

# Analysis of the Effectiveness of G3PARM Algorithm

J.M. Luna, J.R. Romero and S. Ventura

*Knowledge Discovery and Intelligent Systems Research Group  
University of Córdoba, Spain*

HAIS 2010. San Sebastián, Spain. June 2010.

# Outline

- Introduction
- Model description
- Experiments
- Results
- Conclusions
- Future research lines

# Introduction

- Association rules mining (ARM): a Data Mining technique
- Classical algorithms for mining association rules:
  - Apriori
  - FP-Growth
- Genetic Programming: computer programs
- Grammar Guided Genetic Programming (G3P)
- Our proposal is a G3P-based approach for ARM

# Model description

- G3PARM = G3P Association Rules Mining
- Designed for the extraction of association rules
- Based on the use of a context-free grammar
- Each individual represents a rule
- An auxiliary population:
  - Exceed a minimum threshold for two different measures
  - Maximum auxiliary population size

$G = (\Sigma_N, \Sigma_T, P, Rule)$  with:

$\Sigma_N = \{Rule, Antecedent, Consequent, Comparison, Categorical Comparator, Categorical Attribute Comparison \}$

$\Sigma_T = \{AND, "!=" , "=", "name", "value" \}$

$P = \{Rule = Antecedent, Consequent ;$

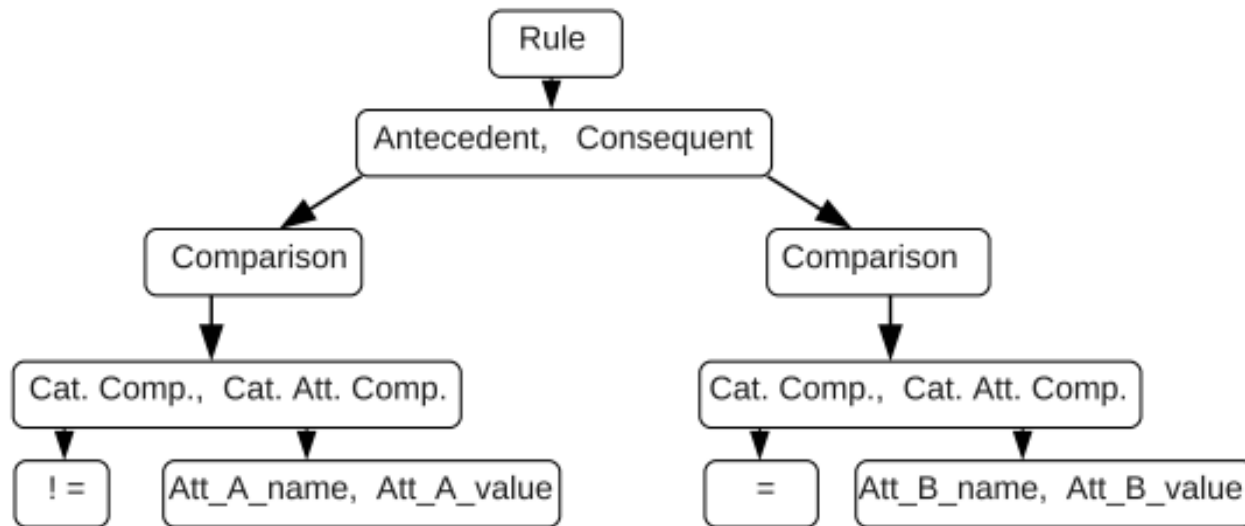
Antecedent = Comparison | AND, Comparison, Antecedent ;

Consequent = Comparison ;

Comparison = Categorical Comparator, Categorical Attribute Comparison ;

Categorical Comparator = "!=" | "=" ;

