

# A Laravel package for systematic construction of symbolic intelligent systems

Alberto Paparella

Università degli Studi di Ferrara  
Dipartimento di Matematica e Informatica

Relatore  
Prof. **Guido Sciavicco**

Correlatori  
Dott. **Giovanni Pagliarini**  
Prof. **Giacomo Piva**

 **Università  
degli Studi  
di Ferrara**

4 ottobre 2021

# Contents

- 1 Introduction
  - Pitón
- 2 Preliminaries
  - Hierarchical problems
  - Rule-based models
  - Rule-extraction
- 3 How to use the package
  - Install
  - Workflow
  - Configuring the database
  - Training
  - Prediction
- 4 Automated treatment suggestion
  - Problem
  - Database
  - Experiment results
- 5 Backup slides

# Pitón

## Pitón

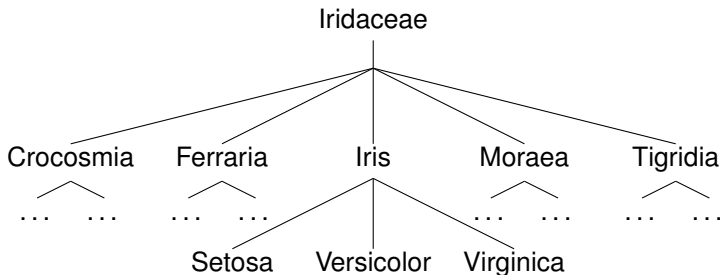
Pitón is a Laravel package which offers machine learning utilities for solving **hierarchically-arranged problems**, where the available data is stored in a **MySQL database**, through the creation of **rule-based models**.

# Setting

The package can be placed in the *supervised learning* area of machine learning. The target of the package is learning from a *dataset* (e.g. a set of data representable in tabular form), making use of an *algorithm* we refer to as *learner*, which synthesizes a *model* that can be used to classify new *instances*.

Specifically, an algorithm divides the dataset into a *training* set and a *test* set, and it *fits* a model based on *input* and *output* values (e.g. features values and respective classification label in the *train* set); then, it tests the extracted model using only the *input* values of the *test* set and comparing the results to their original classification.

# Hierarchical problems



**Figure:** Hierarchy of problems for Iridaceae classification

# Rule-based models

$$\Gamma = \begin{cases} p_1^1 \wedge \dots \wedge p_n^1 \rightarrow C_1 & \text{else} \\ p_1^2 \wedge \dots \wedge p_n^2 \rightarrow C_2 & \text{else} \\ \dots \\ p_1^n \wedge \dots \wedge p_n^n \rightarrow C_n \end{cases}$$

Figure: Example of a generic rule-based model

# Binary rule-based models

$$\Gamma = \begin{cases} p_1^1 \wedge \dots \wedge p_n^1 \rightarrow \mathcal{C} & \text{else} \\ p_1^2 \wedge \dots \wedge p_n^2 \rightarrow \mathcal{C} & \text{else} \\ \dots \\ p_1^n \wedge \dots \wedge p_n^n \rightarrow \mathcal{C} & \text{else} \\ \neg \mathcal{C} \end{cases}$$

Figure: Example of a generic *binary* rule-based model

# Binary classification

$$\Gamma = \begin{cases} p_1^1 \wedge \dots \wedge p_n^1 \rightarrow C & \text{else} \\ p_1^2 \wedge \dots \wedge p_n^2 \rightarrow C & \text{else} \\ \dots & \\ p_1^n \wedge \dots \wedge p_n^n \rightarrow C & \text{else} \\ \neg C & \end{cases}$$

**Figure:** Example of a generic *binary* rule-based model

Classification with binary rule-based models:

- 1 try the first rule, if it activates predict  $C$ , if not try the following rule, and so on;
- 2 if all rules fail, then predict  $\neg C$ , such as the instance is not of class  $C$ .



**Università  
degli Studi  
di Ferrara**

# Binary classification

$$\Gamma = \begin{cases} p_1^1 \wedge \dots \wedge p_n^1 \rightarrow \mathcal{C} & \text{else} \\ p_1^2 \wedge \dots \wedge p_n^2 \rightarrow \mathcal{C} & \text{else} \\ \dots & \\ p_1^n \wedge \dots \wedge p_n^n \rightarrow \mathcal{C} & \text{else} \\ \neg \mathcal{C} & \end{cases}$$

**Figure:** Example of a generic *binary* rule-based model

Classification with binary rule-based models:

- 1 try the first rule, if it activates predict  $\mathcal{C}$ , if not try the following rule, and so on;
- 2 if all rules fail, then predict  $\neg \mathcal{C}$ , such as the instance is not of class  $\mathcal{C}$ .



