

# DOES THE MODEL THINK AS WE EXPECT? EXPLORING ML MODEL LOGIC AND TRUSTWORTHINESS THROUGH DECISION RULES

---

Natalia Andrienko, Gennady Andrienko, Bahavathy Kathirgamanathan

Human-Centred AI

# INTRODUCTION

- **Key question:** ML models can be accurate, but do they reason as we expect?
- **Why This Matters:**
  - Trust in ML models is not just about accuracy – it's about understanding *why* they make decisions.
  - A model may produce correct predictions while relying on reasoning that differs from human logic.
  - This misalignment can affect model adoption, interpretation, and decision-making in critical applications.
- **Our goal:** support exploring model alignment with human expectations using visual analytics.
  - We consider models represented by systems of decision rules.
- **Focus of this talk:** What insights we gain about trustworthiness when analyzing rule-based ML models.

# RUNNING EXAMPLE: CLASSIFICATION MODEL FOR COVID-19 PREDICTION

- Model was developed using a dataset including daily counts and trips for 52 regions during the COVID-19 pandemic period, normalized by population.
- Data aggregated weekly over 64 weeks, excluding initial outbreak phase.
- Discretized data into four levels for disease incidence and population mobility. Low mobility levels indicate restrictions.
- Focus on interdependencies between disease incidence and mobility levels.
  - Increases in disease incidence may lead to reduced mobility through restrictions, which subsequently contribute to a decrease in disease levels.
  - Conversely, relaxed mobility restrictions may result in increased disease incidence.
  - The effects may become noticeable after a delay.

# FEATURES USED FOR DERIVING PREDICTIONS

- Temporal features related to COVID-19 levels and mobility trends over six weeks preceding a target event:
  - **COVID-19 Features:** Weekly categorical indicators (Week6\_Covid to Week1\_Covid) with values c1, c2, c3, and c4, representing increasing severity levels.
  - **Mobility Features:** Weekly categorical indicators (Week6\_Mobility to Week1\_Mobility) with values m1, m2, m3, and m4, representing mobility levels from low (lockdown) to normal.
- The number of days passed since the start of the pandemic monitoring.
- **Target Class:** The outcome variable categorizing the event into one of the four classes (c1 to c4).
- The categorical values were encoded by numbers from 1 to 4.



# REPRESENTATION OF RULES IN A TABLE

Predicted class or value (for regression)

features

Id	Tree...	Tr...	W...	Class	or...	N conditions	Rule	Days_pass...	Week6_Co...	Week5_Co...	Week4_Co...	Week3_Co...	Week2_C...	Week6_Mob...	Week5_M...	Week4_Mobi...	Week3_Mo...	Week2_Mobility	Week1_Mobil...
1	0	3	1	2	0	2		45..171					0..0						
2	0	3	1	1	1	3		45..171					1..1		0..0				
3	0	3	1	3	2	4		45..171			0..0		1..1		1..1				
4	0	3	1	4	3	5		45..171			0..0		1..1		2..2			0..2	
5	0	3	1	3	4	6		45..171			0..0		1..1		2..2		0..2	3..3	
6	0	3	1	1	5	6		45..154			0..0		1..1		2..2		3..3	3..3	
7	0	3	1	3	6	6		154..171			0..0		1..1		2..2		3..3	3..3	
8	0	3	1	1	7	5		45..171					1..1		1..1				0..1
9	0	3	1	1	8	6		45..171			1..3		1..1		1..1				1..3
10	0	3	1	3	9	6		45..171			1..3		1..1	2..2	1..1				1..3
11	0	3	1	1	10	4		45..171					2..2	0..2				0..1	
12	0	3	1	4	11	4		45..171					2..2	0..2			1..3		
13	0	3	1	2	12	4		45..171	0..2				2..2	3..3					
14	0	3	1	1	13	4		45..171	3..3				2..2	3..3					
15	0	3	1	3	14	2		45..171					3..3	3..3					
16	0	3	1	3	15	4		172..434			0..0	0..2			0..2				
17	0	3	1	2	16	7		172..178	0..0		0..0	0..0	0..0		3..3			0..2	
18	0	3	1	2	17	7		172..178	0..0		0..0	0..0	0..0		3..3			3..3	
19	0	3	1	2	18	6		179..434	0..0		0..0	0..0	0..0		3..3				
20	0	3	1	3	19	6		172..434	0..0		0..0	0..0	1..3		3..3				
21	0	3	1	3	20	5		172..434	0..0		0..0	1..2			3..3				
22	0	3	1	3	21	5		172..434	1..3			0..2			3..3				
23	0	3	1	4	22	5		172..434	0..0			1..1	0..2						0..2
24	0	3	1	4	23	6		172..304	1..3			1..1	0..2						0..2
25	0	3	1	3	24	7		305..434	1..3			1..1	0..2			0..2	0..2		0..2
26	0	3	1	2	25	7		305..434	1..3			1..1	0..2			0..2	3..3		0..2
27	0	3	1	3	26	6		172..434	1..3			1..1	0..2			3..3			0..2
28	0	3	1	4	27	6		172..360	0..2			2..2	0..2					0..1	0..2
29	0	3	1	3	28	6		361..434	0..2			2..2	0..2					0..1	0..2
30	0	3	1	4	29	6		172..259	0..2			2..2	0..2					2..2	0..2
31	0	3	1	3	30	8		259..434	0..2			2..2	0..2			0..1	0..1	2..2	0..2
32	0	3	1	4	31	8		259..434	0..2			2..2	0..2			0..1	2..2	2..2	0..2

Conditions  
Bars represent intervals of feature values  
Texts can be hidden

Graphical representations of rules

# GRAPHICAL REPRESENTATION OF A RULE

Rules explorer v.23.01.2025 15:10

R0

Id	Tree...	Tr...	W...	Class	or...	N conditions	Rule	Days_pass...	Week6_Co...	Week5_Co...	Week4_Co...	Week3_Co...	Week2_C...	Week6_Mob...	Week5_M...	Week4_Mobi...	Week3_Mo...	Week2_Mobility	Week1_Mobil...
1	0	3	1	2	0	2		45..171					0.0						
2	0	3	1	1	1	3		45..171					1.1		0.0				
3	0	3	1	3	2	4		45..171			0.0		1.1		1.1				
4	0	3	1	4	3	5		45..171			0.0		1.1		2..3			0.2	
5	0	3	1	3	4	6		45..171			0.0		1.1		2..3			0.2	0.2
6	0	3	1	1	5	6		45..154			0.0		1.1		2..3			0.2	0.2
7	0	3	1	3	6	6		154..171			0.0		1.1		2..3			0.2	0.2
8	0	3	1	1	7	5		45..171			1..3		1.1		1..3				0.1
9	0	3	1	1	8	6		45..171			1..3								1..3
10	0	3	1	3	9	6		45..171			1..3								1..3
11	0	3	1	1	10	4		45..171										0.1	
12	0	3	1	4	11	4		45..171										1..3	
13	0	3	1	2	12	4		45..171	0.2										
14	0	3	1	1	13	4		45..171	3..3										
15	0	3	1	3	14	2		45..171											
16	0	3	1	3	15	4		172..434			0.								
17	0	3	1	2	16	7		172..178	0.0		0.							0.2	
18	0	3	1	2	17	7		172..178	0.0		0.							3..3	
19	0	3	1	2	18	6		179..434	0.0		0.								
20	0	3	1	3	19	6		172..434	0.0		0.								
21	0	3	1	3	20	5		172..434	0.0		0.								
22	0	3	1	3	21	5		172..434	1..3		0.								
23	0	3	1	4	22	5		172..434	0.0		1.								0.2
24	0	3	1	4	23	6		172..304	1..3		1.					0.2			0.2
25	0	3	1	3	24	7		305..434	1..3		1.					0.2			0.2
26	0	3	1	2	25	7		305..434	1..3		1.					0.2			0.2
27	0	3	1	3	26	6		172..434	1..3		1.					3..3			0.2
28	0	3	1	4	27	6		172..360	0.2		2.							0.1	0.2
29	0	3	1	3	28	6		361..434	0.2		2.							0.1	0.2
30	0	3	1	4	29	6		172..259	0.2		2..2							2..2	0.2
31	0	3	1	3	30	8		259..434	0.2		2..2						0.1	0.1	2..2
32	0	3	1	4	31	8		259..434	0.2		2..2						0.1	2..3	2..2

