

Challenges in Explainable Artificial Intelligence for Industry 4.0

Szymon Bobek

AIRA seminar
13 January 2022



<https://geist.re>



The presentation covers several different works founded from the PACMEL (NCN 2018/27/Z/ST6/03392) and XPM (NCN UMO-2020/02/Y/ST6/00070) projects funded by the National Science Centre, Poland under CHIST-ERA program .

GEIST (<https://geist.re>)

Group for Engineering of Intelligent Systems and Technologies

Welcome to GEIST Research Group Webpage!



GEIST is a research group that includes [number of senior and postdocoral researchers](#) as well as PhD and master students. The

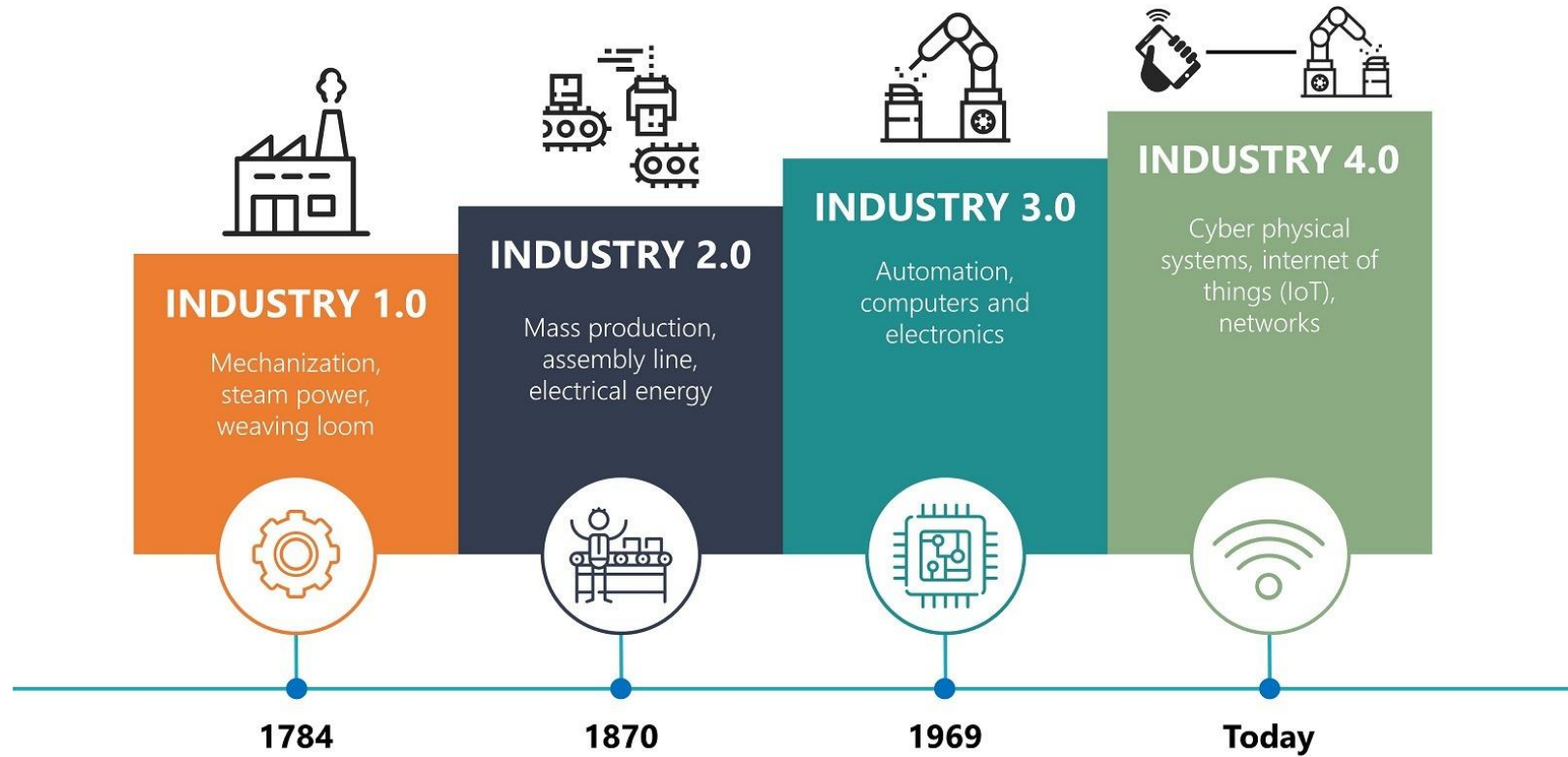
work of the group is coordinated by [Grzegorz J. Nalepa](#). GEIST researchers mostly work at the [Jagiellonian University \(UJ.edu.pl\)](#) as well as the [AGH University of Science and Technology \(AGH.edu.pl\)](#) in Kraków, Poland.

The group is active in the general area of intelligent systems. We work in Explainable AI (XAI), Knowledge and Software Engineering (KE/SE), Business Intelligence (BI), Ambient Intelligence (Aml), and Affective Computing (AfC), (see the group's [research profile](#)).

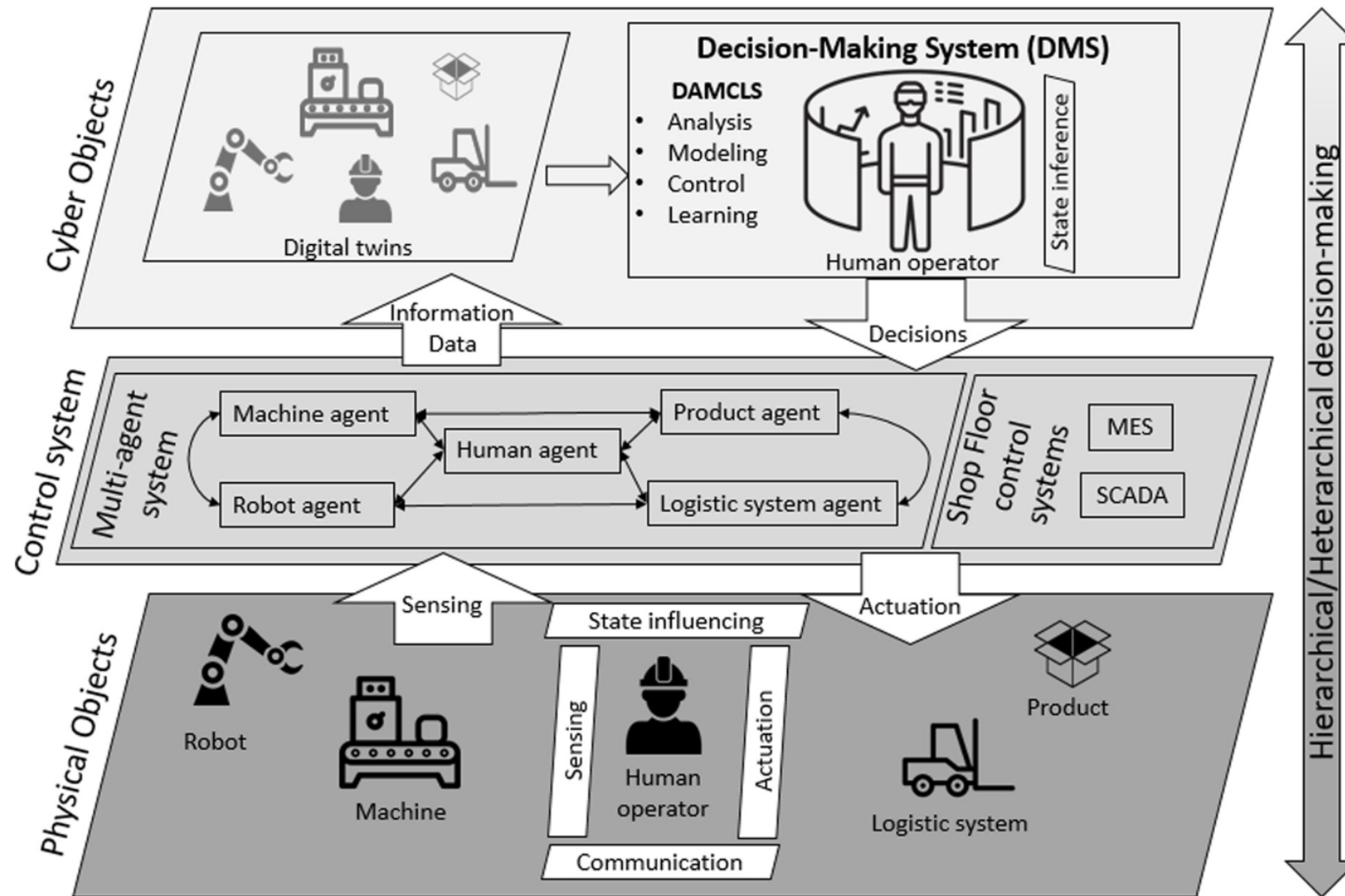
GEIST has been involved in number of [projects](#). For more information see [recent activity](#), [publications](#) and [software](#).



Industry 4.0



Let us have an Industry 4.0 factory!



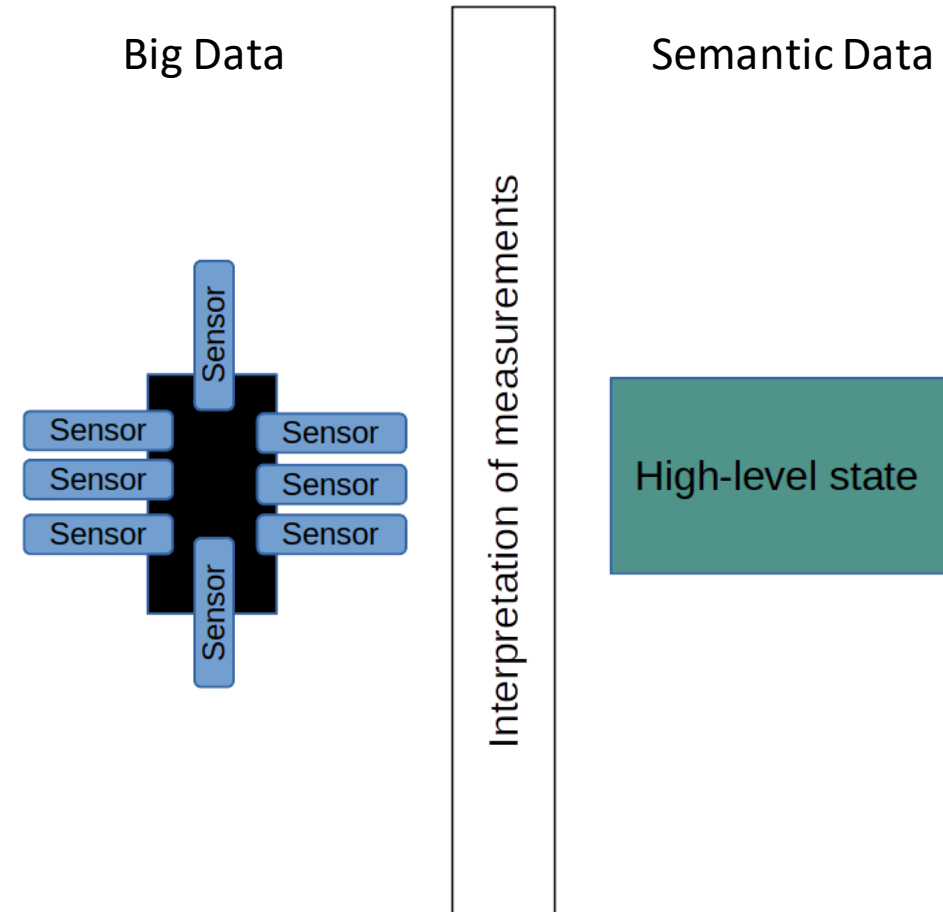
Source: Cimini, C.; Pirola, F.; Pinto, R.; Cavalieri, S. A human-in-the-loop manufacturing control architecture for the next generation of production systems. *Journal of Manufacturing Systems* 2020, 54, 258–271. doi:<https://doi.org/10.1016/j.jmsy.2020.01.002>.





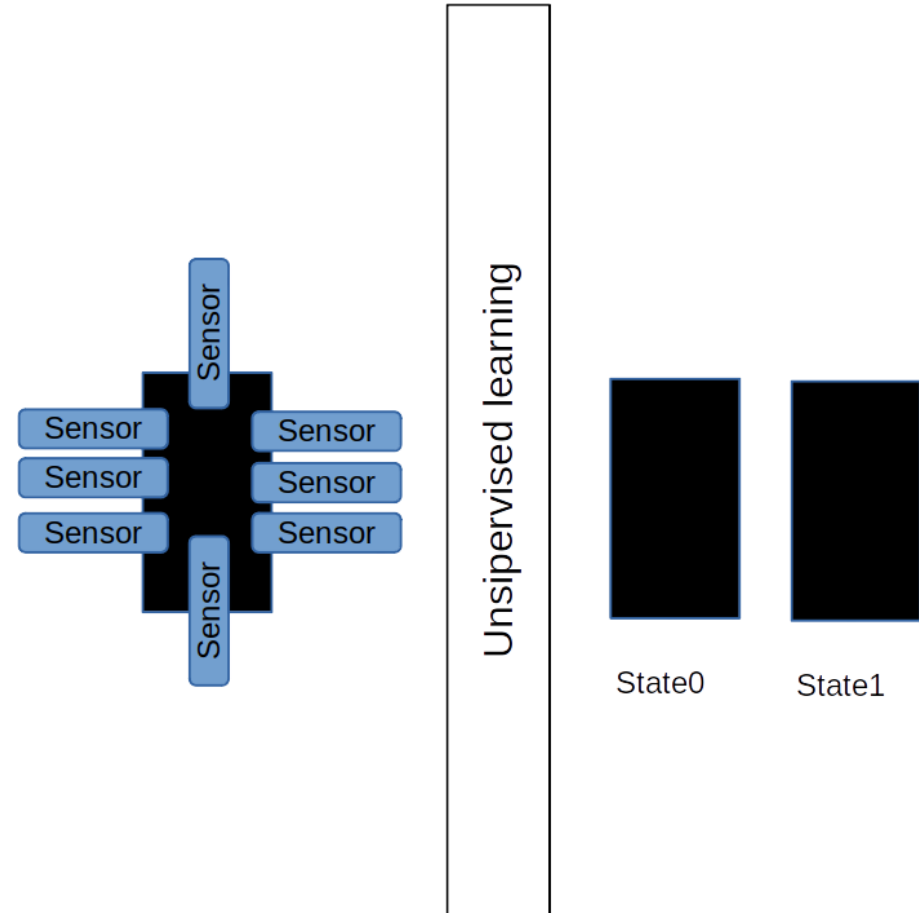
Knowledge discovery from industrial logs

- PACMEL Project (<http://pacmel.geist.re>)
- Industry 4.0: everything is measured
- Low-level measurements to higher-level states
- Expert-defined states vs. Automatically discovered states
 - Data comes with no labels
 - Expert may be too general, or too specific
 - There is a lot of measurements



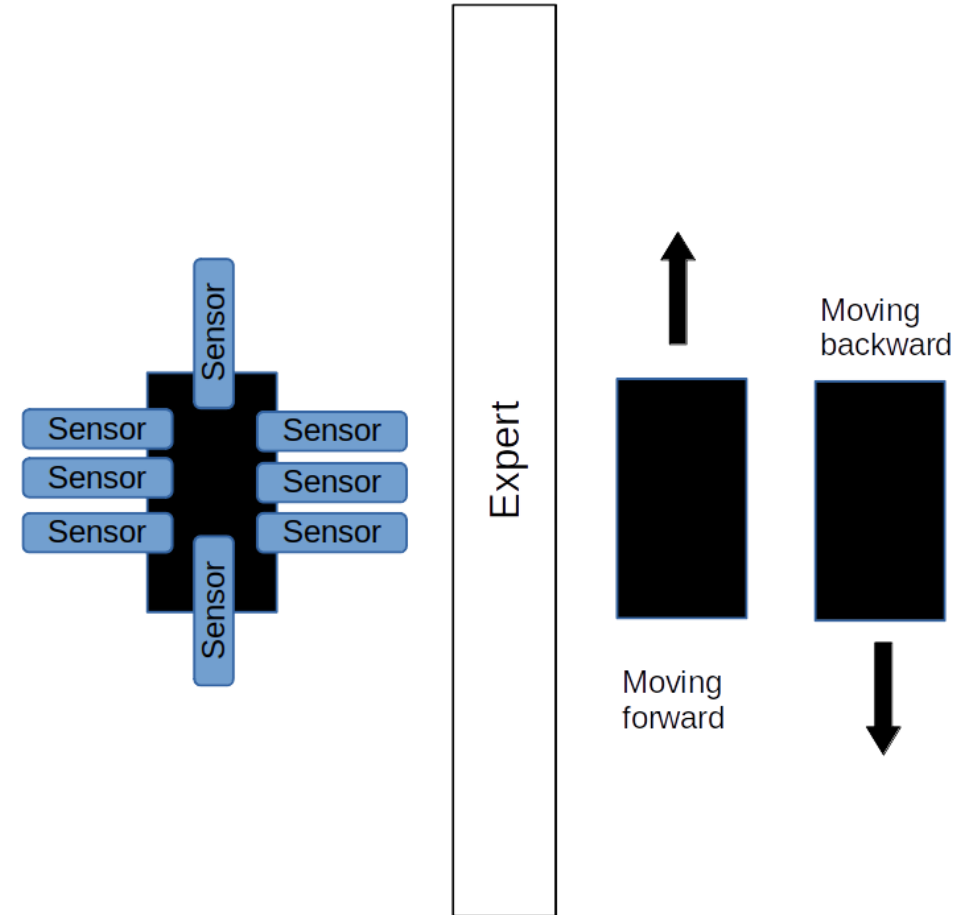
Knowledge discovery from industrial logs

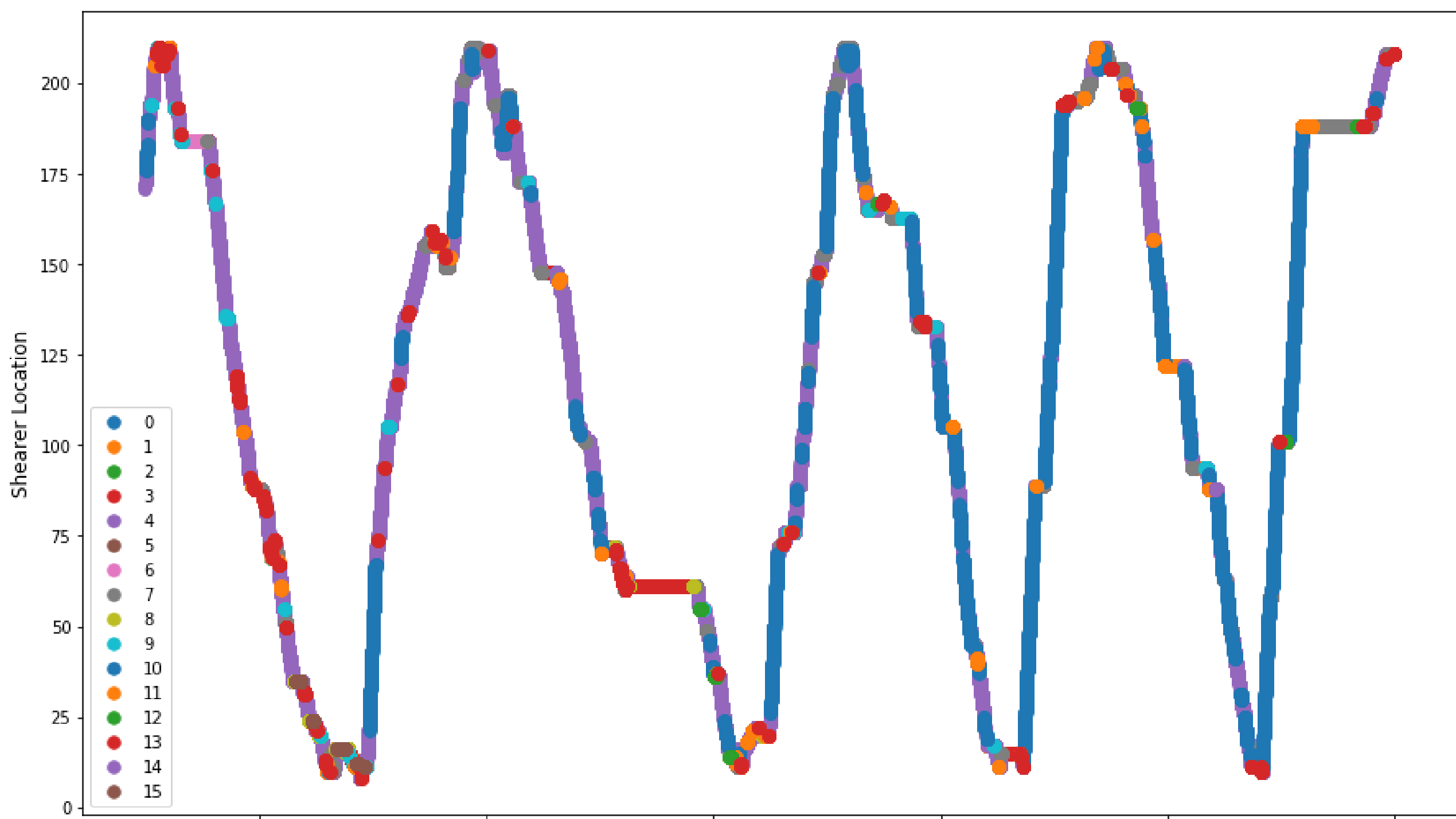
- PACMEL Project (<http://pacmel.geist.re>)
- Industry 4.0: everything is measured
- Low-level measurements to higher-level states
- Expert-defined states vs. Automatically discovered states
 - Data comes with no labels
 - Expert may be too general, or too specific
 - There is a lot of measurements



