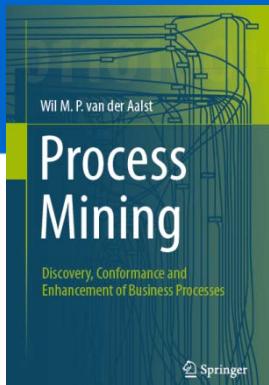


*Process Mining: Data Science in Action*

# Applying Decision Trees

prof.dr.ir. Wil van der Aalst  
[www.processmining.org](http://www.processmining.org)

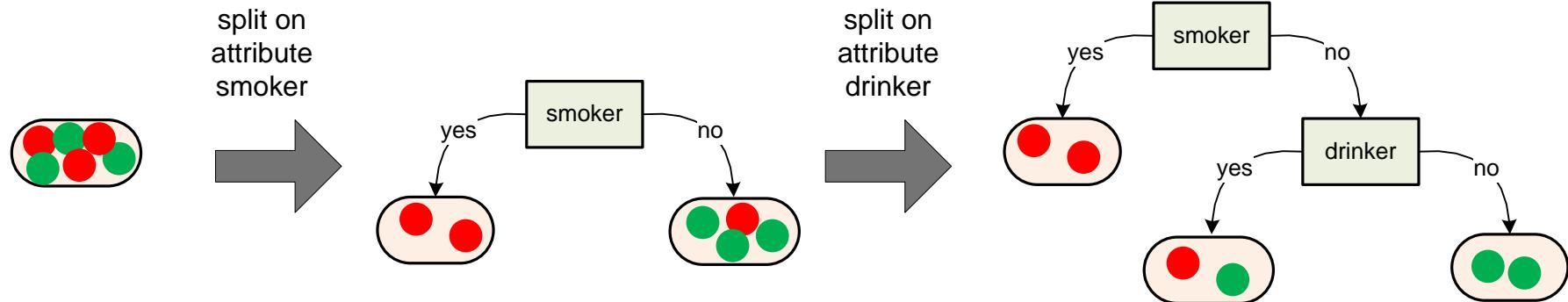


TU/e

Technische Universiteit  
**Eindhoven**  
University of Technology

Where innovation starts

# Decision tree learning



$$E = - \sum_{i=1}^k p_i \log_2(p_i)$$

**Iteratively reduce the overall level of uncertainty (entropy) using label splitting until no significant information gain is possible.**

# Example: 160 students (100 pass, 60 fail)



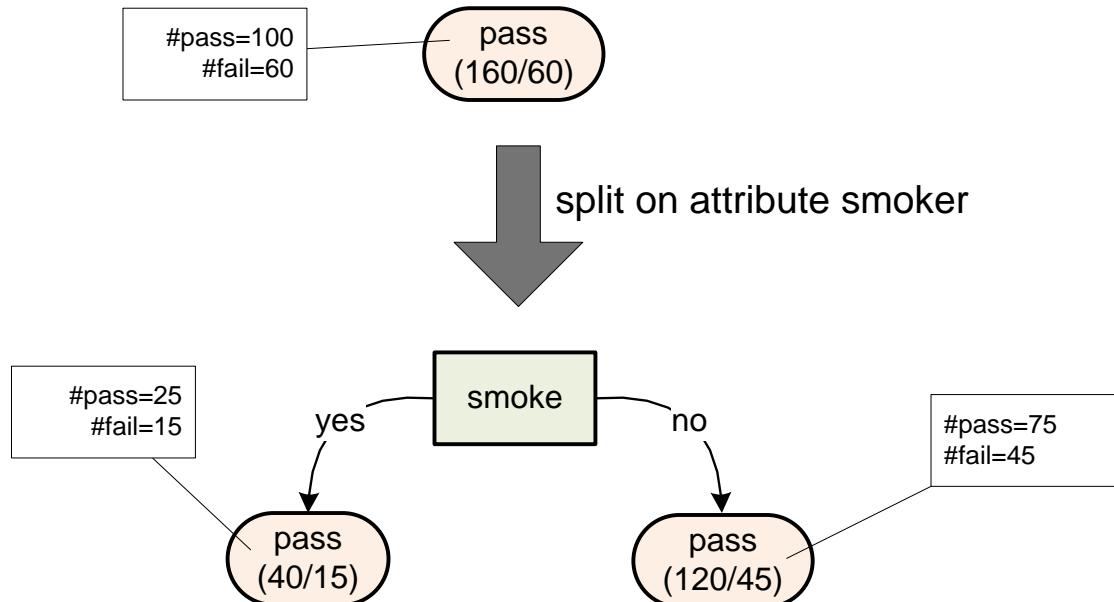
What matters?

