

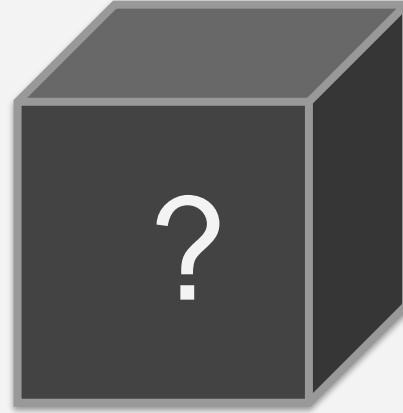
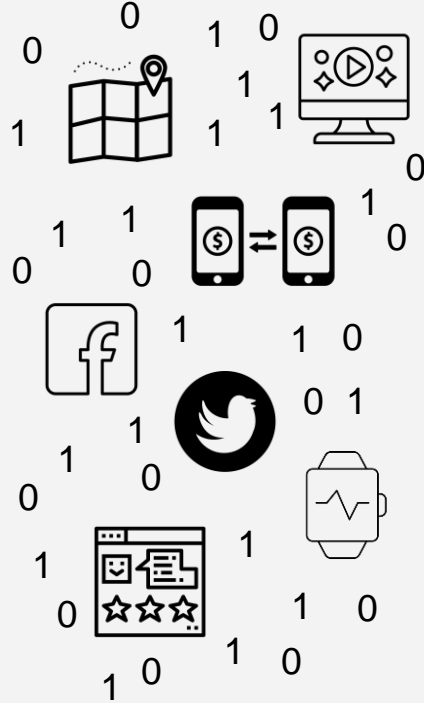
**(Binary)  
Prediction**

### **Black box model**

- ⇒ Thousands of coefficients
- ⇒ Nonlinear techniques

**Behavioral and textual data  
(High-dimensional & sparse)**

# "Explainable AI" or "Interpretable ML"



**(Binary)  
Prediction**

## **Black box model**

- ⇒ Thousands of coefficients
- ⇒ Nonlinear techniques

**Behavioral and textual data  
(High-dimensional & sparse)**

# ABOUT ME



YANOU RAMON



PhD student at University of Antwerp  
Applied Data Mining Group – Prof. dr. ir. David Martens



Towards interpretable classification models built  
from high-dimensional, sparse data



Master's degree in Business Engineering  
University of Antwerp

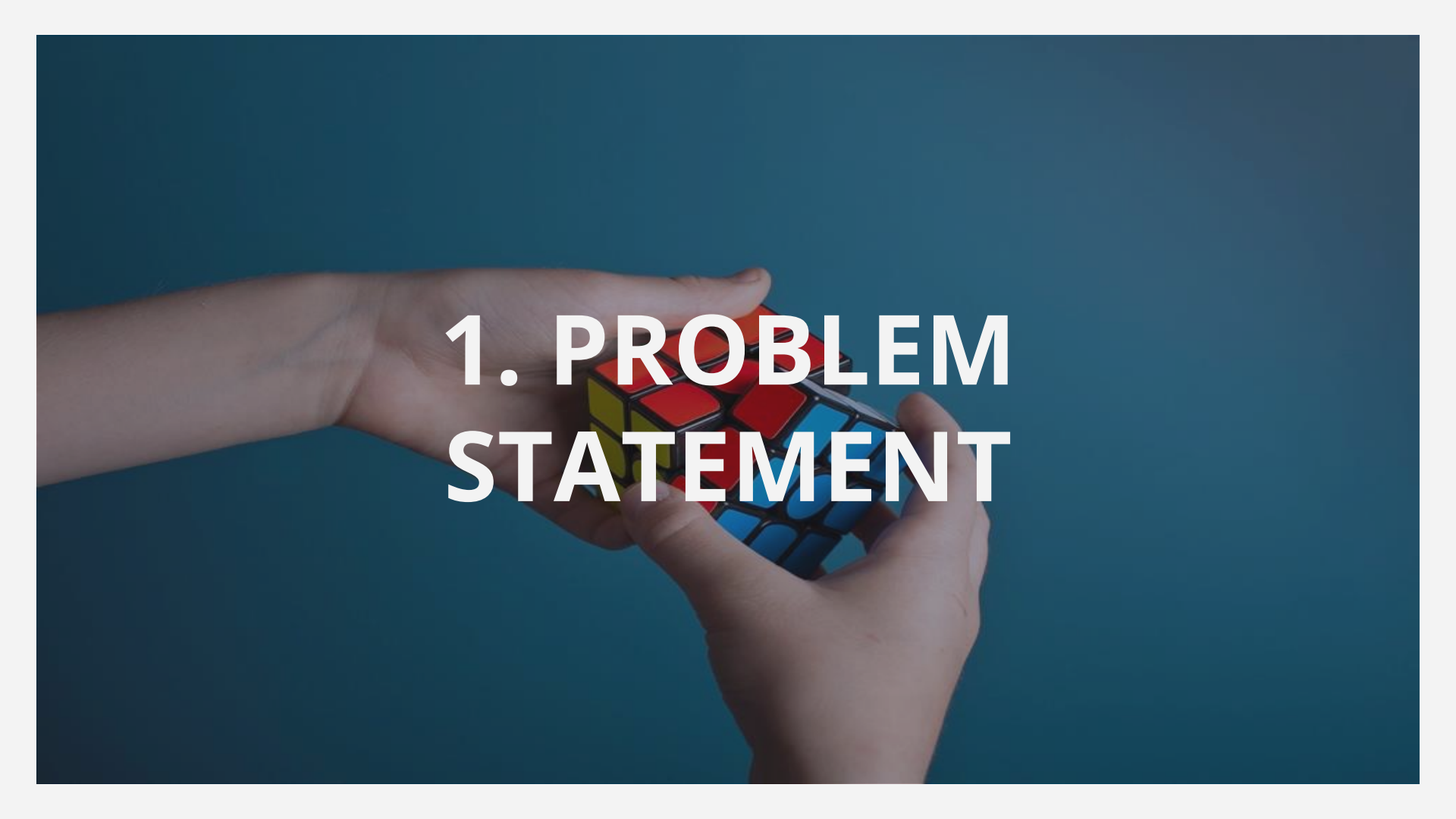


(Big) Data Mining, Artificial Intelligence,  
programming etc.



# FIRST RESEARCH PROJECT

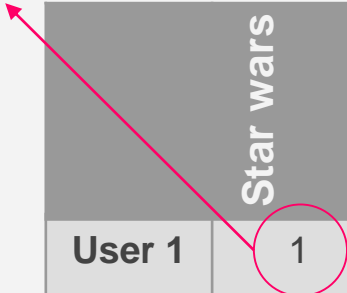
# INSTANCE-LEVEL EXPLANATION ALGORITHMS ON BEHAVIORAL AND TEXTUAL DATA: A COUNTERFACTUAL-ORIENTED COMPARISON

A person's hands are shown holding a Rubik's cube, which is partially solved. The background is a solid dark blue. The text "1. PROBLEM STATEMENT" is overlaid in white, bold, sans-serif font.