Knowledge discovery from industrial logs

- PACMEL Project (<http://pacmel.geist.re>)
- Industry 4.0: everything is measured
- Low-level measurements to higher-level states
- Expert-defined states vs. Automatically discovered states
 - Data comes with no labels
 - Expert may be too general, or too specific
 - There is a lot of measurements

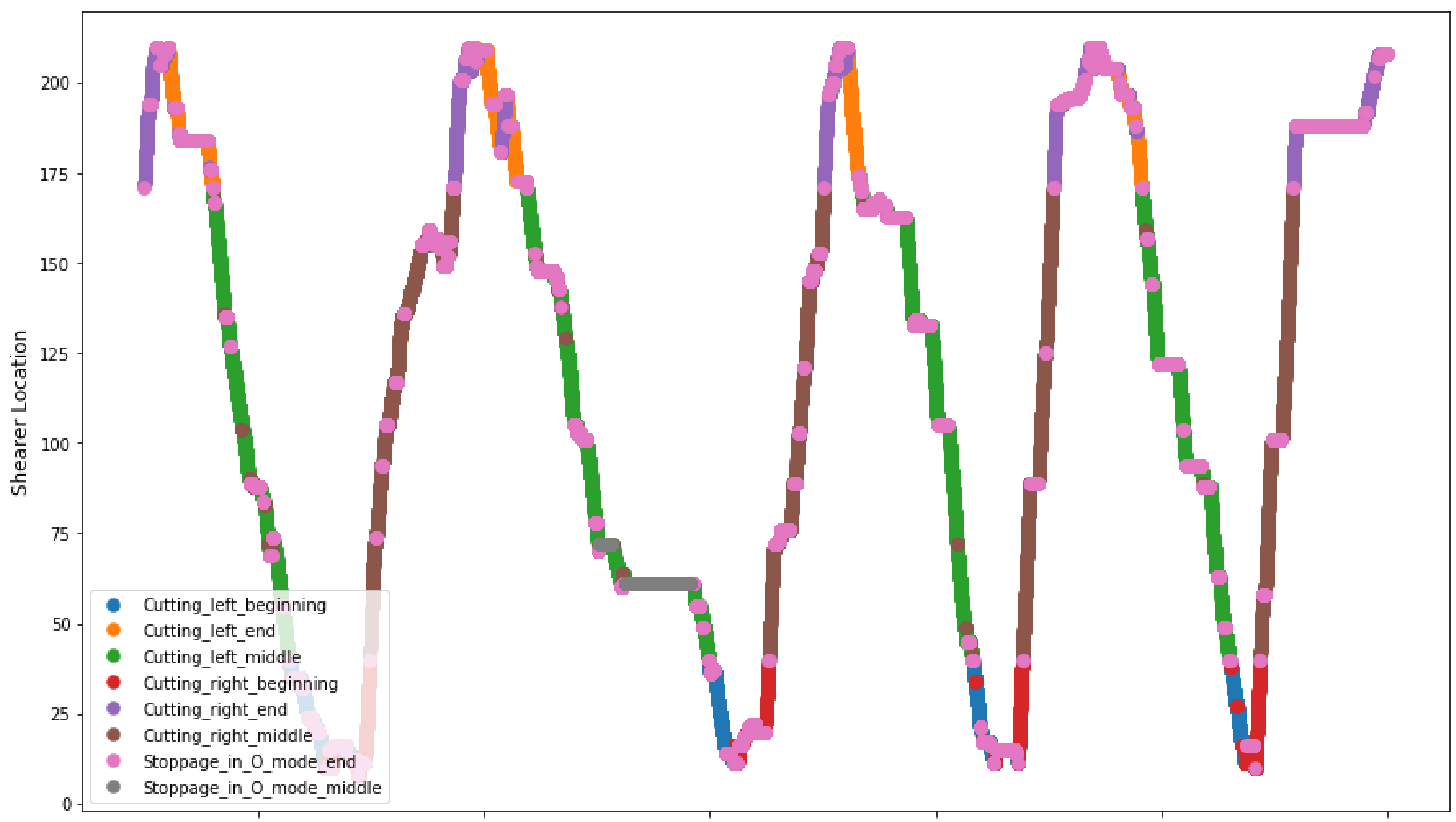




Theoretical states are given by the expert

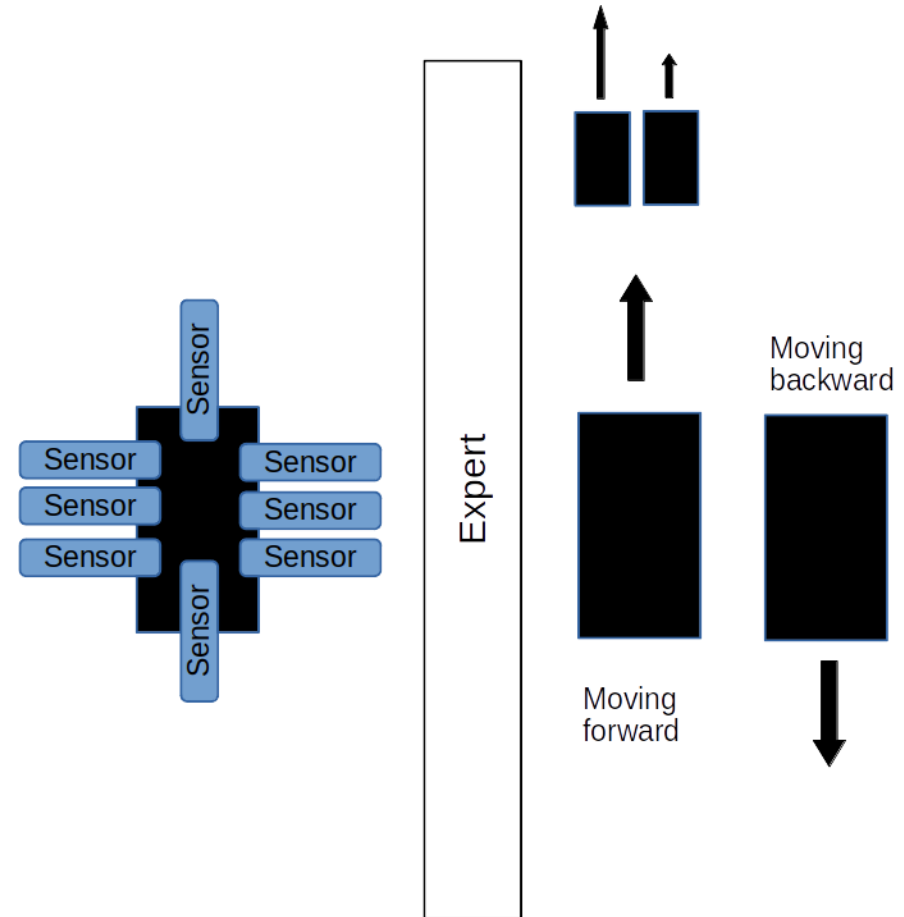
					State	#
= 0	> 0	= 0	= 1	= any	set movingLeft	1
= 0	> 0	= 0	= 0	= any	set movingRight	1
> 0	> 0	= any	= 1	< 40	set cuttingLeftBeginning	1
> 0	> 0	= any	= 1	∈ [40 .. 180]	set cuttingLeftMiddle	1
> 0	> 0	= any	= 1	>= 180	set cuttingLeftEnd	1
> 0	> 0	= any	= 0	< 40	set cuttingRightBeginning	1
> 0	> 0	= any	= 0	∈ [40 .. 180]	set cuttingRightMiddle	1
> 0	> 0	= any	= 0	> 180	set cuttingRightEnd	1
= any	= any	= any	≠ any	< 40	set stoppageInOModeBeginn...	1
= any	= any	= any	≠ any	< 180	set stoppageInOModeMiddle	1
= any	= any	= any	≠ any	>= 180	set stoppageInOModeEnd	1

expertLabeling



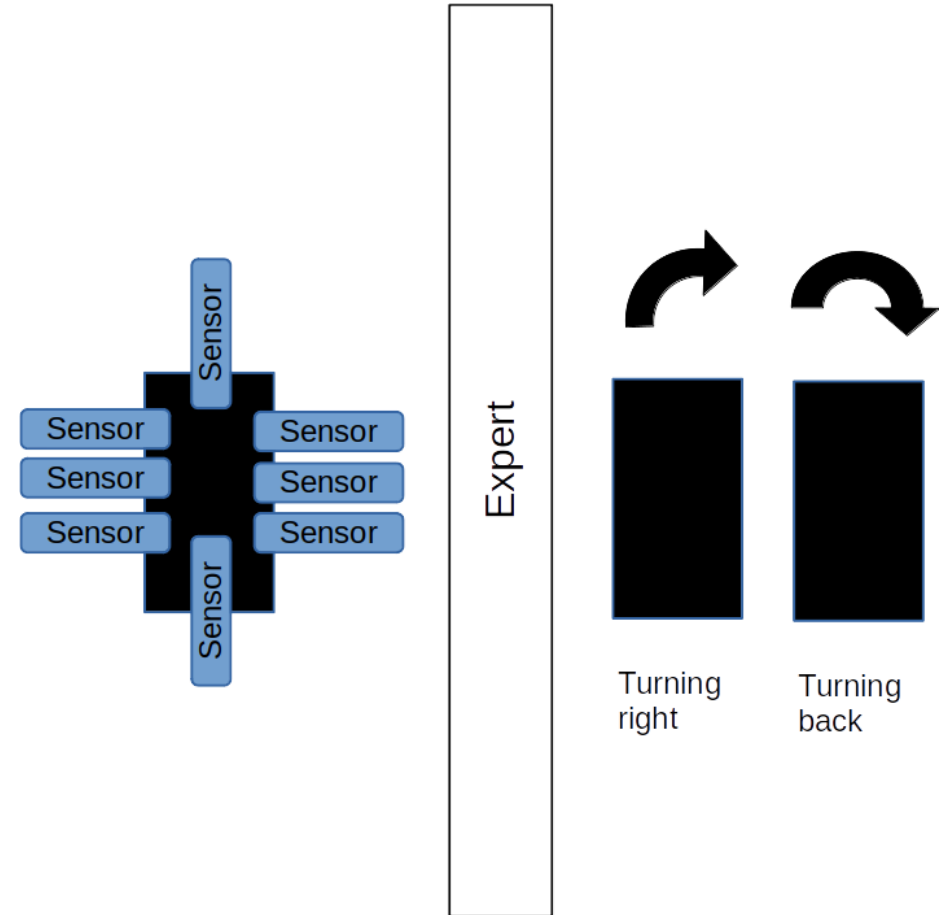
Knowledge discovery from industrial logs

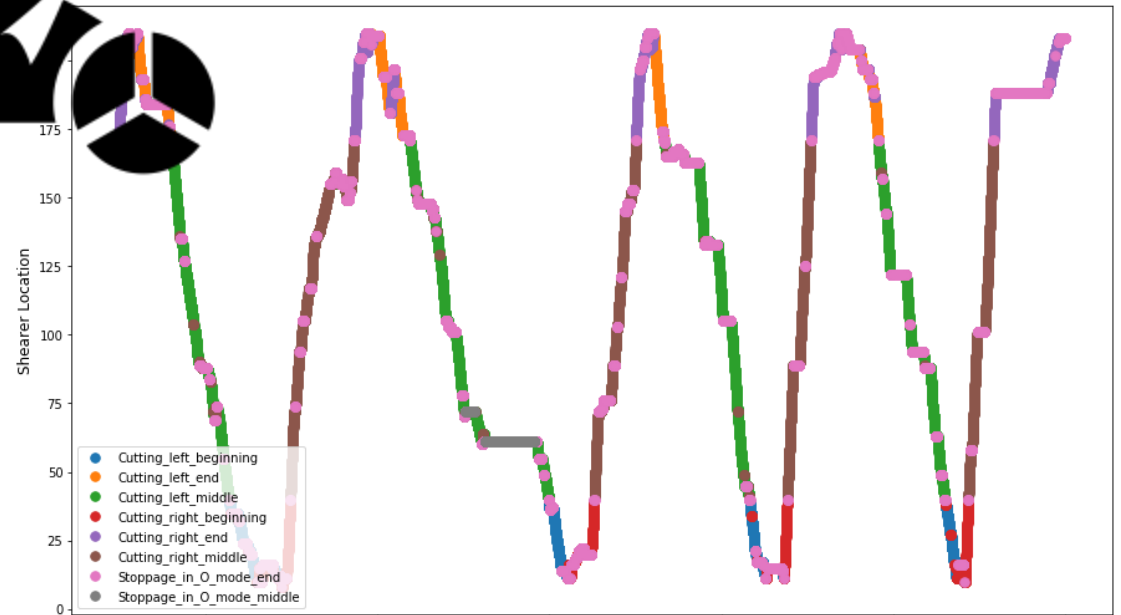
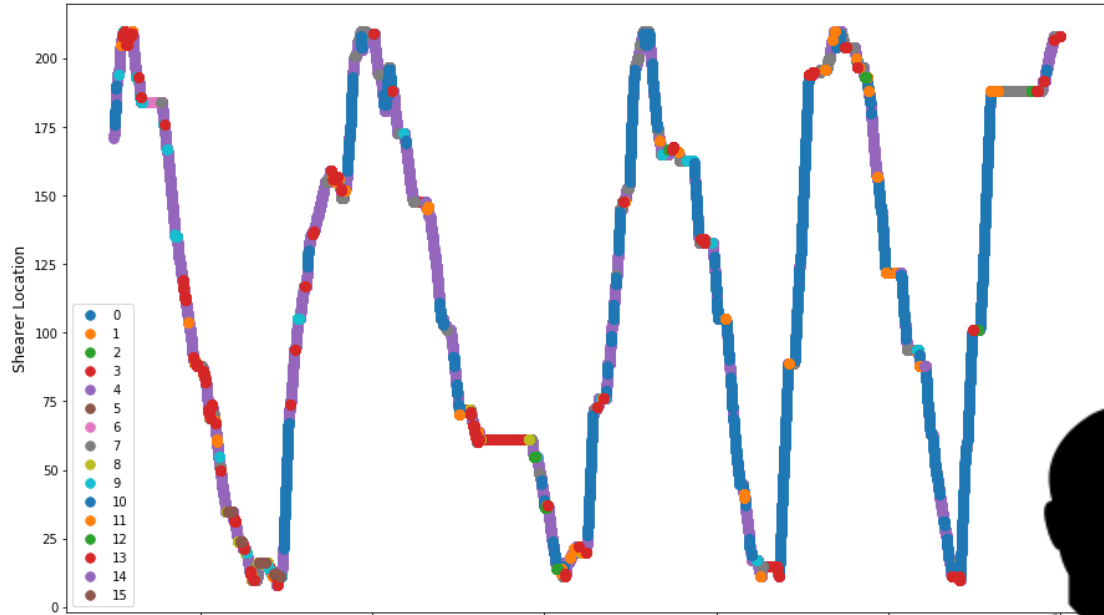
- PACMEL Project (<http://pacmel.geist.re>)
- Industry 4.0: everything is measured
- Low-level measurements to higher-level states
- Expert-defined states vs. Automatically discovered states
 - Data comes with no labels
 - Expert may be too general, or too specific
 - There is a lot of measurements



Knowledge discovery from industrial logs

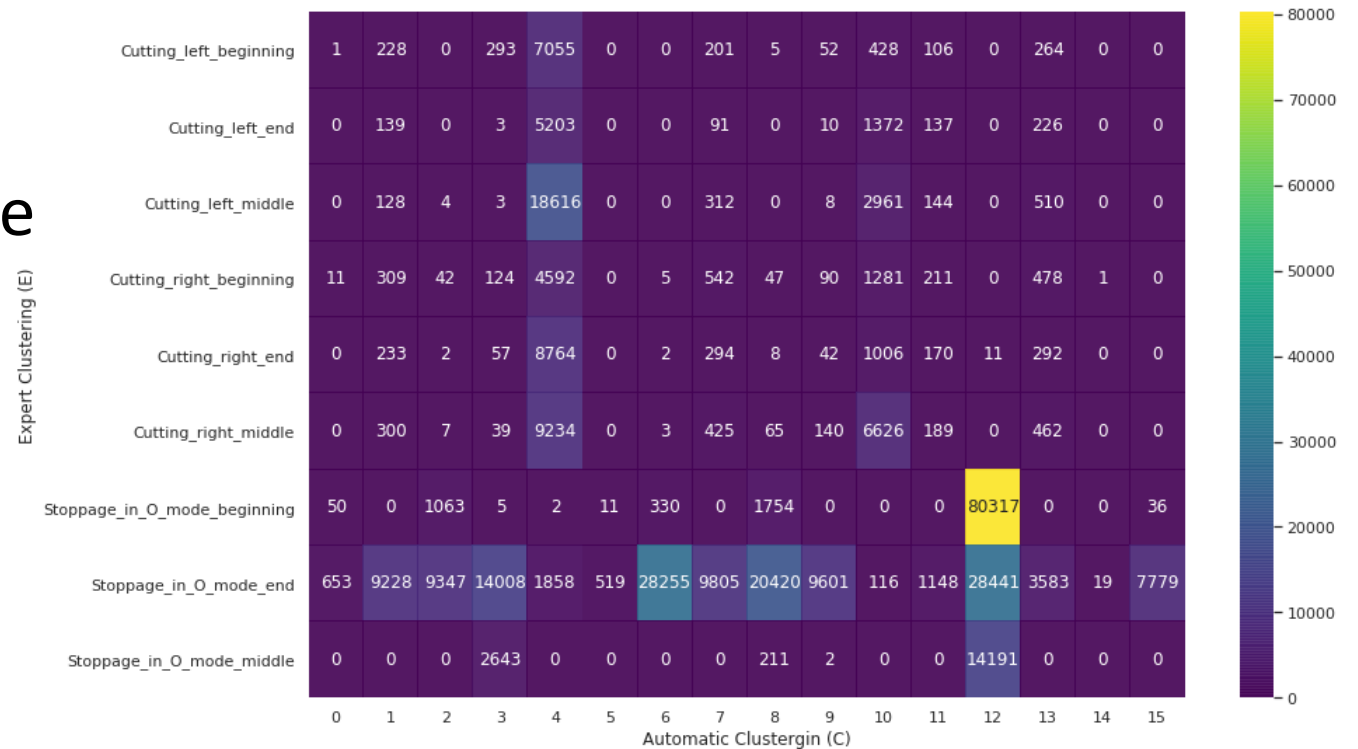
- PACMEL Project (<http://pacmel.geist.re>)
- Industry 4.0: everything is measured
- Low-level measurements to higher-level states
- Expert-defined states vs. Automatically discovered states
 - Data comes with no labels
 - Expert may be too general, or too specific
 - There is a lot of measurements