Rule 8 Weight = 1  
Class 1

Feature	min	from	to	max
Days_passed	45.0000	- inf	171	431.0000
Week4_Covid	0.0000	1	+ inf	3.0000
Week2_Covid	0.0000	1	1	3.0000
Week5_Mobility	0.0000	1	+ inf	3.0000
Week1_Mobility	0.0000	- inf	1	3.0000

Original distinct rules or explanations (7173), N conditions (43128), Total uses (7173)

# CHALLENGES IN UNDERSTANDING MODEL LOGIC

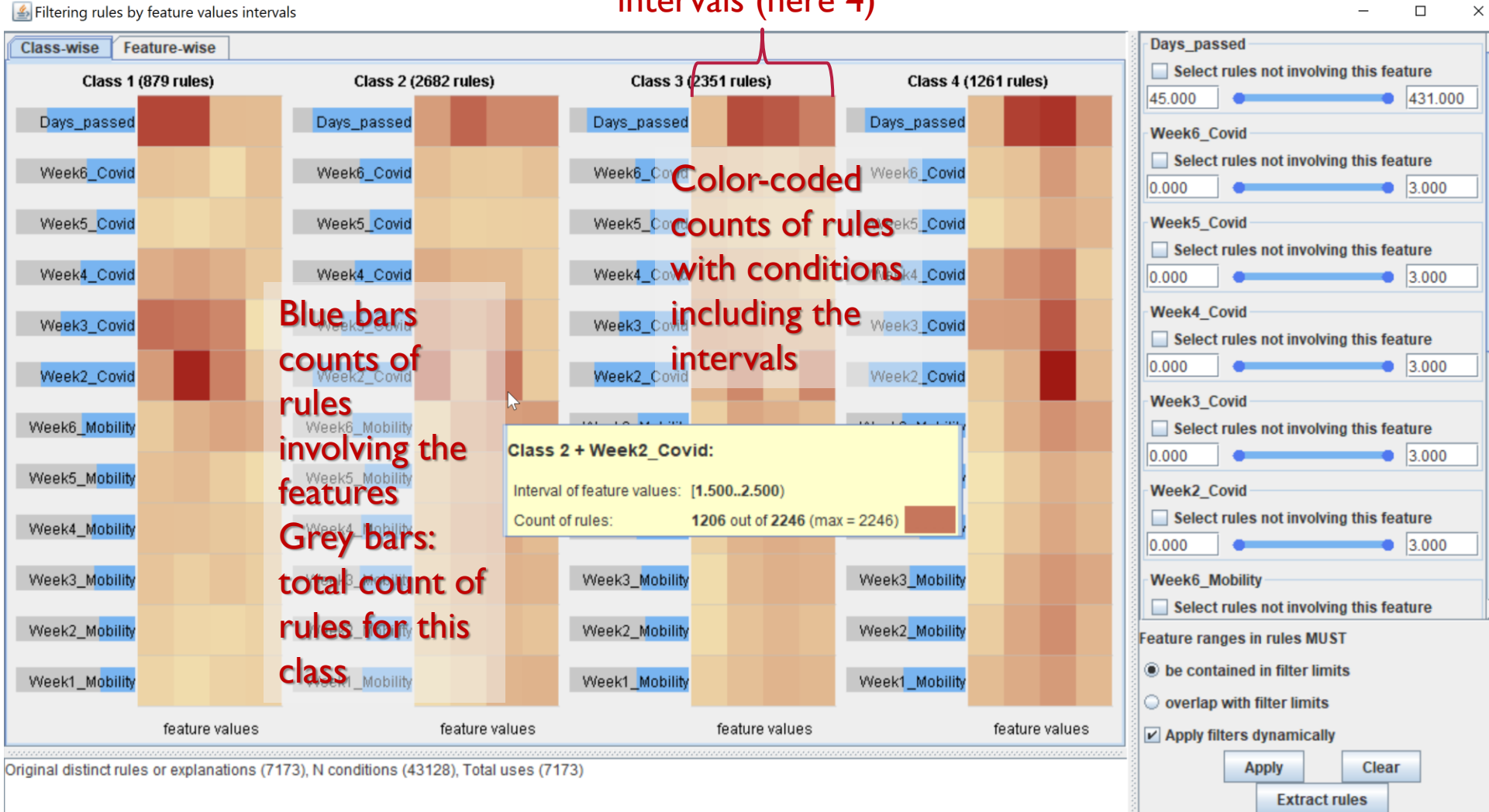
- **We can read and understand the conditions of each rule** >> Representation of a model by rules should allow us to inspect model reasoning, but...
- **The rules are too many** >> detailed examination of individual rules is impractical.
- **The features interact and work jointly** >> investigation of the impacts of individual features on the predictions is insufficient.
- **Our visual analytics solutions:**
  - Provide an overview
  - Enable querying and selection
  - Aggregate and generalize

This research builds on previous work described in paper Adilova, L., Kamp, M., Andrienko, G. and Andrienko, N. **“Re-interpreting rules interpretability”**. *International Journal of Data Science and Analytics* (2023). doi:10.1007/s41060-023-00398-5.

# DISTRIBUTION OF FEATURE VALUE INTERVALS

## Class-wise view

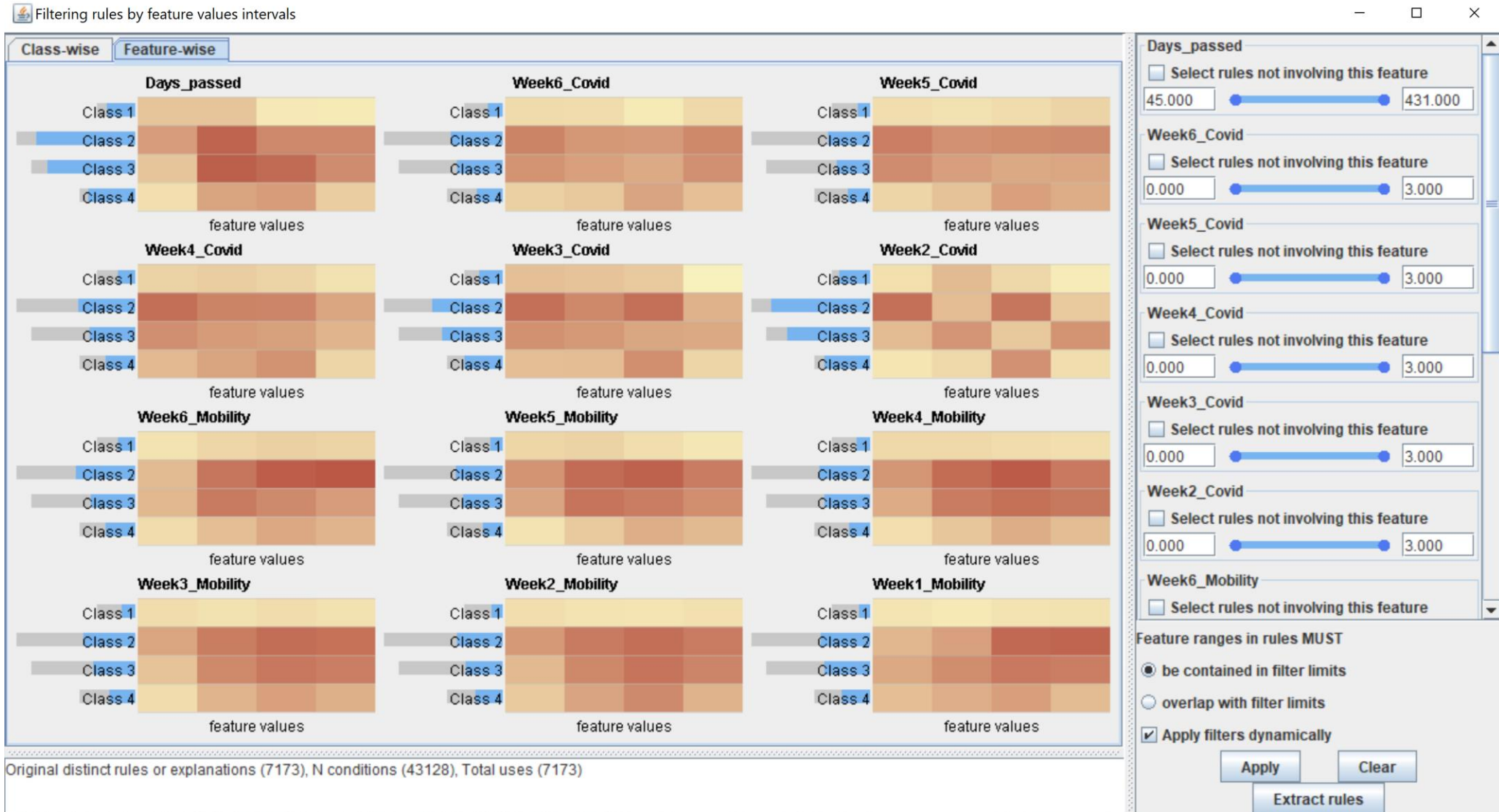
Chosen number of intervals (here 4)