Categorical Attribute Comparison = "name", "value" ;}

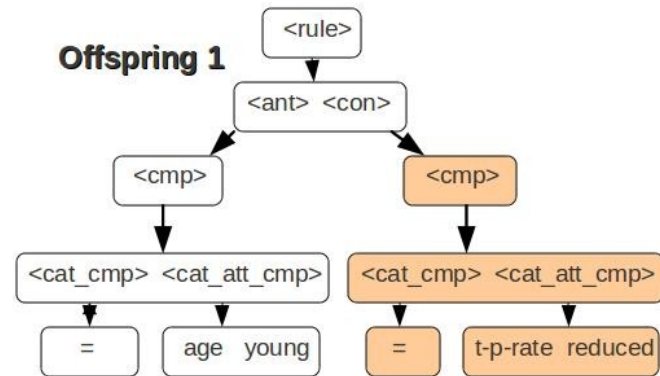
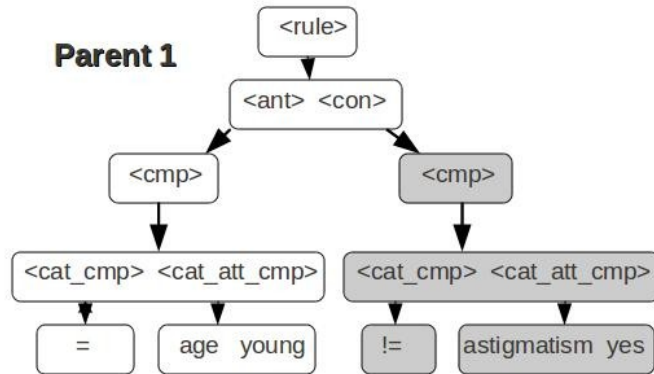


Phenotype:

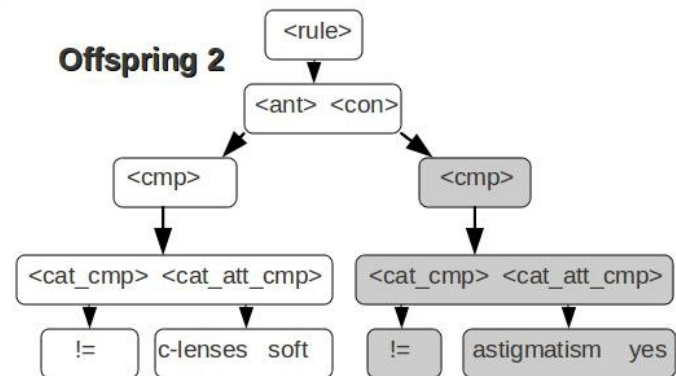
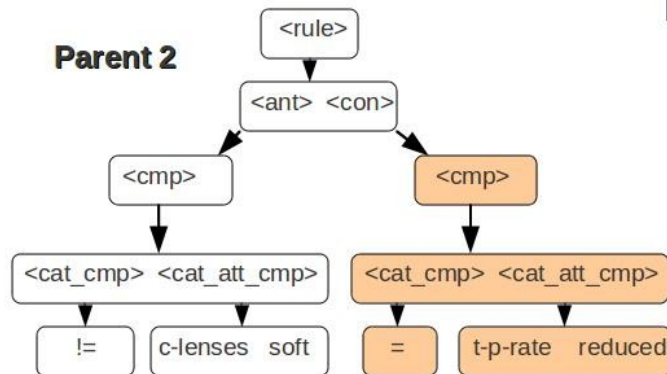
$(\neq \text{Att\_A\_name Att\_A\_value}) \rightarrow (= \text{Att\_B\_name Att\_B\_value})$