# RIPPERk

To build a rule:

- 1 Randomly partition all the examples which have not been covered by any rule yet into two subsets: a *growing set* and a *pruning set*.
- 2 Grow a rule by greedily adding conditions until the rule reaches a confidence of 100% using only the growing set.
- 3 To prevent growing set *overfitting*, immediately prune the rule, deleting some conditions based on a pruning criterion, using the pruning data.
- 4 Lastly, all the positive and negative examples covered by the rule must be removed; then repeat from the step 1.
- 5 RIPPERk stops adding rules when there are no more positive examples left or when a rule has an unacceptably large error rate, or when the last rule added is too *complicated* according to some criterion.



# RIPPERk

To build a rule:

- 1 Randomly partition all the examples which have not been covered by any rule yet into two subsets: a *growing set* and a *pruning set*.
- 2 Grow a rule by greedily adding conditions until the rule reaches a confidence of 100% using only the growing set.
- 3 To prevent growing set *overfitting*, immediately prune the rule, deleting some conditions based on a pruning criterion, using the pruning data.
- 4 Lastly, all the positive and negative examples covered by the rule must be removed; then repeat from the step 1.
- 5 RIPPERk stops adding rules when there are no more positive examples left or when a rule has an unacceptably large error rate, or when the last rule added is too *complicated* according to some criterion.



# RIPPERk

To build a rule:

- 1 Randomly partition all the examples which have not been covered by any rule yet into two subsets: a *growing set* and a *pruning set*.
- 2 Grow a rule by greedily adding conditions until the rule reaches a confidence of 100% using only the growing set.
- 3 To prevent growing set *overfitting*, immediately prune the rule, deleting some conditions based on a pruning criterion, using the pruning data.
- 4 Lastly, all the positive and negative examples covered by the rule must be removed; then repeat from the step 1.
- 5 RIPPERk stops adding rules when there are no more positive examples left or when a rule has an unacceptably large error rate, or when the last rule added is too *complicated* according to some criterion.



# RIPPERk

To build a rule:

- 1 Randomly partition all the examples which have not been covered by any rule yet into two subsets: a *growing set* and a *pruning set*.
- 2 Grow a rule by greedily adding conditions until the rule reaches a confidence of 100% using only the growing set.
- 3 To prevent growing set *overfitting*, immediately prune the rule, deleting some conditions based on a pruning criterion, using the pruning data.
- 4 Lastly, all the positive and negative examples covered by the rule must be removed; then repeat from the step 1.
- 5 RIPPERk stops adding rules when there are no more positive examples left or when a rule has an unacceptably large error rate, or when the last rule added is too *complicated* according to some criterion.



# RIPPERk

To build a rule:

- 1 Randomly partition all the examples which have not been covered by any rule yet into two subsets: a *growing set* and a *pruning set*.
- 2 Grow a rule by greedily adding conditions until the rule reaches a confidence of 100% using only the growing set.
- 3 To prevent growing set *overfitting*, immediately prune the rule, deleting some conditions based on a pruning criterion, using the pruning data.
- 4 Lastly, all the positive and negative examples covered by the rule must be removed; then repeat from the step 1.
- 5 RIPPERk stops adding rules when there are no more positive examples left or when a rule has an unacceptably large error rate, or when the last rule added is too *complicated* according to some criterion.



# Decision Tree

Example data: hospitalization context where patients are *instances* described by their *age*, *weight* and *gender*. Data are normally collected in tabular form:

# Sample	Age	Weight	Gender	Class label
1	37	70	M	ILL
2	49	81	M	HEALTHY
3	20	55	F	ILL
...	...	...	...	...

By means of learning algorithms, a decision tree classifier can be *trained* from such data:

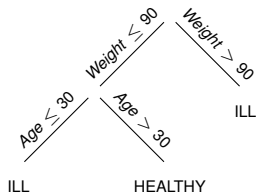


# Decision Tree

Example data: hospitalization context where patients are *instances* described by their *age*, *weight* and *gender*. Data are normally collected in tabular form:

