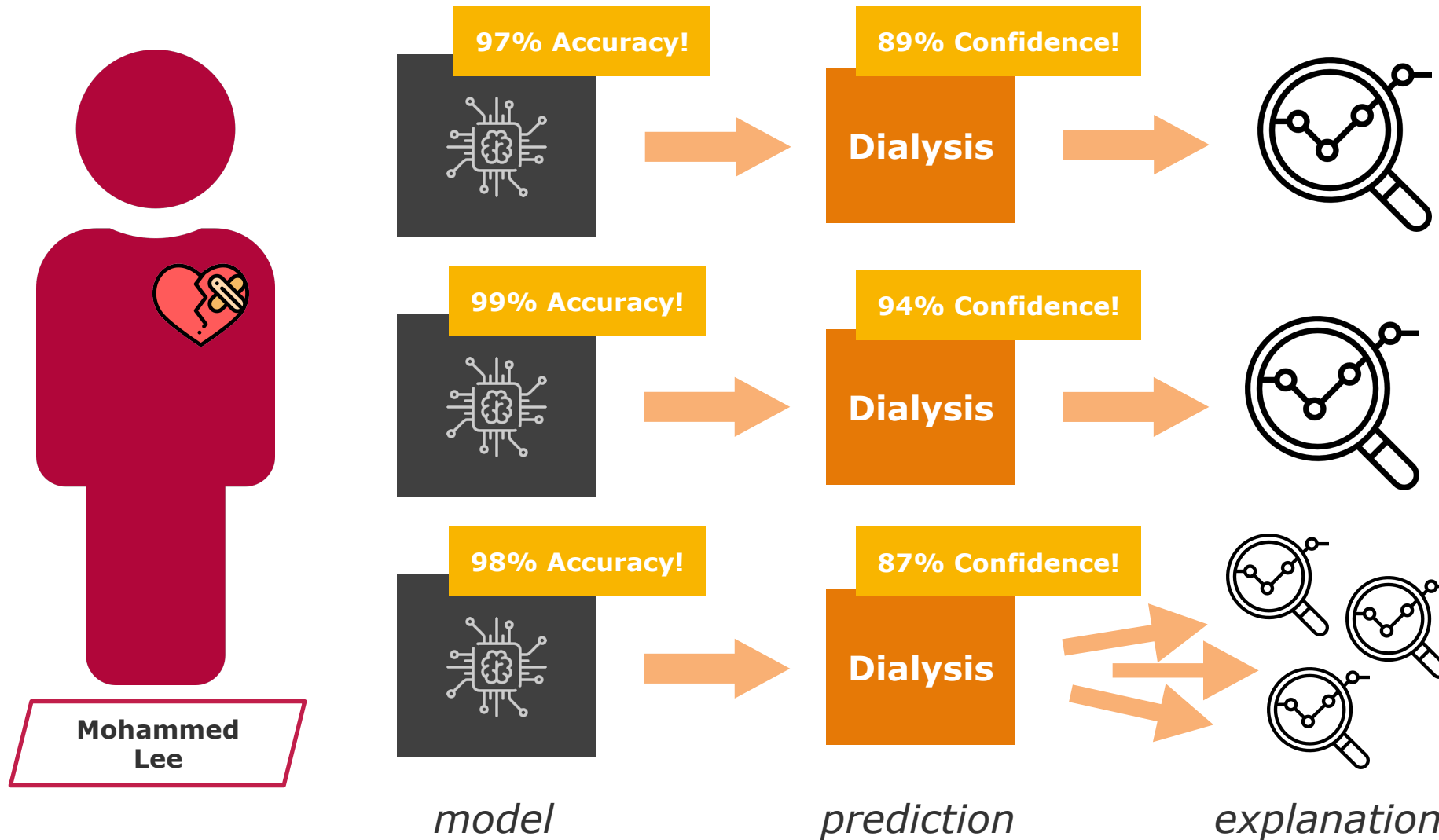
Recap: Use Case – Acute Kidney Injury



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 4

Recap: Use Case – Therapy of Acute Kidney Injury

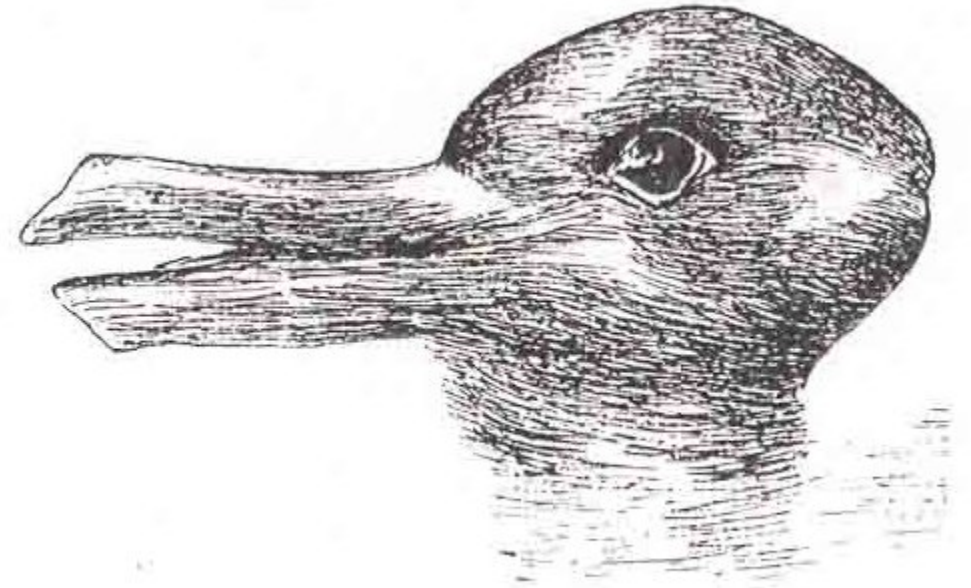


Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 5

Methods

- Building a Clinical Prediction Model
- Applying Interpretability Methods in Detail
- Making Interpretability Available for Domain Experts

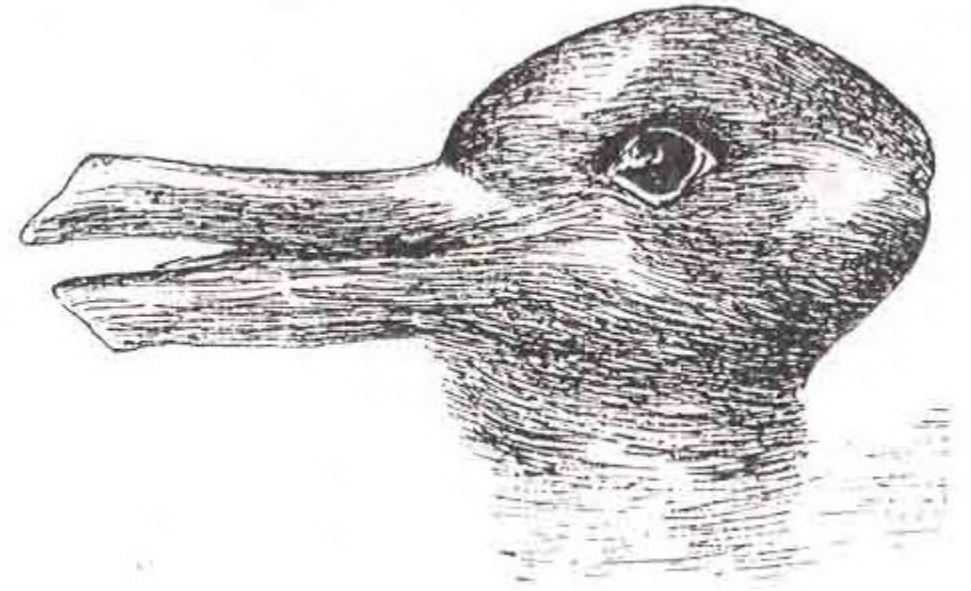


Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **9**

Methods

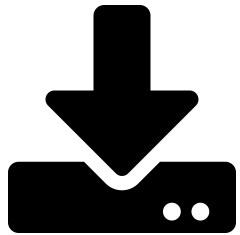
- **Building a Clinical Prediction Model**
- Applying Interpretability Methods in Detail
- Making Interpretability Available for Domain Experts



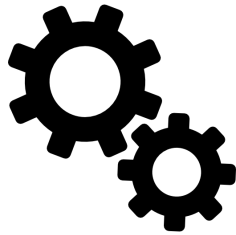
Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **10**

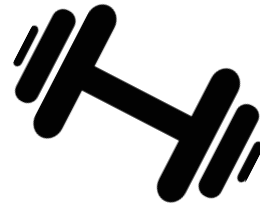
Methods: Building a Clinical Prediction Model



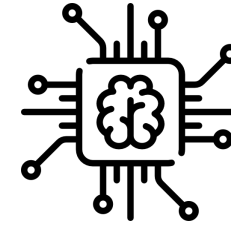
data retrieval



preprocessing



model training

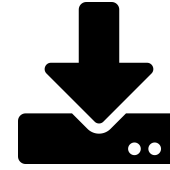


prediction

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **11**

Building a Clinical Prediction Model: Data Retrieval



MIMIC-III Database

LAB EVENTS

- Different lab values
- Flagged
- Timestamp

ICU STAYS

- Start
- End

ICU STAY VITALS (FIRST DAY)