# 1. PROBLEM STATEMENT

# MOVIE VIEWING DATA (MovieLens)

Active feature = "evidence"



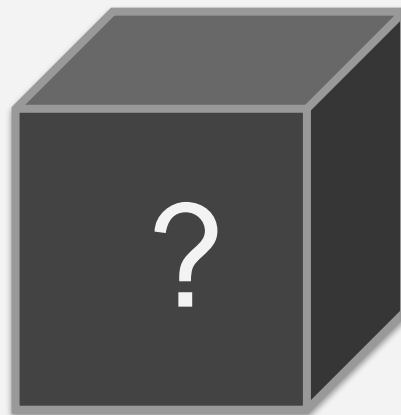
6,040 users

	Star wars	Pearl Harbor	Django	...	Home Alone	Target $\hat{y}$ Gender
User 1	1	0	0	...	1	<i>M</i>
User 2	1	1	0	...	0	<i>F</i>
...	...	...	...	...	...	...
User n	1	0	0	...	0	<i>M</i>

Sparsity  $p = 95,53\%$

	Star wars	Pearl Harbor	Django	...	Home Alone	Target $\hat{y}$ Gender
User 1	1	0	0	...	1	<i>M</i>
User 2	1	1	0	...	0	<i>F</i>
...	...	...	...	...	...	...
User n	1	0	0	...	0	<i>M</i>

**Movie Viewing Data (MovieLens)**



**Black box model**

- ⇒ Thousands of coefficients
- ⇒ Nonlinear techniques



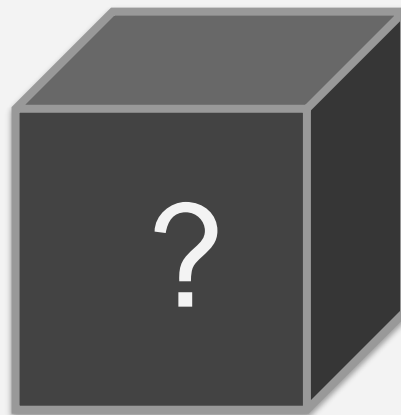
$$\hat{y} = 1 \text{ *if male*}$$

$$\text{else } \hat{y} = 0$$

**Comprehensibility issues**



	Star wars	Pearl Harbor	Django	...	Home Alone	Target $\hat{y}$ Gender
User 1	1	0	0	...	1	<i>M</i>
User 2	1	1	0	...	0	<i>F</i>
...	...	...	...	...	...	...
User n	1	0	0	...	0	<i>M</i>



$$\hat{y} = 1 \text{ if male} \\ \text{else } \hat{y} = 0$$

**Movie Viewing Data (MovieLens)**

**Black box model**

⇒ Thousands of coefficients

⇒ Nonlinear techniques

**INSTANCE-LEVEL EXPLANATIONS:  
Why relevant?**

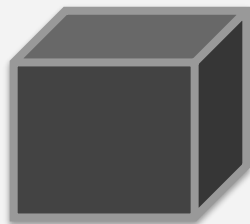
# WHY EXPLANATIONS?

- **Trust and acceptance**

Example: explainable legal document classification (Chhatwal et al., 2019)



Legal documents



Black box model

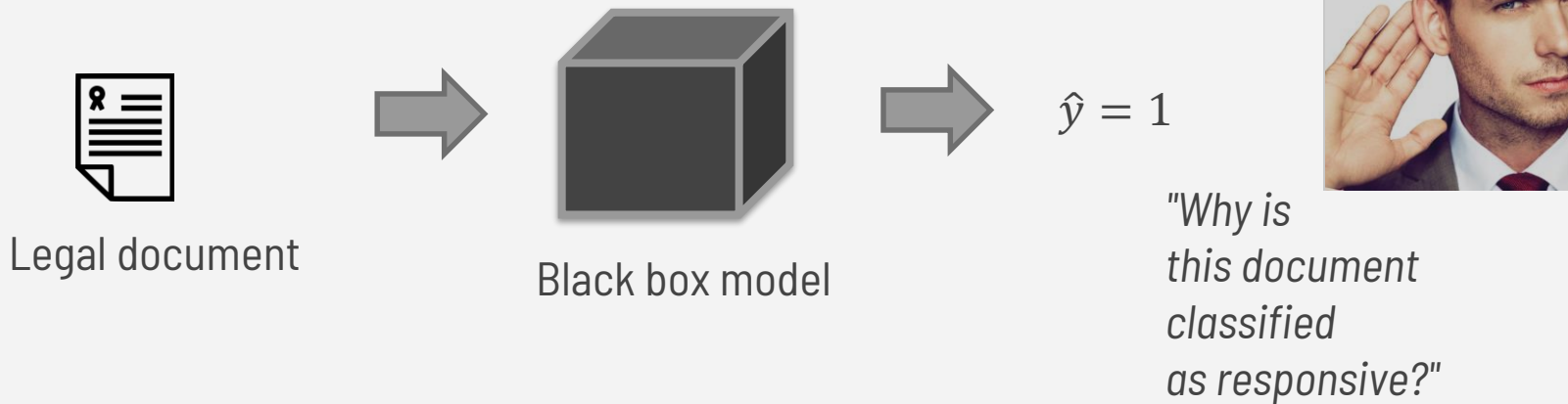


$\hat{y} = 1$  *if responsive*  
*else*  $\hat{y} = 0$

# WHY EXPLANATIONS?

- **Trust and acceptance**

Example: explainable legal document classification (Chhatwal et al., 2019)



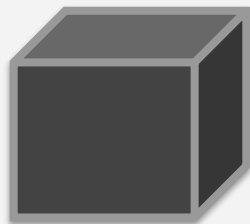
# WHY EXPLANATIONS?

- **Improving the model: explain misclassifications**

Example: objectionable web content detection (Martens & Provost, 2013)



Web pages



Black box model



$\hat{y} = 1$  **if** objectionable  
else  $\hat{y} = 0$

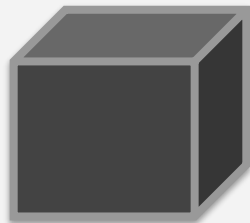
# WHY EXPLANATIONS?

- **Generate insights**

Example: Know your customer (e.g., Hall, 2012; Grossnickle, 2001)



Visited URLs



Black box model



$\hat{y} = 1$  *if interested in product*  
else  $\hat{y} = 0$

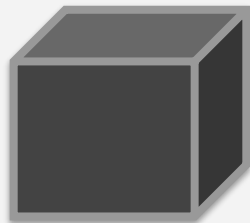
# WHY EXPLANATIONS?

- **Generate insights**

Example: Know your customer (e.g., Hall, 2012; Grossnickle, 2001)



Visited URLs



Black box model



*"Who are we targeting?  
Why are we targeting them?"*

# COUNTERFACTUAL EXPLANATIONS

- Instance-level
- Rule: a minimal set of features such that the predicted class changes when “removing” them (~setting value to zero)
- Comprehensible and concise
- Argued to be the most intuitive and valuable for humans because they are contrastive (*“Why X rather than not-X?”; Miller, 2017*)

# COUNTERFACTUAL EXPLANATIONS

**Example:** gender prediction using movie viewing data



Sam watched 120 movies  
Sam was predicted as 'male'

	Star wars	Pearl Harbor	Django	...	Home Alone	Target $\hat{y}$ Gender
User 1	1	0	0	...	1	<b>M</b>
User 2	1	1	0	...	0	<b>F</b>
...	...	...	...	...	...	...
User n	1	0	0	...	0	<b>M</b>



# COUNTERFACTUAL EXPLANATIONS

**Example:** gender prediction using movie viewing data



Sam watched 120 movies  
Sam was predicted as 'male'

	Star wars	Pearl Harbor	Django	...	Home Alone	Target $\hat{y}$ Gender
User 1	1	0	0	...	1	<b>M</b>
User 2	1	1	0	...	0	<b>F</b>
...	...	...	...	...	...	...
User n	1	0	0	...	0	<b>M</b>

**Why?**

# COUNTERFACTUAL EXPLANATIONS

**Example:** gender prediction using movie viewing data



Sam watched 120 movies

Sam was predicted as 'male'

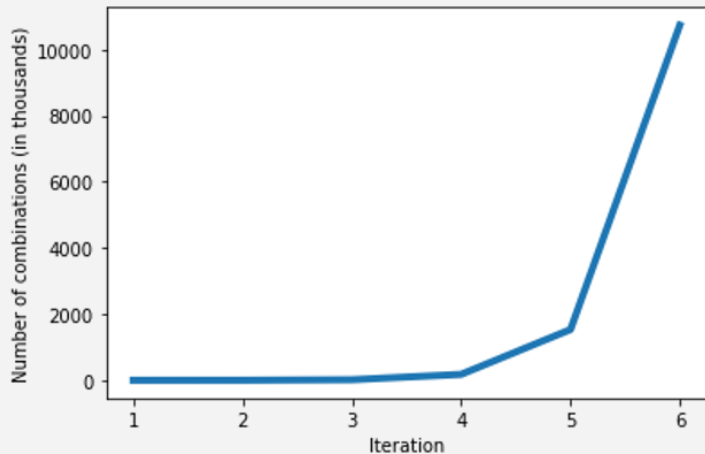
**IF** Sam would not have watched **{Taxi Driver, The Dark Knight, Die Hard, Terminator 2, Now You See Me, Interstellar}**,  
**THEN** the predicted class changes from 'male' to 'female'