# Model description

## Crossover

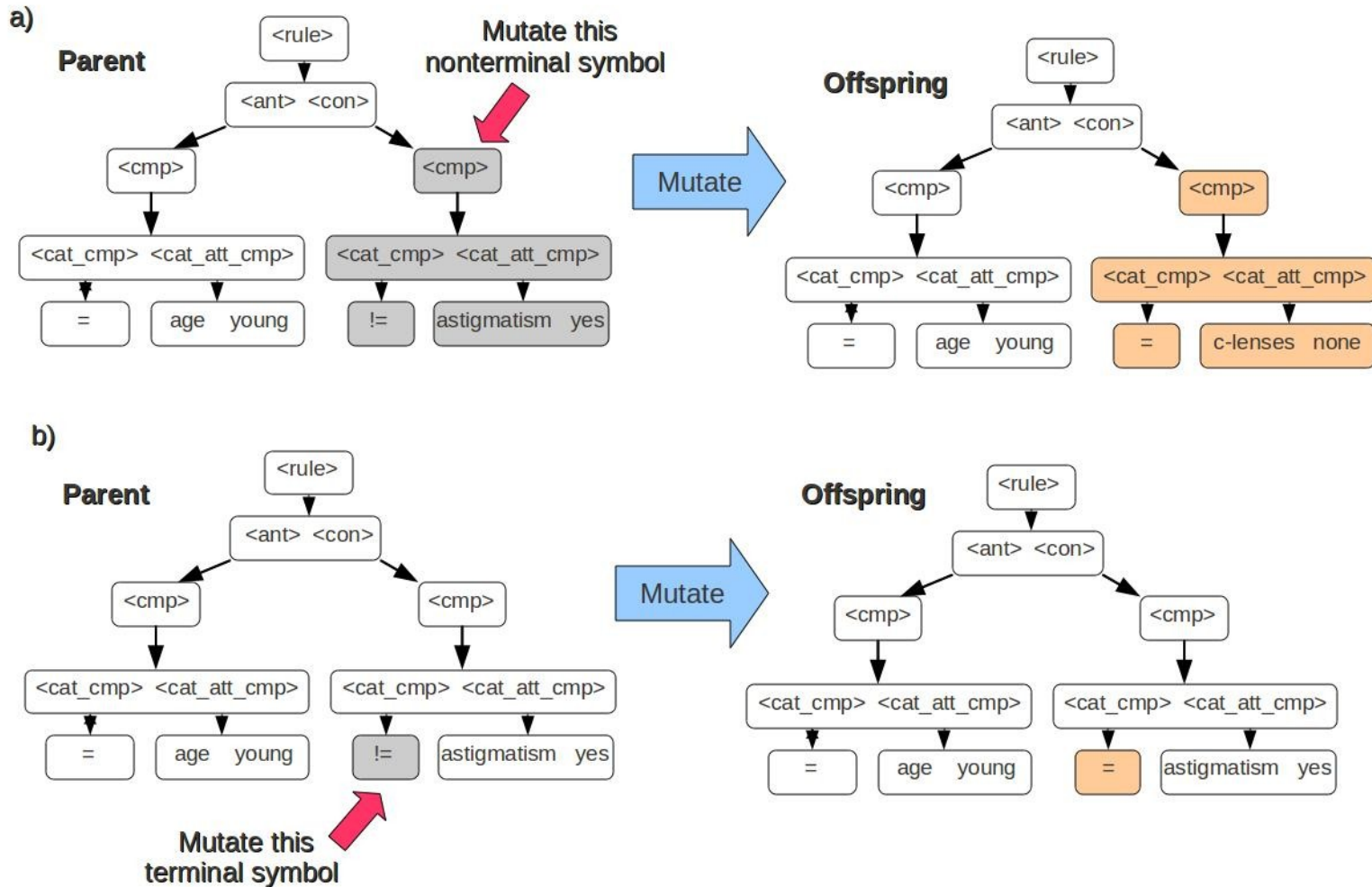


Crossover



# Model description

# Mutation





- G3PARM uses two measures:
  - **Support** measure: proportion of the number of transactions  $T$  including the antecedent  $A$  and the consequent  $C$  in a dataset  $D$ .

$$\text{supp}(A \rightarrow C) = \frac{|\{A \cup C \subseteq T, T \in D\}|}{|D|}$$

- **Confidence** measure: proportion of the number of transactions which include  $A$  and  $C$  in transaction which include  $A$ .

$$\text{conf}(A \rightarrow C) = \frac{|\{A \cup C \subseteq T, T \in D\}|}{|\{A \subseteq T, T \in D\}|}$$

- **Fitness function** = support measure

---

**Algorithm 1** G3PARM algorithm

---

**Require:**  $max\_generations, N$

**Ensure:**  $A$

```
1:  $P \leftarrow random(N)$ 
2:  $A \leftarrow \emptyset$ 
3:  $num\_generations \leftarrow 0$ 
4: while  $num\_generations < max\_generations$  do
5:    $Parents \leftarrow \text{Select parents } (P \cup A)$ 
6:    $Crossover(Parents)$ 
7:    $Mutation(Parents)$ 
8:    $Update\ auxiliary\ population(A, P)$ 
9:    $num\_generations ++$ 
10: end while
11: return  $A$ 
```

---

---

**Algorithm 2** Update auxiliary population

---

**Require:**  $A, P$

**Ensure:**  $A$

```
1:  $A' \leftarrow P \cup A$ 
2:  $Order(A')$ 
3:  $Eliminate\ duplicate(A')$ 
4:  $A \leftarrow Threshold(A')$ 
5: return  $A$ 
```