# Sample	Age	Weight	Gender	Class label
1	37	70	M	ILL
2	49	81	M	HEALTHY
3	20	55	F	ILL
...	...	...	...	...

By means of learning algorithms, a decision tree classifier can be *trained* from such data:



# CART

The algorithm builds a binary tree splitting each node into two child nodes repeatedly using the following steps:

- 1 For each feature (i.e. each attribute) with  $K$  different values, there exist  $K-1$  possible splits; choose for each feature the split that maximizes the splitting criterion.
- 2 Among the best splits from step 1 choose the one which maximizes the splitting criterion.
- 3 Split the node using the best node split from step 2 and repeat from step 1 until a certain stopping criterion is satisfied.



# CART

The algorithm builds a binary tree splitting each node into two child nodes repeatedly using the following steps:

- 1 For each feature (i.e. each attribute) with  $K$  different values, there exist  $K-1$  possible splits; choose for each feature the split that maximizes the splitting criterion.
- 2 Among the best splits from step 1 choose the one which maximizes the splitting criterion.
- 3 Split the node using the best node split from step 2 and repeat from step 1 until a certain stopping criterion is satisfied.



# CART

The algorithm builds a binary tree splitting each node into two child nodes repeatedly using the following steps:

- 1 For each feature (i.e. each attribute) with  $K$  different values, there exist  $K-1$  possible splits; choose for each feature the split that maximizes the splitting criterion.
- 2 Among the best splits from step 1 choose the one which maximizes the splitting criterion.
- 3 Split the node using the best node split from step 2 and repeat from step 1 until a certain stopping criterion is satisfied.



# Install

Add the GitHub project as a repository in the `composer.json` file of your Laravel application as follows:

```
1 "repositories": [  
2     {  
3         "type": "vcs",  
4         "url": "https://github.com/aclai-lab/piton"  
5     }  
6 ]
```

# Install

And then add the package in the require section:

```
1 "require": {  
2     "aclai/piton": "master"  
3 }
```

Finally, run in the terminal the `composer update` command.

# Workflow

- The package reads data directly from a MySQL database following the parameters specified into a *problem* configuration file, which information will be stored into the *Problems* table, and creating a dataset.
- Then, it derives the hierarchy of problems, and creates a new instance in the *Model version* table that will store information about the extracted hierarchy.
- Lastly, it will launch the *training* process for each problem, associating each extracted *rule-based model* to its correspondent problem in the hierarchy; these models are stored in the *Class model* table, while their rules are stored in the *Rules* table.
- Note that the database is dynamic, and therefore the process can be iterated over time (e.g. once a week) with new data.



# Workflow

- The package reads data directly from a MySQL database following the parameters specified into a *problem* configuration file, which information will be stored into the *Problems* table, and creating a dataset.
- Then, it derives the hierarchy of problems, and creates a new instance in the *Model version* table that will store information about the extracted hierarchy.
- Lastly, it will launch the *training* process for each problem, associating each extracted *rule-based model* to its correspondent problem in the hierarchy; these models are stored in the *Class model* table, while their rules are stored in the *Rules* table.
- Note that the database is dynamic, and therefore the process can be iterated over time (e.g. once a week) with new data.



# Workflow

- The package reads data directly from a MySQL database following the parameters specified into a *problem* configuration file, which information will be stored into the *Problems* table, and creating a dataset.
- Then, it derives the hierarchy of problems, and creates a new instance in the *Model version* table that will store information about the extracted hierarchy.
- Lastly, it will launch the *training* process for each problem, associating each extracted *rule-based model* to its correspondent problem in the hierarchy; these models are stored in the *Class model* table, while their rules are stored in the *Rules* table.
- Note that the database is dynamic, and therefore the process can be iterated over time (e.g. once a week) with new data.



# Workflow

- The package reads data directly from a MySQL database following the parameters specified into a *problem* configuration file, which information will be stored into the *Problems* table, and creating a dataset.
- Then, it derives the hierarchy of problems, and creates a new instance in the *Model version* table that will store information about the extracted hierarchy.
- Lastly, it will launch the *training* process for each problem, associating each extracted *rule-based model* to its correspondent problem in the hierarchy; these models are stored in the *Class model* table, while their rules are stored in the *Rules* table.
- Note that the database is dynamic, and therefore the process can be iterated over time (e.g. once a week) with new data.