Aggregated lab values
of first day of ICU stay

Procedure Events

- All procedures in hospital
- Timestamp

Labels

- Labels for classification:
 - Dosage
 - Therapy type

AKI Patients

- Patient master data
- Only patients with AKI

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **13**

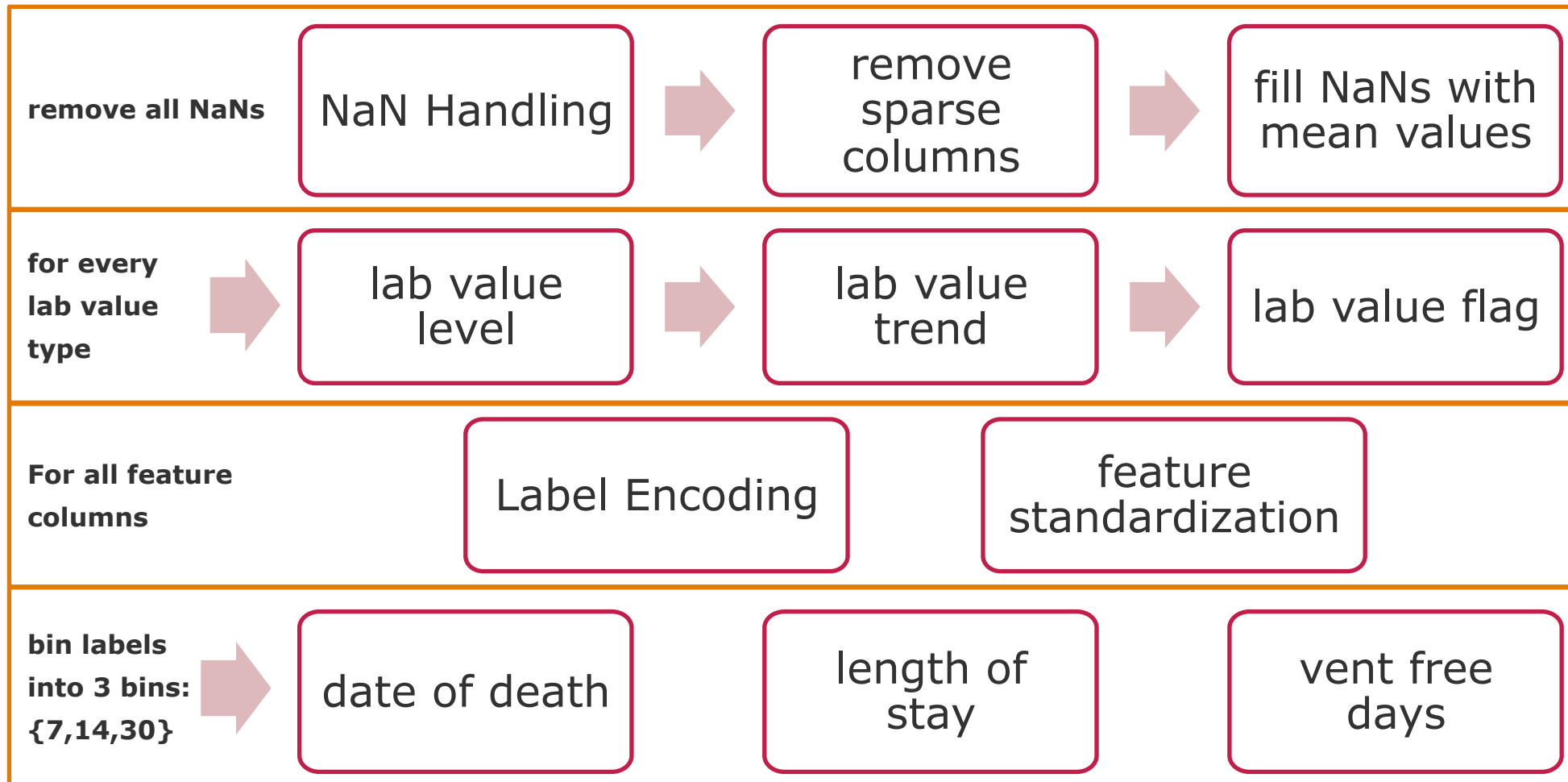
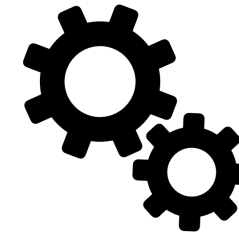
data retrieval

preprocessing

model training

prediction

Building a Clinical Prediction Model: Data Preprocessing



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 14

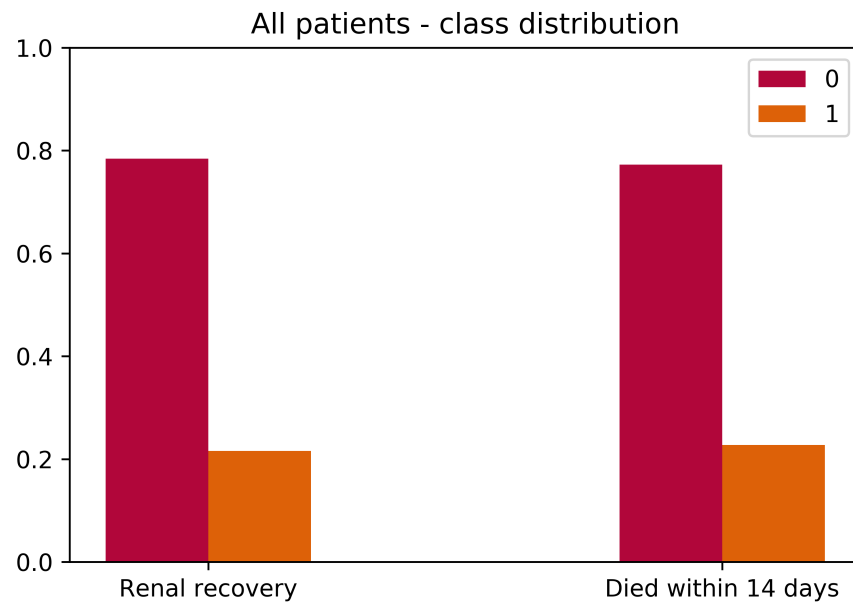
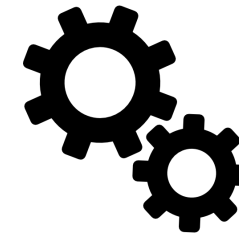
data retrieval

preprocessing

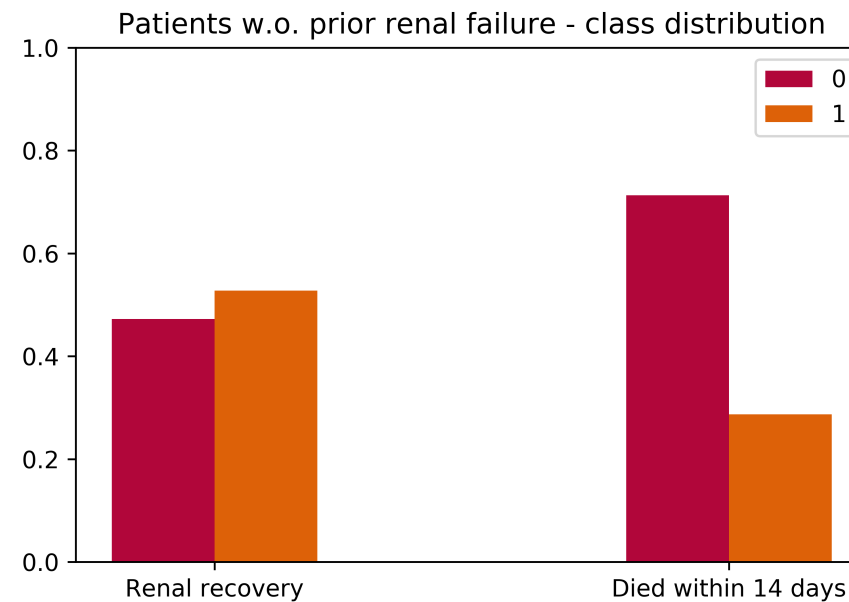
model training

prediction

Building a Clinical Prediction Model: Data Preprocessing – Dataset Characteristics



→ 2945 instances



→ 944 instances

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 16

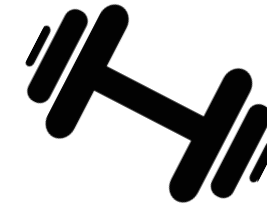
data retrieval

preprocessing

model training

prediction

Building a Clinical Prediction Model: Model Training



Random Parameter Search:

- Randomly pick parameters from specified range
- Create classifier
- 5-fold cross validation
- Evaluate with AUROC score



Trained model with optimal
parameter setting



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **18**

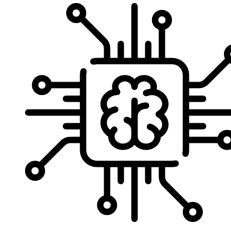
data retrieval

preprocessing

model training

prediction

Building a Clinical Prediction Model: Prediction Patient Outcomes



results_gb_all_0_DIED_14_DAYS.dat

```
{'criterion': 'friedman_mse', 'loss': 'exponential', 'max_depth': 160, 'max_leaf_nodes': 653, 'min_samples_leaf': 38, 'n_estimators': 740}
```

results_gb_all_0_RENAL_RECOVERY.dat

```
{'criterion': 'mse', 'loss': 'exponential', 'max_depth': 77, 'max_leaf_nodes': 202, 'min_samples_leaf': 68, 'n_estimators': 841}
```

results_gb_not_all_0_RENAL_RECOVERY.dat

```
{'criterion': 'friedman_mse', 'loss': 'deviance', 'max_depth': 5, 'max_leaf_nodes': 569, 'min_samples_leaf': 15, 'n_estimators': 903}
```

results_gb_not_all_0_DIED_14_DAYS.dat

```
{'criterion': 'mse', 'loss': 'exponential', 'max_depth': 120, 'max_leaf_nodes': 362, 'min_samples_leaf': 14, 'n_estimators': 165}
```

results_dt_not_all_0_RENAL_RECOVERY.dat

```
{'criterion': 'gini', 'max_depth': 34, 'max_leaf_nodes': 941, 'min_samples_leaf': 6}
```

results_dt_all_0_RENAL_RECOVERY.dat

```
{'criterion': 'gini', 'max_depth': 50, 'max_leaf_nodes': 965, 'min_samples_leaf': 9}
```

results_dt_all_0_DIED_14_DAYS.dat

```
{'criterion': 'gini', 'max_depth': 142, 'max_leaf_nodes': 522, 'min_samples_leaf': 4}
```