---

---

### Algorithm 1 G3PARM algorithm

---

**Require:**  $max\_generations, N$

**Ensure:**  $A$

```
1:  $P \leftarrow random(N)$ 
2:  $A \leftarrow \emptyset$ 
3:  $num\_generations \leftarrow 0$ 
4: while  $num\_generations < max\_generations$  do
5:    $Parents \leftarrow \text{Select parents } (P \cup A)$ 
6:    $Crossover(Parents)$ 
7:    $Mutation(Parents)$ 
8:    $Update\ auxiliary\ population(A, P)$ 
9:    $num\_generations ++$ 
10: end while
11: return  $A$ 
```

---

Rules comprised by the same attributes are considered equals:

1 AND 2 -> 3

2 AND 1 -> 3

---

### Algorithm 2 Update auxiliary population

---

**Require:**  $A, P$

**Ensure:**  $A$

```
1:  $A' \leftarrow P \cup A$ 
2:  $Order(A')$ 
3:  $Eliminate\ duplicate(A')$ 
4:  $A \leftarrow Threshold(A')$ 
5: return  $A$ 
```

---

---

**Algorithm 1** G3PARM algorithm

---

**Require:**  $max\_generations, N$

**Ensure:**  $A$

```
1:  $P \leftarrow random(N)$ 
2:  $A \leftarrow \emptyset$ 
3:  $num\_generations \leftarrow 0$ 
4: while  $num\_generations < max\_generations$  do
5:    $Parents \leftarrow \text{Select parents } (P \cup A)$ 
6:    $Crossover(Parents)$ 
7:    $Mutation(Parents)$ 
8:    $Update\ auxiliary\ population(A, P)$ 
9:    $num\_generations ++$ 
10: end while
11: return  $A$ 
```

---

Minimum threshold for support and confidence measures

---

**Algorithm 2** Update auxiliary population

---

**Require:**  $A, P$

**Ensure:**  $A$

```
1:  $A' \leftarrow P \cup A$ 
2:  $Order(A')$ 
3:  $Eliminate\ duplicate(A')$ 
4:  $A \leftarrow Threshold(A')$ 
5: return  $A$ 
```

---

- To **evaluate** the usefulness of **G3PARM**, several experiments have been carried out on different **datasets**:
  - Credit-g: 1000 instances and 21 attributes
  - HH: 22784 instances and 17 attributes
  - Mushroom: 8124 instances and 23 attributes
  - Segment: 1500 instances and 20 attributes
  - Sonar: 208 instances and 36 attributes
  - Soybean: 683 instances and 36 attributes
  - Wisconsin Breast Cancer: 683 instances and 11 attributes
- **Discretization** of datasets with numerical attributes.

- **G3PARM:**
  - Population size: 50
  - Crossover probability: 70%
  - Mutation probability: 10%
  - Maximum derivation numbers: 24
  - External population size: 20
  - Support threshold: 70%
  - Confidence threshold: 90%
- **Apriori and FP-Growth:**
  - Support threshold: 70%
  - Confidence threshold: 90%