Controls for interactive filtering

# DISTRIBUTION OF FEATURE VALUE INTERVALS

## Feature-wise view





# INTERACTIVE FILTERING

Filtering rules by feature values intervals

Class-wise Feature-wise

	Class 1 (102 rules)	Class 2 (911 rules)	Class 3 (765 rules)	Class 4 (697 rules)
Days_passed				
Week6_Covid				
Week5_Covid				
Week4_Covid				
Week3_Covid				
Week2_Covid				
Week6_Mobility				
Week5_Mobility				
Week4_Mobility				
Week3_Mobility				
Week2_Mobility				
Week1_Mobility				

feature values feature values feature values feature values

Days\_passed  
 Select rules not involving this feature  
45.000 [slider] 431.000

Week6\_Covid  
 Select rules not involving this feature  
0.000 [slider] 3.000

Week5\_Covid  
 Select rules not involving this feature  
0.000 [slider] 3.000

Week4\_Covid  
 Select rules not involving this feature  
0.000 [slider] 3.000

Week3\_Covid  
 Select rules not involving this feature  
0.000 [slider] 3.000

Week2\_Covid  
 Select rules not involving this feature  
1.425 [slider] 3.000

Week6\_Mobility  
 Select rules not involving this feature

Feature ranges in rules MUST

- be contained in filter limits
- overlap with filter limits

Apply filters dynamically

Apply Clear Extract rules

2475 rules selected by applying filter conditions 1) Week2\_Covid: from 1.425 to +infinity

10

# INTERACTIVE FILTERING

Filtering rules by feature values intervals

Class-wise Feature-wise

	Class 1 (13 rules)	Class 2 (260 rules)	Class 3 (231 rules)	Class 4 (241 rules)
Days_passed	High	High	High	High
Week6_Covid	Low	Low	Low	Low
Week5_Covid	Low	Low	Low	Low
Week4_Covid	Low	Low	Low	Low
Week3_Covid	High	High	High	High
Week2_Covid	High	High	High	High
Week6_Mobility	Low	Low	Low	Low
Week5_Mobility	Low	Low	Low	Low
Week4_Mobility	Low	Low	Low	Low
Week3_Mobility	Low	Low	Low	Low
Week2_Mobility	Low	Low	Low	Low
Week1_Mobility	Low	Low	Low	Low

feature values feature values feature values feature values

745 rules selected by applying filter conditions 1) Week3\_Covid: from 1.125 to +infinity 2) Week2\_Covid: from 1.425 to +infinity

Days\_passed  
 Select rules not involving this feature  
45.000 431.000

Week6\_Covid  
 Select rules not involving this feature  
0.000 3.000

Week5\_Covid  
 Select rules not involving this feature  
0.000 3.000

Week4\_Covid  
 Select rules not involving this feature  
0.000 3.000

Week3\_Covid  
 Select rules not involving this feature  
1.125 3.000

Week2\_Covid  
 Select rules not involving this feature  
1.425 3.000

Week6\_Mobility  
 Select rules not involving this feature

Feature ranges in rules MUST  
 be contained in filter limits  
 overlap with filter limits  
 Apply filters dynamically

Apply Clear Extract rules

# INTERACTIVE FILTERING

Filtering rules by feature values intervals

Class-wise Feature-wise

	Class 1 (2 rules)	Class 2 (67 rules)	Class 3 (19 rules)	Class 4 (95 rules)
Days_passed	Red	Orange	Orange	Orange
Week6_Covid	Dark Red	Orange	Orange	Orange
Week5_Covid	White	Orange	White	Orange
Week4_Covid	Orange	Dark Red	Dark Red	Dark Red
Week3_Covid	White	Dark Red	Orange	White
Week2_Covid	White	Dark Red	Orange	Dark Red
Week6_Mobility	Dark Red	Orange	Orange	Orange
Week5_Mobility	Orange	Orange	Orange	Orange
Week4_Mobility	White	Orange	Orange	Orange
Week3_Mobility	Orange	Orange	Orange	Orange
Week2_Mobility	Orange	Orange	Orange	Orange
Week1_Mobility	White	Orange	Orange	Orange

feature values feature values feature values feature values

183 rules selected by applying filter conditions 1) Week4\_Covid: from -infinity to 2.475 2) Week3\_Covid: from 1.125 to +infinity 3) Week2\_Covid: from 1.425 to +infinity

Days\_passed  
 Select rules not involving this feature  
45.000 [Slider] 431.000

Week6\_Covid  
 Select rules not involving this feature  
0.000 [Slider] 3.000

Week5\_Covid  
 Select rules not involving this feature  
0.000 [Slider] 3.000

Week4\_Covid  
 Select rules not involving this feature  
0.000 [Slider] 2.475

Week3\_Covid  
 Select rules not involving this feature  
1.125 [Slider] 3.000

Week2\_Covid  
 Select rules not involving this feature  
1.425 [Slider] 3.000

Week6\_Mobility  
 Select rules not involving this feature

Feature ranges in rules MUST  
 be contained in filter limits  
 overlap with filter limits  
 Apply filters dynamically

Apply Clear Extract rules



# INTERACTIVE FILTERING

Filtering rules by feature values intervals

Class-wise | Feature-wise

	Class 1 (0 rules)	Class 2 (1 rules)	Class 3 (3 rules)	Class 4 (7 rules)
Days_passed		Days_passed	Days_passed	Days_passed
Week6_Covid		Week6_Covid	Week6_Covid	Week6_Covid
Week5_Covid		Week5_Covid	Week5_Covid	Week5_Covid
Week4_Covid		Week4_Covid	Week4_Covid	Week4_Covid
Week3_Covid		Week3_Covid	Week3_Covid	Week3_Covid
Week2_Covid		Week2_Covid	Week2_Covid	Week2_Covid
Week6_Mobility		Week6_Mobility	Week6_Mobility	Week6_Mobility
Week5_Mobility		Week5_Mobility	Week5_Mobility	Week5_Mobility
Week4_Mobility		Week4_Mobility	Week4_Mobility	Week4_Mobility
Week3_Mobility		Week3_Mobility	Week3_Mobility	Week3_Mobility
Week2_Mobility		Week2_Mobility	Week2_Mobility	Week2_Mobility
Week1_Mobility		Week1_Mobility	Week1_Mobility	Week1_Mobility

feature values      feature values      feature values      feature values

11 rules selected by applying filter conditions 1) Week4\_Covid: from -infinity to 1.899 2) Week3\_Covid: from 1.125 to +infinity 3) Week2\_Covid: from 1.425 to +infinity

**Days\_passed**

 Select rules not involving this feature  
 45.000 

 431.000

**Week6\_Covid**

 Select rules not involving this feature  
 0.000 

 3.000

**Week5\_Covid**

 Select rules not involving this feature  
 0.000 

 3.000

**Week4\_Covid**

 Select rules not involving this feature  
 0.000 

 1.899

**Week3\_Covid**

 Select rules not involving this feature  
 1.125 

 3.000

**Week2\_Covid**

 Select rules not involving this feature  
 1.425 

 3.000

**Week6\_Mobility**

 Select rules not involving this feature

**Feature ranges in rules MUST**

be contained in filter limits

overlap with filter limits

Apply filters dynamically

# RELATIONSHIPS BETWEEN FEATURES



# COMBINED EFFECTS OF 2 OR MORE FEATURES

Filtering rules by feature values intervals

Class-wise Feature-wise

Class 1 (13 rules) Class 2 (260 rules) Class 3 (231 rules) Class 4 (241 rules)