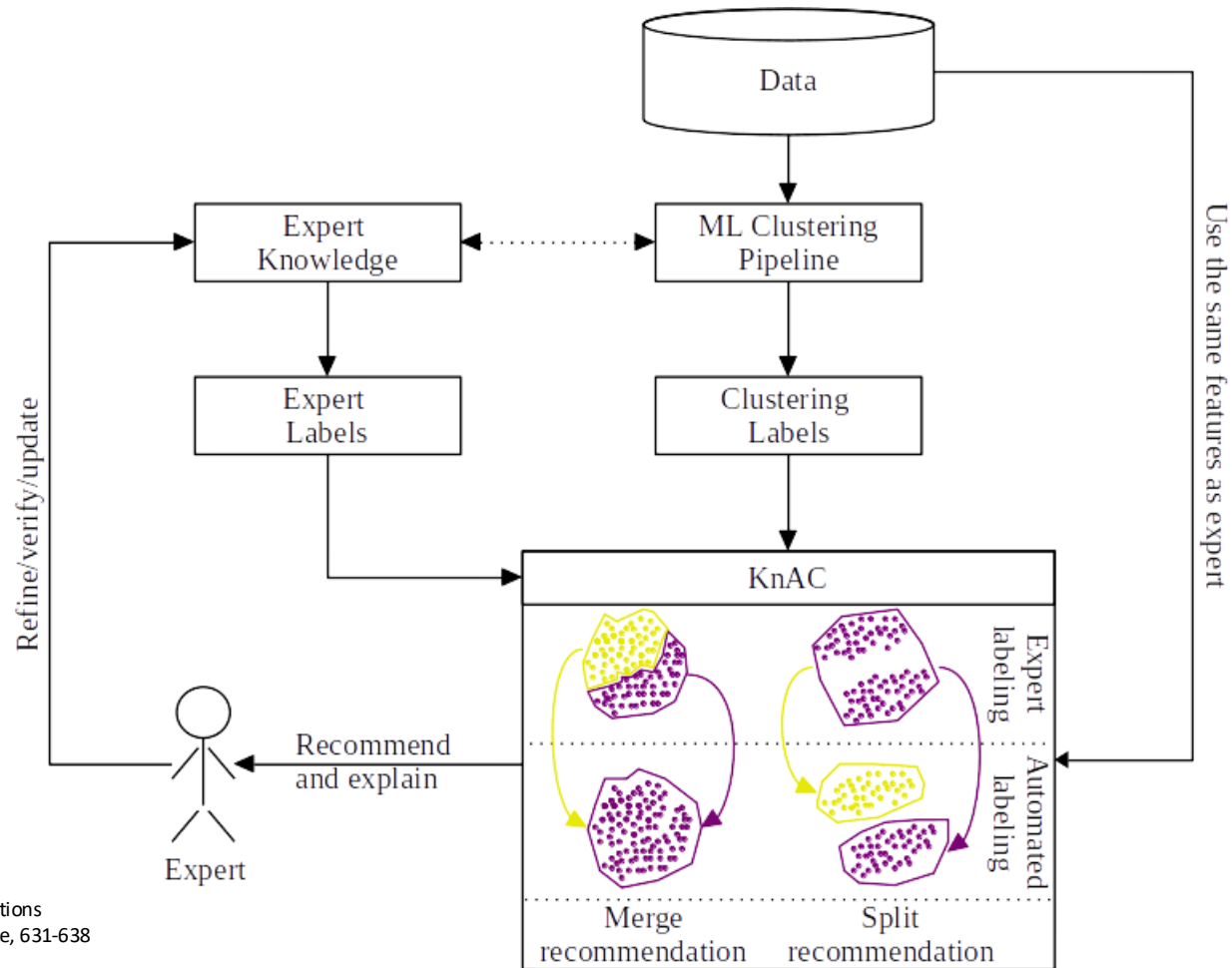


How to confront expert and automatic labelling?

- Analysis of each of the states separately via contingency matrix
- Adjusted rand score
- Adjusted mutual info score
- Homogeneity
- Consistency
- V-measure
- ...



Knowledge Augmented Clustering (KnAC)



Augmenting Automatic Clustering with Expert Knowledge and Explanations
S. Bobek, G. J. Nalepa, International Conference on Computational Science, 631-638

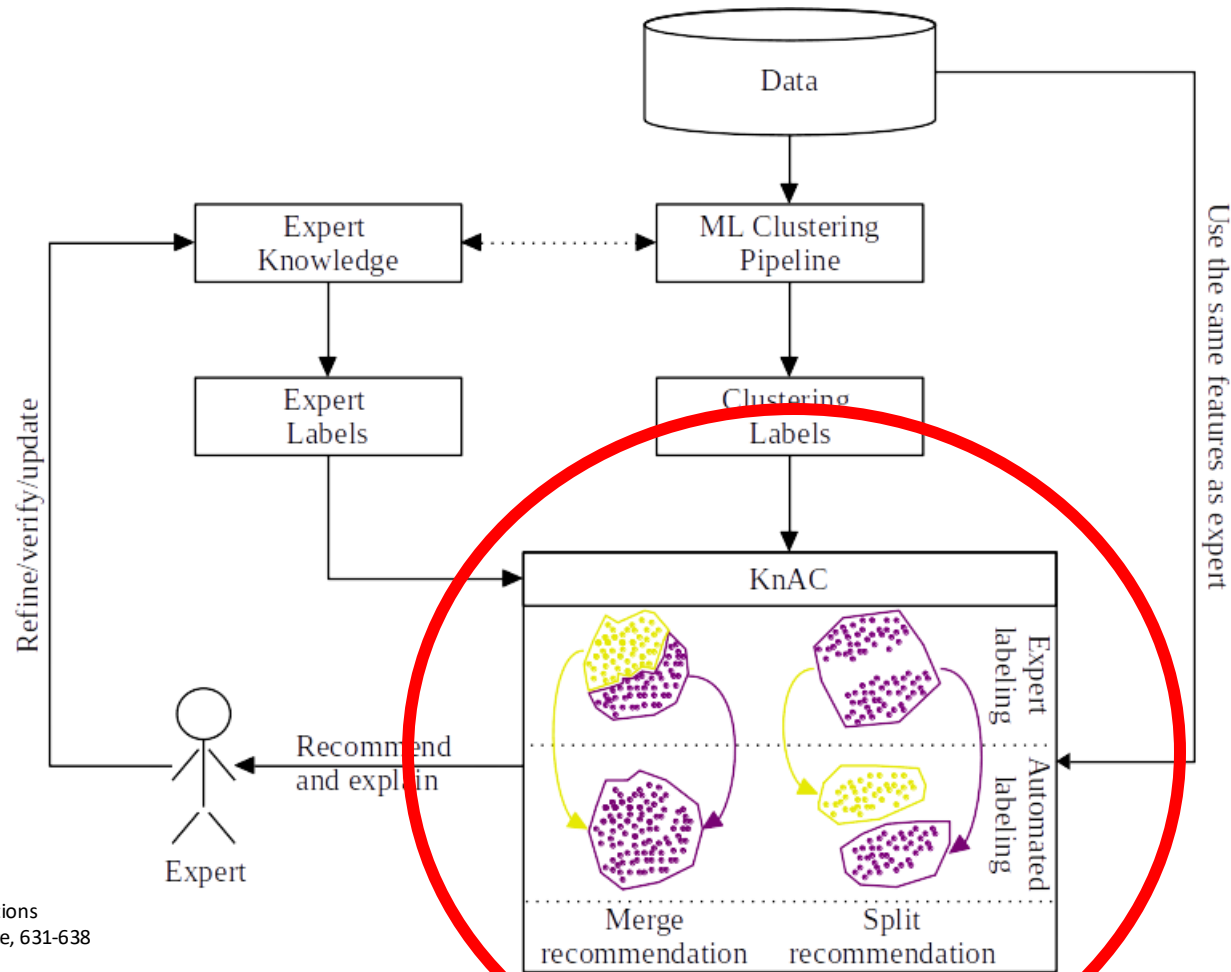
KnAC: an approach for enhancing cluster analysis with background knowledge and explanations

S. Bobek, M. Kuk, J. Brzegowski, E. Brzychczy, and G. J. Nalepa.

ArXiv: <https://arxiv.org/abs/2112.08759>

<https://github.com/sbobek/knac>

Knowledge Augmented Clustering (KnAC)



Augmenting Automatic Clustering with Expert Knowledge and Explanations
S Bobek, GJ Nalepa, International Conference on Computational Science, 631-638

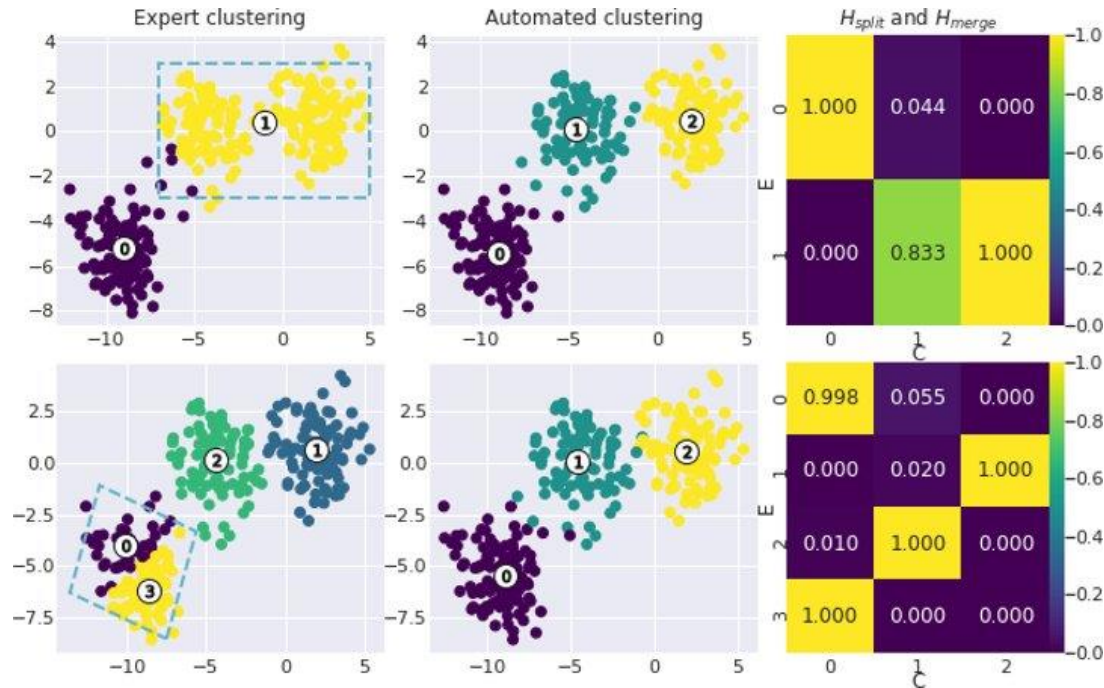
KnAC: an approach for enhancing cluster analysis with background knowledge and explanations

S. Bobek, M. Kuk, J. Brzegowski, E. Brzychczy, and G. J. Nalepa.

ArXiv: <https://arxiv.org/abs/2112.08759>

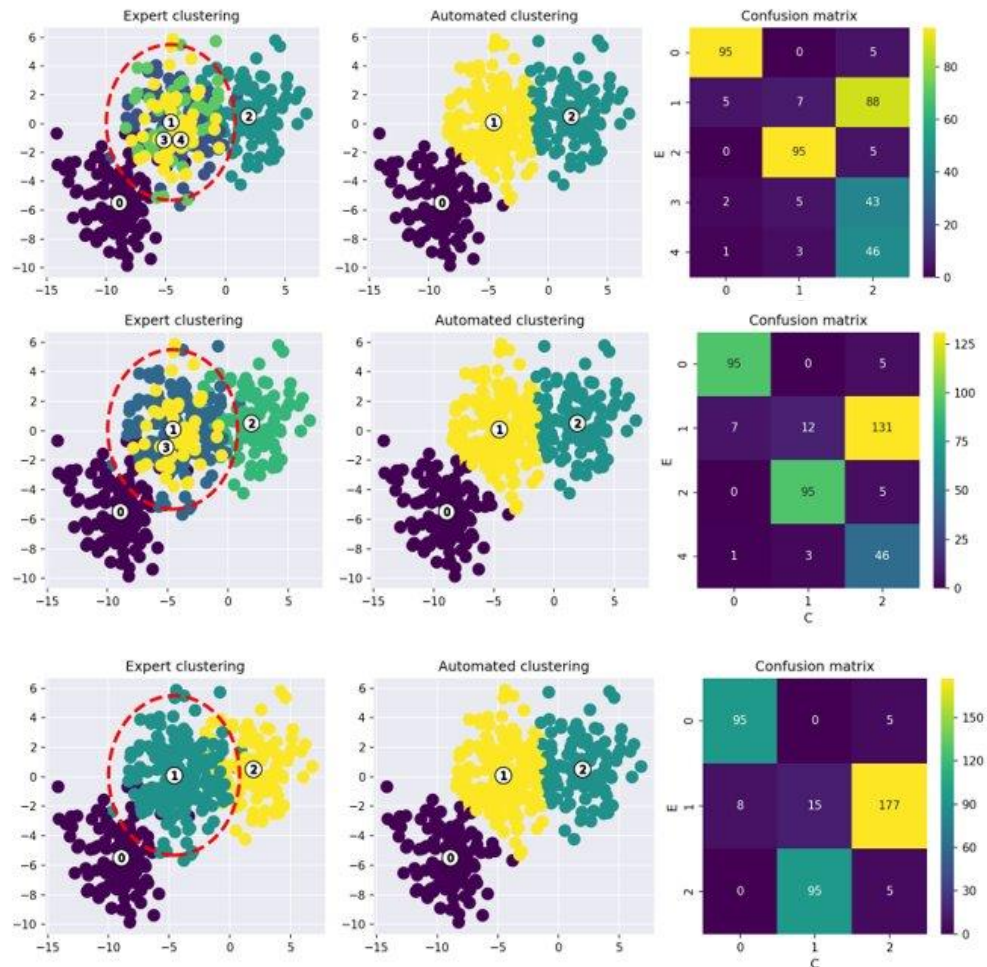
<https://github.com/sbobek/knac>

Expert labels vs. Clustering labels



- Split expert clusters into more specific ones
- Merge expert clusters that seem to be similar
- It is an iterative approach

Expert labels vs. Clustering labels



- Split expert clusters into more specific ones
- Merge expert clusters that seem to be similar
- It is an iterative approach

Splitting expert cluster

- We calculated entropy of each cluster distribution with respect to expert labels
- We scaled rows of distribution matrix to deal with different sized expert clusters
- We divided normalized matrix with entropy values
- The split confidence was calculated by averaging each row of such matrix

	C1	C2	C3	C4
Expert Cluster 1	400	0	1	..
Expert Cluster 2	1000	1000	1000	..
Expert Cluster 3	200	0	3	..
Expert Cluster 4	600	0	10	..

$$H_{i,j}^{split} = \frac{M_{i,j}}{\|M_i\|_2 \left[\frac{H(M_i)}{\log_2(\|E\|)} + 1 \right]}$$

Splitting expert cluster

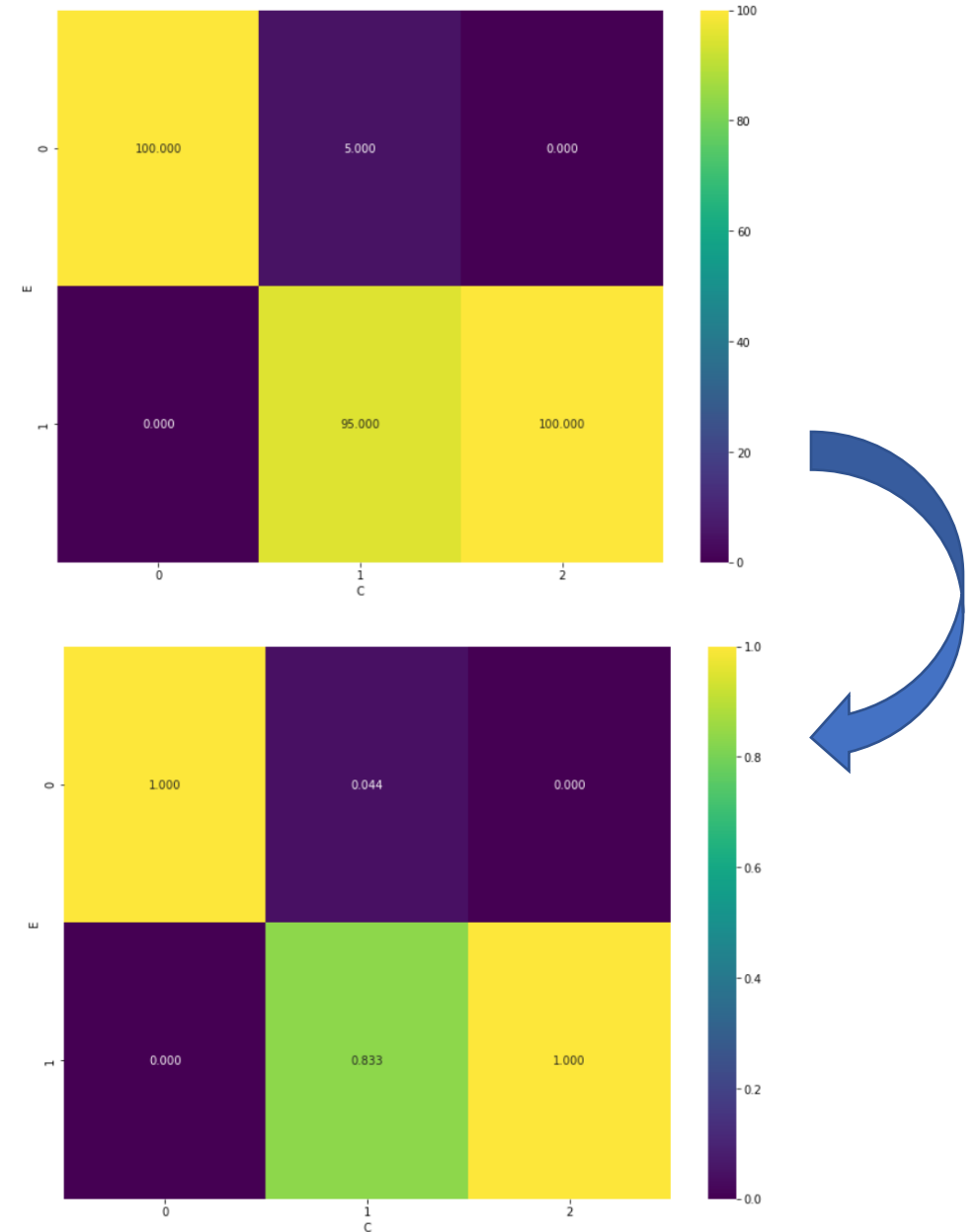
- We calculated entropy of each cluster distribution with respect to expert labels
- We scaled rows of distribution matrix to deal with different sized expert clusters
- We divided normalized matrix with entropy values
- The split confidence was calculated by averaging each row of such matrix

	C1	C2	C3	C4
Expert Cluster 1	400	0	1	..
Expert Cluster 2	1000	0	1000	0
Expert Cluster 3	0	0	3	..
Expert Cluster 4	0	60	0	60

$$H_{i,j}^{split} = \frac{M_{i,j}}{\|M_i\|_2 \left[\frac{H(M_i)}{\log_2(\|E\|)} + 1 \right]}$$

Splitting expert cluster

- We calculated entropy of each cluster distribution with respect to expert labels
- We scaled rows of distribution matrix to deal with different sized expert clusters
- We divided scaled matrix with entropy values
- The split confidence was calculated by averaging each row of such matrix



Merging expert cluster

- We calculated l_2 normalized distribution matrix
- We calculated cosine similarity between rows to denote expert clusters that were similarly splitted with automated method