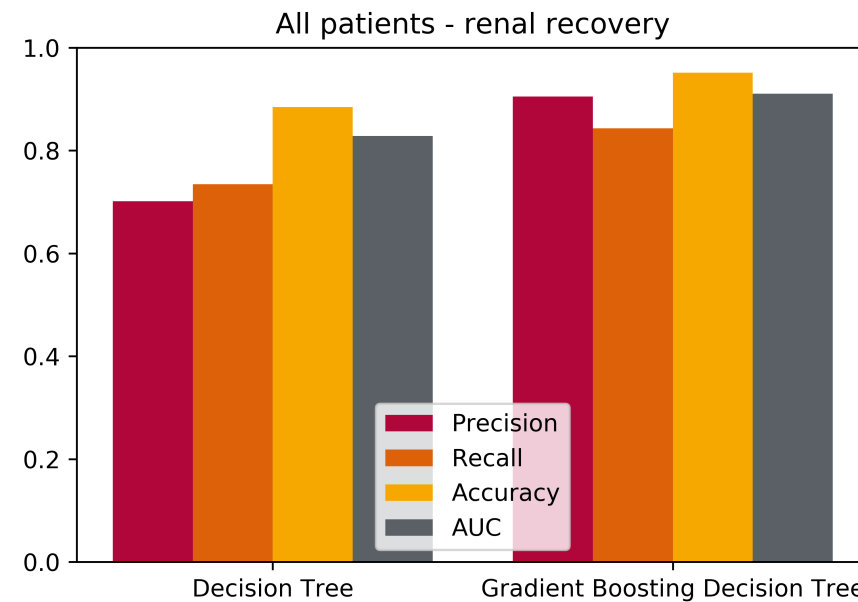
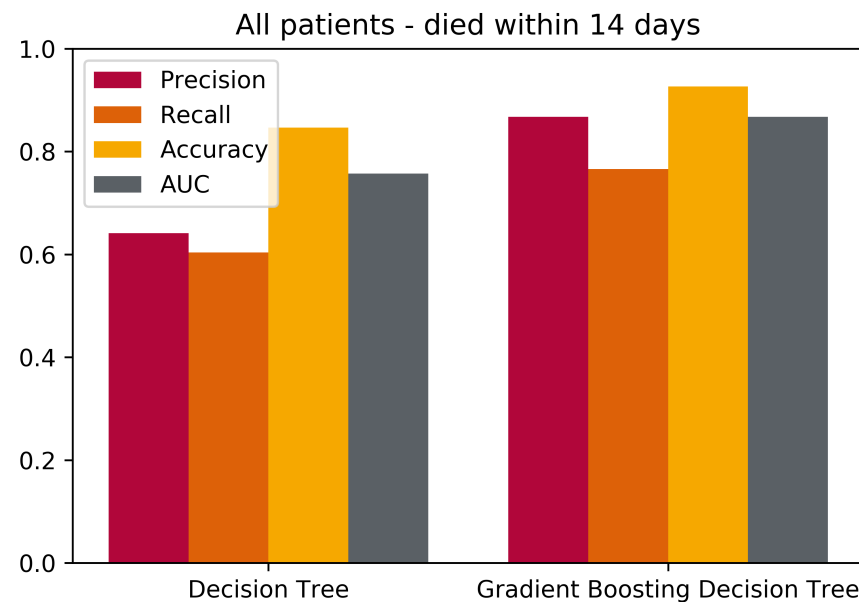
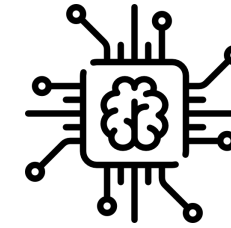
results_dt_not_all_0_DIED_14_DAYS.dat

```
{'criterion': 'gini', 'max_depth': 127, 'max_leaf_nodes': 315, 'min_samples_leaf': 14}
```

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **20**

Building a Clinical Prediction Model: Prediction Patient Outcomes



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **21**

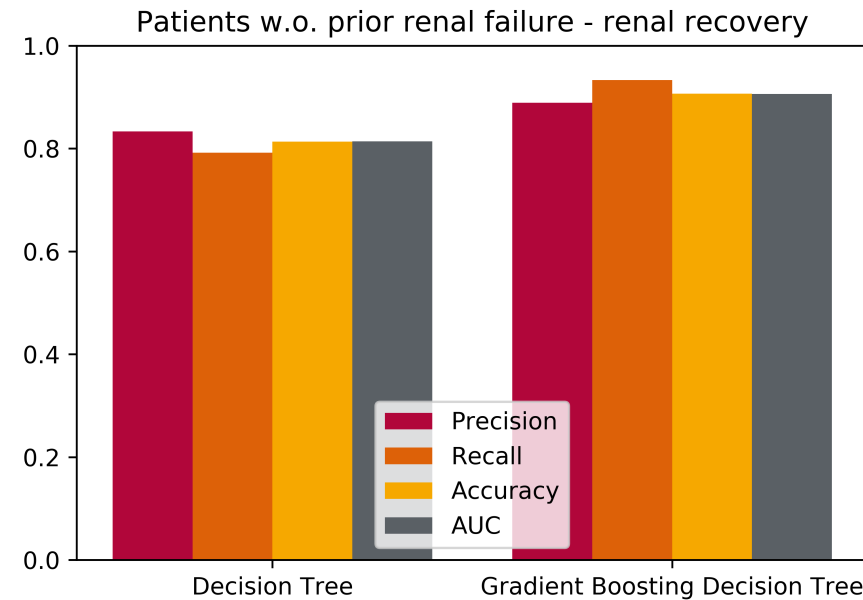
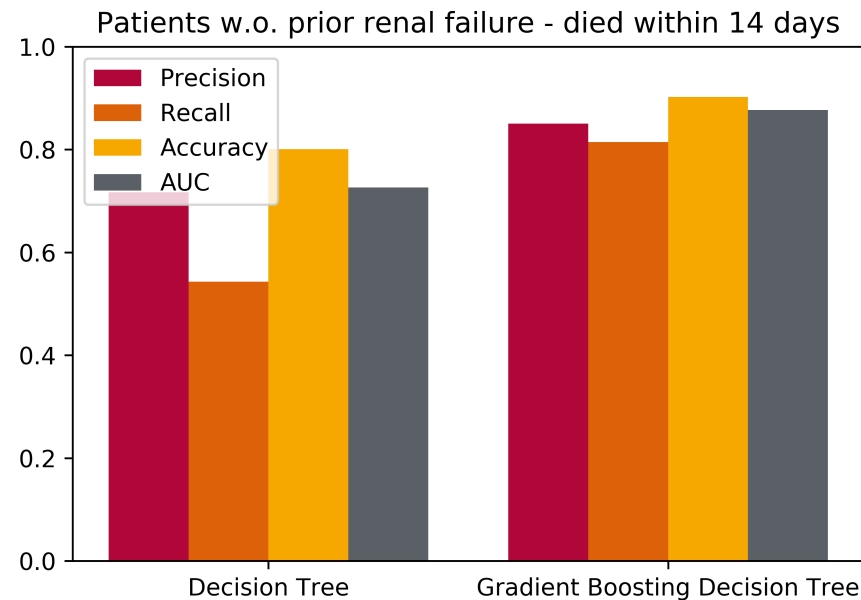
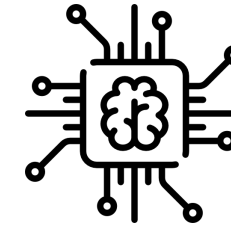
data retrieval

preprocessing

model training

prediction

Building a Clinical Prediction Model: Prediction Patient Outcomes



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 22

data retrieval

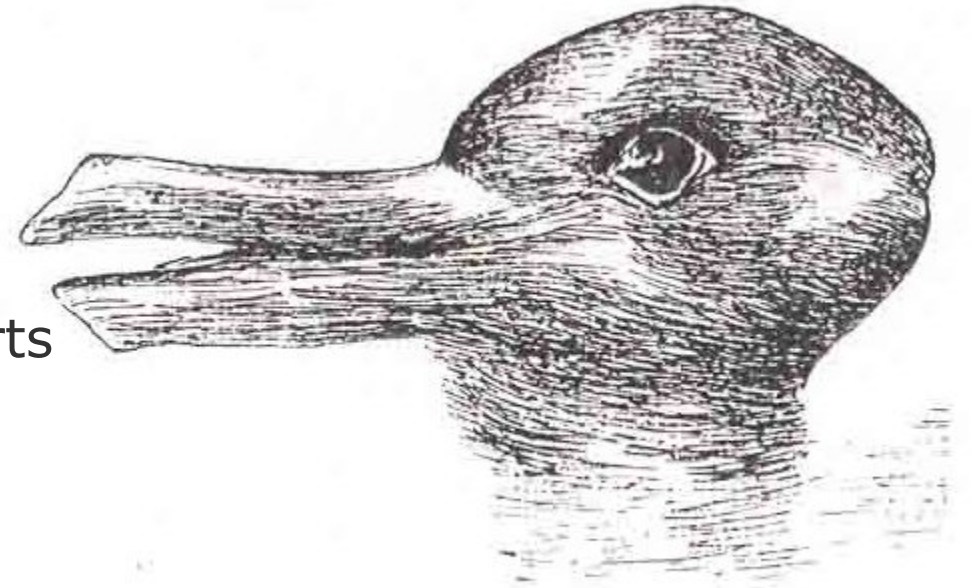
preprocessing

model training

prediction

Methods

- Building a Clinical Prediction Model
- **Applying Interpretability Methods in Detail**
- Making Interpretability Available for Domain Experts



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **24**

Methods: Applying Interpretability Methods in Detail

- Model-based feature importances
- Global Surrogate
- Local Interpretable Model-Agnostic Explanations (LIME)
- Shapley values



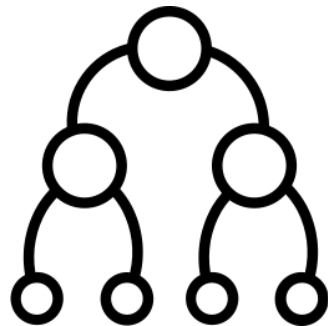
Methods: Applying Interpretability Methods in Detail

- **Model-based feature importances**
- Global Surrogate
- Local Interpretable Model-Agnostic Explanations (LIME)
- Shapley values

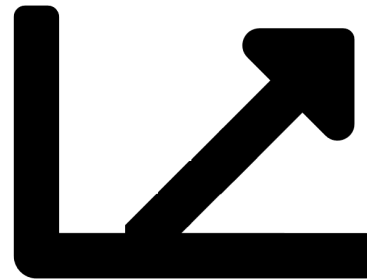


Applying Interpretability Methods in Detail: Model-based Feature Importances

Decision Tree:
= Gini importance



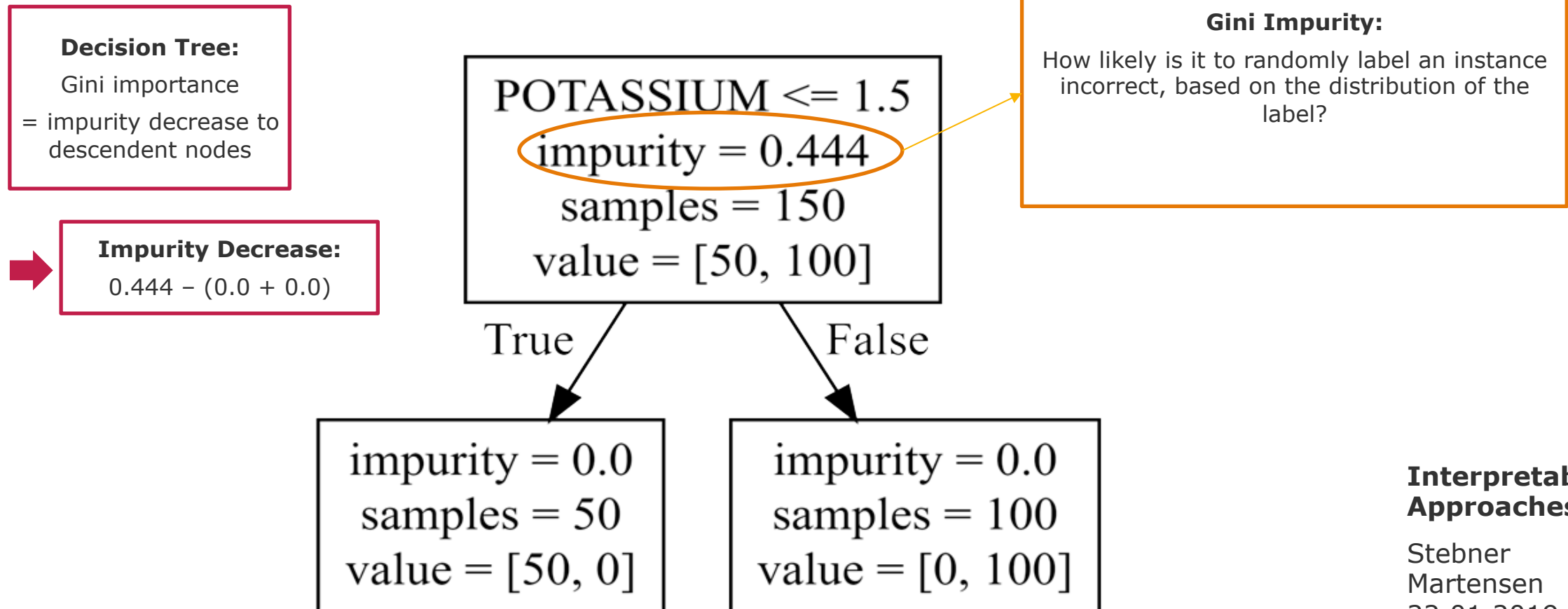
Linear Regression:
Coefficients of linear
function