# WHY COMPLETE SEARCH FAILS

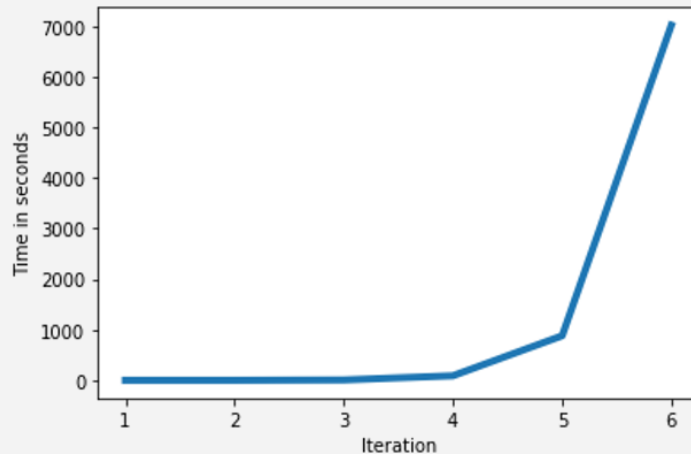
- Start with removing one feature and increase number of features in the subset until the predicted class changes
- Scales exponentially with active features  $m$  and required number of features  $k$  to be removed  
e.g., for an instance with  $m$  features, a combination of  $k$  features requires  $\frac{m!}{(m-k)!k!}$  evaluations

# WHY COMPLETE SEARCH FAILS

**Figure 1: Number of combinations (a) and time elapsed (b) per iteration for an instance with 34 active features and a counterfactual of 6 features (*MovieLens data*)**



(a)



(b)

## 2. COUNTERFACTUAL ALGORITHMS



# ALGORITHMIC ASSUMPTIONS

- **Goal:** find counterfactual explanation as fast and as concise as possible (efficiency-effectiveness tradeoff)
- Model-agnostic
- Max. 30 features in explanation
- Max. 5 minutes to compute explanation

# BEST-FIRST SEARCH (SEDC)

- Explaining document classifications (Martens & Provost, 2013)
- Model-agnostic algorithm SEDC: heuristic best-first search (*lin-SEDC*: linear implementation)
- Optimal for linear models

# BEST-FIRST SEARCH (SEDC)

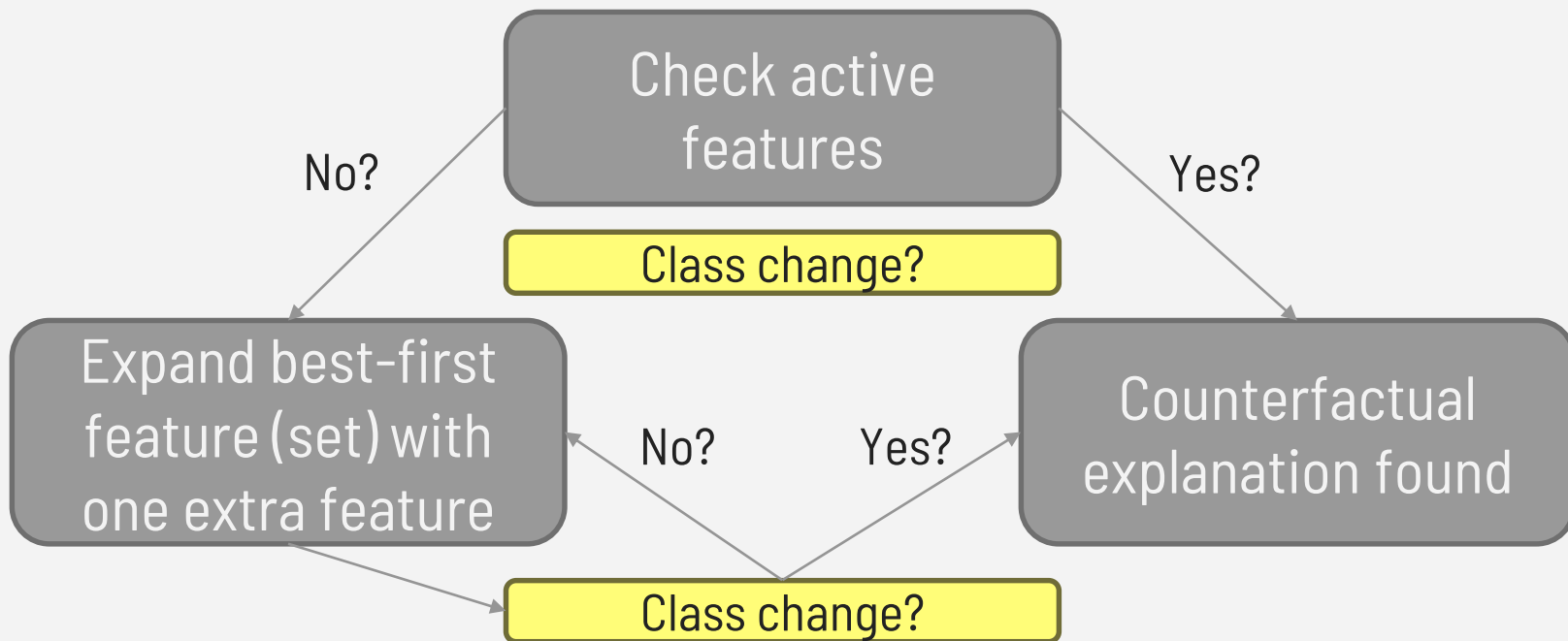
- Explaining document classifications (Martens & Provost, 2013)
- Model-agnostic algorithm SEDC: heuristic best-first search (*lin-SEDC*: linear implementation)
- Optimal for linear models



Implementation on <https://github.com/yramon/edc>



# BEST-FIRST SEARCH (SEDC)



# NOVEL HYBRID ALGORITHMS

## **Additive Feature Attribution methods:**

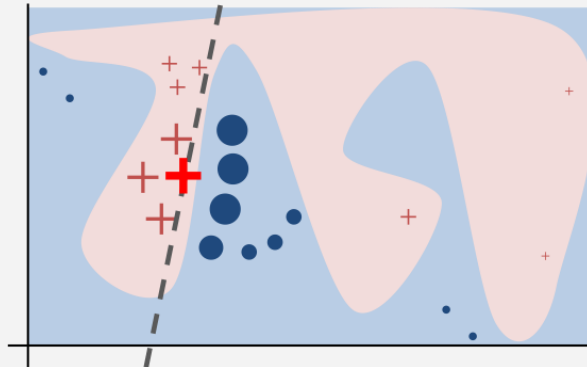
- LIME: Local Model-agnostic Explainer (Ribeiro et al., 2016)
- SHAP: Shapley Additive Explanations (Lundberg et al., 2018)

**Output:** importance-ranked list

# NOVEL HYBRID ALGORITHMS

## LIME / SHAP

- Sparse, linear explanation model
- Approximates original model in neighbourhood of instance
- Perturbed instances