Days\_passed  
Week6\_Covid  
Week5\_Covid  
Week4\_Covid  
Week3\_Covid  
Week2\_Covid  
Week6\_Mobility  
Week5\_Mobility  
Week4\_Mobility  
Week3\_Mobility  
Week2\_Mobility  
Week1\_Mobility

feature values feature values feature values feature values

745 rules selected by applying filter conditions 1) Week3\_Covid: from 1.125 to +infinity 2) Week2\_Covid: from 1.425 to +infinity

Days\_passed  
 Select rules not involving this feature  
45.000 431.000

Week6\_Covid  
 Select rules not involving this feature  
0.000 3.000

Week5\_Covid  
 Select rules not involving this feature  
0.000 3.000

Week4\_Covid  
 Select rules not involving this feature  
0.000 3.000

Week3\_Covid  
 Select rules not involving this feature  
1.125 3.000

Week2\_Covid  
 Select rules not involving this feature  
1.425 3.000

Week6\_Mobility  
 Select rules not involving this feature

Feature ranges in rules MUST

be contained in filter limits  
 overlap with filter limits

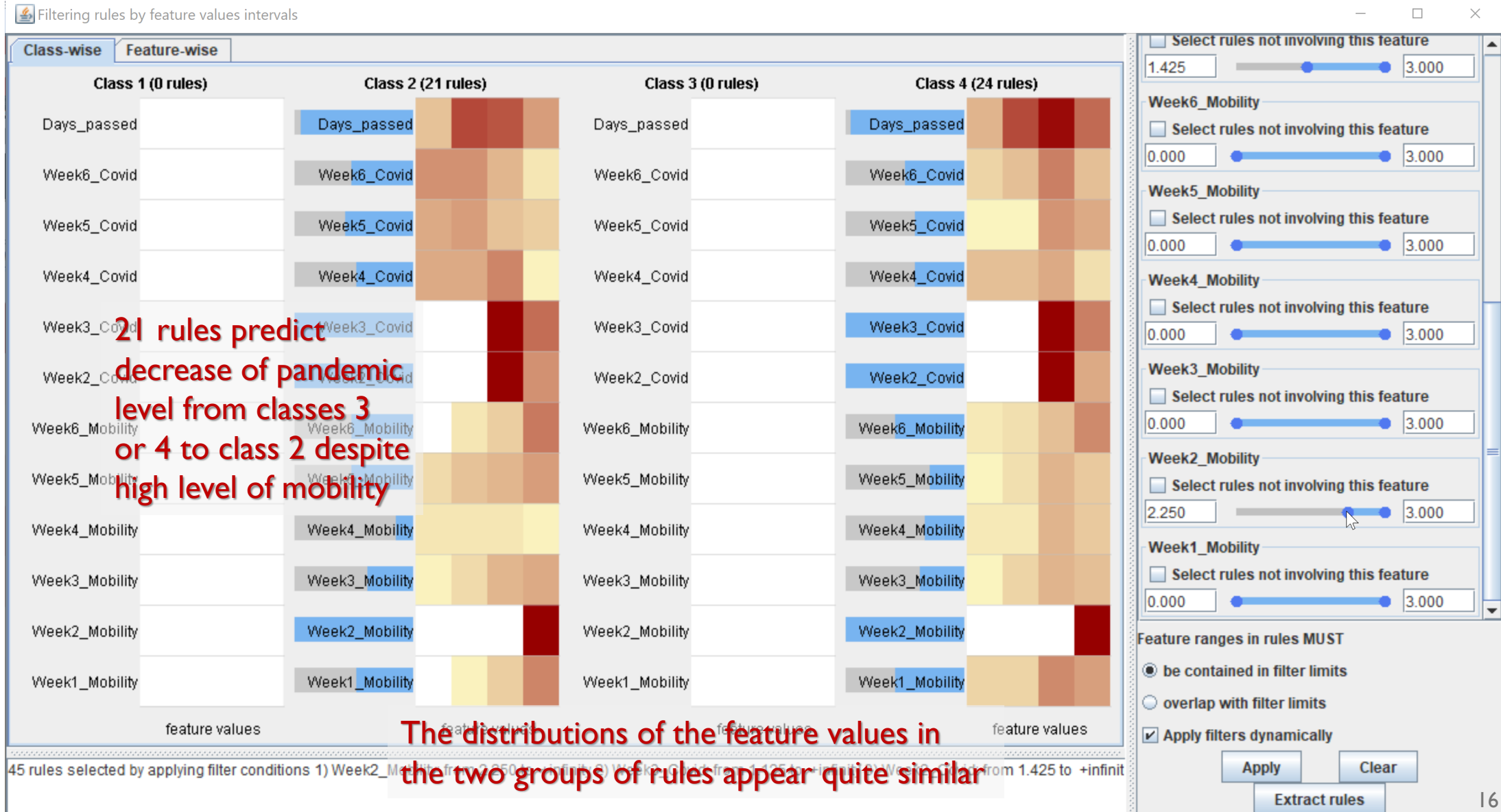
Apply filters dynamically

Apply Clear

Extract rules

15

# SURPRISING FINDINGS



# EXTRACTED RULES SATISFYING THE FILTER

Rules explorer v.30.01.2025 17:20

R0 R0.1

Id	Tree Id	Weight	Class	or...	c...	N con...	Rule	Days_pa...	Week6_...	Week5_...	Week4_...	Week3_...	Week2_...	Week6_...	Week5_...	Week4_...	Week3_...	Week2_M...	Week1_...
561	7	2	1	2	0	-1	7												
562	7	2	1	4	1	-1	7												
6590	91	1	1	2	2	-1	8												
780	10	3	1	4	3	-1	9												
788	10	3	1	2	4	-1	8												
789	10	3	1	4	5	-1	8												
1022	13	3	1	4	6	-1	8												
1098	14	0	1	4	7	-1	6												
1503	20	3	1	4	8	-1	10												
1833	25	2	1	4	9	-1	6												
5632	78	1	1	4	10	-1	5												
4489	62	3	1	2	11	-1	6												
5634	78	1	1	2	12	-1	5												
5633	78	1	1	4	13	-1	5												
2323	32	1	1	2	14	-1	10												
2450	34	1	1	2	15	-1	7												
2451	34	1	1	4	16	-1	7												
4487	62	3	1	4	17	-1	6												
2752	38	1	1	4	18	-1	6												
2753	38	1	1	2	19	-1	6												
5560	77	3	1	2	20	-1	8												
3736	51	2	1	4	21	-1	9												
3933	54	3	1	2	22	-1	9												
3935	54	3	1	4	23	-1	9												
4441	61	2	1	4	24	-1	5												
4445	61	2	1	4	25	-1	9												
4446	61	2	1	2	26	-1	9												
4710	65	2	1	2	27	-1	7												
4864	67	2	1	2	28	-1	6												
4935	68	0	1	4	29	-1	10												
4939	68	0	1	2	30	-1	8												
5135	71	0	1	4	31	-1	8												
5139	71	0	1	4	32	-1	8												
5621	78	1	1	2	33	-1	8												

45 rules satisfying filter conditions 1) Week2\_Mobility: from 0.50 to 1.50

The extracted rules are ordered by similarity (by applying OPTICS algorithm).  
We see that differences between rules predicting different classes may be very subtle.  
The laborious process of comparison can be facilitated by aggregating the rules.

# ITERATIVE AGGREGATION AND GENERALISATION OF RULES

This final rule was derived from 5 original rules.