**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **27**

Applying Interpretability Methods in Detail: Model-based Feature Importances



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **28**

Applying Interpretability Methods in Detail: Model-based Feature Importances

Advantages:

- + Highly compressed, global insight
- + Availability

Disadvantages:

- Faithfulness linked to the error of the model
- Understandability for lay person
- Definition differs per model type

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **29**

Methods: Applying Interpretability Methods in Detail

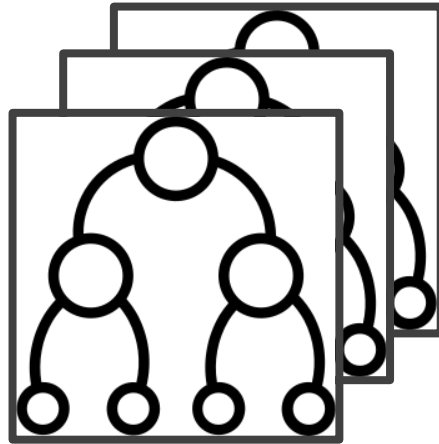
- Model-based feature importances
- **Global Surrogate**
- Local Interpretable Model-Agnostic Explanations (LIME)
- Shapley values



Applying Interpretability Methods in Detail: Global Surrogate

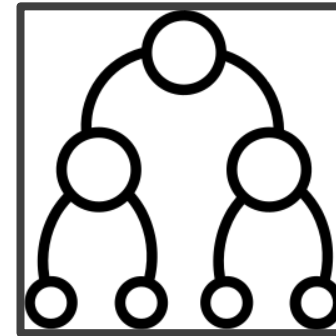
IDEA:

Approximate complicated model output with simpler model



Random forest classifier

Predictions: [0, 1, 0, 1, 1, 0]



Decision Tree (Surrogate)

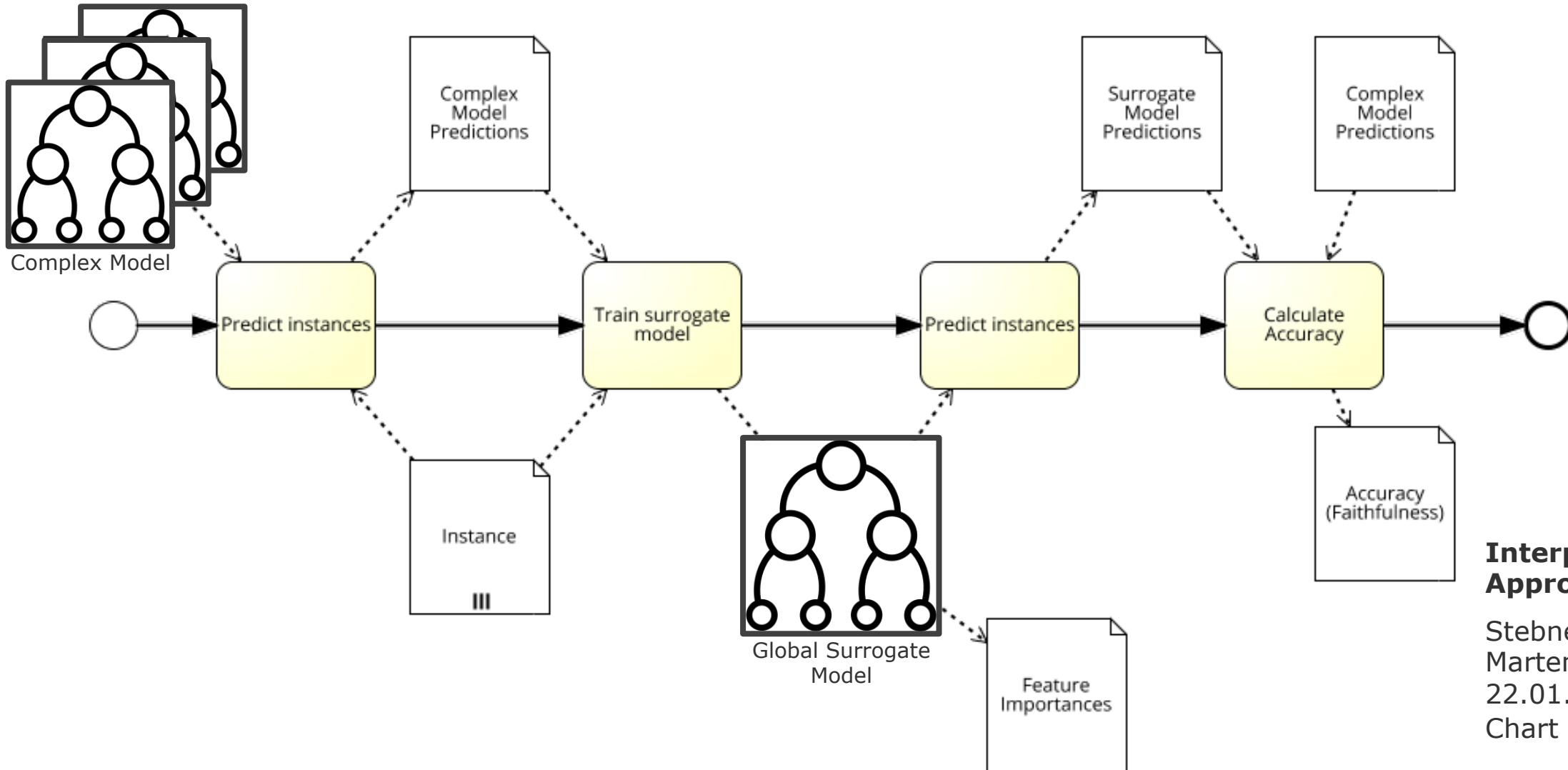
Predictions: [0, 0, 0, 1, 1, 0]

→ 83.33 % accuracy

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **31**

Applying Interpretability Methods in Detail: Global Surrogate



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 32

Applying Interpretability Methods in Detail: Global Surrogate

Advantages:

- + Applicable to any original model (model-agnostic)
- + Surrogate models are “arguably” intuitive
- + Approximation easily measurable

Disadvantages:

- Conclusions about model and not data
- Close for one subset of data, divergent for another?
- Intrinsically interpretable models?

Interpretability Approaches

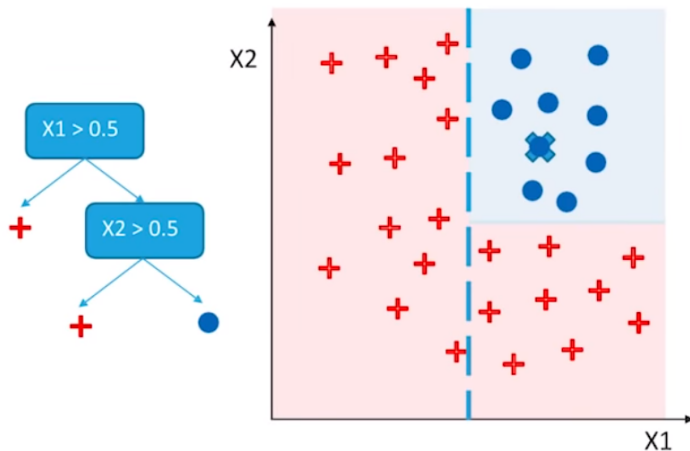
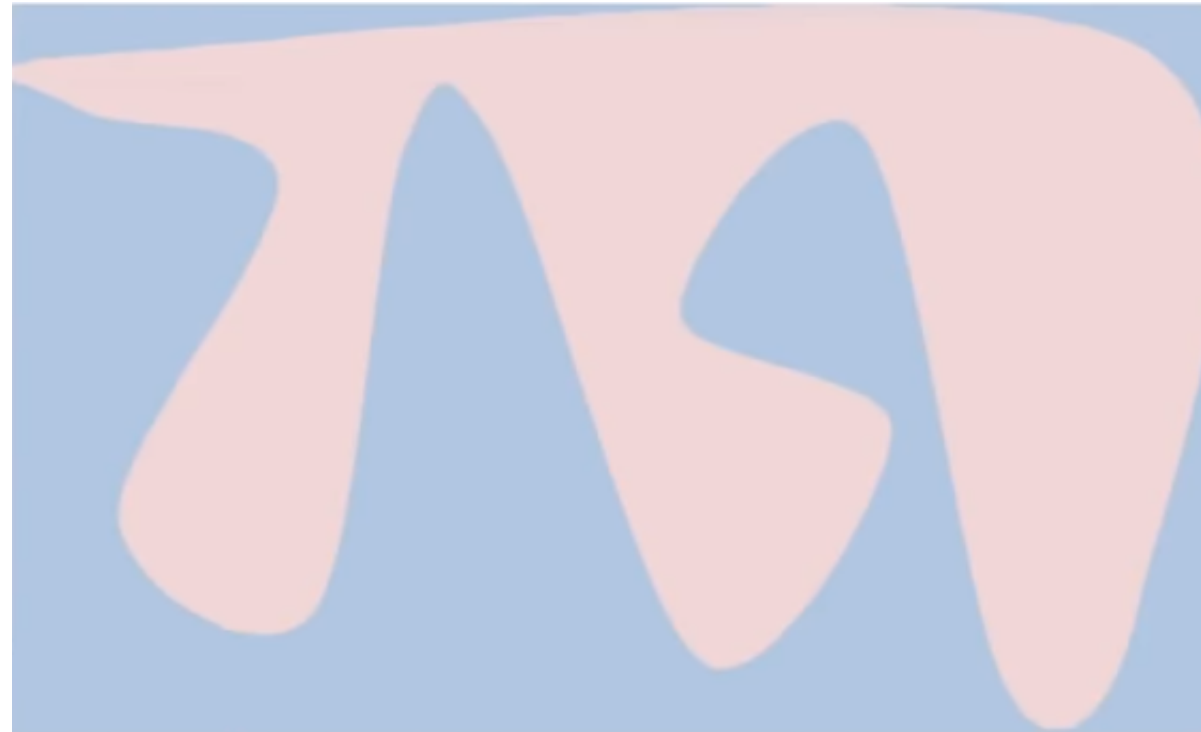
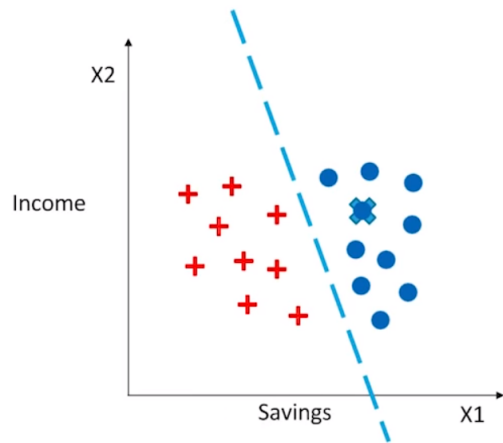
Stebner
Martensen
22.01.2019
Chart **33**