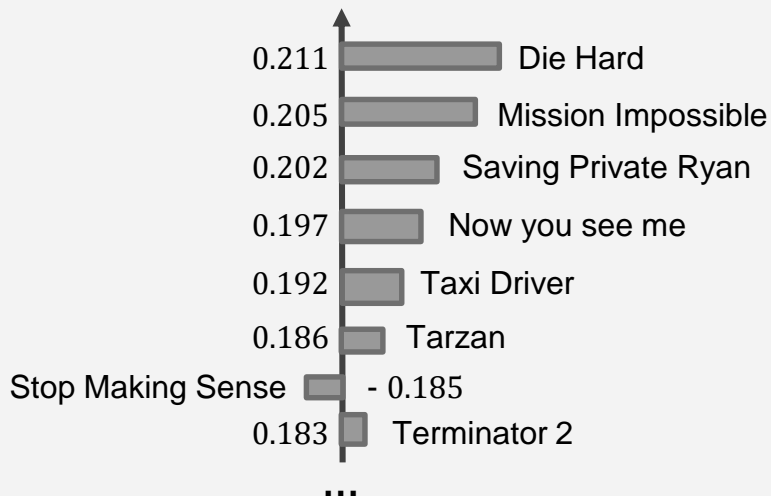


*Ribeiro et al., 2016*

# NOVEL HYBRID ALGORITHMS

## LIME / SHAP

**Example:** gender prediction using movie viewing data



# NOVEL HYBRID ALGORITHMS

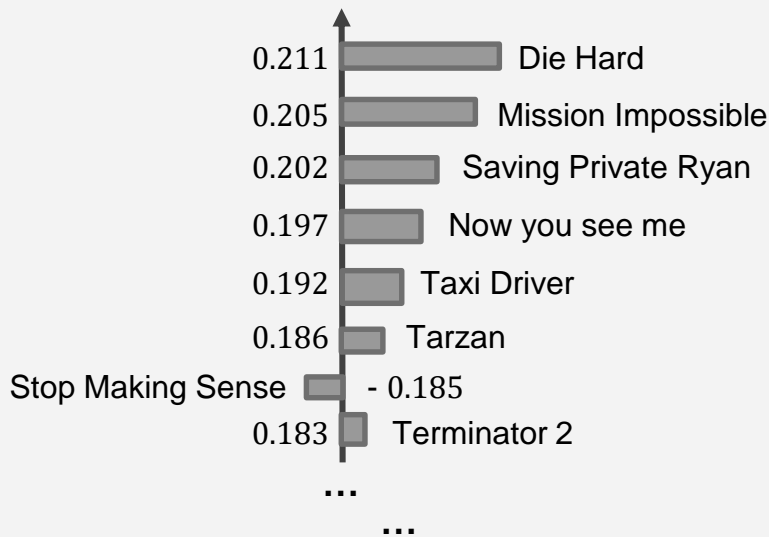
**Originality:** importance rankings may be “intelligent” starting point for efficiently searching counterfactuals

⇒ Novel algorithms: **LIME-C** and **SHAP-C**

# NOVEL HYBRID ALGORITHMS

## LIME-C / SHAP-C

**Example:** gender prediction using movie viewing data



Remove features with positive importance weight until the class changes

# CONTRIBUTIONS

- Two novel model-agnostic algorithms (LIME-C / SHAP-C)
- Define quantitative evaluation criteria
- Evaluate performance against existing SEDC algorithm and make practical recommendations

A background image featuring a glass hourglass with blue sand on a desk. In the background, a laptop and a pair of glasses are visible, all under a warm, orange-toned light. The text '3. EXPERIMENTAL SETUP' is overlaid in the center in a bold, white, sans-serif font.