# Pitón database

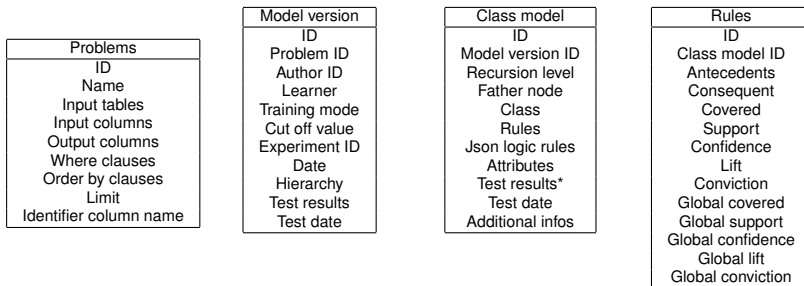


Figure: Database structure

# Configuring the database

First, one should add the following new connection in the `connections` array inside the `config/database.php` file:

```
1 'piton_connection' => [  
2     'driver' => env('DB_CONNECTION_PITON'),  
3     'host' => env('DB_HOST_PITON', '127.0.0.1'),  
4     'port' => env('DB_PORT_PITON', '3306'),  
5     'database' => env('DB_DATABASE_PITON', 'forge'),  
6     'username' => env('DB_USERNAME_PITON', 'forge'),  
7     'password' => env('DB_PASSWORD_PITON', ''),  
8     'unix_socket' => '',  
9     'charset' => 'utf8mb4',  
10    'collation' => 'utf8mb4_unicode_ci',  
11    'prefix' => '',  
12    'prefix_indexes' => true,  
13    'strict' => true,  
14    'engine' => null,  
15 ]
```

# Configuring the database

Then, one should add the following to the `.env` file of the project, to specify how to access this database:

```
1 DB_CONNECTION_PITON=mysql
2 DB_HOST_PITON=127.0.0.1
3 DB_PORT_PITON=3306
4 DB_DATABASE_PITON=<your_piton_database>
5 DB_USERNAME_PITON=<your_mysql_username>
6 DB_PASSWORD_PITON=<your_mysql_password>
```

After the database has been configured, it is finally possible to launch the `php artisan migrate` command from the command line to populate it.



# Training

- 1 Publish the *problem-config* file via the following command:  

```
php artisan vendor:publish --tag=problem-config
```
- 2 Fill it, and rename it after the problem to be solved.
- 3 Publish the configuration file for the specified learner (*between `prip`, `sklearn_cart`, `wittgenstein_irep`, `wittgenstein_ripper`*):  

```
php artisan vendor:publish --tag=<learner_name>-config
```
- 4 Run the `piton:update_models` command, which accepts as parameters a *problem* name, an *author ID*, the name of the learner to be used and, eventually, the specific algorithm to be used for the training.

# Training

- 1 Publish the *problem-config* file via the following command:  

```
php artisan vendor:publish --tag=problem-config
```
- 2 Fill it, and rename it after the problem to be solved.
- 3 Publish the configuration file for the specified learner (*between `prip`, `sklearn_cart`, `wittgenstein_irep`, `wittgenstein_ripper`*):  

```
php artisan vendor:publish --tag=<learner_name>-config
```
- 4 Run the `piton:update_models` command, which accepts as parameters a *problem* name, an *author ID*, the name of the learner to be used and, eventually, the specific algorithm to be used for the training.

Introduction	Install
Preliminaries	Workflow
How to use the package	Configuring the database
Automated treatment suggestion	Training
Backup slides	Prediction

# Training

- 1 Publish the *problem-config* file via the following command:  

```
php artisan vendor:publish --tag=problem-config
```
- 2 Fill it, and rename it after the problem to be solved.
- 3 Publish the configuration file for the specified learner (*between `prip`, `sklearn_cart`, `wittgenstein_irep`, `wittgenstein_ripperk`*):  

```
php artisan vendor:publish --tag=<learner_name>-config
```
- 4 Run the `piton:update_models` command, which accepts as parameters a *problem* name, an *author ID*, the name of the learner to be used and, eventually, the specific algorithm to be used for the training.