Methods: Applying Interpretability Methods in Detail

- Model-based feature importances
- Global Surrogate
- **Local Interpretable Model-Agnostic Explanations (LIME)**
- Shapley values



Applying Interpretability Methods in Detail: LIME

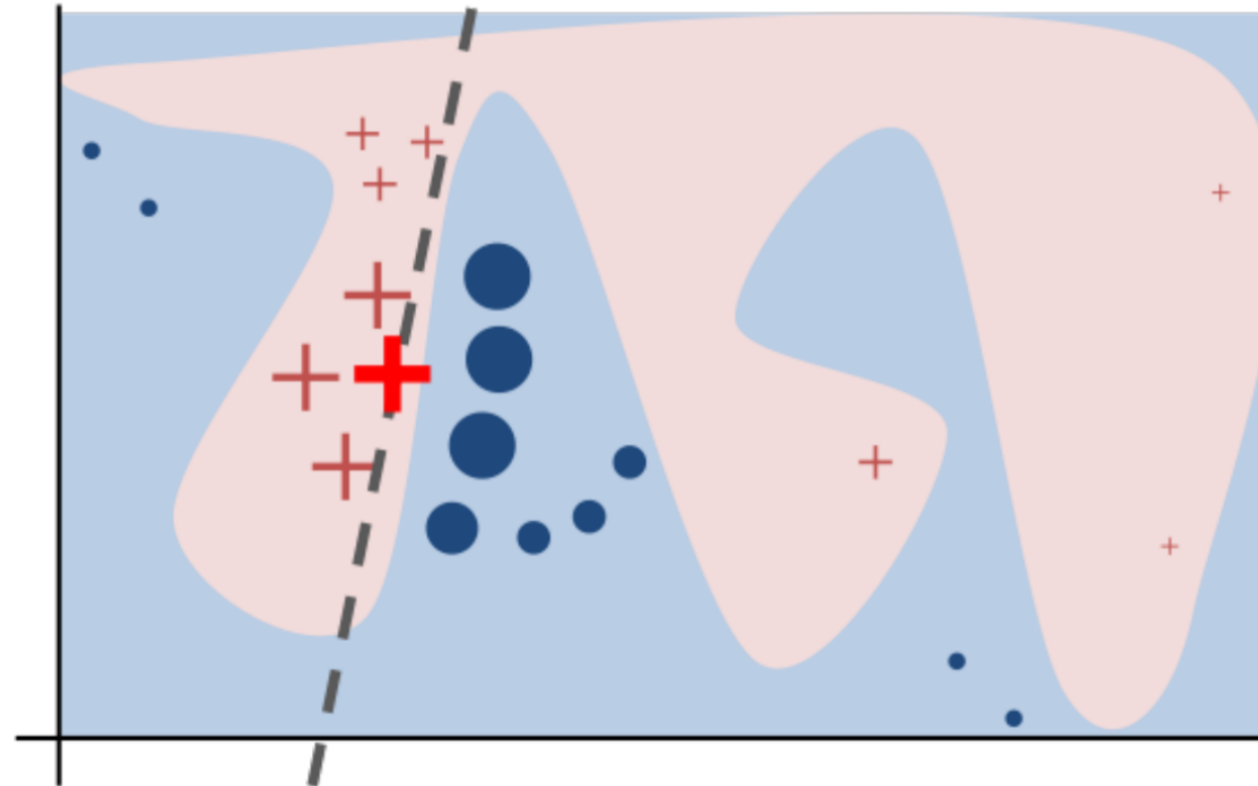


Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **35**

Applying Interpretability Methods in Detail: LIME

1. *Perturbate data*
2. *Compute proximity*
3. *Make predictions*
4. *(Select features)*
5. *Fit a simple model*
6. *Extract explanations*
(feature weights)



**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **36**

- Select a model family and train the model

Fidelity-Interpretability Trade-off

$$\mathcal{L}(f, g, \pi_x)$$

Unfaithfulness of the model

$$\Omega(g)$$

Complexity of the model

$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

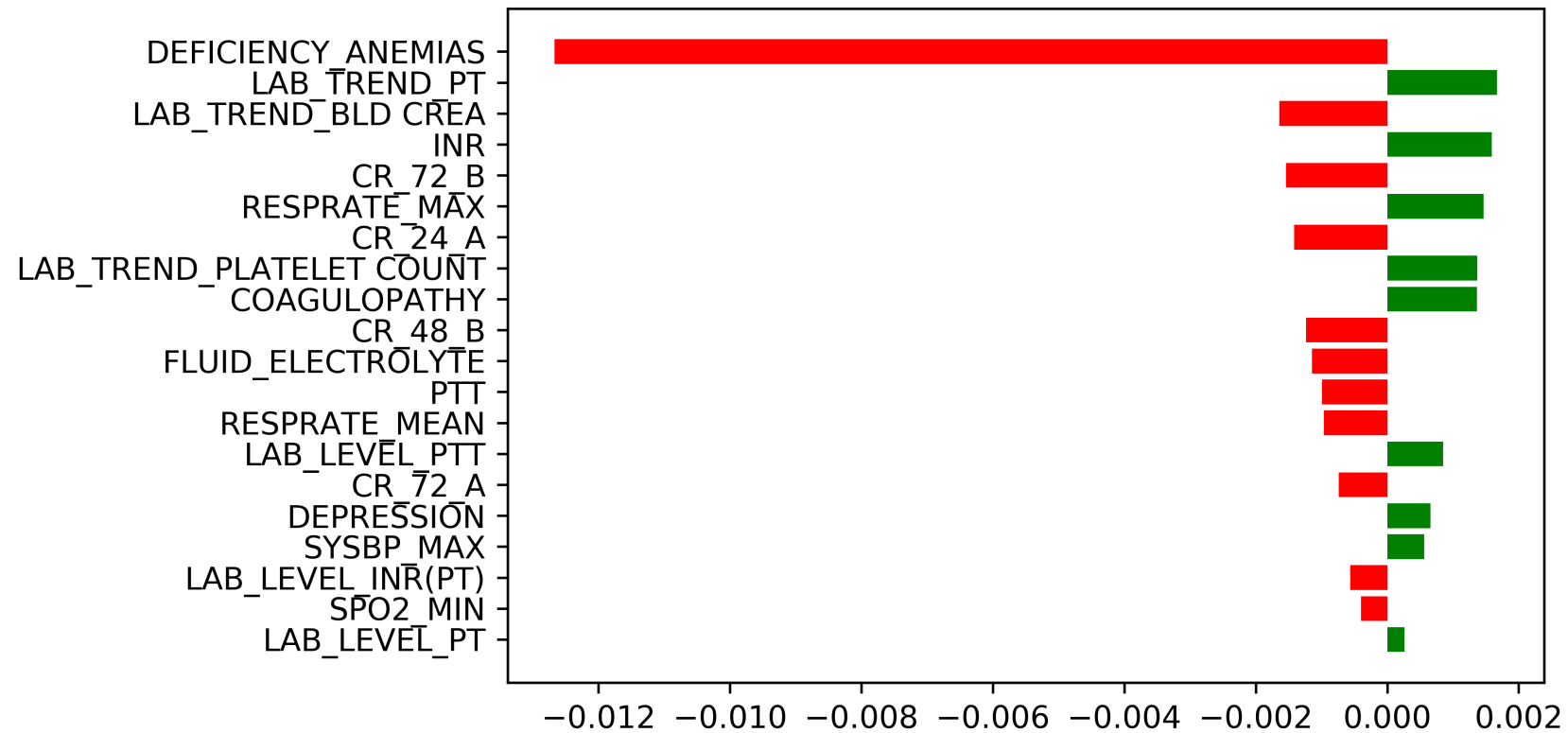
- Extract explanations (e.g. model weights)

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **38**

Applying Interpretability Methods in Detail: LIME

Local explanation for class 1



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 39

Why Submodular Pick?