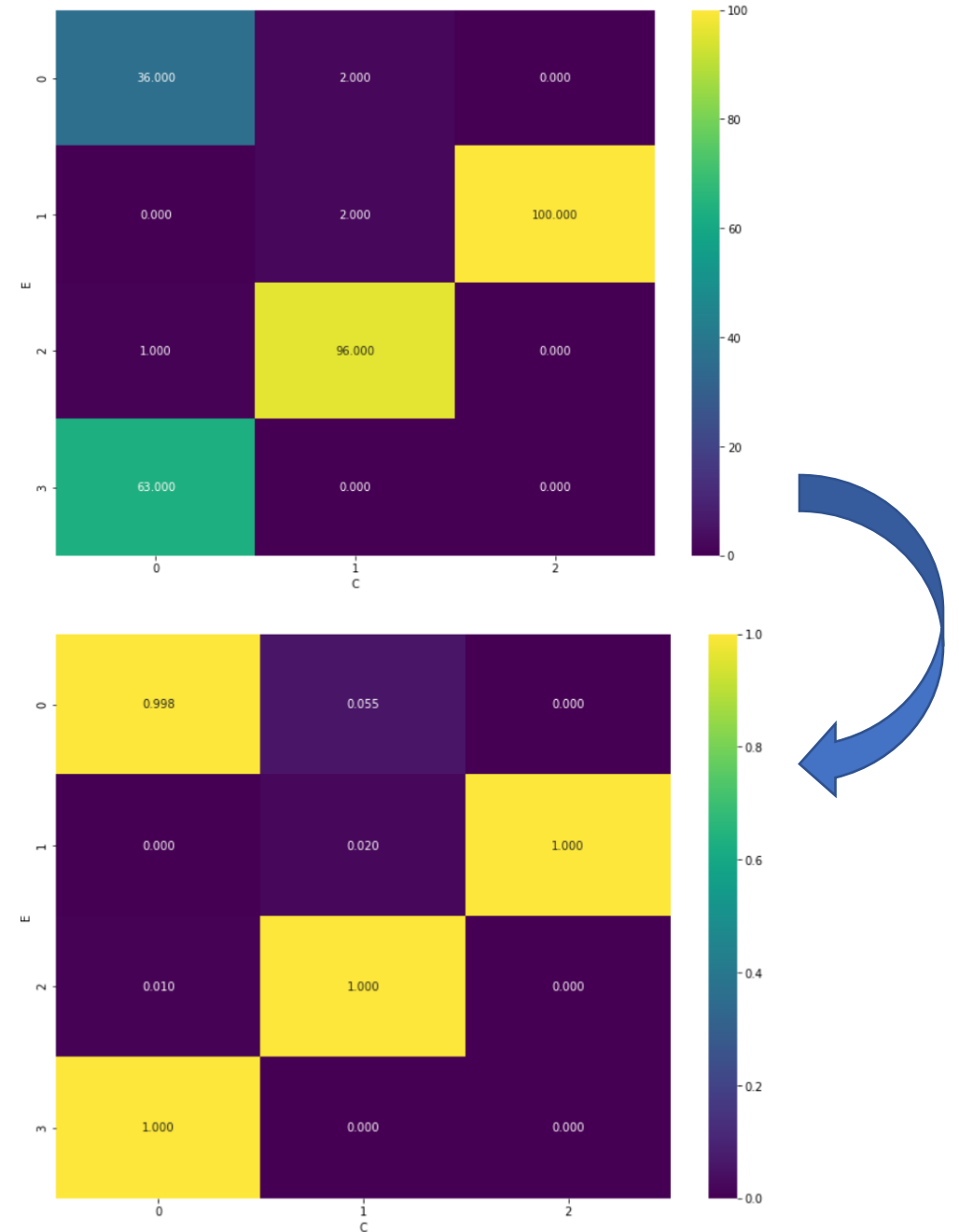
	C1	C2	C3	C4
Expert Cluster 1	400	0	1	..
Expert Cluster 2	1000	0	1000	0
Expert Cluster 3	0	0	3	..
Expert Cluster 4	60	0	60	0

$$H_{i,j}^{merge} = \frac{M_{i,j}}{\|M_i\|_2}$$

Expert Clusters 2 and 4 are similar in their distribution in automated clustering

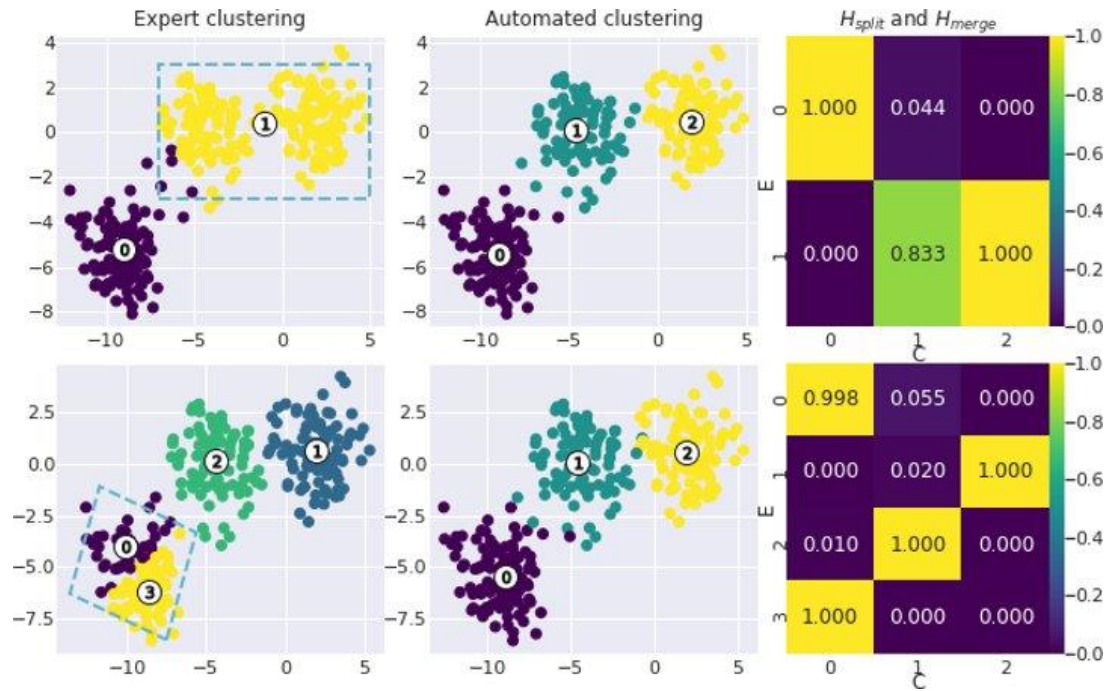
Merging expert cluster

- We calculated l_2 normalized distribution matrix
- We calculated cosine similarity between rows to denote expert clusters that were similarly splitted with automated method



Results - splits

$$C_i^{split} = \left\{ c_j \in H_i^{split} : \frac{c_j}{1 - \lambda^s} > \epsilon_s \right\}$$



Decrease in silhouette score between splitted clusters

$$Conf(C_i^{split}) = \left\{ (1 - \lambda^s)c_j + \lambda^s(S^{dec}(C_i^{split})) : c_j \in C_i^{split} \right\}$$

Assuming $\lambda^s = 0.1$

SPLIT EXPERT CLUSTER

E_1

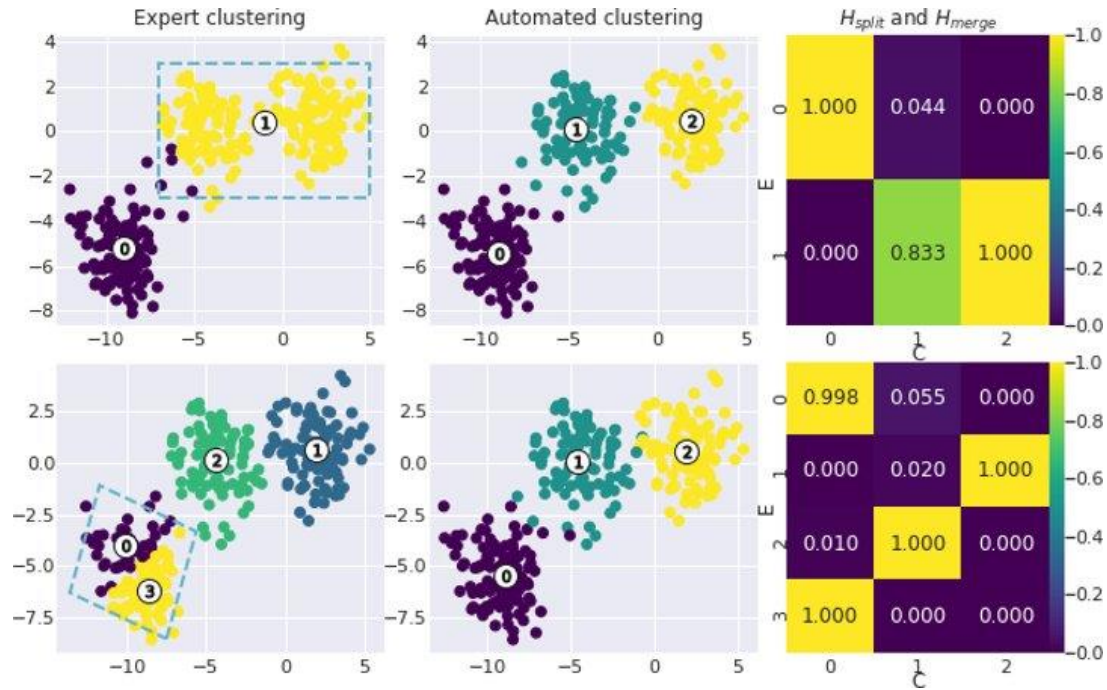
INTO CLUSTERS

[(C_1, C_2)]

(Confidence 0.87)

Results - merges

Linkage distance between merged clusters

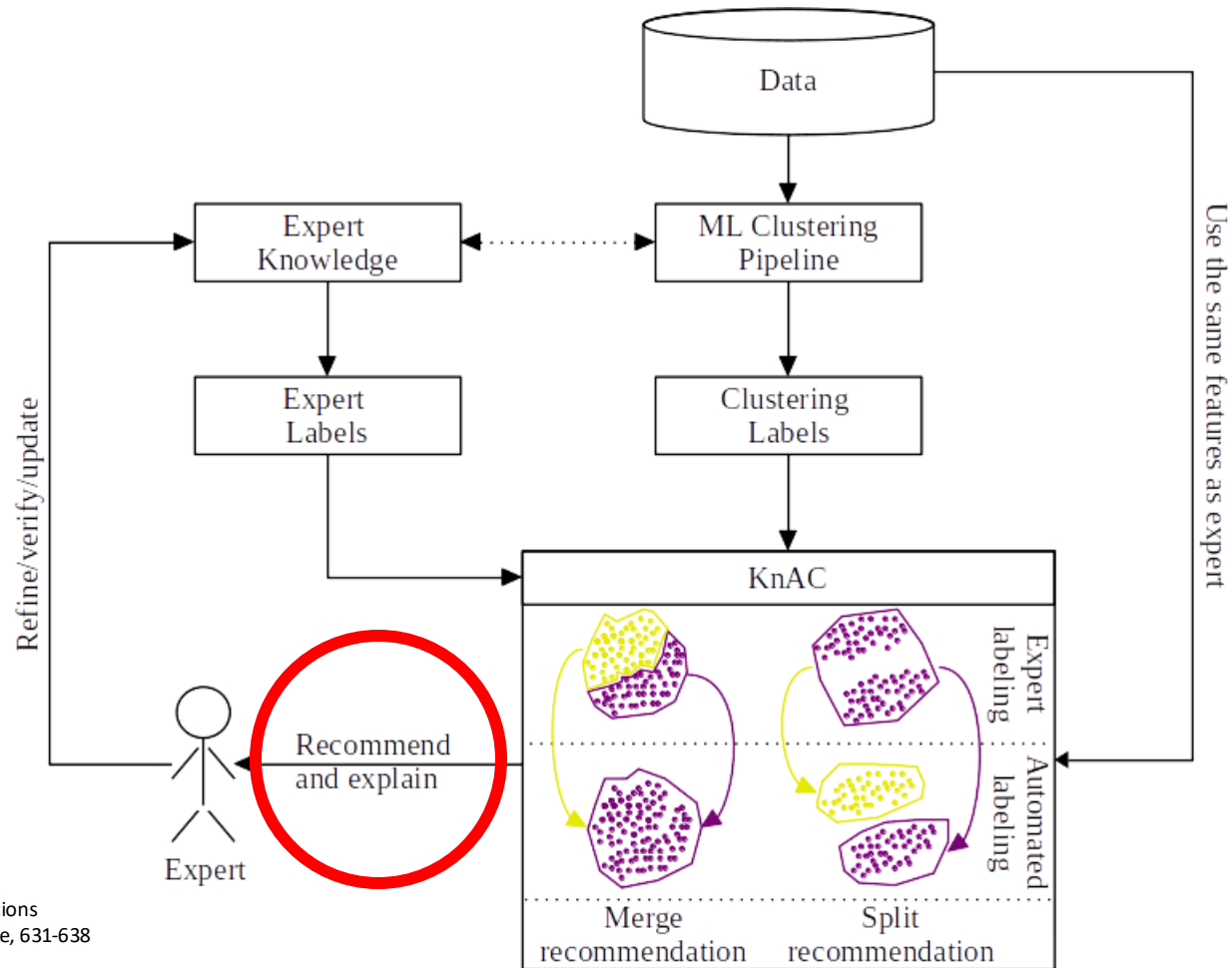


$$C_{j,k}^{merge} = \left\{ E \ni E_j, E_k : (1 - \lambda^m) H_{j,k}^{sim} + \lambda^m (1 - D_{j,k}^{linkage}) > \epsilon_m \right\}$$

Assuming $\lambda^m = 0.2$

MERGE
 EXPERT CLUSTER E_0
 WITH
 EXPERT CLUSTER E_3
 INTO
 CLUSTER C_0 # (Confidence 0.98)

Knowledge Augmented Clustering (KnAC)



Augmenting Automatic Clustering with Expert Knowledge and Explanations
S Bobek, GJ Nalepa, International Conference on Computational Science, 631-638

KnAC: an approach for enhancing cluster analysis with background knowledge and explanations

S. Bobek, M. Kuk, J. Brzegowski, E. Brzychczy, and G. J. Nalepa.

ArXiv: <https://arxiv.org/abs/2112.08759>

<https://github.com/sbobek/knac>

eXplainable Artificial Intelligence (XAI)

To explain an event is to provide some information about its causal history.

In an act of explaining, someone who is in possession of some information about the causal history of some event - explanatory information, I shall call it - tries to convey it to someone else. - David Lewis

Different approaches

- Intelligibility of the system
- Interpretability of models
- Explainability of ML models

DARPA. Broad agency announcement – explainable artificial intelligence (XAI). DARPA-BAA-16-53, August 2016. <https://www.darpa.mil/program/explainable-artificial-intelligence>

C. Molnar. Interpretable Machine Learning. E-book – Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License, 2019. <https://christophm.github.io/interpretable-ml-book/>

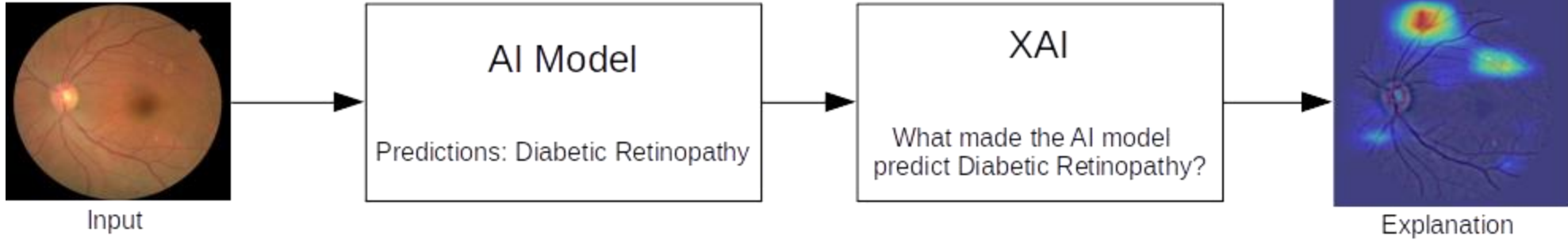
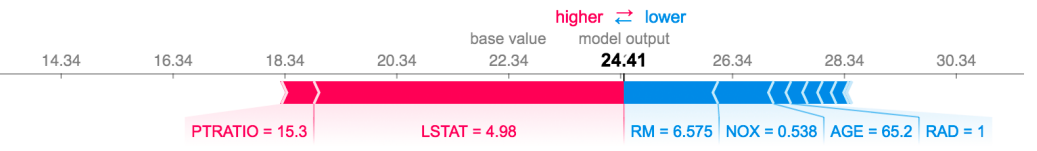
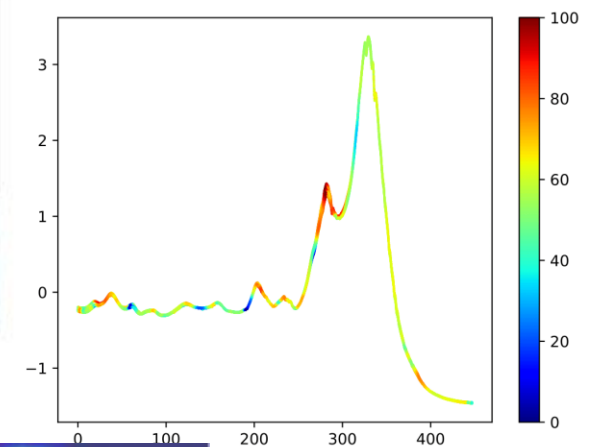
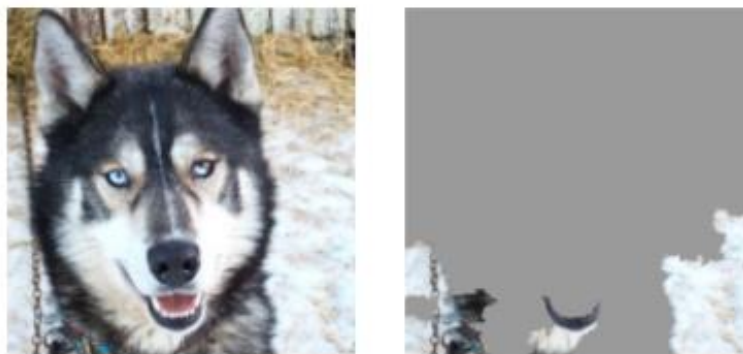
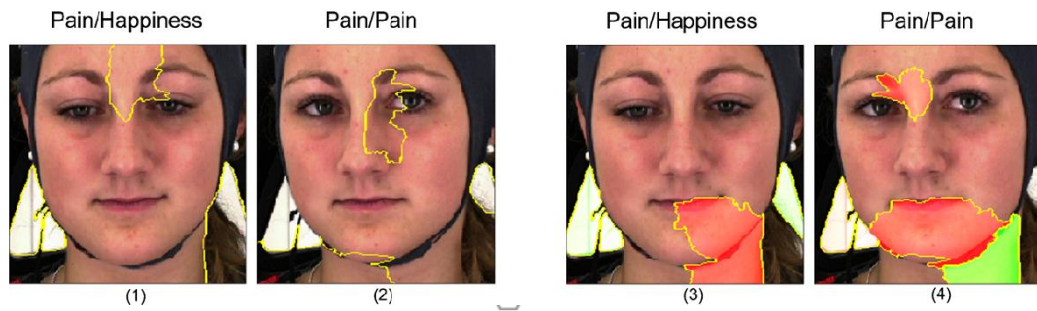
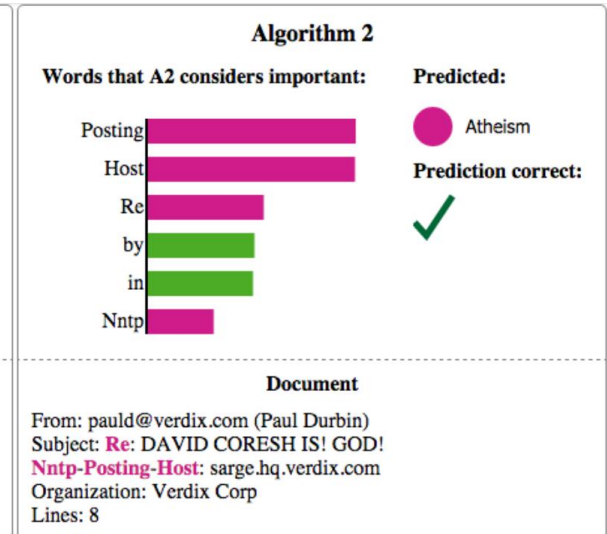
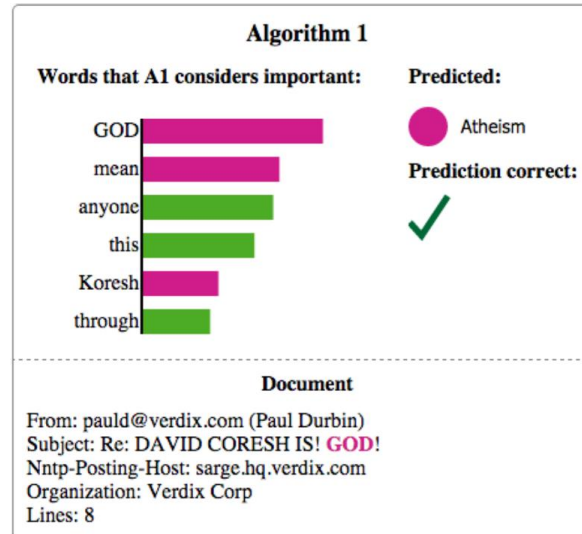
Old topic