# Results

## Results obtained

Dataset	Average_support			Average_confidence			%Instances		
	(1)	(2)	(3)	(1)	(2)	(3)	(1)	(2)	(3)
<i>CreditEqFre10</i>	0.780	0.709	<b>0.850</b>	<b>0.941</b>	0.855	0.939	0.987	0.987	<b>1.000</b>
<i>CreditEqFre5</i>	0.780	0.709	<b>0.850</b>	0.941	0.855	<b>0.953</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid10</i>	0.780	0.709	<b>0.892</b>	0.941	0.855	<b>0.965</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid5</i>	0.773	0.709	<b>0.858</b>	0.942	0.863	<b>0.961</b>	0.989	0.989	<b>1.000</b>
<i>HHEqFre10</i>	None	None	<b>0.803</b>	None	None	<b>0.913</b>	None	None	<b>1.000</b>
<i>HHEqFreq5</i>	None	None	<b>0.740</b>	None	None	<b>0.909</b>	None	None	<b>0.997</b>
<i>HHEqWid10</i>	0.761	0.761	<b>0.922</b>	0.950	0.950	<b>0.986</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>HHEqWid5</i>	0.765	0.765	<b>0.902</b>	0.955	0.955	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>Mushroom</i>	0.824	0.817	<b>0.890</b>	0.968	0.960	<b>0.978</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqFre10</i>	<b>0.876</b>	<b>0.876</b>	0.813	<b>0.975</b>	<b>0.975</b>	0.926	0.996	0.996	<b>1.000</b>
<i>SegmentEqFre5</i>	<b>0.876</b>	<b>0.876</b>	0.817	<b>0.975</b>	<b>0.975</b>	0.974	0.996	0.996	<b>1.000</b>
<i>SegmentEqWid10</i>	0.815	0.815	<b>0.884</b>	0.968	0.968	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqWid5</i>	0.860	0.860	<b>0.882</b>	0.964	0.964	<b>0.969</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SonarEqFre10</i>	None	None	<b>0.782</b>	None	None	<b>0.909</b>	None	None	<b>1.000</b>
<i>SonarEqFre5</i>	None	None	<b>0.583</b>	None	None	<b>0.731</b>	None	None	<b>0.626</b>
<i>SonarEqWid10</i>	None	None	<b>0.958</b>	None	None	<b>0.887</b>	None	None	<b>1.000</b>
<i>SonarEqWid5</i>	0.747	None	<b>0.835</b>	0.942	None	<b>0.947</b>	0.846	None	<b>1.000</b>
<i>Soybean</i>	0.778	0.722	<b>0.822</b>	0.950	0.953	<b>0.957</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>WBCEqFre10</i>	None	None	<b>0.875</b>	None	None	<b>0.958</b>	None	None	<b>1.000</b>
<i>WBCEqFre5</i>	None	None	<b>0.806</b>	None	None	<b>0.928</b>	None	None	<b>1.000</b>
<i>WBCEqWid10</i>	0.821	0.821	<b>0.900</b>	<b>0.996</b>	<b>0.996</b>	0.971	0.821	0.821	<b>1.000</b>
<i>WBCEqWid5</i>	<b>0.872</b>	<b>0.872</b>	0.864	<b>0.996</b>	<b>0.996</b>	0.956	0.872	0.872	<b>1.000</b>
<b>Ranking</b>	2.204	2.522	<b>1.272</b>	2.159	2.431	<b>1.409</b>	2.340	2.386	<b>1.272</b>

- (1) Apriori
- (2) FP-Growth
- (3) G3PARM

# Results

## Results obtained

Dataset	Average_support			Average_confidence			%Instances		
	(1)	(2)	(3)	(1)	(2)	(3)	(1)	(2)	(3)
<i>CreditEqFre10</i>	0.780	0.709	<b>0.850</b>	<b>0.941</b>	0.855	0.939	0.987	0.987	<b>1.000</b>
<i>CreditEqFre5</i>	0.780	0.709	<b>0.850</b>	0.941	0.855	<b>0.953</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid10</i>	0.780	0.709	<b>0.892</b>	0.941	0.855	<b>0.965</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid5</i>	0.773	0.709	<b>0.858</b>	0.942	0.863	<b>0.961</b>	0.989	0.989	<b>1.000</b>
<i>HHEqFre10</i>	None	None	<b>0.803</b>	None	None	<b>0.913</b>	None	None	<b>1.000</b>
<i>HHEqFreq5</i>	None	None	<b>0.740</b>	None	None	<b>0.909</b>	None	None	<b>0.997</b>
<i>HHEqWid10</i>	0.761	0.761	<b>0.922</b>	0.950	0.950	<b>0.986</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>HHEqWid5</i>	0.765	0.765	<b>0.902</b>	0.955	0.955	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>Mushroom</i>	0.824	0.817	<b>0.890</b>	0.968	0.960	<b>0.978</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqFre10</i>	<b>0.876</b>	<b>0.876</b>	0.813	<b>0.975</b>	<b>0.975</b>	0.926	0.996	0.996	<b>1.000</b>
<i>SegmentEqFre5</i>	<b>0.876</b>	<b>0.876</b>	0.817	<b>0.975</b>	<b>0.975</b>	0.974	0.996	0.996	<b>1.000</b>
<i>SegmentEqWid10</i>	0.815	0.815	<b>0.884</b>	0.968	0.968	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqWid5</i>	0.860	0.860	<b>0.882</b>	0.964	0.964	<b>0.969</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SonarEqFre10</i>	None	None	<b>0.782</b>	None	None	<b>0.909</b>	None	None	<b>1.000</b>
<i>SonarEqFre5</i>	None	None	<b>0.583</b>	None	None	<b>0.731</b>	None	None	<b>0.626</b>
<i>SonarEqWid10</i>	None	None	<b>0.958</b>	None	None	<b>0.887</b>	None	None	<b>1.000</b>
<i>SonarEqWid5</i>	0.747	None	<b>0.835</b>	0.942	None	<b>0.947</b>	0.846	None	<b>1.000</b>
<i>Soybean</i>	0.778	0.722	<b>0.822</b>	0.950	0.953	<b>0.957</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>WBCEqFre10</i>	None	None	<b>0.875</b>	None	None	<b>0.958</b>	None	None	<b>1.000</b>
<i>WBCEqFre5</i>	None	None	<b>0.806</b>	None	None	<b>0.928</b>	None	None	<b>1.000</b>
<i>WBCEqWid10</i>	0.821	0.821	<b>0.900</b>	<b>0.996</b>	<b>0.996</b>	0.971	0.821	0.821	<b>1.000</b>
<i>WBCEqWid5</i>	<b>0.872</b>	<b>0.872</b>	0.864	<b>0.996</b>	<b>0.996</b>	0.956	0.872	0.872	<b>1.000</b>
<b>Ranking</b>	2.204	2.522	<b>1.272</b>	2.159	2.431	<b>1.409</b>	2.340	2.386	<b>1.272</b>