# **3. EXPERIMENTAL SETUP**



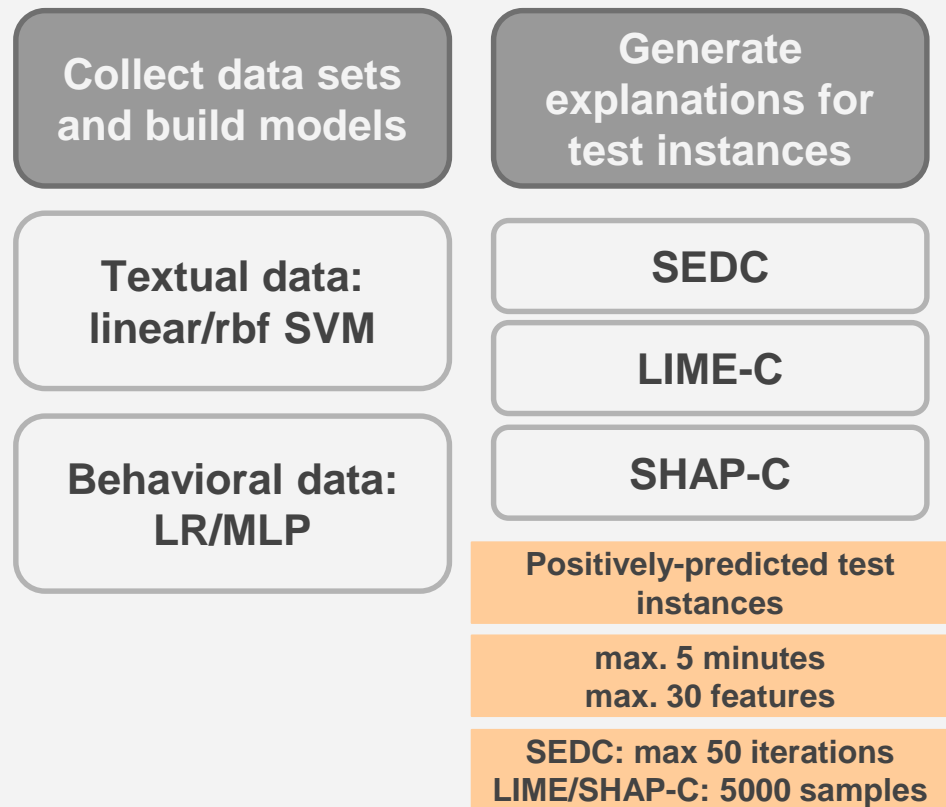
Collect data sets  
and build models

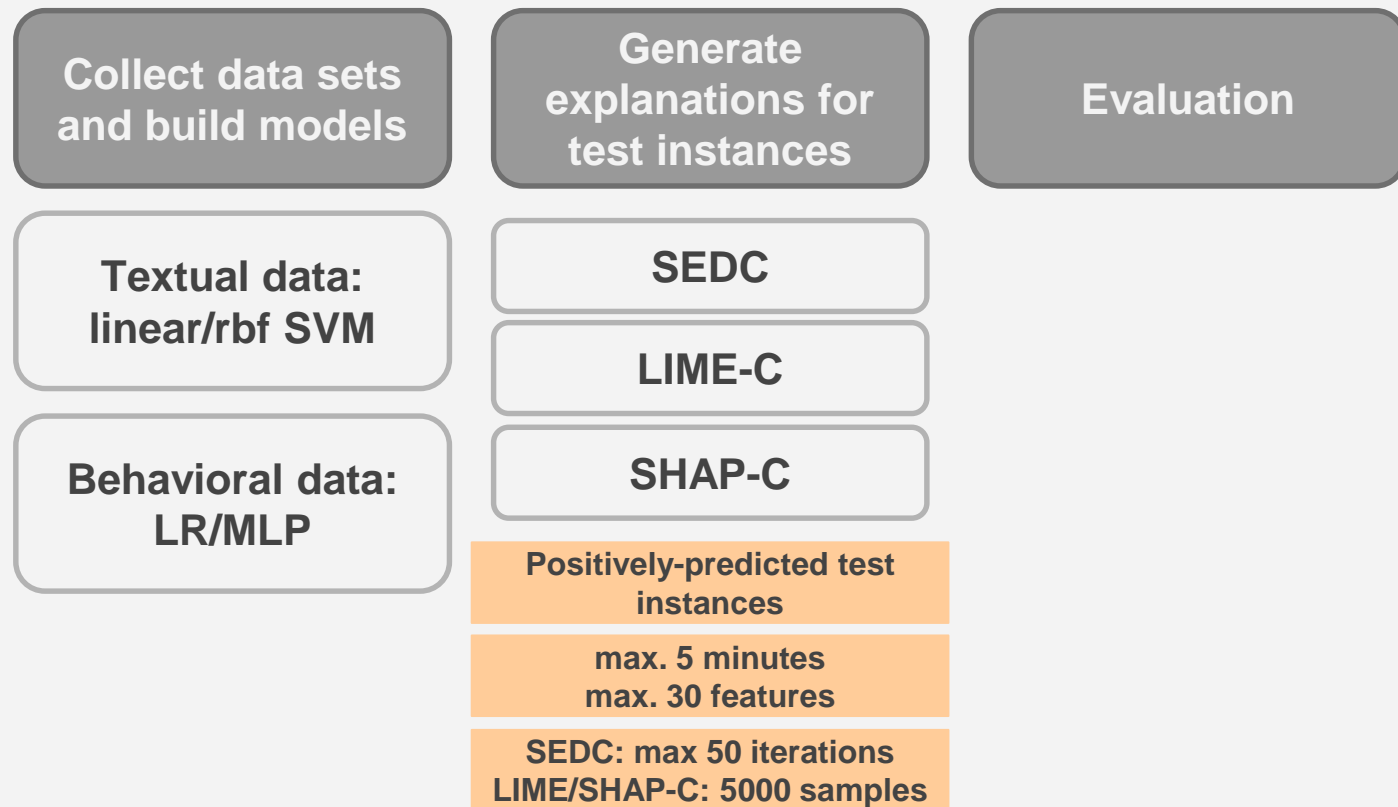
Textual data:  
linear/rbf SVM

Behavioral data:  
LR/MLP

**Table 1: Data sets and characteristics**

Dataset	Type	Target	Instances	Features	$b$	$p$	Test set (%)	$\hat{m}_{lin}$	$\hat{m}_{nonlin}$	ref
Flickr*	B	comments	100,000	190,991	36.91%	99.99%	20,000 (20%)	2.02	2.96	[38]
Ecommerce*	B	gender	15,000	21,880	21.98%	99.99%	3,000 (15%)	2.60	2.67	[3]
Airline*	T	sentiment	14,640	5,183	16.14%	99.82%	2,928 (15%)	7.81	8.21	[2]
Twitter	T	topic	6,090	4,569	9.15%	99.74%	1,218 (10%)	9.52	9.35	[5]
Fraud*	B	fraudulent	858,131	107,345	6.4e-5%	99.99%	171,627(1%)	11.83	14.09	n.a.
YahooMovies*	B	gender	7,642	11,915	28.87%	99.76%	1,529 (20%)	25.24	25.00	[6]
TaFeng*	B	age	31,640	23,719	45.23%	99.90%	6,328 (15%)	44.32	37.24	[23]
KDD2015*	B	dropout	120,542	4,835	20.71%	99.67%	24,109 (20%)	49.01	46.40	[4]
20news	T	atheism	18,846	41,356	4.24%	99.84%	3,770 (5%)	67.96	62.77	[1]
Movielens.100k	B	gender	943	1,682	28.95%	93.69%	189 (25%)	68.73	73.42	[21]
Facebook*	B	gender	386,321	122,924	44.57%	99.94%	77,265 (30%)	83.03	84.55	[9]
Movielens.1m*	B	gender	6,040	3,706	28.29%	95.53%	1,208 (25%)	168.46	153.46	[21]
Libimseti*	B	gender	137,806	166,353	44.53%	99.93%	27,562 (30%)	229.16	226.97	[8]





# EVALUATION CRITERIA

The **goal** is to find a small-sized counterfactual as fast as possible → **tradeoff** between

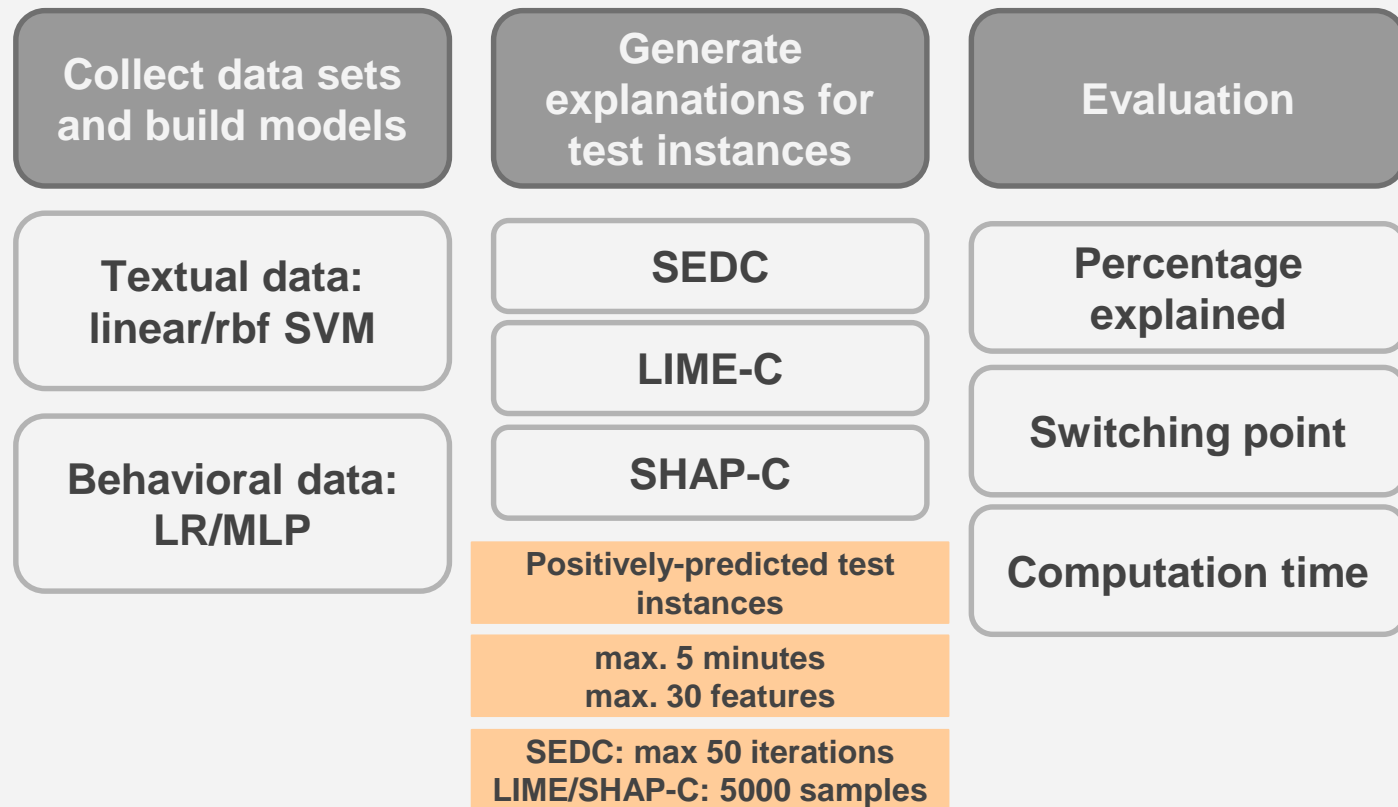
- **Effectiveness**

- Percentage explained

- Switching point: # features in explanation

- **Efficiency**

- Computation time in seconds



A low-angle shot of a basketball player in a black jersey jumping high to shoot a basketball. The player's right arm is extended upwards, holding the ball, while their left arm is also raised. A red basketball hoop is visible at the top of the frame. The background is a bright, overcast sky with soft, white clouds. The text "4. RESULTS & CONCLUSION" is overlaid in the center in a large, white, sans-serif font.