Università  
degli Studi  
di Ferrara

# Training

- 1 Publish the *problem-config* file via the following command:  

```
php artisan vendor:publish --tag=problem-config
```
- 2 Fill it, and rename it after the problem to be solved.
- 3 Publish the configuration file for the specified learner (*between `prip`, `sklearn_cart`, `wittgenstein_irep`, `wittgenstein_ripperk`*):  

```
php artisan vendor:publish --tag=<learner_name>-config
```
- 4 Run the `piton:update_models` command, which accepts as parameters a *problem* name, an *author ID*, the name of the learner to be used and, eventually, the specific algorithm to be used for the training.

# Prediction

Prediction can be done in two different ways:

- Using the `php artisan piton:predict_by_identifier` command.
- Using the `predictByIdentifier()` function of the `DBFit` class. This function accepts as parameters an instance id, and an instance of model version id, to specify which rule-based models to use (i.e. associated with which execution), and return the prediction in json format, so that it can easily be displayed on a web page (even using AJAX).

Example:

```
$lastMV = ModelVersion::orderByDesc('id')->first();  
$dbfit = new DBFit();  
$prediction = $dbfit->predictByIdentifier(1020, [], $lastMV['id']);
```

# Prediction

Prediction can be done in two different ways:

- Using the `php artisan piton:predict_by_identifier` command.
- Using the `predictByIdentifier()` function of the `DBFit` class. This function accepts as parameters an instance id, and an instance of model version id, to specify which rule-based models to use (i.e. associated with which execution), and return the prediction in json format, so that it can easily be displayed on a web page (even using AJAX).

Example:

```
$lastMV = ModelVersion::orderByDesc('id')->first();  
$dbfit = new DBFit();  
$prediction = $dbfit->predictByIdentifier(1020, [], $lastMV['id']);
```

# Problem

Given information about the patient, such as its *anamnesis*, its *menopausal state*, its *densitometry*, we would like to predict which *therapies* and respective *active principles* the physician could suggest the patient.

This can easily be seen as a hierarchy of problems consisting of two levels: the first concerning the *type of the suggested therapy*, the second concerning which *active principle* to suggest.

# Problem

Therapy name	Abbreviation	Active principle	Abbreviation
Hormonal Therapy	$T_{horm}$	MHT (tibolone) MHT (oral) MHT transdermal MHT (TSEC)	<i>tib</i> <i>oral</i> <i>trans</i> <i>tsec</i>
Osteoprotective Therapy	$T_{osteop}$	Alendronate Alendronate + vit D Risendronate Ibandronate Clodronate Raloxifene Bazedoxifene Denomasub Teriparatide Zoledronate	<i>ale</i> <i>ale + vD</i> <i>ris</i> <i>iba</i> <i>clo</i> <i>ral</i> <i>baz</i> <i>den</i> <i>ter</i> <i>zol</i>
Vitamin D Therapy	$T_{vitDth}$	Colecalciferol Calcifediol	<i>colec</i> <i>calci</i>
Vitamin D Supplementation	$S_{vitDsup}$	Colecalciferol Calcifediol	<i>colec</i> <i>calci</i>
Calcium Supplementation	$S_{calsup}$	Carbonated calcium Citratd calcium	<i>carb</i> <i>citr</i>