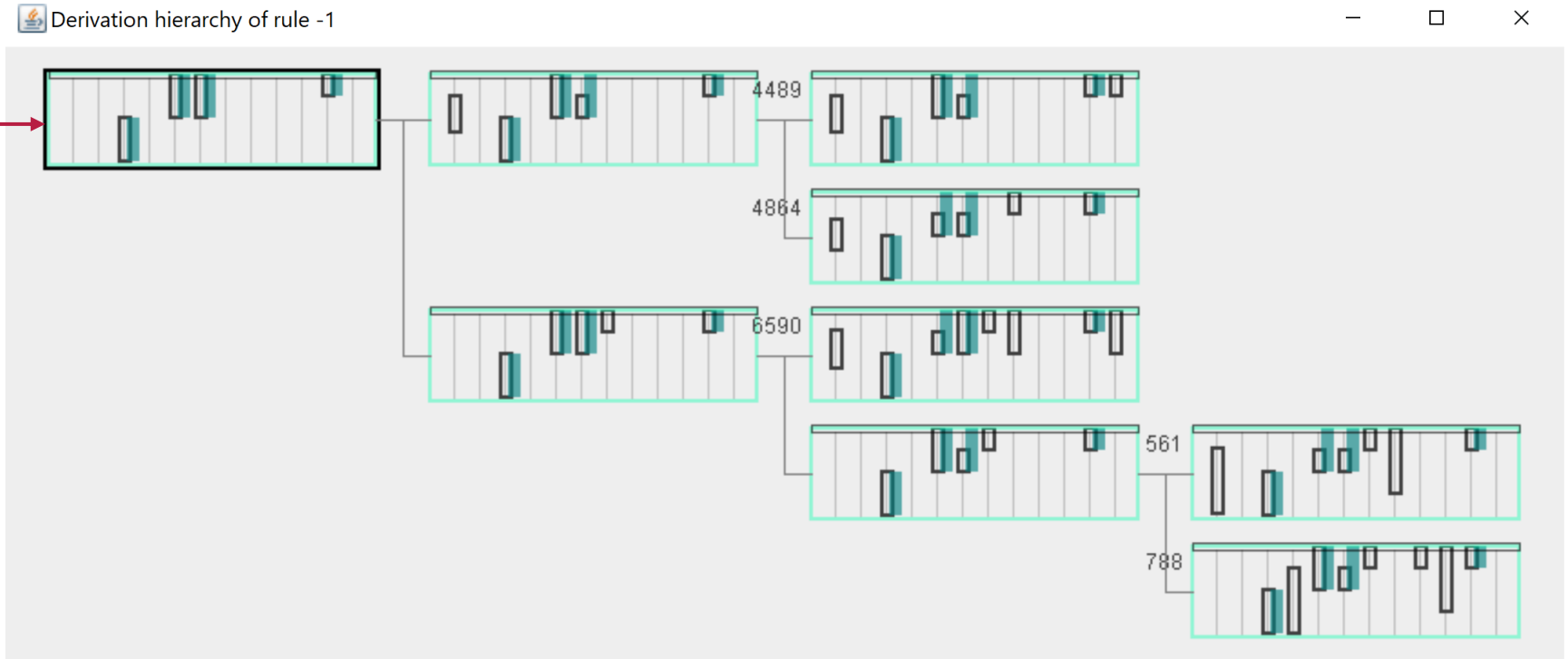
The glyph representing the rule is interactively selected.

The conditions from this rule are represented by blue vertical bars in all glyphs for comparison.

Rule  $R_1$  covers rule  $R_2$  if  $R_1$  is applicable to all instances where  $R_2$  is applicable.

$R_1$  is more general than  $R_2$ .

The operation of uniting 2 or more rules creates a more general rule covering each of the original rules.



Uniting rules predicting the same class:

- Identical conditions remain in the resulting rules.
- Differing intervals are joined by creating a covering interval.
- Conditions with features missing in one of the rules are omitted.



# A RESULT OF AGGREGATING 45 SELECTED RULES

## Coherence threshold = 0.75 (fraction of allowed exceptions)

Rules explorer v.30.01.2025 17:20

R0		R0.1		R0.1.1																					
Id	...	Weight	Class	...	c...	N con...	Rule	N right...	N wro...	Coher...	N unit...	Depth	Days_...	Week...	Week5...	Week...	Week3...	Week...	Week6...	Week5...	Week4...	Week3...	Week2...	Week1...	
183325	2	1	4 1	0	-1	6		1	0	100.00...	0	1													
679094	0	1	4 1	1	-1	7		1	0	100.00...	0	1													
-1 -1 -1		3	4 5	2	-1	5		5	0	100.00...	5	3													
-1 -1 -1		3	2 3	3	-1	5		3	0	75.000...	4	3													
-1 -1 -1		4	4 4	4	-1	4		4	0	100.00...	6	3													
-1 -1  1		2	2 3	5	-1	4		3	0	100.00...	3	2													
-1 -1 -1		2	2 2	6	-1	6		3	0	100.00...	2	2													
-1 -1 -1		2	4 3	7	-1	4		2	0	100.00...	3	2													
-1 -1 -1		2	2 2	8	-1	4		2	0	100.00...	2	2													
-1 -1 -1		5	2 5	9	-1	4		7	0	100.00...	8	4													
-1 -1 -1		2	4 3	10	-1	5		3	0	100.00...	3	2													
-1 -1 3		2	4 2	11	-1	7		2	0	100.00...	2	2													
-1 -1 -1		2	4 2	12	-1	6		2	0	100.00...	2	2													
-1 -1 -1		3	4 3	13	-1	4		6	0	100.00...	4	3													
-1 -1 -1		2	2 2	14	-1	8		2	0	100.00...	2	2													
-1 -1 -1		2	2 2	15	-1	5		3	0	100.00...	2	2													
-1 -1 0		2	2 2	16	-1	7		2	0	100.00...	2	2													

17 aggregated and generalized rules obtained from the set of 45 earlier selected or derived rules using the following parameter settings: min coherence = 0.750; stepwise aggregation starting from 1.000 with step 0.050

These rules have exceptions

Differences between 2 similar rules predicting distinct classes cannot be grounded in domain knowledge.

### Rules covered by one of the “rough” rules:

Id	Tr...	Weight	Class	N uses	or...	N con...	Rule	N ri...	N wr...	Cohe...	N unit...	Dept...	Days_...	Week...	Week5...	Week...	Week3...	Week...	Week6...	Week...	Week4_...	Week...	Week2_...	Week1_...
1	-1	-1	2	3	0	5		3	1	75.00...	4	3												
261	2	1	2	1	1	9		1	0	100.0...	0	0												
368	0	1	2	1	2	8		1	0	100.0...	0	0												
482	1	1	2	1	3	10		1	0	100.0...	0	0												
561	2	1	4	1	4	9		1	0	100.0...	0	0												

Exception (a rule predicting a different class)

# A RESULT OF AGGREGATING 45 SELECTED RULES

Coherence threshold = 0.75 (fraction of allowed exceptions)

Rules explorer v.30.01.2025 17:20

R0		R0.1		R0.1.1																					
Id	...	Weight	Class	...	c...	N con...	Rule	N right...	N wro...	Coher...	N unit...	Depth	Days_...	Week...	Week5...	Week...	Week3...	Week...	Week6...	Week5...	Week4...	Week3...	Week2...	Week1...	
183325	2	1	4 1	0	-1	6		1	0	100.00...	0	1													
679094	0	1	4 1	1	-1	7		1	0	100.00...	0	1													
-1 -1 -1		3	4 5	2	-1	5		5	0	100.00...	5	3													
-1 -1 -1		3	2 3	3	-1	5		3	0	75.000...	4	3													
-1 -1 -1		4	4 4	4	-1	4		4	0	100.00...	6	3													
-1 -1 -1		2	2 3	5	-1	4		3	0	100.00...	3	2													
-1 -1 -1		2	2 2	6	-1	6		3	0	100.00...	2	2													
-1 -1 -1		2	4 3	7	-1	4		2	0	100.00...	3	2													
-1 -1 -1		5	2 5	9	-1	4		7	0	100.00...	8	4													
-1 -1 -1		2	4 3	10	-1	5		3	0	100.00...	3	2													
-1 -1 3		2	4 2	11	-1	7		2	0	100.00...	2	2													
-1 -1 -1		2	4 2	12	-1	6		2	0	100.00...	2	2													
-1 -1 -1		3	4 3	13	-1	4		6	0	75.000...	4	3													
-1 -1 -1		2	2 2	14	-1	8		2	0	100.00...	2	2													
-1 -1 -1		2	2 2	15	-1	5		3	0	100.00...	2	2													
-1 -1 0		2	2 2	16	-1	7		2	0	100.00...	2	2													

17 aggregated and generalized rules obtained from the set of 45 earlier selected or derived rules using the following parameter settings: min coherence = 0.750; stepwise aggregation starting from 1.000 with step 0.050

These rules have exceptions

Differences between 2 similar rules predicting distinct classes cannot be grounded in domain knowledge.