# 4. RESULTS & CONCLUSION

# EFFECTIVENESS

**Table 2: Percentage explained**

Dataset	Linear			Nonlinear		
	SEDC (%)	LIME-C (%)	SHAP-C (%)	SEDC (%)	LIME-C (%)	SHAP-C (%)
Flickr	100	100	100	28.67	28.33	28.67
Ecommerce	100	97.33	100	95.00	97.00	99.67
Airline	100	100	100	100	100	100
Twitter	100	100	100	100	100	100
Fraud	100	100	81.67	100	100	75
YahooMovies	100	100	100	98.67	100	100
TaFeng	100	100	100	93.33	100	100
KDD2015	100	100	100	99.67	100	99.67
20news	99.47	99.47	100	99.47	98.94	100
MovieLens_100k	100	100	100	100	100	100
Facebook	95.67	95.00	95.00	70.33	92.67	89.67
MovieLens_1m	98.67	98.67	98.67	88.33	95.00	95.67
Libimseti	92.67	90.33	88.67	77.00	81.67	72.33
Average	98.96	98.52	97.23	88.49	91.82	89.28
# wins	12	9	10	5	9	9

# EFFECTIVENESS

**Table 2: Percentage explained**

Dataset	Linear			Nonlinear		
	SEDC (%)	LIME-C (%)	SHAP-C (%)	SEDC (%)	LIME-C (%)	SHAP-C (%)
Flickr	<b>100</b>	<b>100</b>	<b>100</b>	<b>28.67</b>	28.33	<b>28.67</b>
Ecommerce	<b>100</b>	97.33	<b>100</b>	<u>95.00</u>	<u>97.00</u>	<b>99.67</b>
Airline	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
Twitter	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
Fraud	<b>100</b>	<b>100</b>	<u>81.67</u>	<b>100</b>	<b>100</b>	<u>75</u>
YahooMovies	<b>100</b>	<b>100</b>	<b>100</b>	98.67	<b>100</b>	<b>100</b>
TaFeng	<b>100</b>	<b>100</b>	<b>100</b>	<u>93.33</u>	<b>100</b>	<b>100</b>
KDD2015	<b>100</b>	<b>100</b>	<b>100</b>	99.67	<b>100</b>	99.67
20news	99.47	99.47	<b>100</b>	99.47	98.94	<b>100</b>
MovieLens_100k	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
Facebook	<b>95.67</b>	95.00	95.00	<u>70.33</u>	<b>92.67</b>	<u>89.67</u>
MovieLens_1m	<b>98.67</b>	<b>98.67</b>	<b>98.67</b>	<u>88.33</u>	95.00	<b>95.67</b>
Libimseti	<b>92.67</b>	<u>90.33</u>	<u>88.67</u>	<u>77.00</u>	<b>81.67</b>	<u>72.33</u>
Average	<b>98.96</b>	98.52	97.23	88.49	<b>91.82</b>	89.28
# wins	<b>12</b>	9	10	5	<b>9</b>	<b>9</b>



# EFFECTIVENESS

**Table 2: Percentage explained**

Dataset	Linear			Nonlinear		
	SEDC (%)	LIME-C (%)	SHAP-C (%)	SEDC (%)	LIME-C (%)	SHAP-C (%)
Flickr	100	100	100	28.67	28.33	28.67
Ecommerce	100	97.33	100	95.00	97.00	99.67
Airline	100	100	100	100	100	100
Twitter	100	100	100	100	100	100
Fraud	100	100	81.67	100	100	75
YahooMovies	100	100	100	98.67	100	100
TaFeng	100	100	100	93.33	100	100
KDD2015	100	100	100	99.67	100	99.67
20news	99.47	99.47	100	99.47	98.94	100
MovieLens_100k	100	100	100	100	100	100
Facebook	95.67	95.00	95.00	70.33	92.67	89.67
MovieLens_1m	98.67	98.67	98.67	88.33	95.00	95.67
Libimseti	92.67	90.33	88.67	77.00	81.67	72.33
Average	98.96	98.52	97.23	88.49	91.82	89.28
# wins	12	9	10	5	9	9

# EFFECTIVENESS

**Table 2: Percentage explained**

Dataset	Linear			Nonlinear		
	SEDC (%)	LIME-C (%)	SHAP-C (%)	SEDC (%)	LIME-C (%)	SHAP-C (%)
Flickr	100	100	100	<b>28.67</b>	28.33	<b>28.67</b>
Ecommerce	100	97.33	100	<u>95.00</u>	<u>97.00</u>	<b>99.67</b>
Airline	100	100	100	100	100	100
Twitter	100	100	100	100	100	100
Fraud	100	100	81.67	100	100	75
YahooMovies	100	100	100	98.67	100	100
TaFeng	100	100	100	<u>93.33</u>	100	100
KDD2015	100	100	100	99.67	100	99.67
20news	99.47	99.47	100	99.47	98.94	100
MovieLens_100k	100	100	100	100	100	100
Facebook	<b>95.67</b>	95.00	95.00	<u>70.33</u>	<b>92.67</b>	<u>89.67</u>
MovieLens_1m	<b>98.67</b>	<b>98.67</b>	<b>98.67</b>	<u>88.33</u>	95.00	<b>95.67</b>
Libimseti	<b>92.67</b>	<u>90.33</u>	<u>88.67</u>	<u>77.00</u>	<b>81.67</b>	<u>72.33</u>
Average	<b>98.96</b>	98.52	97.23	88.49	<b>91.82</b>	89.28
# wins	12	9	10	5	9	9

# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	4(1-13)	3(1-4.4)	3(1-2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(13-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-38.5)
# wins	13	13	13	3	12	11	11	3

# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	1(1-13)	3(1-4.4)	3(1.2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(13-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-38.5)
# wins	13	13	13	3	12	11	11	3

# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	4(1-13)	3(1-4.4)	3(1.2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(13-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-38.5)
# wins	13	13	13	3	12	11	11	3

# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	4(1-13)	3(1-4.4)	3(1.2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(13-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-38.5)
# wins	13	13	13	3	12	11	11	3

# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	4(1-13)	3(1-4.4)	3(1-2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(15-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-58.5)
# wins	13	13	13	3	12	11	11	3

# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	4(1-13)	3(1-4.4)	3(1.2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(13-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-38.5)
# wins	13	13	13	3	12	11	11	3



# EFFECTIVENESS

**Table 3: Switching point in # features (Median + Interquantile range)**

Dataset	Linear				Nonlinear			
	SEDC	LIME-C	SHAP-C	Random	SEDC	LIME-C	SHAP-C	Random
Flickr	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
Ecommerce	1(1-1)	1(1-1)	1(1-1)	1(1-2)	1(1-1)	1(1-1)	1(1-1)	1(1-1)
Airline	1(1-2)	1(1-2)	1(1-2)	2(1-3)	1(1-1)	1(1-1)	1(1-1)	2(1-3)
Twitter	2(1-3)	2(1-3)	2(1-3)	3(2-5)	1(1-1)	1(1-1)	1(1-1)	3(2-5.5)
Fraud	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-1)	1(1-2)
YahooMovies	2(1-4)	2(1-4)	2(1-4)	4(2-7)	1(1-3)	2(1-3)	2(1-3)	4(2-12)
TaFeng	2(1-4)	2(1-4)	2(1-4)	5(3-11)	2(1-8)	2(1-3)	2(1-3.05)	6(3-17)
KDD2015	3(1-7)	3(1-7)	3(1-7)	8.5(3-17.25)	2(1-3)	2(1-3.95)	2(1-4.5)	5(2-9)
20news	2(1-4)	2(1-4)	2(1-4)	11(4-23.5)	1(1-3)	1(1-3)	1(1-3)	8(3-18)
Movielens_100k	2(1-4)	2(1-4)	2(1-4)	5.5(3-10)	2(1-4)	2(1-4)	2(1-4)	5(2-9.25)
Facebook	3(2-8)	3(2-8)	3(2-8)	8(4-20)	1(1-13)	3(1-4.4)	3(1.2-5)	9(4.5-19.5)
Movielens_1m	3(2-7)	3(2-7)	3(2-7)	9(4-19.25)	3(1-5)	3(1-6)	3(1-6)	7(3-14)
Libimseti	3(2-6)	3(2-6.2)	3(2-6.2)	29(13-52)	2(1-5)	4.2(1.8-8.8)	5(2.5-11.2)	19(8-38.5)
# wins	13	13	13	3	12	11	11	3