- Expert systems
- Recommender systems
- Context-aware systems
- Machine learning

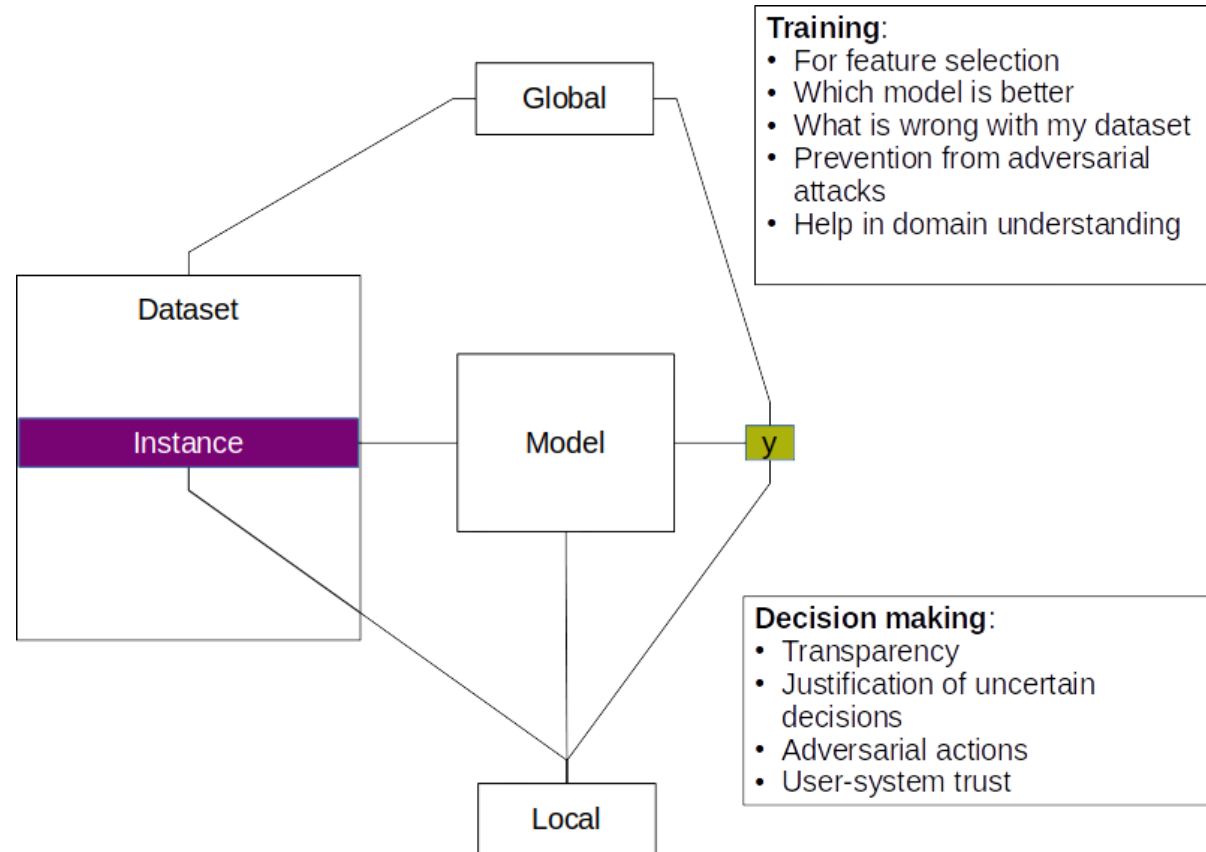
Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell* **1**, 206–215 (2019). <https://doi.org/10.1038/s42256-019-0048-x>

A. Barredo Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Bar-bado, S. Garcia, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82 – 115, 2020.

Examples of XAI



General Goals of XAI



Why XAI is non trivial

In an **act** of explaining, **someone** who is in possession of **some information**

Artificial intelligence

Feature contribution

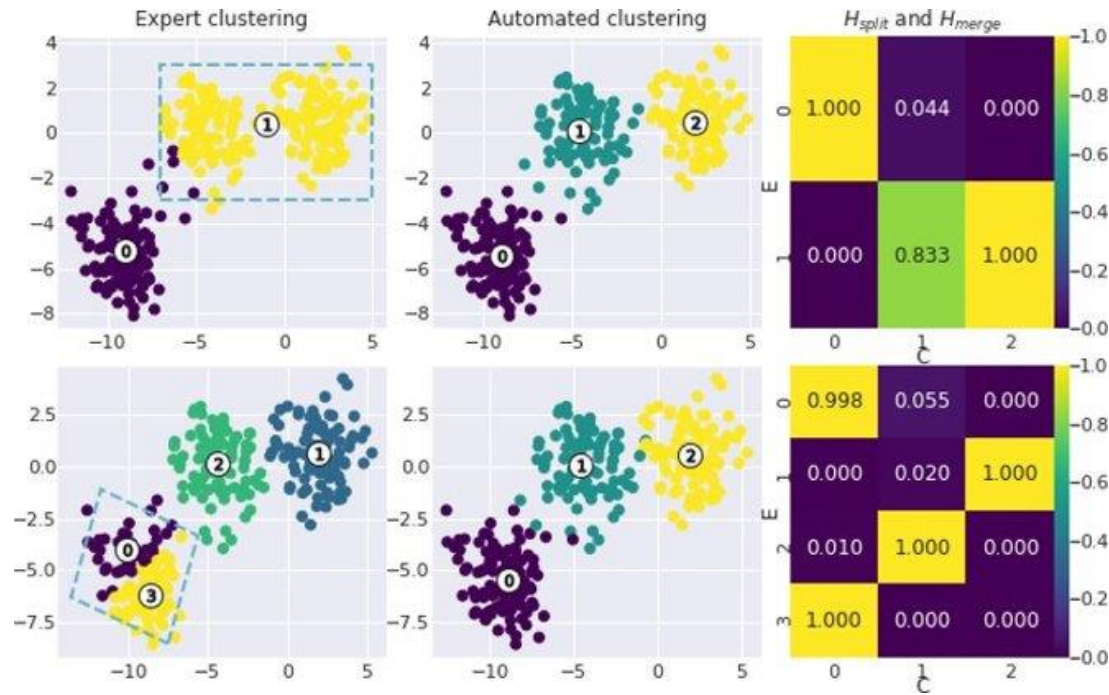
about the **causal history of some event** - explanatory information,

Why input to the model generated
such output

I shall call it - tries to **convey it to someone else.**

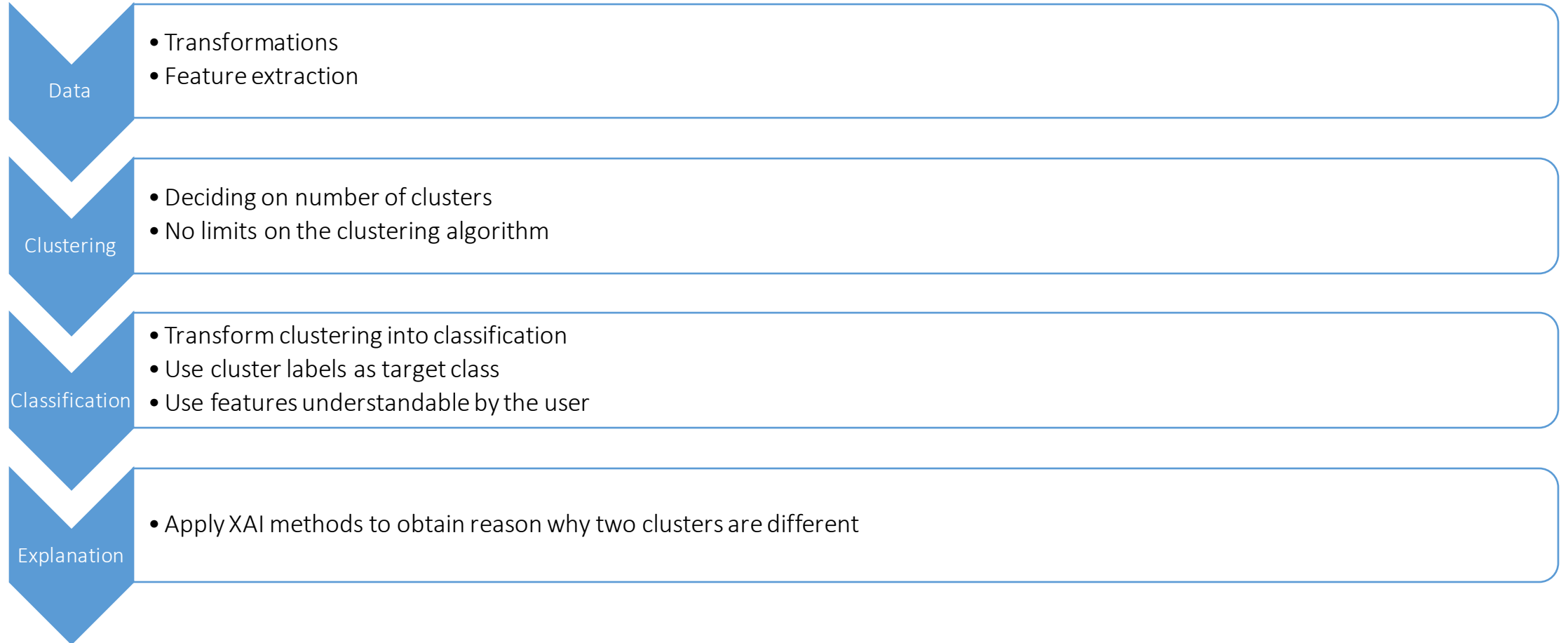
Human

How to use XAI in KnAC?



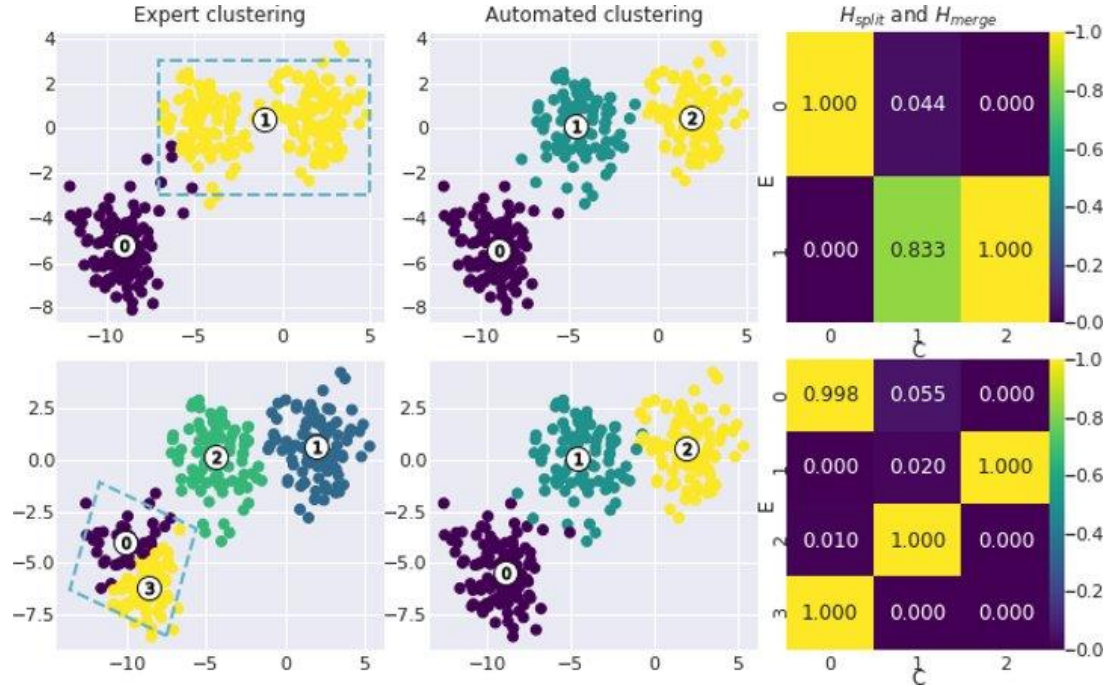
- Split: What makes the two new clusters different from each other to convince expert they are different entities?
- Merge: What makes the two expert clusters different from each other to convince expert that they are the same entity (difference is irrelevant)

From clustering to classification

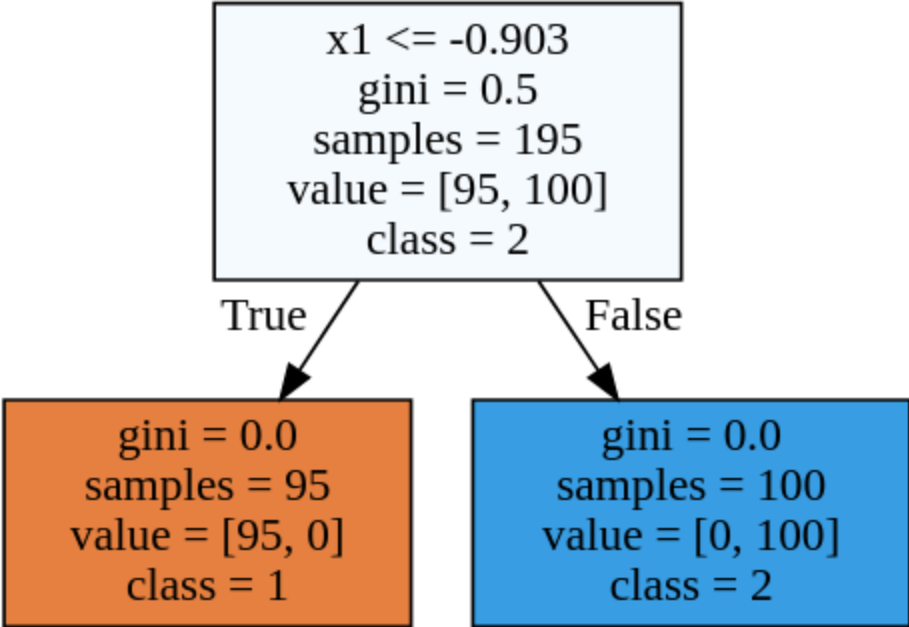


Explanations of splits

SPLIT EXPERT CLUSTER
 E_1
 INTO CLUSTERS
 [(C_1, C_2)]
 (Confidence 0.87)

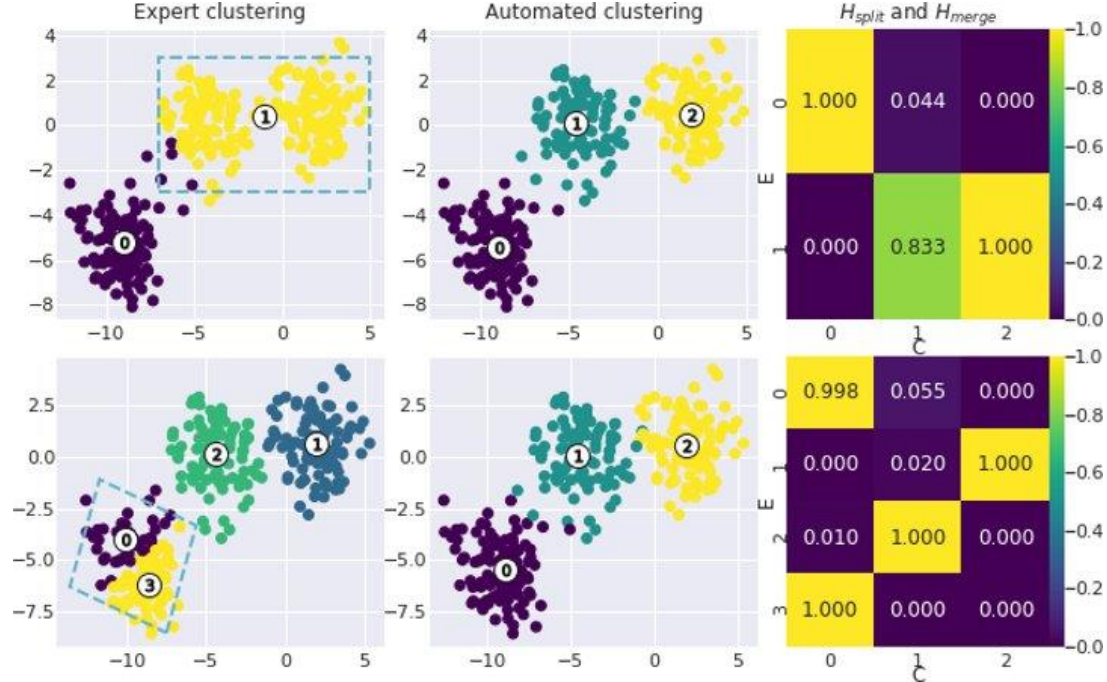


C_1: $x_1 \leq -0.30$ (Precision: 0.99, Coverage: 0.49)
 C_2: $x_1 > -0.30$ (Precision: 1.00, Coverage: 0.49)



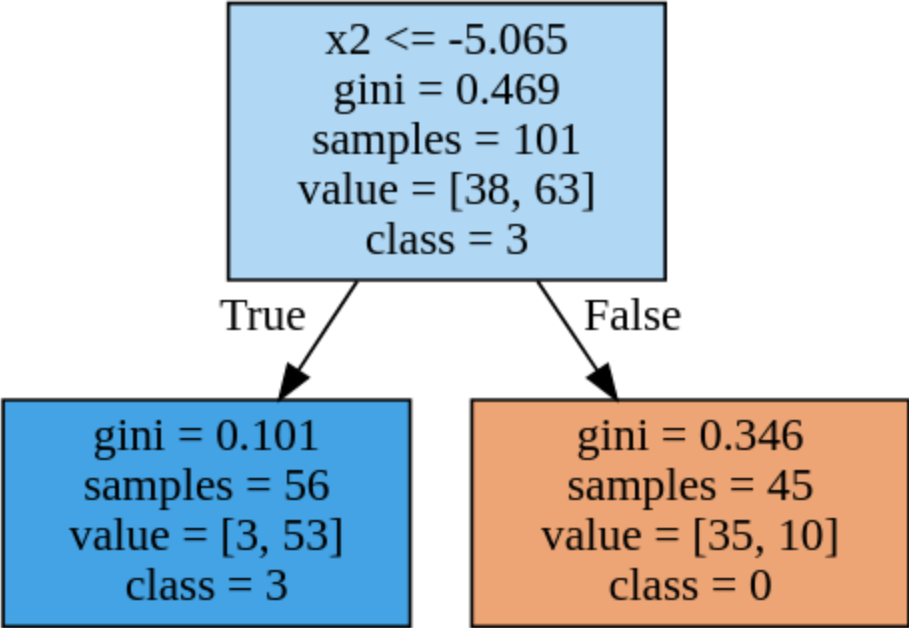
Explanations of merges

MERGE
 EXPERT CLUSTER E_0
 WITH
 EXPERT CLUSTER E_3
 INTO
 CLUSTER C_0 # (Confidence 0.98)

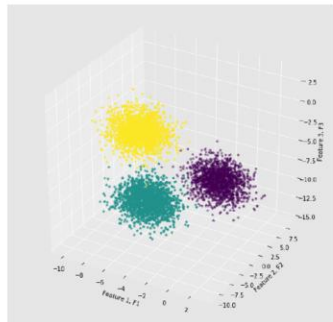


E_0: $x_1 \leq -8.20$ AND $x_2 > -4.34$ (Precision: 1.00, Coverage: 0.07)

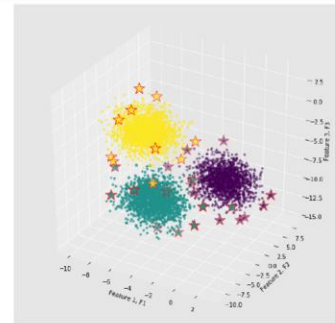
E_1: $x_1 \leq -4.34$ (Precision: 0.90, Coverage: 0.25)



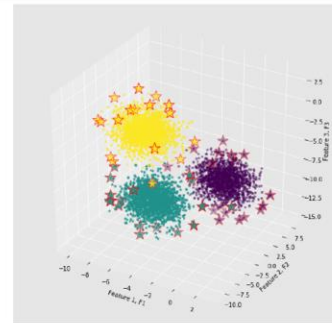
Explainable clusters



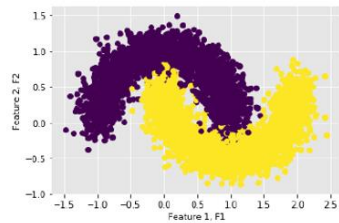
(a) Make blobs 3d dataset visualization.