## Rules covered by one of the "rough" rules:

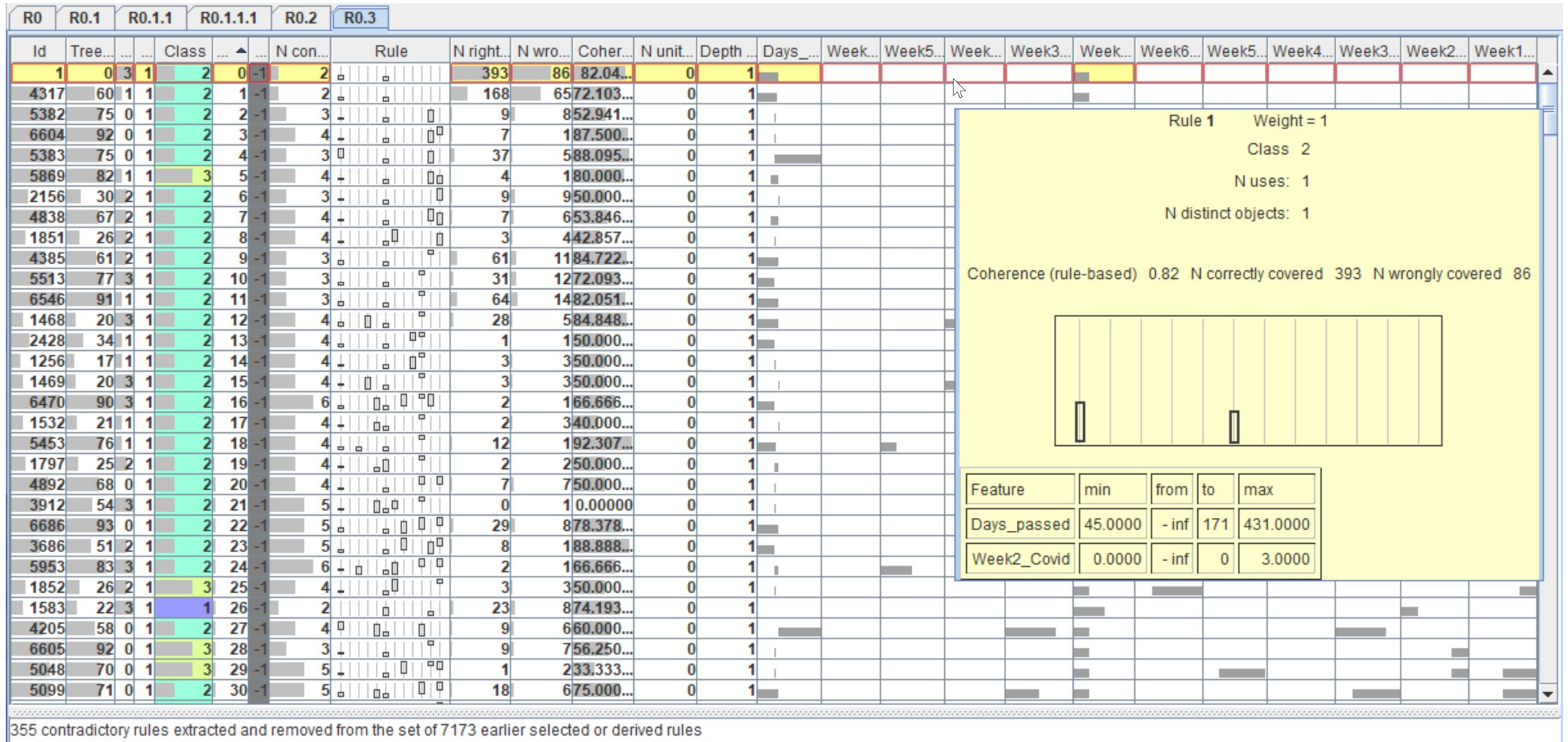
All coverages of rule # 1



Exception (a rule predicting a different class)



# DETECTING AND REMOVING CONTRADICTIONARY RULES



Our model contains 355 contradictory rules. Some of them include only a few conditions. According to the domain and/or commonsense knowledge, these conditions cannot be sufficient for making valid predictions.

# MANY NON-CONTRADICTIONARY RULES ALSO INCLUDE TOO FEW CONDITIONS

Rules explorer v.30.01.2025 17:20

R0		R0.1		R0.1.1		R0.2		R0.3		R0.4		R0.2.1		R0.2.1.1											
Id	Tr...	...	Class	...	N ...	Rule	N right ...	N w...	Coher...	N unite...	Depth ...	Days_p...	Week6...	Week5...	Week4...	Week3...	Week2...	Week6...	Week5...	Week4...	Week3...	Week2...	Week1...		
3213	45	0	1	4	...	2		2	0	100.00...	0	1													
4497	62	3	1	4	...	2		1	0	100.00...	0	1													
15	0	3	1	3	11	2		40	0	100.00...	0	1													
282	3	1	1	3	...	2		78	0	100.00...	0	1													
504	6	2	1	3	...	2		9	0	100.00...	0	1													
579	7	2	1	3	...	2		3	0	100.00...	0	1													
797	10	3	1	3	...	2		113	0	100.00...	0	1													
1251	16	0	1	3	...	2		65	0	100.00...	0	1													
1396	19	0	1	3	...	2		6	0	100.00...	0	1													
1506	20	3	1	3	...	2		47	0	100.00...	0	1													
1593	22	3	1	3	...	2		21	0	100.00...	0	1													
1901	26	2	1	3	...	2		29	0	100.00...	0	1													
4378	60	1	1	3	...	2		42	0	100.00...	0	1													
4581	63	0	1	3	...	2		1	0	100.00...	0	1													
5545	77	3	1	3	...	2		8	0	100.00...	0	1													
5868	82	1	1	3	...	2		14	0	100.00...	0	1													
5945	83	3	1	3	...	2		7	0	100.00...	0	1													
6217	87	0	1	3	...	2		5	0	100.00...	0	1													
7111	99	1	1	3	...	2		8	0	100.00...	0	1													
228	3	1	1	2	...	2		3	0	100.00...	0	1													
1118	15	1	1	2	...	2		5	0	100.00...	0	1													
4491	62	3	1	2	...	2		16	0	100.00...	0	1													
4496	62	3	1	2	...	2		27	0	100.00...	0	1													
4498	62	3	1	2	...	2		31	0	100.00...	0	1													
4889	68	0	1	2	...	2		80	0	100.00...	0	1													
4893	68	0	1	2	...	2		82	0	100.00...	0	1													
5041	70	0	1	2	...	2		5	0	100.00...	0	1													
5859	82	1	1	2	...	2		1	0	100.00...	0	1													
6963	97	1	1	2	...	2		8	0	100.00...	0	1													
229	3	1	1	1	...	2		4	0	100.00...	0	1													

6761 from the 7173 earlier selected or derived rules remaining after removal of 412 contradictory rules.