- (1) Apriori
- (2) FP-Growth
- (3) G3PARM



# Results

## Results obtained

Dataset	Average_support			Average_confidence			%Instances		
	(1)	(2)	(3)	(1)	(2)	(3)	(1)	(2)	(3)
<i>CreditEqFre10</i>	0.780	0.709	<b>0.850</b>	<b>0.941</b>	0.855	0.939	0.987	0.987	<b>1.000</b>
<i>CreditEqFre5</i>	0.780	0.709	<b>0.850</b>	0.941	0.855	<b>0.953</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid10</i>	0.780	0.709	<b>0.892</b>	0.941	0.855	<b>0.965</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid5</i>	0.773	0.709	<b>0.858</b>	0.942	0.863	<b>0.961</b>	0.989	0.989	<b>1.000</b>
<i>HHEqFre10</i>	None	None	<b>0.803</b>	None	None	<b>0.913</b>	None	None	<b>1.000</b>
<i>HHEqFreq5</i>	None	None	<b>0.740</b>	None	None	<b>0.909</b>	None	None	<b>0.997</b>
<i>HHEqWid10</i>	0.761	0.761	<b>0.922</b>	0.950	0.950	<b>0.986</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>HHEqWid5</i>	0.765	0.765	<b>0.902</b>	0.955	0.955	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>Mushroom</i>	0.824	0.817	<b>0.890</b>	0.968	0.960	<b>0.978</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqFre10</i>	<b>0.876</b>	<b>0.876</b>	0.813	<b>0.975</b>	<b>0.975</b>	0.926	0.996	0.996	<b>1.000</b>
<i>SegmentEqFre5</i>	<b>0.876</b>	<b>0.876</b>	0.817	<b>0.975</b>	<b>0.975</b>	0.974	0.996	0.996	<b>1.000</b>
<i>SegmentEqWid10</i>	0.815	0.815	<b>0.884</b>	0.968	0.968	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqWid5</i>	0.860	0.860	<b>0.882</b>	0.964	0.964	<b>0.969</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SonarEqFre10</i>	None	None	<b>0.782</b>	None	None	<b>0.909</b>	None	None	<b>1.000</b>
<i>SonarEqFre5</i>	None	None	<b>0.583</b>	None	None	<b>0.731</b>	None	None	<b>0.626</b>
<i>SonarEqWid10</i>	None	None	<b>0.958</b>	None	None	<b>0.887</b>	None	None	<b>1.000</b>
<i>SonarEqWid5</i>	0.747	None	<b>0.835</b>	0.942	None	<b>0.947</b>	0.846	None	<b>1.000</b>
<i>Soybean</i>	0.778	0.722	<b>0.822</b>	0.950	0.953	<b>0.957</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>WBCEqFre10</i>	None	None	<b>0.875</b>	None	None	<b>0.958</b>	None	None	<b>1.000</b>
<i>WBCEqFre5</i>	None	None	<b>0.806</b>	None	None	<b>0.928</b>	None	None	<b>1.000</b>
<i>WBCEqWid10</i>	0.821	0.821	<b>0.900</b>	<b>0.996</b>	<b>0.996</b>	0.971	0.821	0.821	<b>1.000</b>
<i>WBCEqWid5</i>	<b>0.872</b>	<b>0.872</b>	0.864	<b>0.996</b>	<b>0.996</b>	0.956	0.872	0.872	<b>1.000</b>
Ranking	2.204	2.522	<b>1.272</b>	2.159	2.431	<b>1.409</b>	2.340	2.386	<b>1.272</b>