- LIME is **Local** Interpretable Model Explanations
- Submodular Pick explains model globally by combining local explanations

Parameters:

- # instances (10 percent of dataset)
- # explanations (1 percent of dataset)
- # features (complexity value)

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **40**

Applying Interpretability Methods in Detail: LIME Submodular Pick

1. Select k instances

	f1	f2	f3	f4	f5
Do					
Do					
Do					
Do					
Do					

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **41**

Applying Interpretability Methods in Detail: LIME Submodular Pick

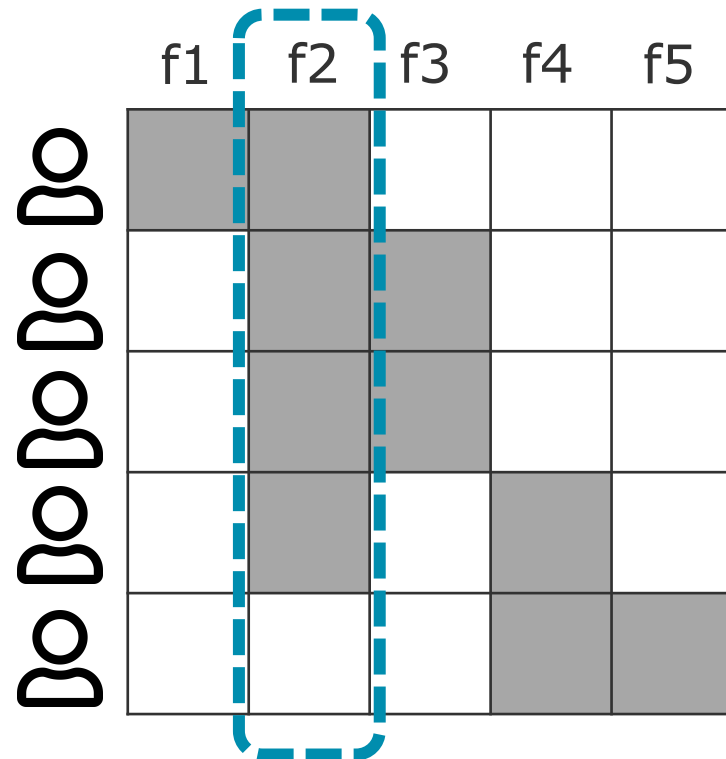
	f1	f2	f3	f4	f5
Do	■	■			
Do		■	■		
Do		■	■		
Do		■		■	
Do				■	■

1. Select k instances
2. Get k local explanations and the important features

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **42**

Applying Interpretability Methods in Detail: LIME Submodular Pick



1. Select k instances
2. Get k local explanations and the important features
3. (f2 has highest importance, because important in 4/5 explanations)

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 43

Applying Interpretability Methods in Detail: LIME Submodular Pick

	f1	f2	f3	f4	f5
Do	Do	Do	Do	Do	Do
Do	Do	Do	Do	Do	Do
Do	Do	Do	Do	Do	Do
Do	Do	Do	Do	Do	Do
Do	Do	Do	Do	Do	Do

1. Select k instances
2. Get k local explanations and the important features
3. (f2 has highest importance, because important in 4/5 explanations)
4. Pick i explanations with highest coverage

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 44

Applying Interpretability Methods in Detail: LIME Submodular Pick

Coverage of an explanation:

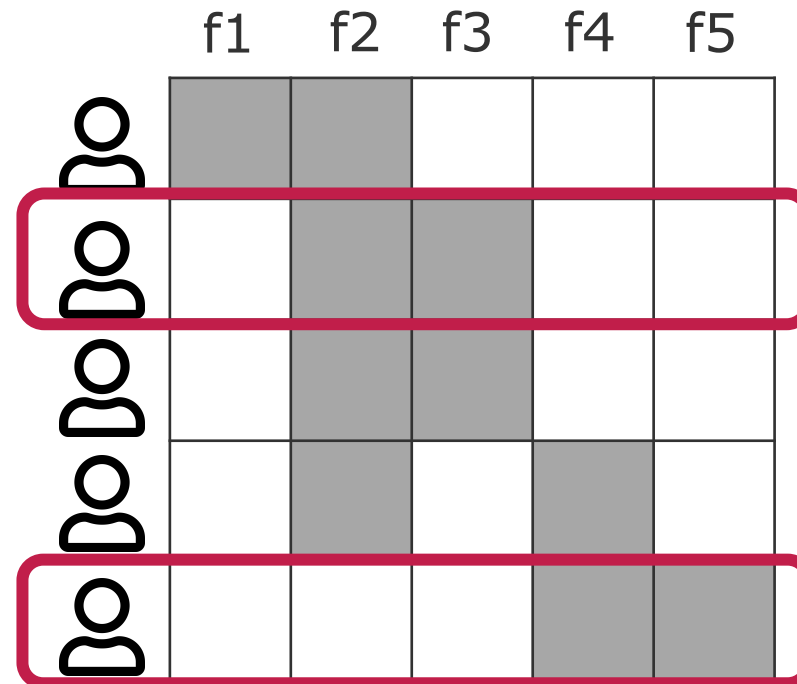
$$c(V, \mathcal{W}, I) = \sum_{j=1}^{d'} \mathbb{1}_{[\exists i \in V: \mathcal{W}_{ij} > 0]} I_j$$

for some set V .

But which V ?

Pick B explanations to
maximize the coverage:

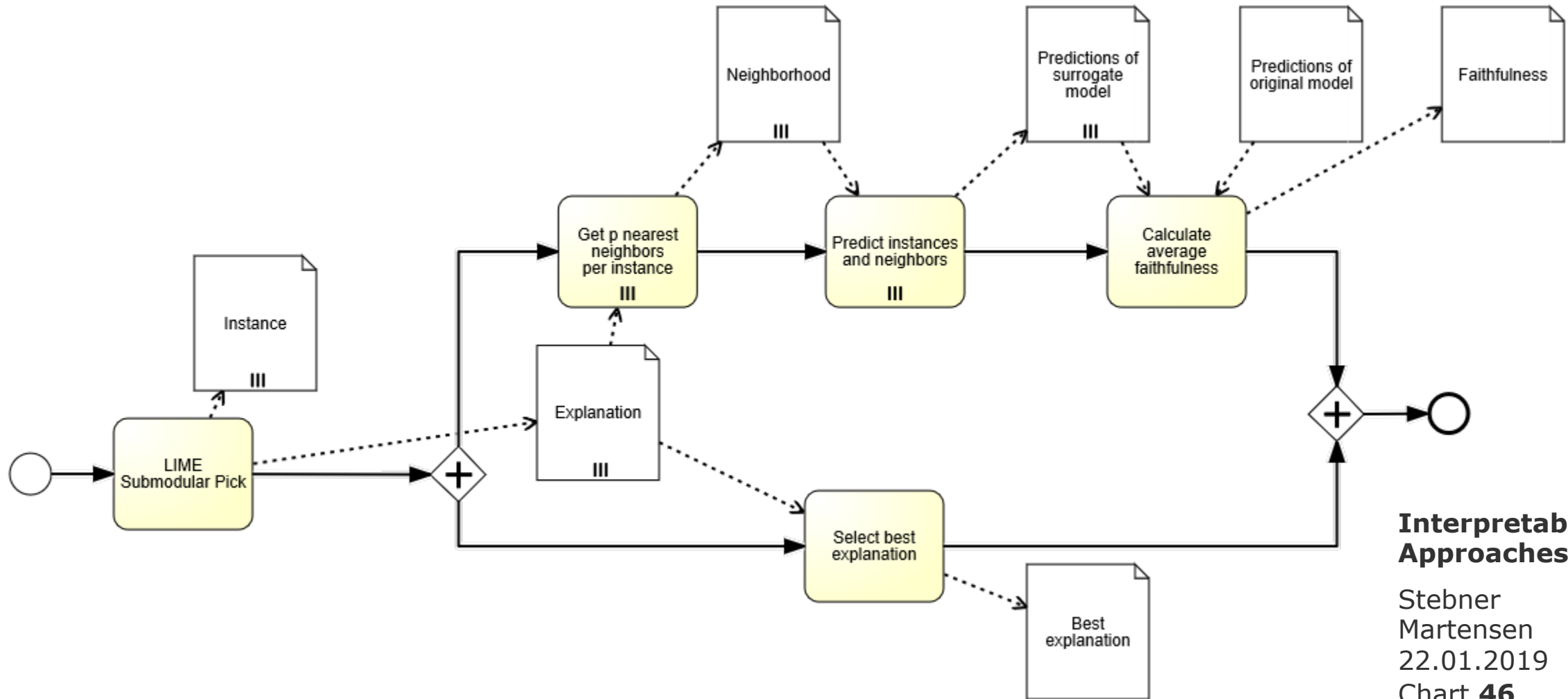
$$Pick(\mathcal{W}, I) = \operatorname{argmax}_{V, |V| \leq B} c(V, \mathcal{W}, I)$$



**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **45**

Applying Interpretability Methods in Detail: LIME Submodular Pick - Evaluation



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 46

Applying Interpretability Methods in Detail: LIME Submodular Pick

Advantages:

- + Not model dependent, based on data!
- + Includes visualization
- + Local and global approach

Disadvantages:

- Requires correct definition of neighborhood
- Submodular pick optimizes coverage, potentially disregards feature interactions
- Instability of model explanations (non-deterministic results)

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **47**

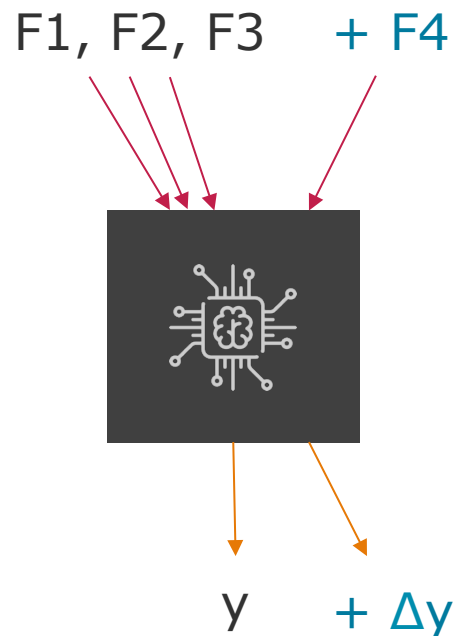
Methods: Applying Interpretability Methods in Detail

- Model-based feature importances
- Global Surrogate
- Local Interpretable Model-Agnostic Explanations (LIME)
- **Shapley values**



Applying Interpretability Methods in Detail: Shapley Values

How much did the feature **contribute** to the model's prediction?



→ Figure out the **marginal contribution** of F4.

$$\varphi_i(x) = f(x_1, \dots, x_n) - E[f(x_1, \dots, X_i, \dots, x_n)]$$

- Example for a simple linear model:

$$f(x_1, \dots, x_n) \approx y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

$$\varphi_i(x) = \beta_i x_i - \beta_i E[X_i]$$

→ care!: it's an additive model
with no feature interactions

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **49**

Applying Interpretability Methods in Detail: Shapley Values

$$\varphi_i(x) = \sum_{Q \subseteq S \setminus \{i\}} \frac{|Q|!(|S| - |Q| - 1)!}{|S|!} (\Delta_{Q \cup \{i\}}(x) - \Delta_Q(x)).$$

What? Some difference.

A Number. Maybe a scaling factor?

Summation over subsets?