**Table:** Possible therapies and respective active principles.

# Problem

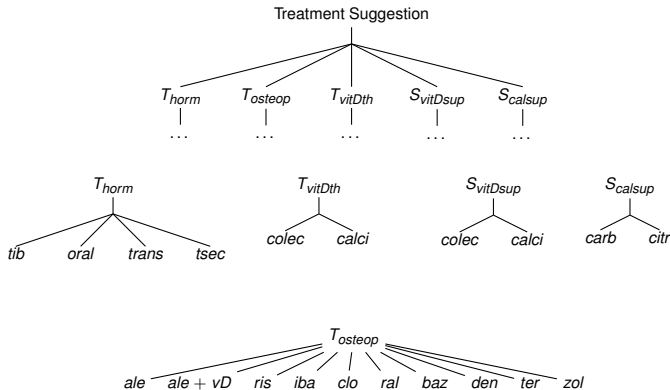


Figure: Possible hierarchy to solve the treatment suggestion problem

# Database

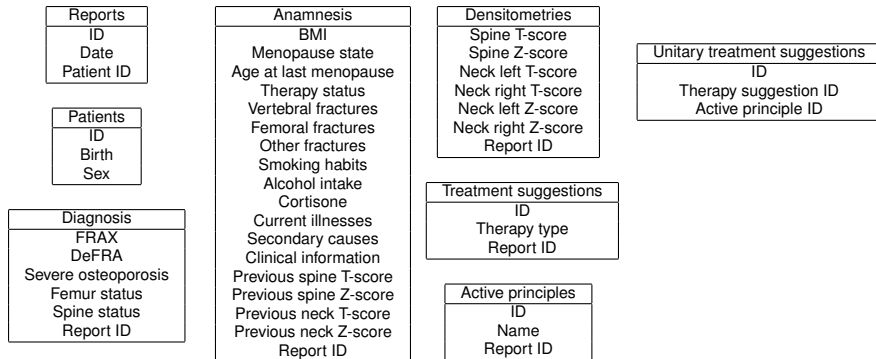


Figure: Sub-set of tables of the CMO database.

# Experiment results

To estimate the effectiveness of the application of the package to the CMO system, a complete *experiment* has been executed, considering the recommendations that have been given from Sept. the 1st, 2018 to Aug. the 31th, 2020.

In this experiment, only three data set have been considered for classifier extraction at the first level and seven at the second level; this is because the others were *unbalanced* (e.g. one class appeared for less than 10% of the instances).

For each data set, 80% of the instances are used as a *training set*, while the 20% stratified most recent records of each data set are used for *testing* purposes.

# Experiment results

For each problem and its corresponding test, we reported the following values:

- *accuracy*, that is, the rate of corrected classification;
- *sensitivity*, that is, the rate of true positives;
- *specificity*, that is, the rate of true negatives;
- *positive predicted value*, that is, the inverse of the false discovery rate;
- *negative predicted value*, that is, the inverse of the false omission rate;
- *F1 score*, that is, the harmonic mean of sensitivity and positive predicted value.

# Experiment results

Classifier	Accuracy	$F_1$	Sensitivity	Specificity	PPV	NPV
$\Gamma_{osteop}$	0.86	0.63	0.55	0.95	0.74	0.88
$\Gamma_{vitDsup}$	0.75	0.76	0.74	0.76	0.79	0.71
$\Gamma_{calsup}$	0.83	0.61	0.47	0.97	0.85	0.82

**Table:** Results of the experiment: global evaluation of first level classifiers.

# Experiment results

Classifier	Accuracy	$F_1$	Sensitivity	Specificity	PPV	NPV
$\Gamma_{osteop}^{ale}$	0.72	0.39	0.27	0.95	0.73	0.72
$\Gamma_{osteop}^{den}$	0.71	0.35	0.25	0.92	0.58	0.73
$\Gamma_{osteop}^{ris}$	0.83	0.12	0.10	0.92	0.14	0.89
$\Gamma_{vitDsup}^{calci}$	0.93	0.65	0.54	0.98	0.83	0.94
$\Gamma_{vitDsup}^{colec}$	0.92	0.96	0.98	0.54	0.94	0.79
$\Gamma_{calsup}^{citr}$	0.62	0.52	0.46	0.75	0.60	0.63
$\Gamma_{calsup}^{carb}$	0.62	0.68	0.75	0.46	0.63	0.60

**Table:** Results of the experiment: global evaluation of second level classifiers.

# Experiment results

For local evaluation, we consider the number of rules of each classifier and their distribution among the following four types, based on *support* and *confidence*:

- I Relevant and reliable, showing  $\text{support} > 0.2$  and  $\text{confidence} > 0.7$ ;
- II Relevant, but unreliable, showing  $\text{support} > 0.2$  but  $\text{confidence} \leq 0.7$ ;
- III Irrelevant, but reliable, showing  $\text{support} \leq 0.2$  but  $\text{confidence} > 0.7$ ;
- IV Irrelevant and unreliable, showing  $\text{support} \leq 0.2$  and  $\text{confidence} \leq 0.7$ .

# Experiment results

Types I and III (which include only reliable rules) are the most common ones, indicating that our approach is relatively stable.