- (1) Apriori
- (2) FP-Growth
- (3) G3PARM

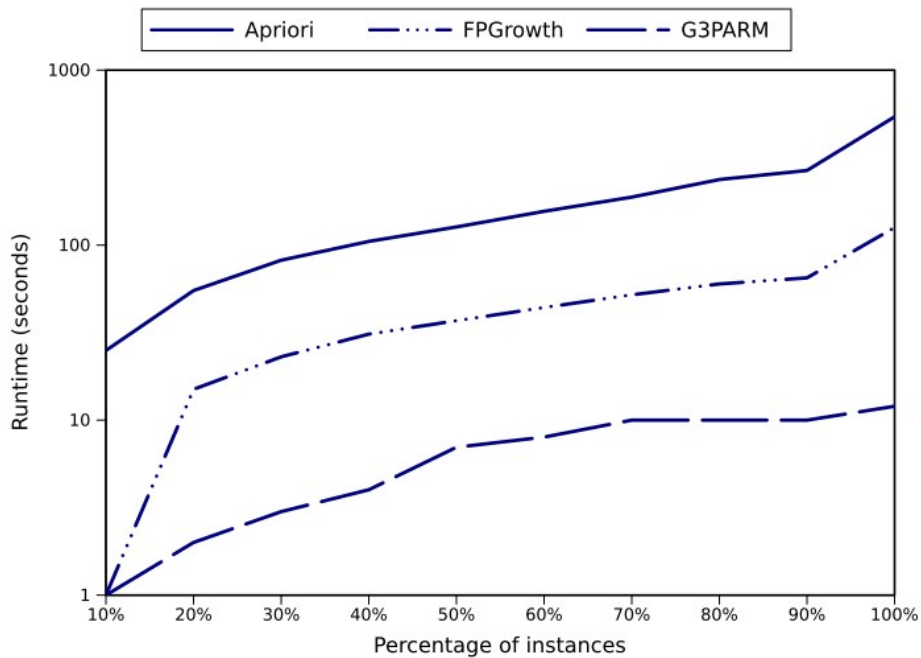
# Results

## Results obtained

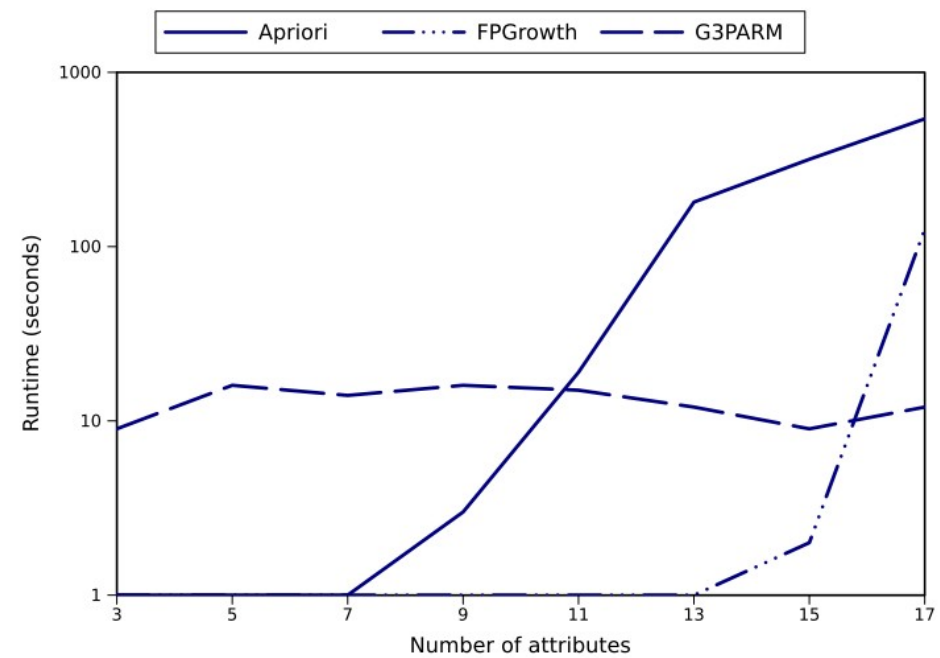
Dataset	Average_support			Average_confidence			%Instances		
	(1)	(2)	(3)	(1)	(2)	(3)	(1)	(2)	(3)
<i>CreditEqFre10</i>	0.780	0.709	<b>0.850</b>	<b>0.941</b>	0.855	0.939	0.987	0.987	<b>1.000</b>
<i>CreditEqFre5</i>	0.780	0.709	<b>0.850</b>	0.941	0.855	<b>0.953</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid10</i>	0.780	0.709	<b>0.892</b>	0.941	0.855	<b>0.965</b>	0.987	0.987	<b>1.000</b>
<i>CreditEqWid5</i>	0.773	0.709	<b>0.858</b>	0.942	0.863	<b>0.961</b>	0.989	0.989	<b>1.000</b>
<i>HHEqFre10</i>	None	None	<b>0.803</b>	None	None	<b>0.913</b>	None	None	<b>1.000</b>
<i>HHEqFreq5</i>	None	None	<b>0.740</b>	None	None	<b>0.909</b>	None	None	<b>0.997</b>
<i>HHEqWid10</i>	0.761	0.761	<b>0.922</b>	0.950	0.950	<b>0.986</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>HHEqWid5</i>	0.765	0.765	<b>0.902</b>	0.955	0.955	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>Mushroom</i>	0.824	0.817	<b>0.890</b>	0.968	0.960	<b>0.978</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqFre10</i>	<b>0.876</b>	<b>0.876</b>	0.813	<b>0.975</b>	<b>0.975</b>	0.926	0.996	0.996	<b>1.000</b>
<i>SegmentEqFre5</i>	<b>0.876</b>	<b>0.876</b>	0.817	<b>0.975</b>	<b>0.975</b>	0.974	0.996	0.996	<b>1.000</b>
<i>SegmentEqWid10</i>	0.815	0.815	<b>0.884</b>	0.968	0.968	<b>0.979</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SegmentEqWid5</i>	0.860	0.860	<b>0.882</b>	0.964	0.964	<b>0.969</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>SonarEqFre10</i>	None	None	<b>0.782</b>	None	None	<b>0.909</b>	None	None	<b>1.000</b>