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **50**

Applying Interpretability Methods in Detail: Shapley Values

$$\varphi_i(x) = \sum_{Q \subseteq S \setminus \{i\}} \frac{|Q|!(|S| - |Q| - 1)!}{|S|!} (\Delta_{Q \cup \{i\}}(x) - \Delta_Q(x)).$$

F1	F2	F3	F4
X			~
	X		~
		X	~
X	X		~
	X	X	~
X		X	~
X	X	X	~

- S is a set of all features
- Q a subset of S not including i

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **51**

Applying Interpretability Methods in Detail: Shapley Values

$$\varphi_i(x) = \sum_{Q \subseteq S \setminus \{i\}} \frac{|Q|!(|S| - |Q| - 1)!}{|S|!} (\Delta_{Q \cup \{i\}}(x) - \Delta_Q(x)).$$

F1	F2	F3	F4
X	X		~

- S is a set of all features
- Q a subset of S not including i

f1	f2	E[F3]	f4
f1	f2	E[F3]	E[F4]

← Feature values with i

← Feature values without i

$$f_Q(x) = \mathbb{E}[f | X_i = x_i, \forall i \in Q]$$

$$\Delta_Q(x) = f_Q(x) - f_{\emptyset}(x)$$

$$f_{\emptyset}(x) = \mathbb{E}[f]$$

$$\Delta_{Q \cup \{i\}}(x) - \Delta_Q(x)$$

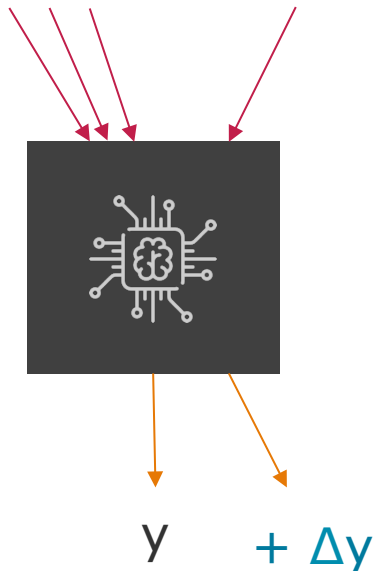
Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **52**

Applying Interpretability Methods in Detail: Shapley Values

$$\varphi_i(x) = \sum_{Q \subseteq S \setminus \{i\}} \underbrace{\frac{|Q|!(|S| - |Q| - 1)!}{|S|!}}_{\text{weight}} (\Delta_{Q \cup \{i\}}(x) - \Delta_Q(x)).$$

F1, F2, F3 + F4 + F6, F5, F7, ...



- $|Q|!$ -many possible rearrangements
- $(|S| - |Q| - 1)!$ -many possibilities to arrange features following i

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **53**

Applying Interpretability Methods in Detail: Shapley Values

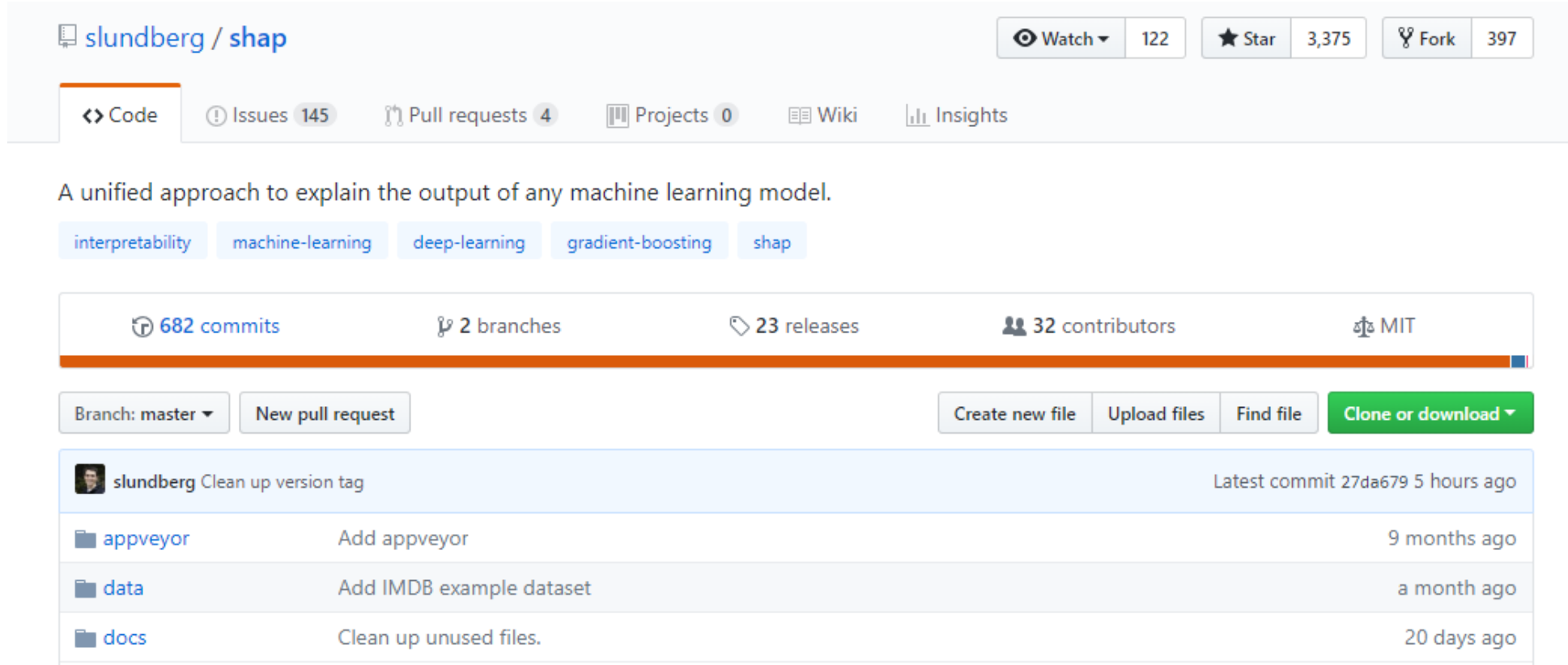
Some unique properties:

- **Efficiency**
 - Contributions add up to the difference of prediction and expectation
- **Symmetry**
 - Same value for same contributions
- **Dummy Feature**
 - Non-contributing features have value 0
- **Additivity**
 - Multi-model predictions (e.g. random forest) can be analyzed

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **54**

Applying Interpretability Methods in Detail: SHAP



slundberg / shap

Watch 122 Star 3,375 Fork 397





Code Issues 145 Pull requests 4 Projects 0 Wiki Insights

A unified approach to explain the output of any machine learning model.

interpretability machine-learning deep-learning gradient-boosting shap

682 commits 2 branches 23 releases 32 contributors MIT

Branch: master New pull request Create new file Upload files Find file Clone or download

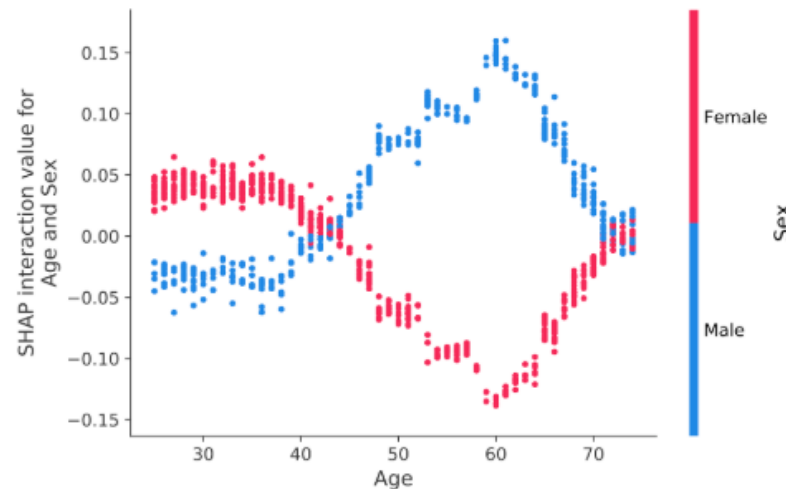
 slundberg	Clean up version tag	Latest commit 27da679 5 hours ago
 appveyor	Add appveyor	9 months ago
 data	Add IMDB example dataset	a month ago
 docs	Clean up unused files.	20 days ago

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **55**

Applying Interpretability Methods in Detail: SHAP

Customer A:



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 56

Applying Interpretability Methods in Detail: Shapley values and SHAP

Advantages:

- + Contrastive explanations (with respect to the expectation)
- + Applicable for whole dataset, subset or single instance
- + Solid foundation from game theory

Disadvantages:

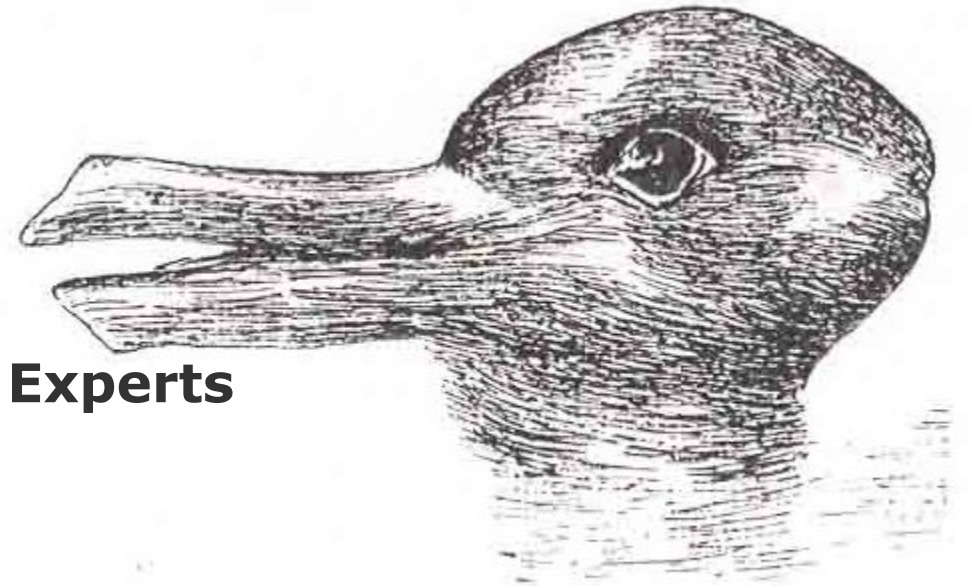
- Exponential computational complexity
- Always returns all features
- No prediction model

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **57**

Methods

- Building a Clinical Prediction Model
- Applying Interpretability Methods in Detail
- **Making Interpretability Available for Domain Experts**



**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **58**

Methods: Making Interpretability Available for Domain Experts

Requirements:

- Compare different Interpretability Method outputs for one CPM
- Rank interpretability models
- Faithfulness-Complexity tradeoff

Visualizations:

- Feature Importances
- Complexity-Faithfulness-Graph

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **59**

Results

- Feature Importances
- Complexity-Faithfulness-Graph
- Clinical Hypotheses



Results: Feature Importances

Feature	Model-based A.	LIME	Linear Surrogate Model	Tree Surrogate Model	SHAP
Age					
Platelets					
Blood Gas					
...					