# The rules predicting decrease of pandemic level despite high mobility remain after removing the contradictions



# RULES MISSING ALL MOBILITY FEATURES

R0		R0.1		R0.1.1		R0.1.1.1		R0.2		R0.3		R0.4	
Id	Tre...	Wei...	Class	or...	N conditions	Rule	Days_passed	Week6_Covid	Week5_Covid	Week4_Covid	Week3_Covid	Week2_Covid	
1	0	3	1	2	0	1	2						
4317	60	1	1	2	1	1	2						
4889	68	0	1	2	2	1	2						
4893	68	0	1	2	3	1	2						
6606	92	0	1	2	4	1	2						
5859	82	1	1	2	5	1	2						
2706	38	1	1	2	6	1	4						
1530	21	1	1	2	7	1	3						
5860	82	1	1	1	8	1	2						
1677	23	1	1	1	9	1	3						
6602	92	0	1	2	10	1	2						
2170	30	2	1	1	11	1	2						
5946	83	3	1	1	12	1	2						
5863	82	1	1	2	13	1	3						
15	0	3	1	3	14	1	2						
1935	27	0	1	3	15	1	2						
5945	83	3	1	3	16	1	2						
797	10	3	1	3	17	1	2						
5868	82	1	1	3	18	1	2						
2189	30	2	1	3	19	1	2						
2611	36	2	1	3	20	1	2						
1844	25	2	1	3	21	1	1						
4373	60	1	1	3	22	1	2						
5928	82	1	1	3	23	1	2						
7026	97	1	1	3	24	1	1						
282	3	1	1	3	25	1	2						
1969	27	0	1	3	26	1	3						
2682	37	1	1	3	27	1	2						
3060	42	1	1	3	28	1	3						

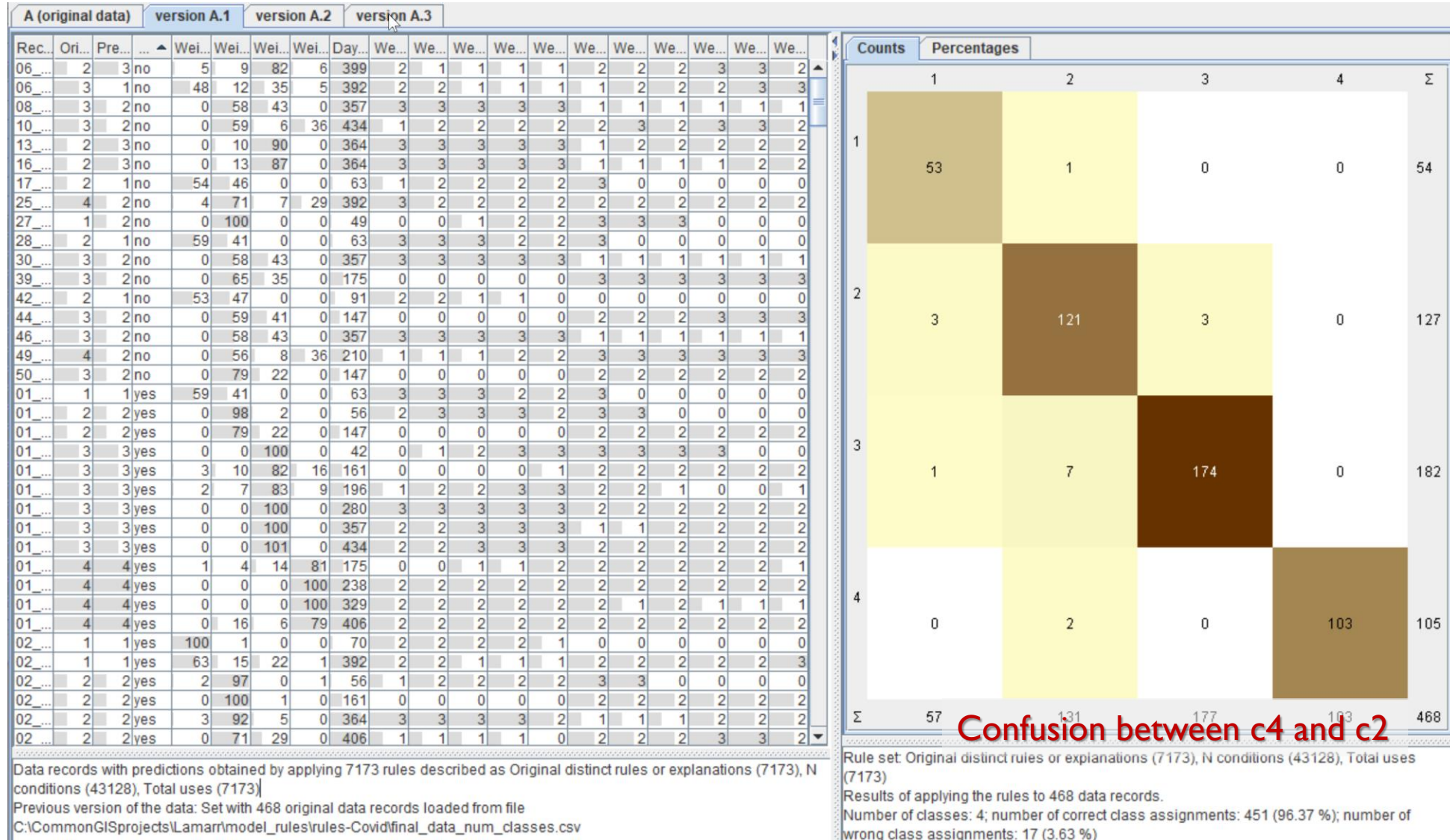
Our model contains 232 rules that do not involve any of the mobility features. There are also 46 rules missing all COVID features. This does not align with the initial goal to predict the impact of mobility levels on COVID-19 development depending on the prior temporal context.



# APPLICATION OF RULES TO DATA

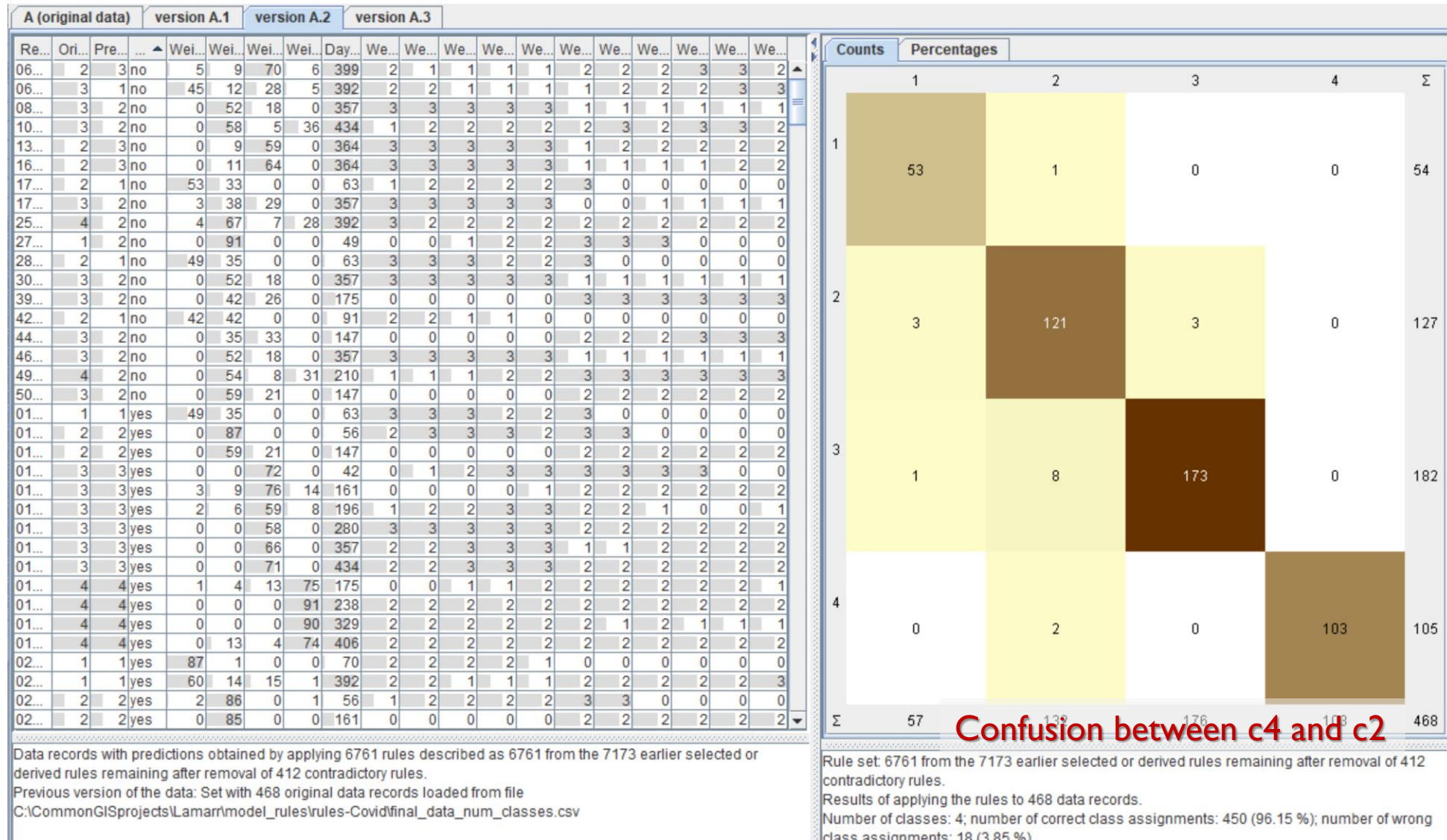
## Full set of 7173 original rules

17 misclassified instances (3.63%)