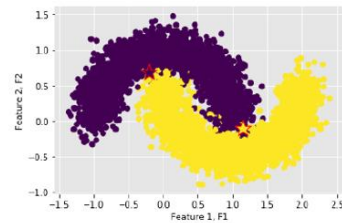
(b) Make blobs 3d dataset - KDTree query describing method.



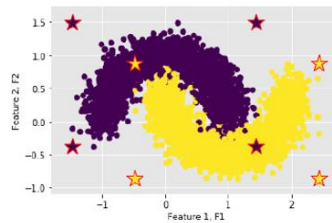
(c) Make blobs 3d dataset - Isolation Forest describing method.



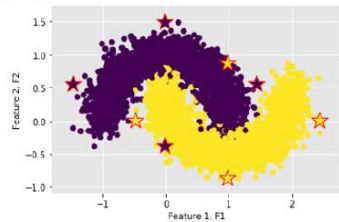
(a) Make moons dataset clusters visualization.



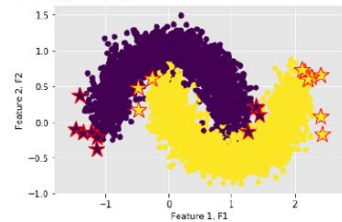
(b) Make moons dataset - K-medoids describing method.



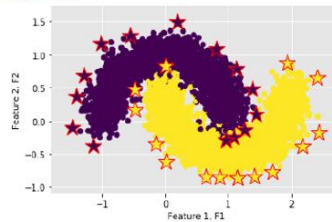
(c) Make moons dataset - Corners describing method.



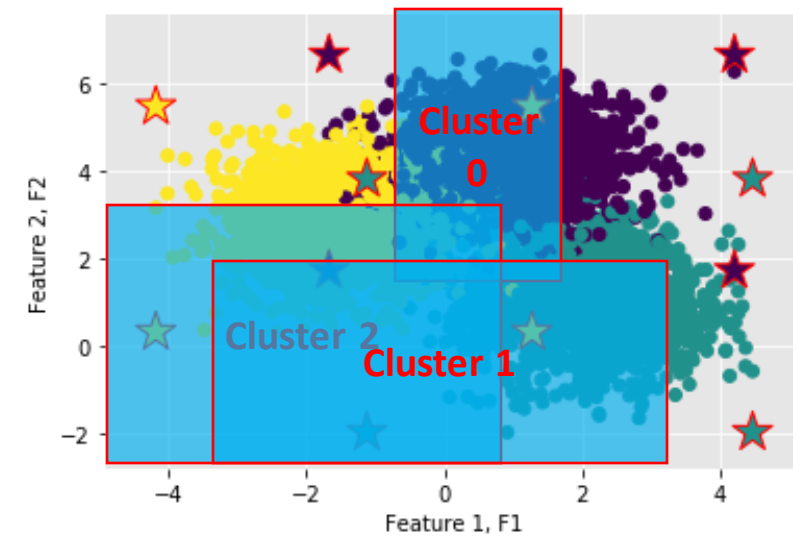
(d) Make moons dataset - Middle points describing method.



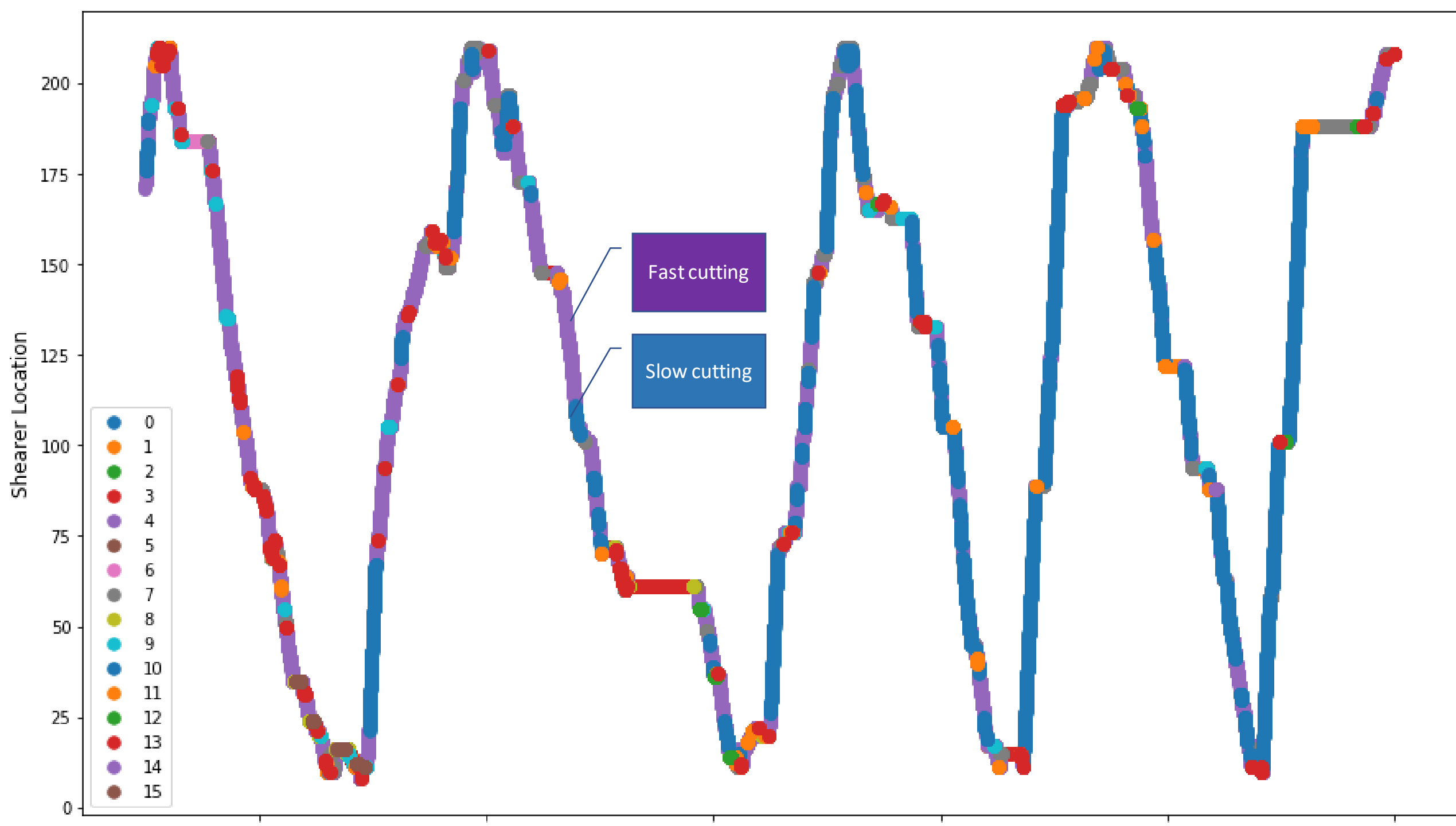
(e) Make moons dataset - Maximum distance describing method.



(f) Make moons dataset - Alpha shape describing method.

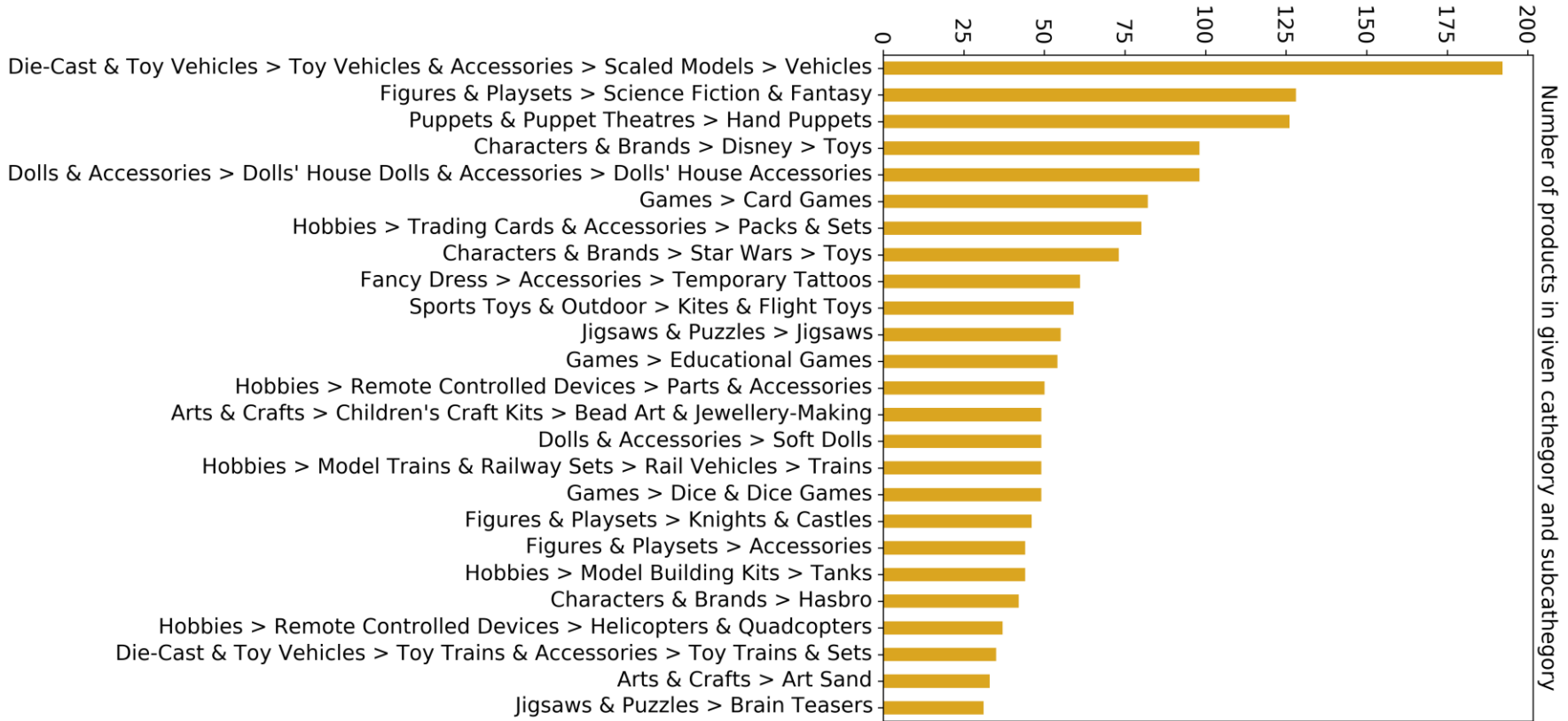


Rule no.	Rule	Cluster	Certainty
1	$F1 > 0.68$ and $F2 > 2.99$	0	0.48
2	$0.68 < F1 \leq 1.77$ and $F2 > 1.64$	0	0.64
3	$-1.14 < F1 \leq 1.77$ and $F2 > 1.64$	0	0.54
4	$F1 > 0.68$ and $F2 \leq 2.99$	1	0.44
5	$F1 > -1.14$ and $F2 \leq 1.64$	1	0.68
6	$F1 \leq -1.14$	2	0.25
7	$F1 \leq 0.68$ and $F2 \leq 2.99$	2	0.43



E-commerce and coal mine

Product to category



Chuggington is an action-packed contemporary animated train series for pre-schoolers that follows the exciting adventures of three young trainees: Wilson, Brewster and Koko. In each energetic, vibrant episode, the trainees ride the rails through the world of Chuggington, exploring many locations and taking on exciting challenges that test their courage, speed and determination. With the help support and guidance of the more experienced Chuggers, they learn positive values, including respect and loyalty, and new skills such as teamwork and patience, empowering them to be the best trainees they can be. Box Contains 1 x Chuggington Train ”

E-commerce and coal mine

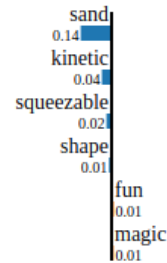
See the results at online tutorial: <https://github.com/sbobek/knac>

Justification for Arts & Crafts > Art Sand
[('sand', -0.14062779721924024), ('kinetic', -0.04427780308500065), ('squeezable', -0.023407082689875347), ('shape', -0.011500639821694475), ('fun', 0.0103043090148735)

Prediction probabilities



Arts & Crafts > Art Sand & Brands > Disney



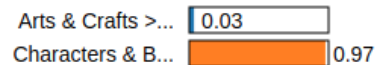
Text with highlighted words

have hours of fun with this magic sand playset creating brilliant sand shapes or create your own sculptures the magic sand is the squeezable sand where you can feel the fun pack it pull it shape it and love it motion sand is so incredible you can put it down it kinetic meaning it sticks to itself and not to you it easy to shape and mould and flows through your fingers like slow moving liquid but leaves them completely dry kinetic sand stimulates children creative skills allowing them to create anything they can imagine it never dries out and is gluten free this soft and stretchy sand easily cleans up while delivering non stop fun it squeezable sand you can put down for ages years and over

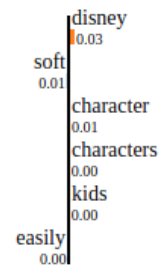
Justification for Characters & Brands > Disney > Toys

[('disney', 0.027092040859663918), ('soft', -0.008068585098757642), ('character', 0.006040355286770795), ('characters', 0.004645787316197593), ('kids', 0.0043490668605)

Prediction probabilities



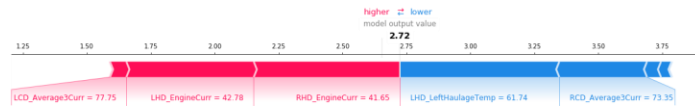
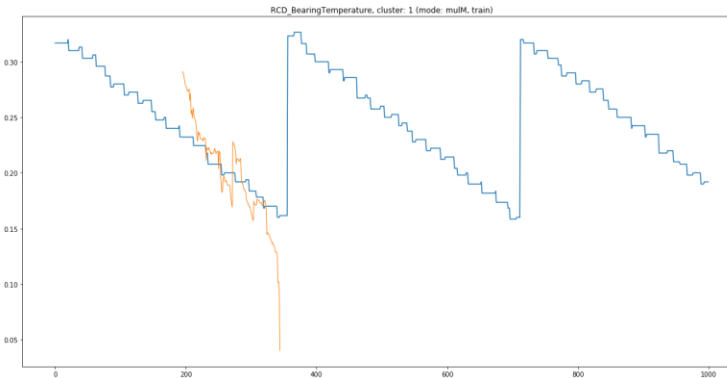
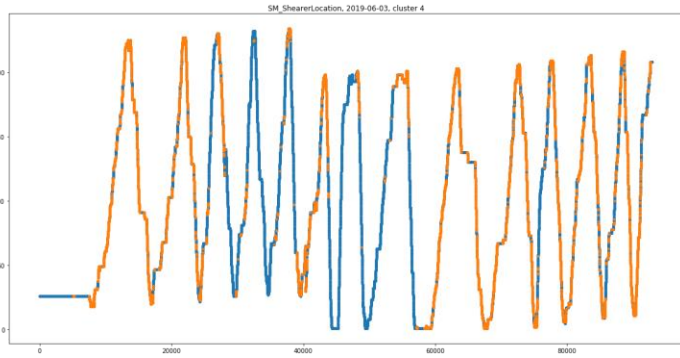
Arts & Crafts > Art Sand & Brands > Disney



Text with highlighted words

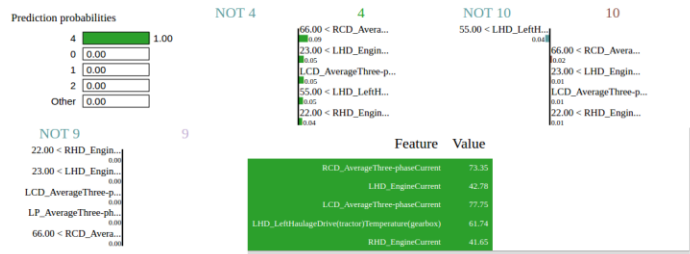
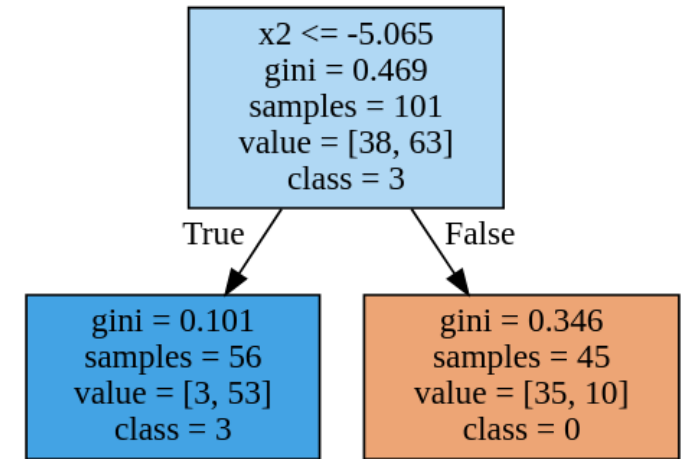
product description whether at home on the road or in the air your favourite disney character can provide great company and comfort these soft colourful cushions can be easily transformed into disney character soft toy by simply opening and closing the velcro loved by children of all ages these classic disney characters will keep kids entertained for hours and when sleepy just rest your head on the cushion and dream away all our disney character cushions are washable please read washing label for further instructions box contains x

How to explain? Which explanation should we trust?