# EFFICIENCY

**Table 4: Computation time in seconds (Median + Interquantile range)**

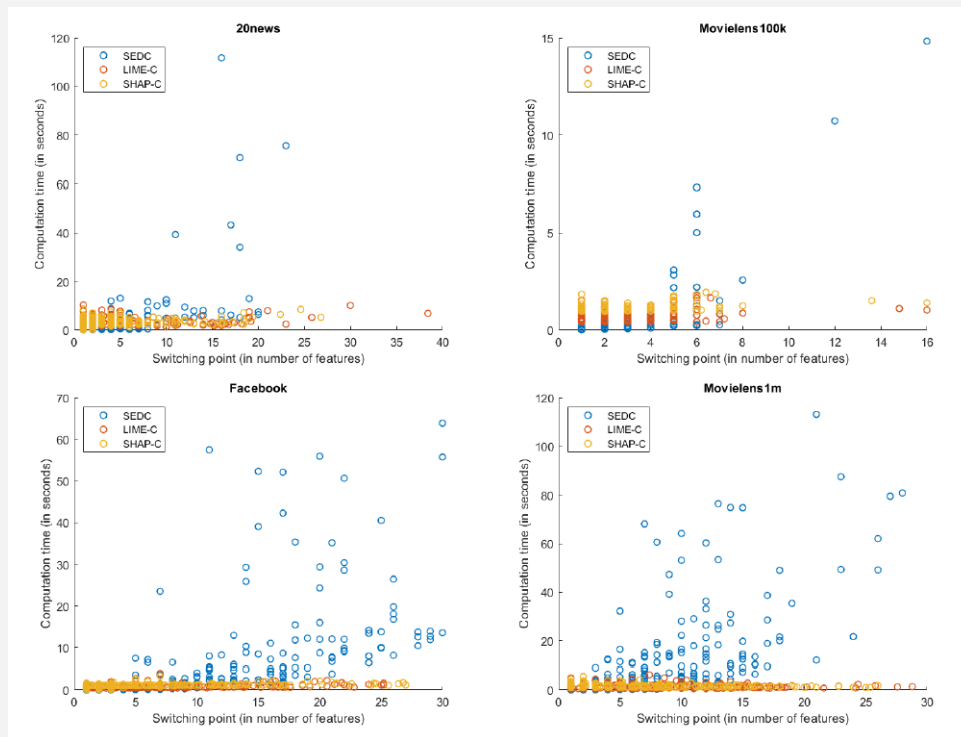
Dataset	Linear			Nonlinear		
	SEDC	LIME-C	SHAP-C	SEDC	LIME-C	SHAP-C
Flickr	<b>0.01(0.00-0.02)</b>	0.34(0.33 – 0.35)	0.08(0.08 – 0.08)	<b>0.02(0.00-0.02)</b>	0.39(0.39 – 0.42)	0.12(0.09 – 0.25)
Ecommerce	<b>0.02(0.00-0.02)</b>	0.34(0.33 – 0.36)	<b>0.02(0.02-0.03)</b>	<b>0.02(0.00-0.02)</b>	0.39(0.38 – 0.41)	0.03(0.03 – 0.03)
Airline	<b>0.02(0.02-0.02)</b>	0.94(0.81 – 1.08)	0.09(0.03 – 0.60)	<b>0.02(0.02-0.02)</b>	1.35(1.17 – 1.51)	0.13(0.04 – 0.82)
Twitter	<b>0.03(0.02-0.05)</b>	0.61(0.56 – 0.64)	0.18(0.06 – 0.46)	<b>0.02(0.01-0.02)</b>	0.67(0.63 – 0.69)	0.15(0.06 – 0.47)
Fraud	<b>0.01(0.00-0.02)</b>	0.38(0.36 – 0.39)	0.07(0.06 – 0.08)	<b>0.01(0.01-0.01)</b>	0.43(0.42 – 0.44)	0.09(0.07 – 0.17)
YahooMovies	<b>0.03(0.02-0.08)</b>	0.44(0.43 – 0.49)	0.96(0.90 – 1.00)	<b>0.06(0.03-0.20)</b>	0.82(0.79 – 0.85)	1.35(1.28 – 1.39)
TaFeng	<b>0.05(0.02-0.22)</b>	0.50(0.45 – 0.59)	1.03(0.97 – 1.08)	<b>0.04(0.02-0.40)</b>	0.51(0.46 – 0.59)	1.01(0.95 – 1.06)
KDD2015	<b>0.11(0.02-0.79)</b>	0.52(0.47 – 0.61)	1.04(0.99 – 1.09)	<b>0.14(0.04-0.56)</b>	0.84(0.78 – 0.94)	1.37(1.31 – 1.45)
20news	<b>0.19(0.05-1.34)</b>	3.12(2.09 – 4.18)	3.65(2.74 – 4.49)	<b>0.09(0.03-0.68)</b>	2.16(1.49 – 2.95)	2.53(1.99 – 3.09)
Movielens_100k	<b>0.06(0.03-0.30)</b>	0.49(0.44 – 0.69)	0.87(0.83 – 1.04)	<b>0.09(0.04-0.35)</b>	0.55(0.50 – 0.83)	1.10(1.02 – 1.27)
Facebook	<b>0.12(0.03-1.17)</b>	0.55(0.46 – 0.75)	1.11(1.04 – 1.23)	<b>0.19(0.02-2.20)</b>	0.51(0.46 – 0.59)	1.06(1.00 – 1.12)
Movielens_1m	<b>0.37(0.06-3.09)</b>	0.74(0.52 – 1.21)	1.21(1.05 – 1.53)	<b>0.39(0.07-1.56)</b>	0.76(0.59 – 1.12)	1.29(1.16 – 1.54)
Libimseti	<b>0.36(0.14-2.26)</b>	1.07(0.92 – 1.38)	1.37(1.27 – 1.52)	<b>0.39(0.09-1.56)</b>	1.02(0.91 – 1.23)	1.42(1.35 – 1.53)
# wins	13	0	1	13	0	0

# EFFICIENCY

**Table 4: Computation time in seconds (Median + Interquantile range)**

Dataset	Linear			Nonlinear		
	SEDC	LIME-C	SHAP-C	SEDC	LIME-C	SHAP-C
Flickr	<b>0.01(0.00-0.02)</b>	0.34(0.33 – 0.35)	0.08(0.08 – 0.08)	<b>0.02(0.00-0.02)</b>	0.39(0.39 – 0.42)	0.12(0.09 – 0.25)
Ecommerce	<b>0.02(0.00-0.02)</b>	0.34(0.33 – 0.36)	<b>0.02(0.02-0.03)</b>	<b>0.02(0.00-0.02)</b>	0.39(0.38 – 0.41)	0.03(0.03 – 0.03)
Airline	<b>0.02(0.02-0.02)</b>	0.94(0.81 – 1.08)	0.09(0.03 – 0.60)	<b>0.02(0.02-0.02)</b>	1.35(1.17 – 1.51)	0.13(0.04 – 0.82)
Twitter	<b>0.03(0.02-0.05)</b>	0.61(0.56 – 0.64)	0.18(0.06 – 0.46)	<b>0.02(0.01-0.02)</b>	0.67(0.63 – 0.69)	0.15(0.06 – 0.47)
Fraud	<b>0.01(0.00-0.02)</b>	0.38(0.36 – 0.39)	0.07(0.06 – 0.08)	<b>0.01(0.01-0.01)</b>	0.43(0.42 – 0.44)	0.09(0.07 – 0.17)
YahooMovies	<b>0.03(0.02-0.08)</b>	0.44(0.43 – 0.49)	0.96(0.90 – 1.00)	<b>0.06(0.03-0.20)</b>	0.82(0.79 – 0.85)	1.35(1.28 – 1.39)
TaFeng	<b>0.05(0.02-0.22)</b>	0.50(0.45 – 0.59)	1.03(0.97 – 1.08)	<b>0.04(0.02-0.40)</b>	0.51(0.46 – 0.59)	1.01(0.95 – 1.06)
KDD2015	<b>0.11(0.02-0.79)</b>	0.52(0.47 – 0.61)	1.04(0.99 – 1.09)	<b>0.14(0.04-0.56)</b>	0.84(0.78 – 0.94)	1.37(1.31 – 1.45)
20news	<b>0.19(0.05-1.34)</b>	3.12(2.09 – 4.18)	3.65(2.74 – 4.49)	<b>0.09(0.03-0.68)</b>	2.16(1.49 – 2.95)	2.53(1.99 – 3.09)
Movielens_100k	<b>0.06(0.03-0.30)</b>	0.49(0.44 – 0.69)	0.87(0.83 – 1.04)	<b>0.09(0.04-0.35)</b>	0.55(0.50 – 0.83)	1.10(1.02 – 1.27)
Facebook	<b>0.12(0.03-1.17)</b>	0.55(0.46 – 0.75)	1.11(1.04 – 1.23)	<b>0.19(0.02-2.20)</b>	0.51(0.46 – 0.59)	1.06(1.00 – 1.12)
Movielens_1m	<b>0.37(0.06-3.09)</b>	0.74(0.52 – 1.21)	1.21(1.05 – 1.53)	<b>0.39(0.07-1.56)</b>	0.76(0.59 – 1.12)	1.29(1.16 – 1.54)
Libimseti	<b>0.36(0.14-2.26)</b>	1.07(0.92 – 1.38)	1.37(1.27 – 1.52)	<b>0.39(0.09-1.56)</b>	1.02(0.91 – 1.23)	1.42(1.35 – 1.53)
# wins	13	0	1	13	0	0