# APPLICATION OF RULES TO DATA

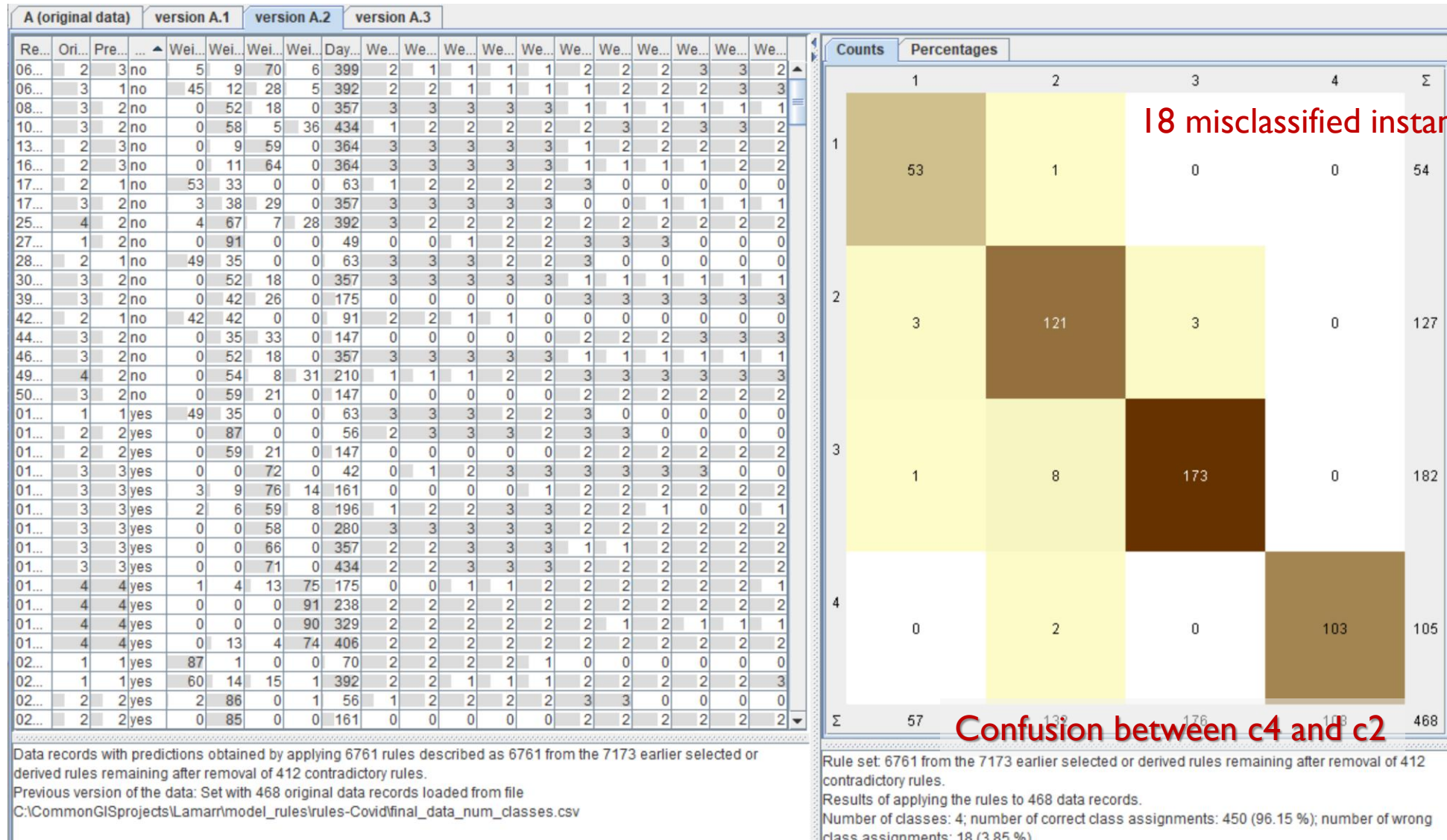
Subset of 676 | non-contradictory original rules 18 misclassified instances (3.85%)





# APPLICATION OF RULES TO DATA

6600 non-contradictory rules after removing the rules missing mobility features



# WHERE THE CONFUSION HAPPENS

Here the pandemic level remains stable at c2 while the level of mobility is reduced (m2). The prediction of class c2 rather than c4 appears reasonable.

Record ID: 25\_c4\_7\_15/03/2021\_02/05/2021  
Row: 242

Column	Value
Original Class/Value	4.0
Predicted Class/Value	2
Match?	no
Weight 1	4
Weight 2	62
Weight 3	7
Weight 4	28
Days_passed	392
Week6_Covid	3
Week5_Covid	2
Week4_Covid	2
Week3_Covid	2
Week2_Covid	2
Week6_Mobility	2
Week5_Mobility	2
Week4_Mobility	2
Week3_Mobility	2
Week2_Mobility	2
Week1_Mobility	2

Record ID: 49\_c4\_10\_14/09/2020\_22/11/2020  
Row: 457

Column	Value
Original Class/Value	4.0
Predicted Class/Value	2
Match?	no
Weight 1	0
Weight 2	54
Weight 3	8
Weight 4	29
Days_passed	210
Week6_Covid	1
Week5_Covid	1
Week4_Covid	1
Week3_Covid	2
Week2_Covid	2
Week6_Mobility	3
Week5_Mobility	3
Week4_Mobility	3
Week3_Mobility	3
Week2_Mobility	3
Week1_Mobility	3

Here the pandemic level slightly increased from c1 to c2 while the mobility level remains relatively high (m3). Here an increase of the pandemic level would be expected.

Evidently, the 21 rules predicting class c2 when mobility in week -2 is high are not responsible for these confusions. Hence, there were training data instances corresponding to these rules.



# WHAT WE HAVE LEARNT ABOUT THE MODEL

- **Unwanted behaviour**: making predictions based on insufficient (too few) conditions
- **Unwanted behaviour**: making predictions while ignoring mobility features or prior pandemic levels
- **Unwanted property**: contradictions among the rules – multiple rules predicting different classes can be applied to the same instances
  - Removal of the contradictions, rules ignoring mobility or pandemic features, and remaining rules with less than 3 conditions only slightly (by 0.22%) decreases the model accuracy for an available test dataset. This could be acceptable for the sake of improving model logic.
- **Unwanted property**: some rules are not justifiable by domain logic or common sense.
  - However, they seem to be in accord with the training and test data. Hence, there are real cases contradicting the logical expectations.

# GENERAL INSIGHTS: DIFFERENCES BETWEEN ML AND HUMAN REASONING

- A model may reach correct conclusions but may not "think" like a human.
- Some rules may lack domain-relevant conditions yet still function well.
- Adjusting or filtering rules based on human logic might not significantly affect accuracy—suggesting redundancy or alternative reasoning paths.

# CONCLUSIONS

- Trustworthiness is not just about accuracy—it's about understanding *why* the model makes decisions.
- How should we balance human logic vs. data-driven inferences when interpreting and explaining models?
- **Open question:** Should ML models be adjusted to align better with human reasoning, even if accuracy does not improve or may even slightly degrade?
  - If so, how can we incorporate domain knowledge and human logic at the stage of model training?