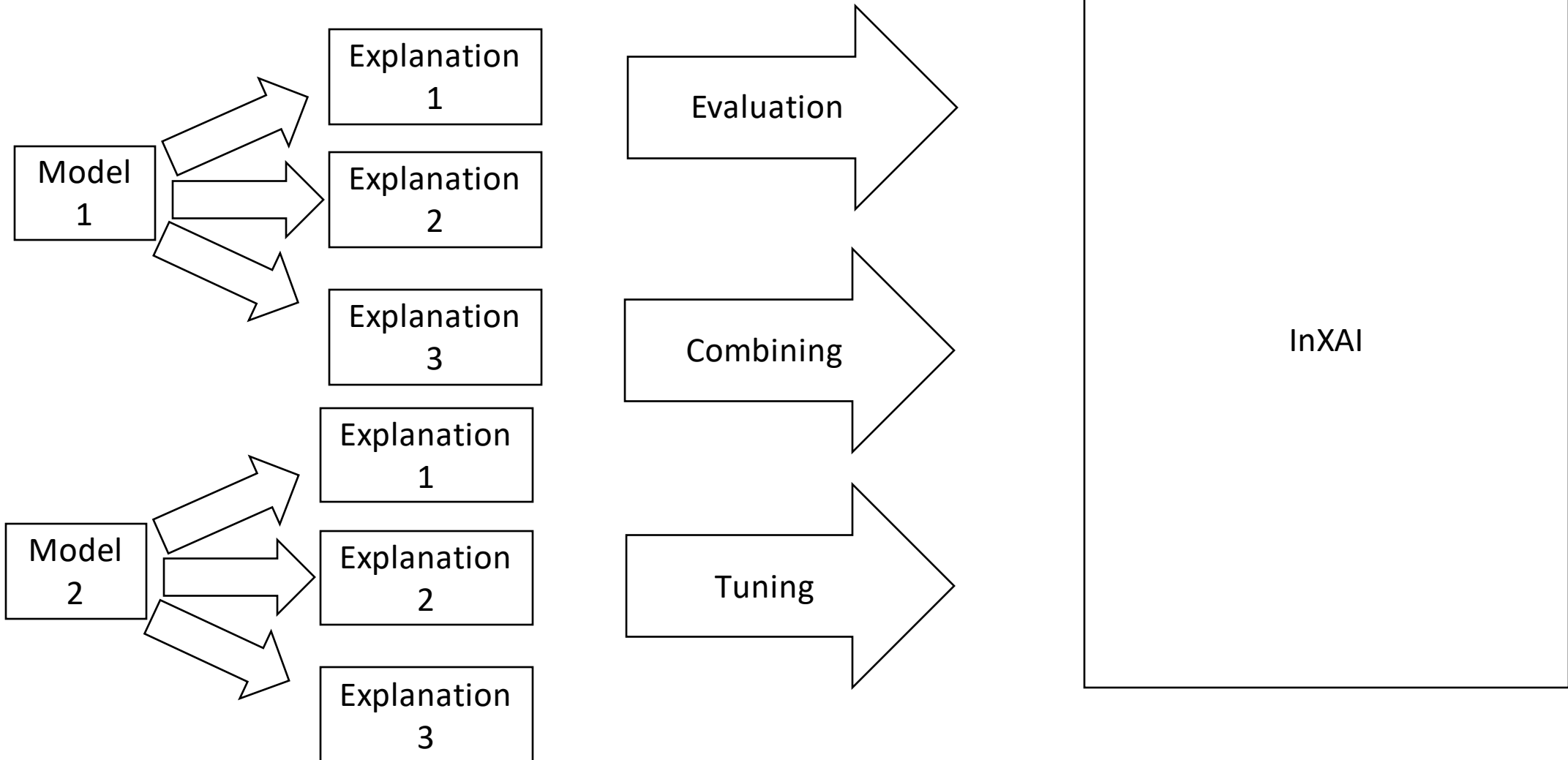


ANCHOR

SM_ShearerSpeed <= 2.00 AND
 RCD_AverageThree-phaseCurrent_LOW OR
 RCD_AverageThree-phaseCurrent_MEDIUM OR
 RCD_AverageThree-phaseCurrent_HIGH OR
 RCD_AverageThree-phaseCurrent_OVERLOAD) AND
 (RHD_EngineCurrent_IDLE OR
 RHD_EngineCurrent_LOW) AND
 (LCD_AverageThree-phaseCurrent_IDLE OR
 LCD_AverageThree-phaseCurrent_LOW) AND
 RHD_RightHaulageDrive(tractor)Temperature(gearbox) > 0.00 AND
 (LHD_EngineCurrent_IDLE OR LHD_EngineCurrent_LOW) AND
 43.00 < LA_LeftArmTemperature <= 52.00 AND
 LHD_LeftHaulageDrive(tractor)Temperature(gearbox) > 65.00 AND
 RA_RightArmTemperature <= 54.00 AND LP_AverageThree-phaseCurrent <= 4.00 AND
 57.00 < RCD_BearingTemperature <= 64.00

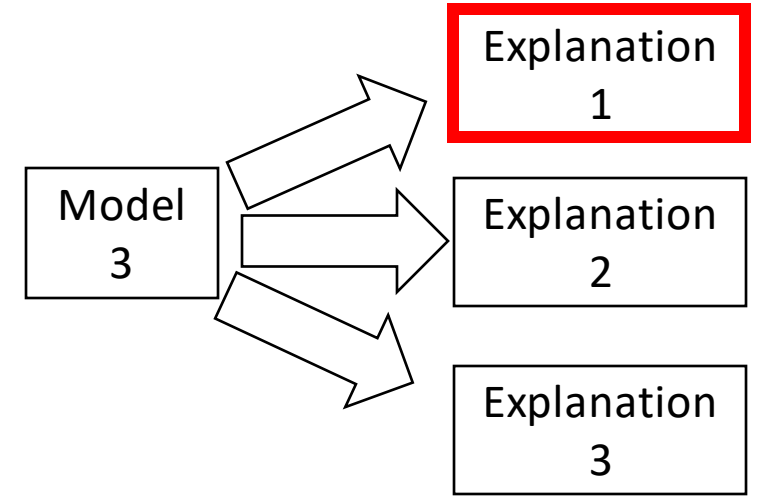
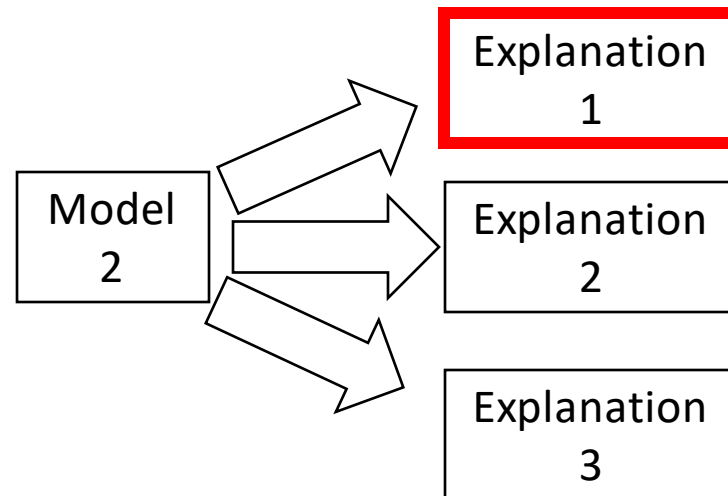
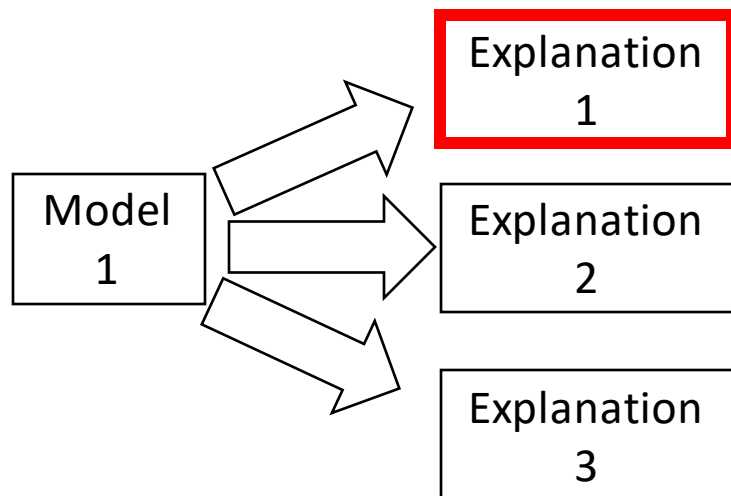
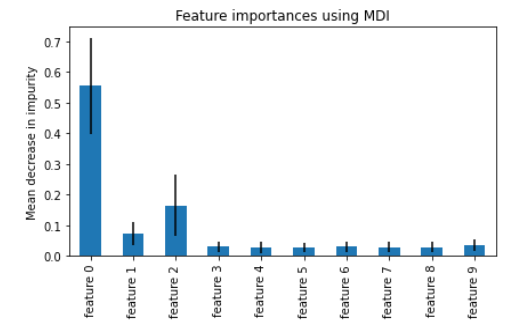
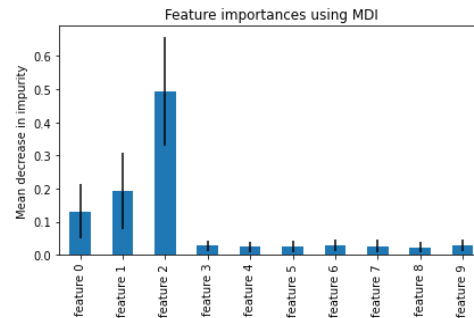
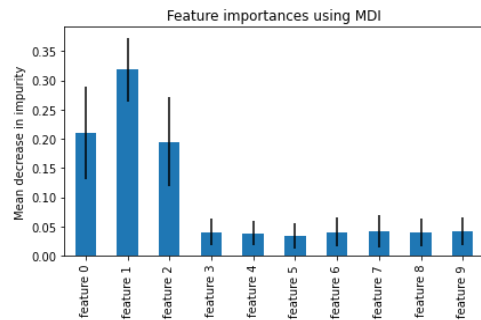


Intelligible XAI (InXAI)



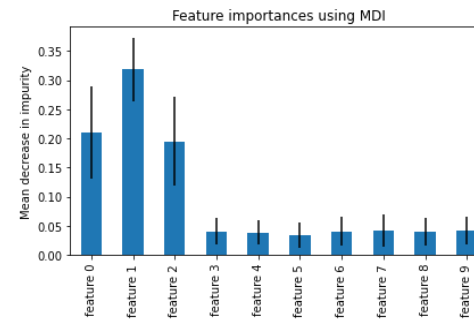
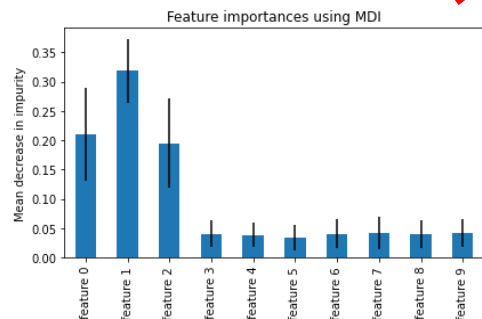
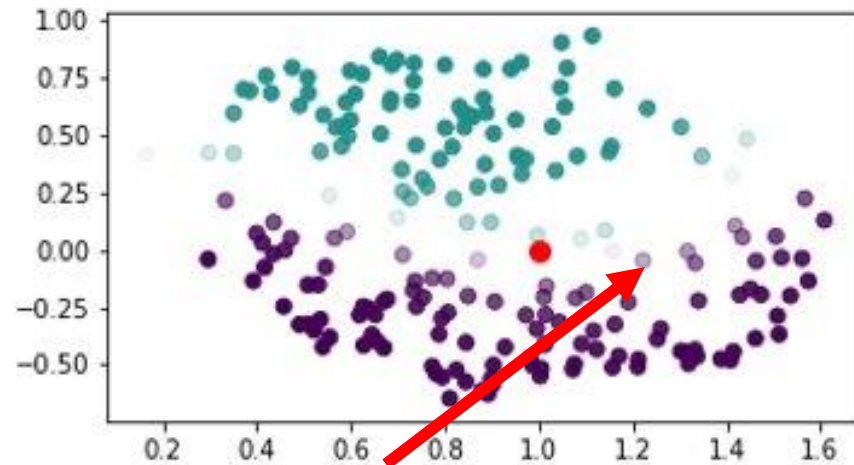
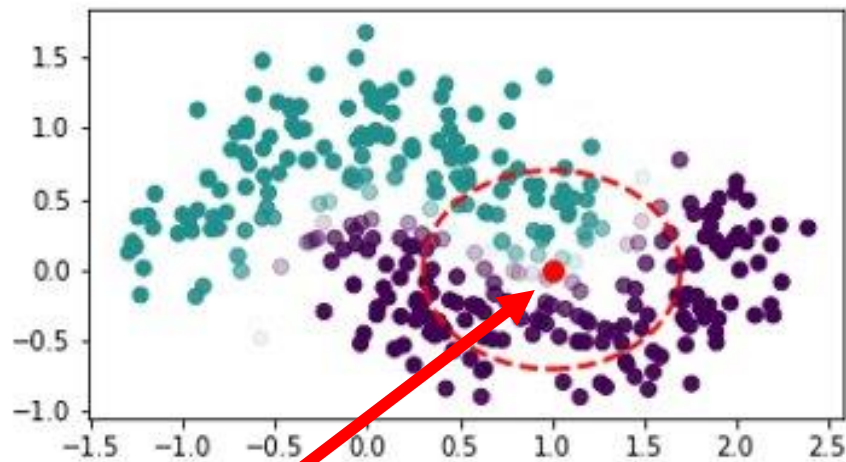
Consistency between explanations for different models (or explainers)

$$C(\Phi^{e \rightarrow m_1}, \Phi^{e \rightarrow m_2}, \dots, \Phi^{e \rightarrow m_n}) = \frac{1}{\max_{a,b \in m_1, m_2, \dots, m_n} \|\Phi_j^{e \rightarrow m_a} - \Phi_j^{e \rightarrow m_b}\|_2 + 1}$$

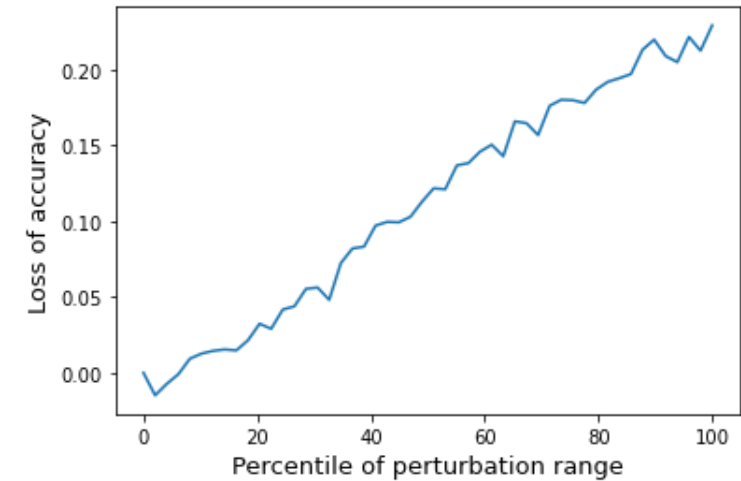
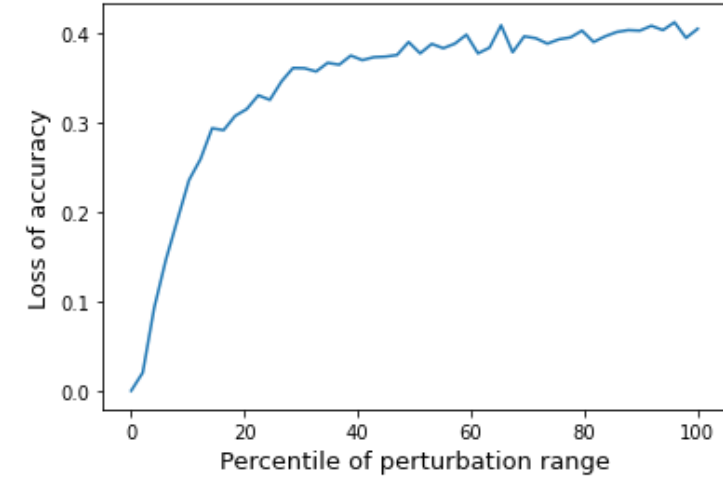
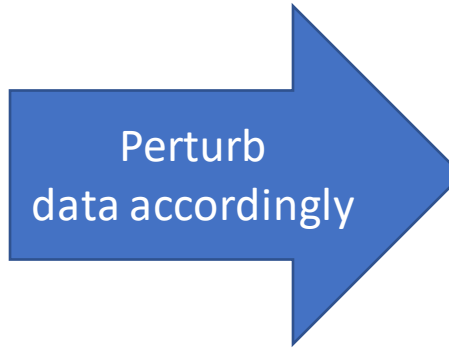
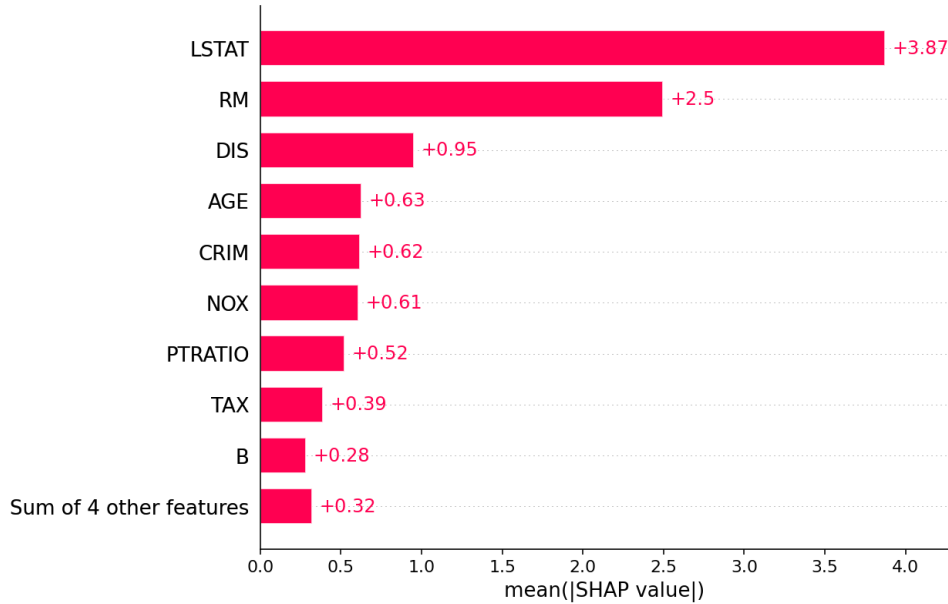


Stability of explanations for similar instances

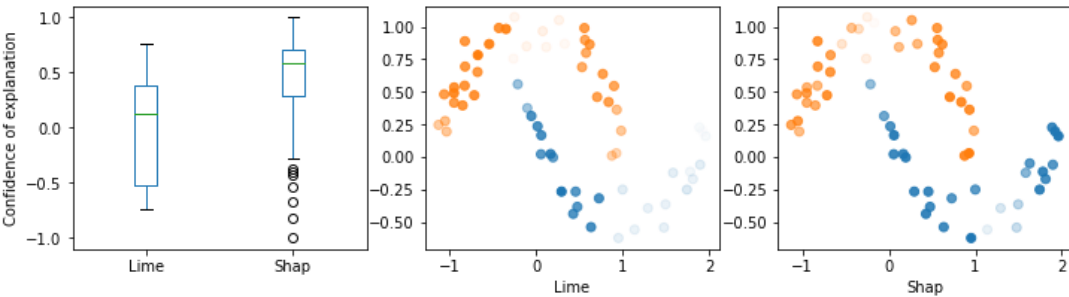
$$\hat{L}(\Phi^{e \rightarrow m}, X) = \max_{x_j \in N_\epsilon(x_i)} \frac{\|x_i - x_j\|_2}{\|\Phi_i^{e \rightarrow m} - \Phi_j^{e \rightarrow m}\|_2 + 1}$$



Quality Loss (AUCx)