Classifier	#	#I	#II	#III	#IV
$\Gamma_{osteop}$	6	1	0	5	0
$\Gamma_{vitDsup}$	6	2	2	2	0
$\Gamma_{calsup}$	4	1	0	2	1

**Table:** Results of the experiment: local evaluation of first level classifiers.

# Thank you for your attention

# Thank you for your attention



# CART

## Pruning

Then, in order to enhance the generalization of the resulting decision tree, pruning is applied:

- 1 Split randomly training data into  $N$  folds.
- 2 Select a pruning level for the tree (level 0 equals to the full tree).
- 3 Use  $N - 1$  folds to create  $N - 1$  new pruned trees and estimate the error on the  $N$ th fold.
- 4 Repeat from step 2 until all pruning levels are used.
- 5 Find the smallest error and use the pruning level assigned to it.
- 6 Until the pruning level is reached, remove all the leafs in the lowest tree level and assign the decision class (e.g. the class with the higher number of cases covered by the node) to their parent node.





# CART

## Pruning

Then, in order to enhance the generalization of the resulting decision tree, pruning is applied:

- 1 Split randomly training data into  $N$  folds.
- 2 Select a pruning level for the tree (level 0 equals to the full tree).
- 3 Use  $N - 1$  folds to create  $N - 1$  new pruned trees and estimate the error on the  $N$ th fold.
- 4 Repeat from step 2 until all pruning levels are used.
- 5 Find the smallest error and use the pruning level assigned to it.
- 6 Until the pruning level is reached, remove all the leafs in the lowest tree level and assign the decision class (e.g. the class with the higher number of cases covered by the node) to their parent node.





# CART

## Pruning

Then, in order to enhance the generalization of the resulting decision tree, pruning is applied:

- 1 Split randomly training data into  $N$  folds.
- 2 Select a pruning level for the tree (level 0 equals to the full tree).
- 3 Use  $N - 1$  folds to create  $N - 1$  new pruned trees and estimate the error on the  $N$ th fold.
- 4 Repeat from step 2 until all pruning levels are used.
- 5 Find the smallest error and use the pruning level assigned to it.
- 6 Until the pruning level is reached, remove all the leafs in the lowest tree level and assign the decision class (e.g. the class with the higher number of cases covered by the node) to their parent node.





# CART

## Pruning

Then, in order to enhance the generalization of the resulting decision tree, pruning is applied:

- 1 Split randomly training data into  $N$  folds.
- 2 Select a pruning level for the tree (level 0 equals to the full tree).
- 3 Use  $N - 1$  folds to create  $N - 1$  new pruned trees and estimate the error on the  $N$ th fold.
- 4 Repeat from step 2 until all pruning levels are used.
- 5 Find the smallest error and use the pruning level assigned to it.
- 6 Until the pruning level is reached, remove all the leafs in the lowest tree level and assign the decision class (e.g. the class with the higher number of cases covered by the node) to their parent node.





# CART

## Pruning

Then, in order to enhance the generalization of the resulting decision tree, pruning is applied:

- 1 Split randomly training data into  $N$  folds.
- 2 Select a pruning level for the tree (level 0 equals to the full tree).
- 3 Use  $N - 1$  folds to create  $N - 1$  new pruned trees and estimate the error on the  $N$ th fold.
- 4 Repeat from step 2 until all pruning levels are used.
- 5 Find the smallest error and use the pruning level assigned to it.
- 6 Until the pruning level is reached, remove all the leafs in the lowest tree level and assign the decision class (e.g. the class with the higher number of cases covered by the node) to their parent node.



# CART

## Pruning

Then, in order to enhance the generalization of the resulting decision tree, pruning is applied:

- 1 Split randomly training data into  $N$  folds.
- 2 Select a pruning level for the tree (level 0 equals to the full tree).
- 3 Use  $N - 1$  folds to create  $N - 1$  new pruned trees and estimate the error on the  $N$ th fold.
- 4 Repeat from step 2 until all pruning levels are used.
- 5 Find the smallest error and use the pruning level assigned to it.
- 6 Until the pruning level is reached, remove all the leafs in the lowest tree level and assign the decision class (e.g. the class with the higher number of cases covered by the node) to their parent node.



# Experiment results

Classifier	#	#I	#II	#III	#IV
$\Gamma_{osteop}^{ale}$	2	0	2	0	0
$\Gamma_{osteop}^{den}$	4	0	1	1	2
$\Gamma_{osteop}^{ris}$	2	1	0	0	1
$\Gamma_{vitDsup}^{calci}$	5	1	0	2	2
$\Gamma_{vitDsup}^{colec}$	5	1	0	2	2
$\Gamma_{calsup}^{citr}$	3	0	2	1	0
$\Gamma_{calsup}^{carb}$	3	0	2	1	0

**Table:** Results of the experiment: local evaluation of second level classifiers.