# EFFICIENCY: time vs switching point

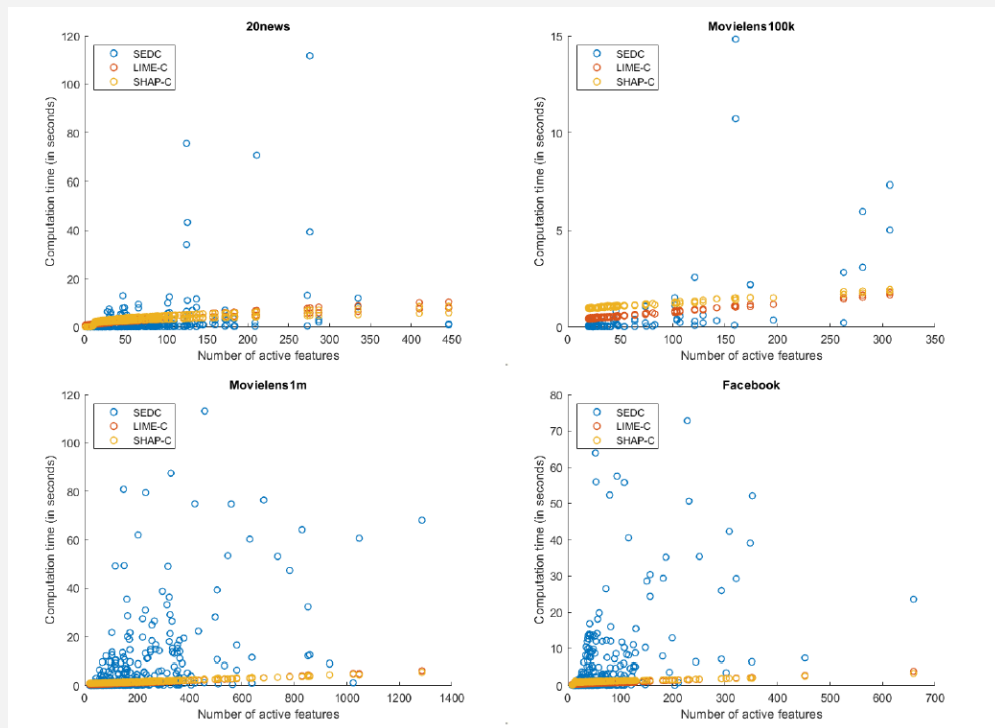


# EFFICIENCY

**Table 4: Computation time in seconds (Median + Interquantile range)**

Dataset	Linear			Nonlinear		
	SEDC	LIME-C	SHAP-C	SEDC	LIME-C	SHAP-C
Flickr	<b>0.01(0.00-0.02)</b>	0.34(0.33 – 0.35)	0.08(0.08 – 0.08)	<b>0.02(0.00-0.02)</b>	0.39(0.39 – 0.42)	0.12(0.09 – 0.25)
Ecommerce	<b>0.02(0.00-0.02)</b>	0.34(0.33 – 0.36)	<b>0.02(0.02-0.03)</b>	<b>0.02(0.00-0.02)</b>	0.39(0.38 – 0.41)	0.03(0.03 – 0.03)
Airline	<b>0.02(0.02-0.02)</b>	0.94(0.81 – 1.08)	0.09(0.03 – 0.60)	<b>0.02(0.02-0.02)</b>	1.35(1.17 – 1.51)	0.13(0.04 – 0.82)
Twitter	<b>0.03(0.02-0.05)</b>	0.61(0.56 – 0.64)	0.18(0.06 – 0.46)	<b>0.02(0.01-0.02)</b>	0.67(0.63 – 0.69)	0.15(0.06 – 0.47)
Fraud	<b>0.01(0.00-0.02)</b>	0.38(0.36 – 0.39)	0.07(0.06 – 0.08)	<b>0.01(0.01-0.01)</b>	0.43(0.42 – 0.44)	0.09(0.07 – 0.17)
YahooMovies	<b>0.03(0.02-0.08)</b>	0.44(0.43 – 0.49)	0.96(0.90 – 1.00)	<b>0.06(0.03-0.20)</b>	0.82(0.79 – 0.85)	1.35(1.28 – 1.39)
TaFeng	<b>0.05(0.02-0.22)</b>	0.50(0.45 – 0.59)	1.03(0.97 – 1.08)	<b>0.04(0.02-0.40)</b>	0.51(0.46 – 0.59)	1.01(0.95 – 1.06)
KDD2015	<b>0.11(0.02-0.79)</b>	0.52(0.47 – 0.61)	1.04(0.99 – 1.09)	<b>0.14(0.04-0.56)</b>	0.84(0.78 – 0.94)	1.37(1.31 – 1.45)
20news	<b>0.19(0.05-1.34)</b>	3.12(2.09 – 4.18)	3.65(2.74 – 4.49)	<b>0.09(0.03-0.68)</b>	2.16(1.49 – 2.95)	2.53(1.99 – 3.09)
Movielens_100k	<b>0.06(0.03-0.30)</b>	0.49(0.44 – 0.69)	0.87(0.83 – 1.04)	<b>0.09(0.04-0.35)</b>	0.55(0.50 – 0.83)	1.10(1.02 – 1.27)
Facebook	<b>0.12(0.03-1.17)</b>	0.55(0.46 – 0.75)	1.11(1.04 – 1.23)	<b>0.19(0.02-2.20)</b>	0.51(0.46 – 0.59)	1.06(1.00 – 1.12)
Movielens_1m	<b>0.37(0.06-3.09)</b>	0.74(0.52 – 1.21)	1.21(1.05 – 1.53)	<b>0.39(0.07-1.56)</b>	0.76(0.59 – 1.12)	1.29(1.16 – 1.54)
Libimseti	<b>0.36(0.14-2.26)</b>	1.07(0.92 – 1.38)	1.37(1.27 – 1.52)	<b>0.39(0.09-1.56)</b>	1.02(0.91 – 1.23)	1.42(1.35 – 1.53)
# wins	13	0	1	13	0	0

# EFFICIENCY: time vs active features



# CONCLUSION

- **SEDC** most efficient and effective for small instances, *however*
    - computation time very sensitive to switching point
    - flaw in heuristic best-first for some nonlinear models
  - **SHAP-C** overall good performance, *however*
    - problems with highly unbalanced data
    - computation time more sensitive to # active features than LIME-C
- ⇒ **LIME-C** seems most favourable search algorithm: best tradeoff
- low computation times
  - least sensitive to switching point and # active features
  - stable performance in terms of effectiveness criteria



# 5. FURTHER RESEARCH



# FURTHER RESEARCH

- More data sets and models
- Study efficiency-effectiveness tradeoff of the algorithms
- Evaluate other hybrid algorithms
- Other objectives of the algorithm



**NEXT PROJECT?**

# GLOBAL EXPLANATIONS

- Approximate original model?
- Rule extraction?



# THANKS!

Further questions?

+32 497 901 304

[yanou.ramon@uantwerp.be](mailto:yanou.ramon@uantwerp.be)

[www.linkedin.com/in/yanou-ramon](https://www.linkedin.com/in/yanou-ramon)

[www.applieddatamining.com](https://www.applieddatamining.com)