Ensemble explanations

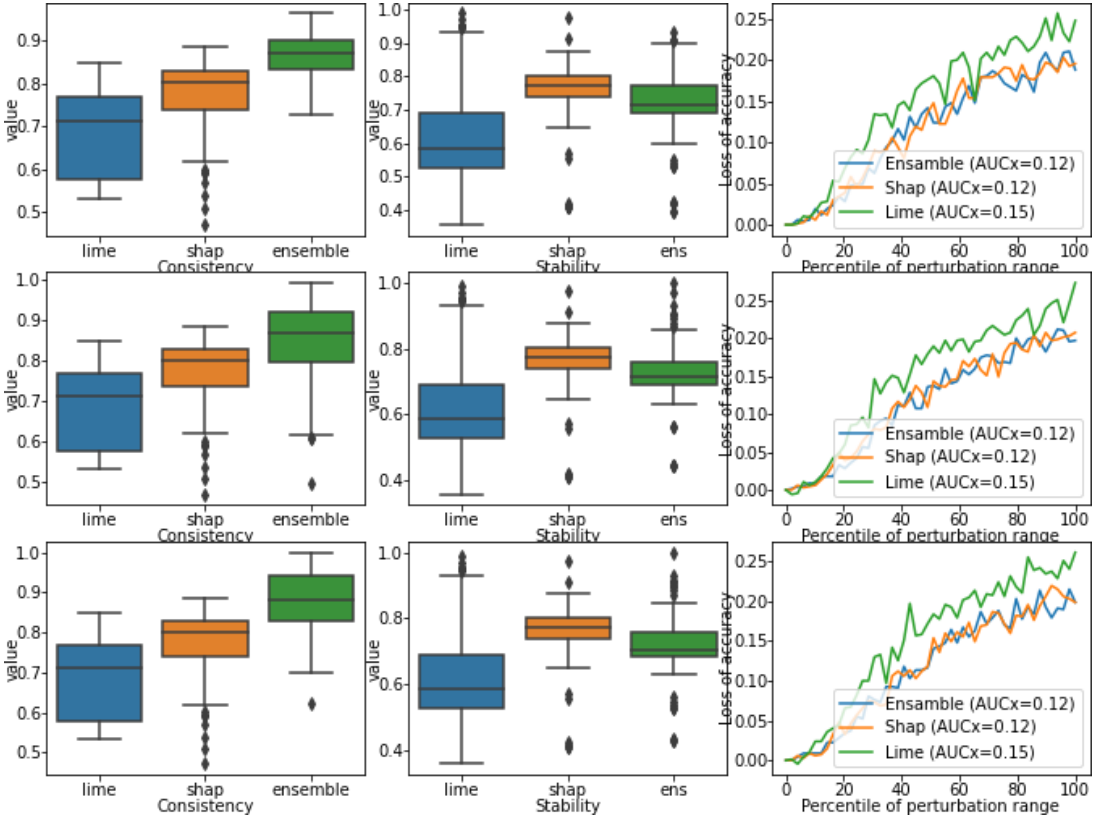


$$C(\Phi^{e \rightarrow m_1}, \Phi^{e \rightarrow m_2}, \dots, \Phi^{e \rightarrow m_n}) = \frac{1}{\max_{a,b \in m_1, m_2, \dots, m_n} \|\Phi_j^{e \rightarrow m_a} - \Phi_j^{e \rightarrow m_b}\|_2 + 1}$$

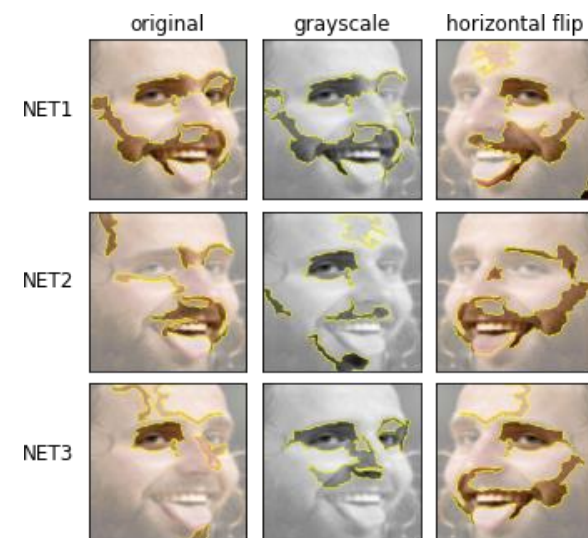
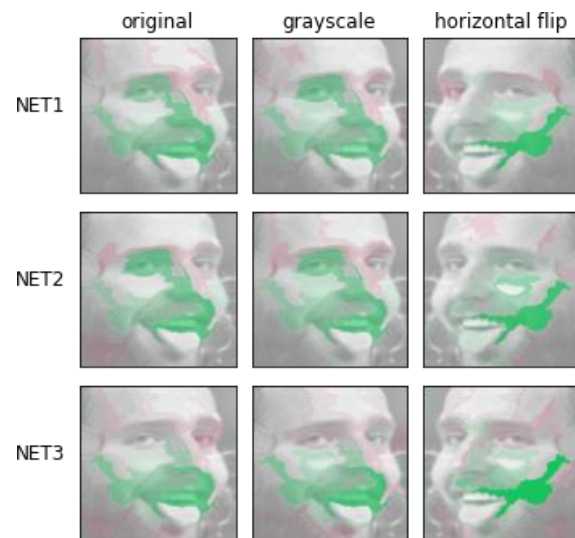
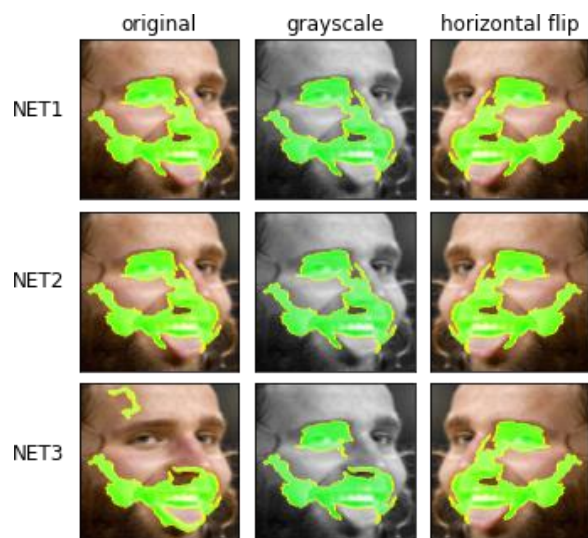
$$\hat{L}(\Phi^{e \rightarrow m}, X) = \max_{x_j \in N_\epsilon(x_i)} \frac{\|x_i - x_j\|_2}{\|\Phi_i^{e \rightarrow m} - \Phi_j^{e \rightarrow m}\|_2 + 1}$$

$$ES(M, w) = \sum w_i \cdot M_i$$

$$\Phi^{ens} = \frac{ES(M, w) \cdot [\gamma_1 \Phi^{e_1}, \gamma_2 \Phi^{e_2}, \dots, \gamma_n \Phi^{e_n}]}{\sum_{i=1}^n ES_i(M, w)}$$

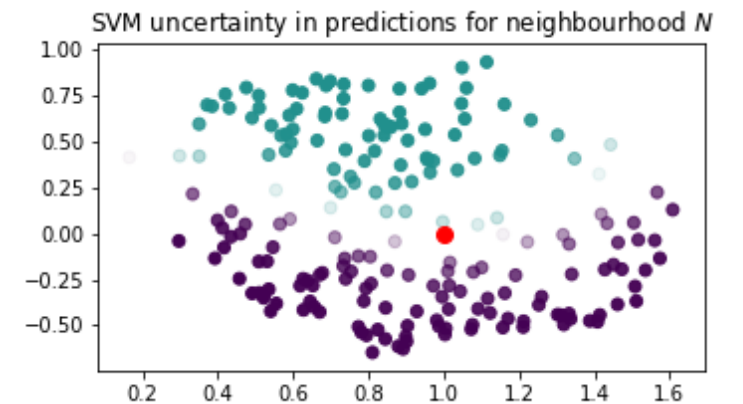
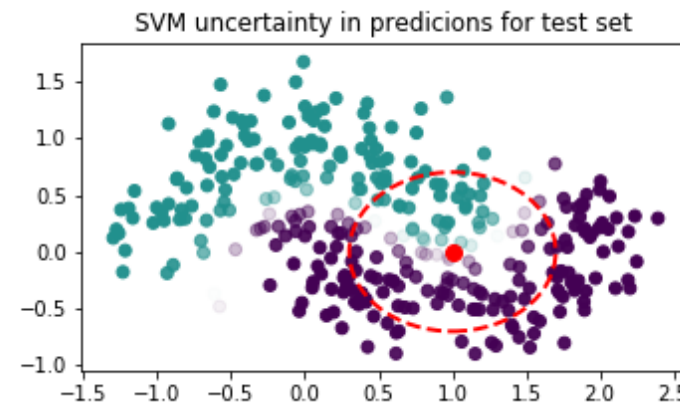
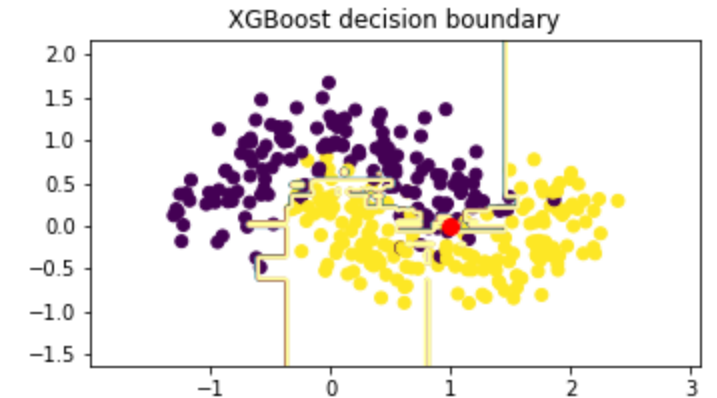
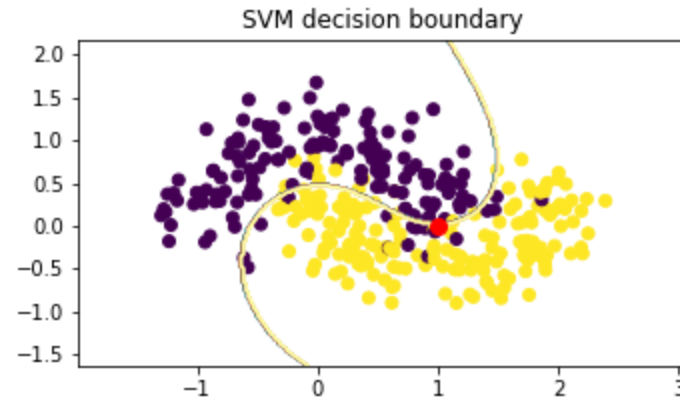


Ensemble explanations



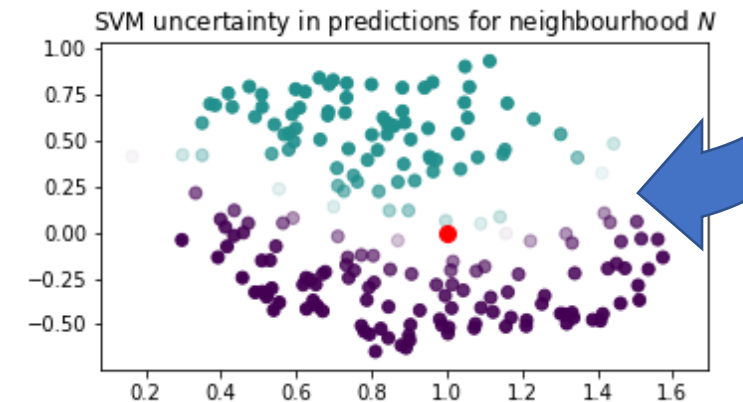
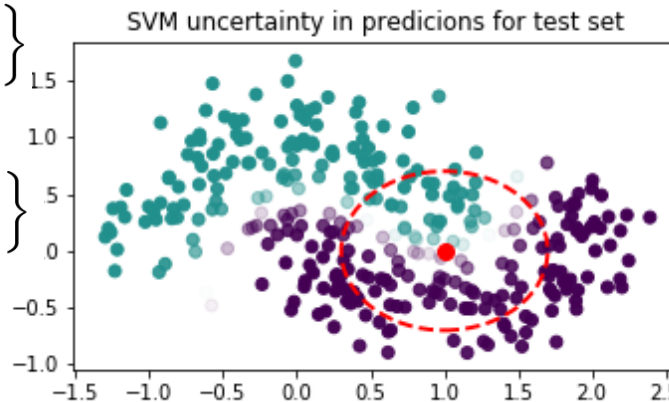
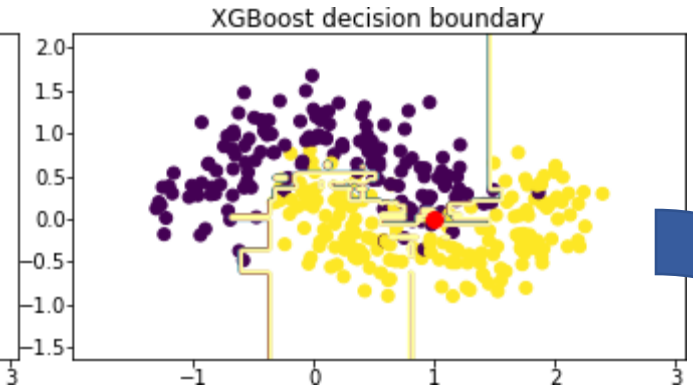
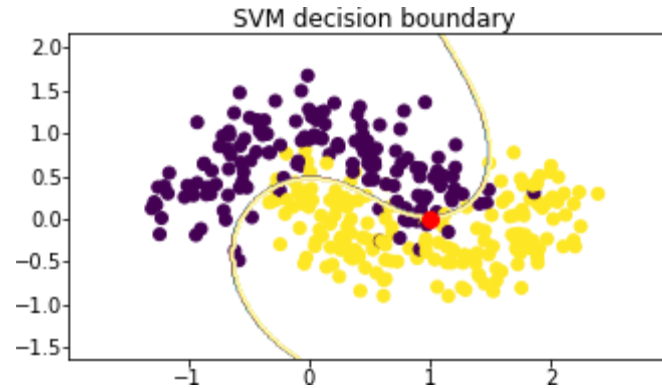
Local Uncertain Explanations

- Machine learning model we explain is **uncertain**
- **Target variable** for training local interpretable model is **uncertain**
- **Neighbourhood** (training data for local explainable model) is **uncertain**



Local Uncertain Explanations

- We use neighbourhood as uncertain dataset
- Instead of probability calculated as frequency, we average the probability obtained from ML model
- We modify Information Gain split criterion to use these measure and build decision tree



$$N(x^{(i)}, K) = \left\{ x^{(k)} \in X : d(x^{(i)}, x^{(k)}) \leq D_i^{(K)} \right\}$$

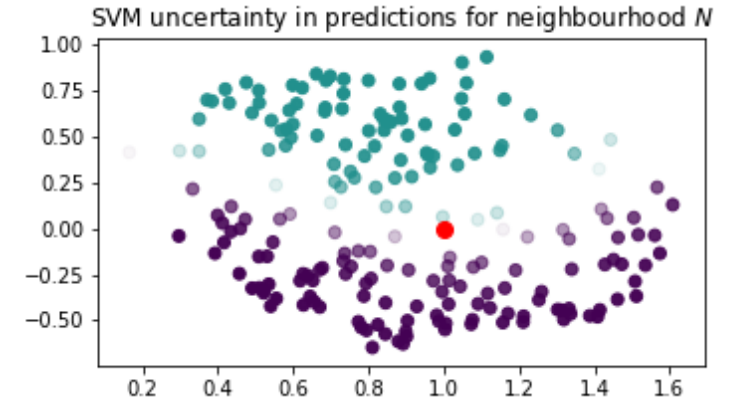
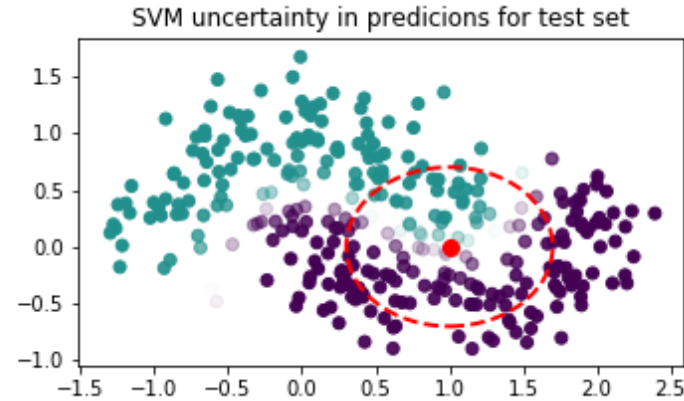
$$D_i = \left\{ d(x^{(i)}, x^{(1)}), d(x^{(i)}, x^{(2)}), \dots, d(x^{(i)}, x^{(m)}) \right\}$$

Prediction

S

Local Uncertain eXplanations (LUX)

- We use neighbourhood as uncertain dataset
- Instead of probability calculated as frequency, we average the probability obtained from ML model
- We modify Information Gain split criterion to use these measure and build decision tree

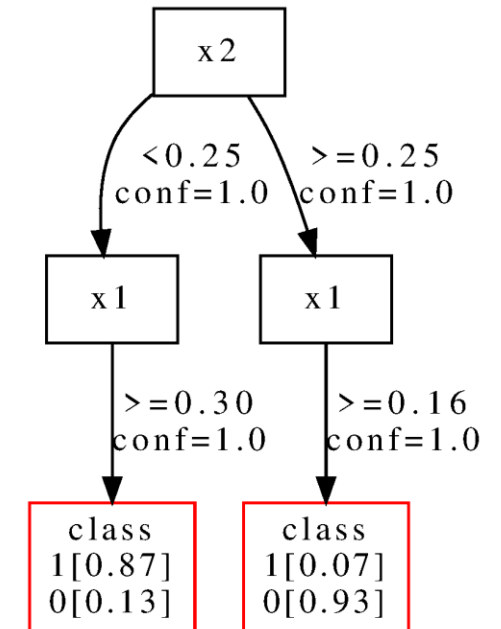


$$H^U(X) = - \sum_{v \in \text{Domain}(C)} P_{total}(C = v) \log_2 P_{total}(C = v)$$

$$P_{total}(A_i = v_j^i) = \frac{1}{k} \sum_{X_t \ni P_j=1 \dots n} P_j(A_i = v_j^i)$$

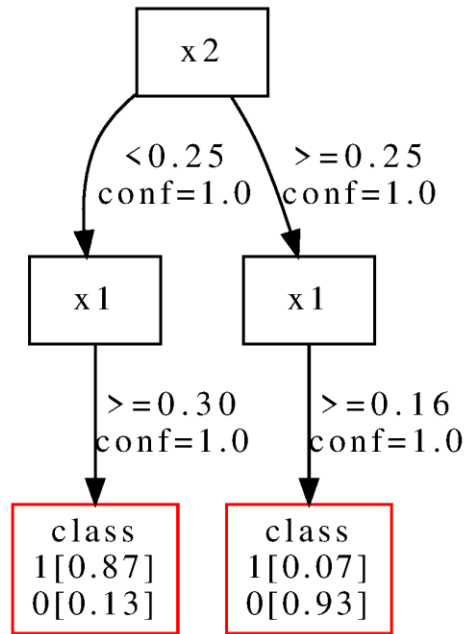
$$\text{Gain}^U(A) = H^U(X) - \sum_{v \in \text{Domain}(A)} P_{total}(A = v) H^U(X_v)$$

```
@relation lux
@attribute x1 @REAL
@attribute x2 @REAL
@attribute class {1,0}
@data
0.94,0.01,1[0.48]
0.87,-0.04,1[0.64]
1.02,-0.16,1[0.78]
1.14,0.08,1[0.37]
1.01,-0.21,1[0.83]
1.10,-0.19,1[0.81]
0.80,-0.13,1[0.81]
0.91,-0.23,1[0.87]
0.77,-0.12,1[0.83]
1.01,-0.28,1[0.89]
0.97,-0.28,1[0.89]
...
```



<https://github.com/sbobek/lux>

Local Uncertain Explanations (LUX)



- Translate branches to rules (XTT2 format)
- XTT2 format is a **knowledge** representation that is
 - **extensible** (HWEd editor),
 - **formalized** (ALSV(FD) logic),
 - **executable** (HeaRTDrtoid engine)
- **Get the uncertainty of an explanation that is transferred from uncertainty of data and model prediction**

x1	x2	class	#
≥ 0.30	< 0.25	set 1	0,8
≥ 0.16	≥ 0.25	set 0	0,9

tree

Add condition Add decision Add rule

Summary

- Knowledge Augmented Clustering (KnAC)
- Local Uncertain eXplanbations (LUX)
- Intelligible XAI (InXAI)
- Technology needs to be human-centric
- Explanations are important for unsupervised methods (KnAC/Explainable clusters)
- *The truth is out there*

Open Challenges in XAI for (not only) Industry 4.0

- Mediating explanations between human and XAI system.
 - Explanation is an act of conveying knowledge
 - Technology needs to be human-centric. Good explanation does not always mean useful or understandable
- Defining mediatable information granules via human-in-the-loop conceptualization.
 - Semantic gap between XAI and different explanation addressees (stakeholders)
- Multi-faced continuous assessment of quality of explanations.
 - Why should I trust... your explanation
 - Correlation does not mean causation

Thank you for your attention!

Give us a feedback @ <https://github.com/sbobek/knac>