<i>SonarEqFre5</i>	None	None	<b>0.583</b>	None	None	<b>0.731</b>	None	None	<b>0.626</b>
<i>SonarEqWid10</i>	None	None	<b>0.958</b>	None	None	<b>0.887</b>	None	None	<b>1.000</b>
<i>SonarEqWid5</i>	0.747	None	<b>0.835</b>	0.942	None	<b>0.947</b>	0.846	None	<b>1.000</b>
<i>Soybean</i>	0.778	0.722	<b>0.822</b>	0.950	0.953	<b>0.957</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
<i>WBCEqFre10</i>	None	None	<b>0.875</b>	None	None	<b>0.958</b>	None	None	<b>1.000</b>
<i>WBCEqFre5</i>	None	None	<b>0.806</b>	None	None	<b>0.928</b>	None	None	<b>1.000</b>
<i>WBCEqWid10</i>	0.821	0.821	<b>0.900</b>	<b>0.996</b>	<b>0.996</b>	0.971	0.821	0.821	<b>1.000</b>
<i>WBCEqWid5</i>	<b>0.872</b>	<b>0.872</b>	0.864	<b>0.996</b>	<b>0.996</b>	0.956	0.872	0.872	<b>1.000</b>
Ranking	2.204	2.522	<b>1.272</b>	2.159	2.431	<b>1.409</b>	2.340	2.386	<b>1.272</b>

- (1) Apriori
- (2) FP-Growth
- (3) G3PARM

## Different number of instances



## Different number of attributes



# Conclusions

- Novel **G3P-based** algorithm for mining association rules
- **G3PARM** obtains rules with:
  - High support
  - High confidence
  - High representative rules
- **G3PARM** scales quite linearly as we increase up the dataset size and the number of attributes

# Future research lines

- Use of numerical attributes
- Modify our approach to work with infrequent pattern
- Multiobjective version using support and confidence as objectives to be optimized



# Thanks!



## Analysis of the Effectiveness of G3PARM Algorithm

J.M. Luna, J.R. Romero and S. Ventura

*Knowledge Discovery and Intelligent Systems Research Group  
University of Córdoba, Spain*

HAIS 2010. San Sebastián, Spain. June 2010.