- Comparing interpretability methods output for every feature
- Filter, sort, threshold, ... operations
- (Weighted) average

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **61**

Results: Feature Importances

- Feature Importances ordered by maximal importance

Feature	Model-based importances	LIME	Linear Surrogate Model	Tree Surrogate Model	SHAP
Lab Flag PT			0.4024		
Lab Flag INR(PT)			-0.3983		
Deficiency Anemias			0.2606		
AIDS			-0.2547		
Lab Level Hematocrit (Calculated)			0.2515		
GFR_72	0.1618	0.0607		0.2127	0.1618
Lab Flag Bilirubin	0.0678			0.1610	
CR_72	0.0440	0.0397		0.1530	0.0440
Lactate	0.4508			0.1050	0.4508

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **62**

Results: Feature Importances

- Feature Importances ordered by occurrences (if occurred more than once)

Feature	Model-based importances	LIME	Linear Surrogate Model	Tree Surrogate Model	SHAP
GFR_72	0.1618	0.0607		0.2127	0.1618
CR_72	0.0440	0.0397		0.1529	0.0440
Lactate	0.0451			0.1050	0.0451
Lab Flag Bilirubin	0.0678			0.1610	
Bicarbonate	0.0295			0.0192	

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **64**

Results: Complexity-Faithfulness-Graph

Complexity – Faithfulness – Tradeoff:

Complexity \sim Faithfulness
Complexity \sim 1 / Interpretability



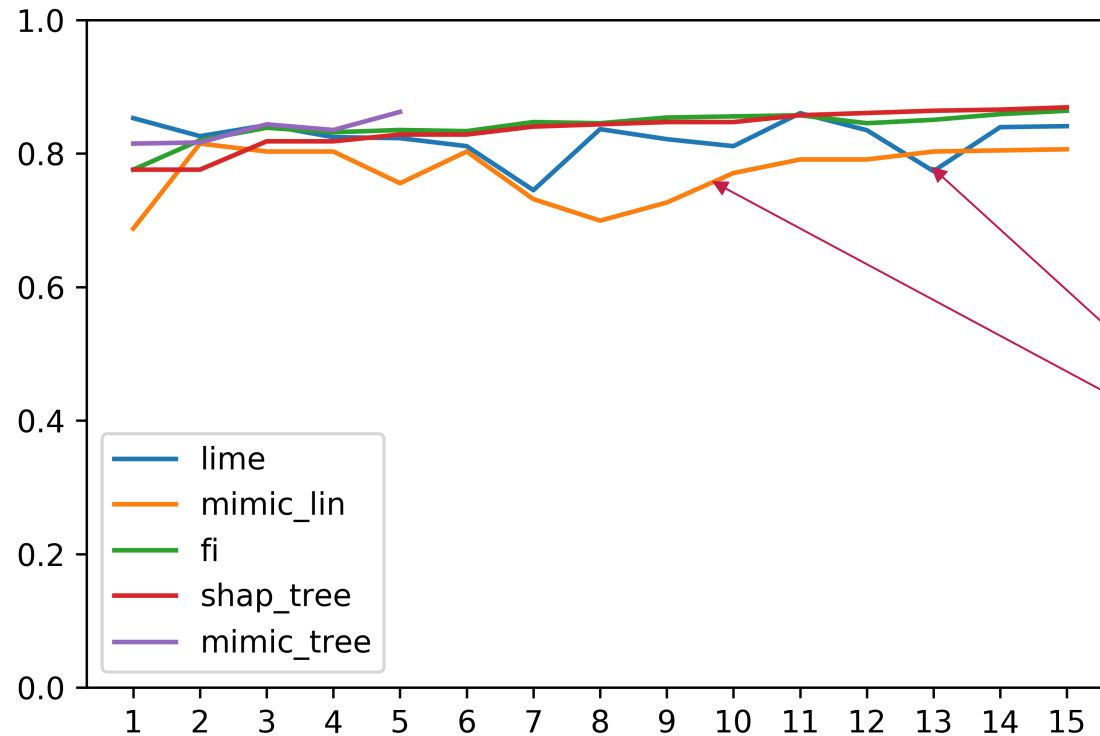
Increased complexity \rightarrow increase in faithfulness
Increased complexity \rightarrow decrease in interpretability

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **65**

Results: Complexity-Faithfulness-Graph

Complexity – Faithfulness – Tradeoff:



**Why does this not
increase
monotonically?**

→ Maybe showing
incompetence of linear
models for complex
relations

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **66**

Results:

Tentative Clinical Hypotheses

Glomerular Filtration Rate 72h before procedure:

- Flow rate of filtered fluid through the kidney
- Known as indicator of kidney function

Creatinine Clearance Rate 72h before procedure:

- Volume of blood plasma cleared of creatinine per unit time
- AKI is defined as increase of CR over baseline

Do these patients have higher chances of survival/recovery because their AKI is detected earlier?

Bilirubin:

- Product of breakdown of red blood cells

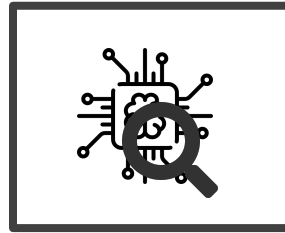
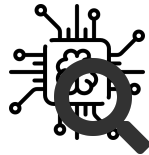
Platelets / Thrombocytes:

- First responders to sites of damages in the body

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **67**

Outlook: Next Steps



Python Notebooks for
Interpretability

Interpretability as a Service

Clinical Predictive Model

Clinical Hypotheses

**Evaluation of Clinical
Hypotheses with Charité**

**Interpretability
Approaches**

Stebner
Martensen
22.01.2019
Chart **68**

Visions & Objectives: Contribution

VISION 1

Find and validate medical hypotheses regarding mortality and recovery of AKI

- ✓ Train CPM
- ✓ Predict patient outcomes
- ✓ Gather interpretations
- ✓ Derive and evaluate clinical hypotheses

VISION 2

Make interpretations of CPMs available to physicians

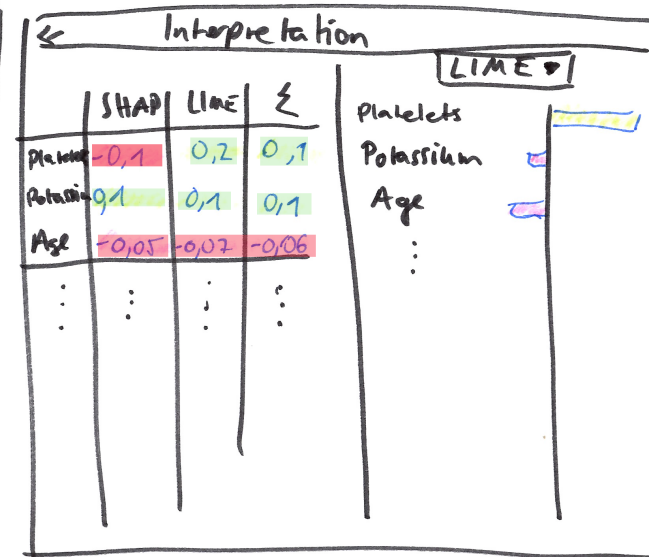
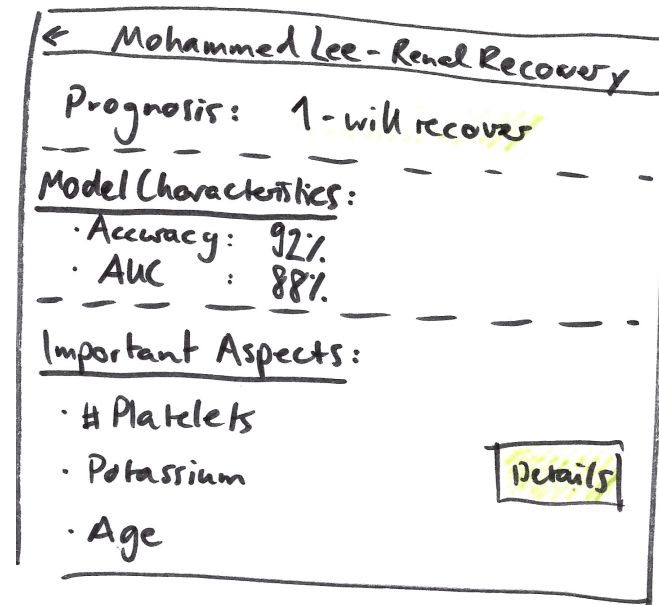
- ✓ Interpret any CPM
- ✓ Make interpretations comparable side-by-side
- ✓ Show complexity-faithfulness tradeoff

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 69

Outlook: Future Work

- Evaluate Approach with different cohorts (Heidelberg database, different disease)
- Patient Predictor and Diagnosis Explainer (UI)



Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart 70

Sources

Duck-Rabbit-Illusion: https://en.wikipedia.org/wiki/Ambiguous_image#/media/File:Duck-Rabbit_illusion.jpg

Sherlock: <https://images.fineartamerica.com/images-medium-large-5/sherlock-holmes-c1905-granger.jpg>

Cardiopulmonary Bypass:

https://upload.wikimedia.org/wikipedia/commons/thumb/2/24/Blausen_0468_Heart-Lung_Machine.png/300px-Blausen_0468_Heart-Lung_Machine.png

Injured Kidney:

https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcQ4kVzdKHZ81KazmyE9YXLQvvqp9iF00PI56PfPI0MOV_Fxorw1aA

Error Plane:

<https://image.slidesharecdn.com/navdeepmlinov0117-171102184007/95/ideas-on-machine-learning-interpretability-9-638.jpg?cb=1509648095>

Icons by Fontawesome (<https://fontawesome.com/license>) and by Freepik, Appzgear, Pixel perfect , phatplus & Eucalyp on <https://flaticon.com>

LIME:

<https://www.slideshare.net/0xdata/interpretable-machine-learning-using-lime-framework-kasia-kulma-phd-data-scientist>

Feature Importances from sci-kit learn: <https://github.com/scikit-learn/scikit-learn/blob/master/sklearn/tree/tree.py>

LIME Paper: Ribeiro et al. "Why Should I Trust You?" Explaining the Predictions of Any Classifier (ACL Proceedings 2016)

Interpretable Method (Dis-)Advantages: Molnar, C. (2018). Interpretable Machine Learning. Retrieved from <https://christophm.github.io/interpretable-ml-book/>

Evaluating Interpretability: Explaining Explanations: An Approach to Evaluating Interpretability of Machine Learning. Leilani Gilpin et al. Massachusetts Institute of Technology. 2018.



Interpretability Approaches applied to Clinical Predictive Modeling

Trends in Bioinformatics
Intermediate Presentation

Discussion

Possible Questions:

- How should we normalize the importances, so that they are actually comparable?
- As a patient, in how much level of detail would you expect your doctor to explain Machine Learning results?
- As a physician, how do you want to be trained for interpretable models?

Interpretability Approaches

Stebner
Martensen
22.01.2019
Chart **73**