attending  
lectures?

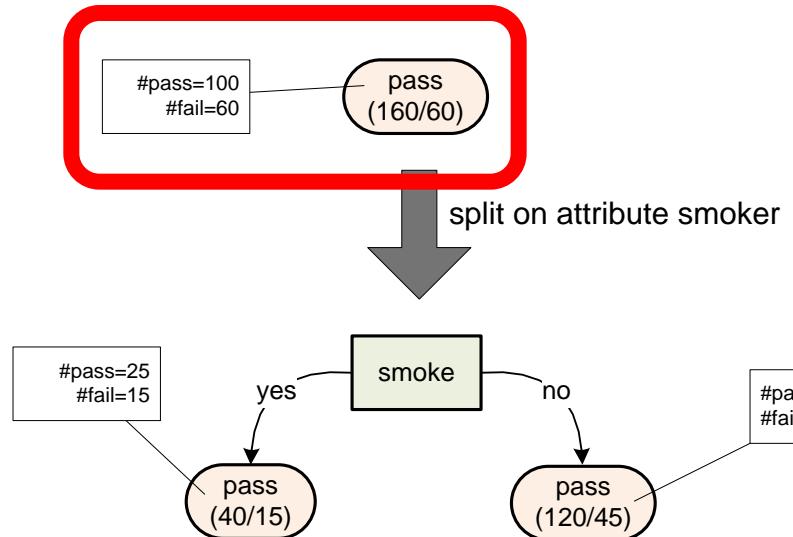
gender?

smoking?

# Question: What is the information gain?

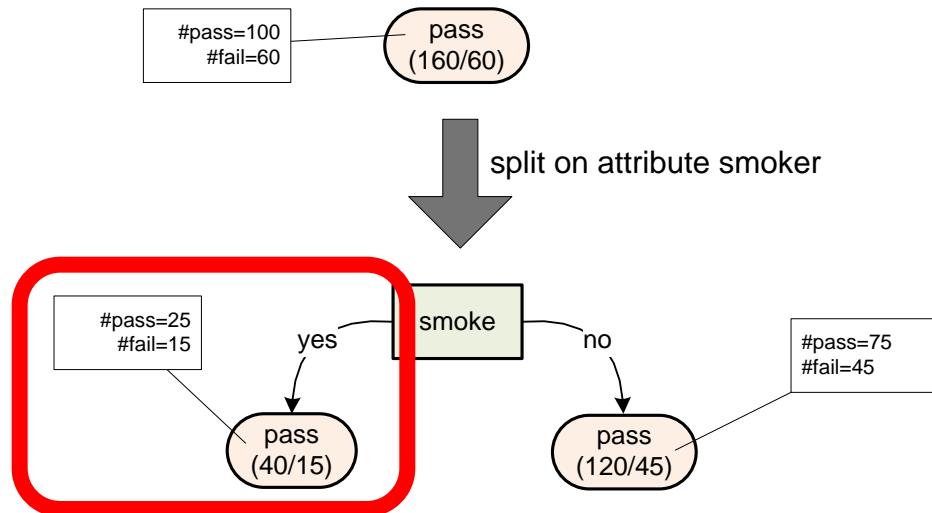


# Answer: Entropy of root node



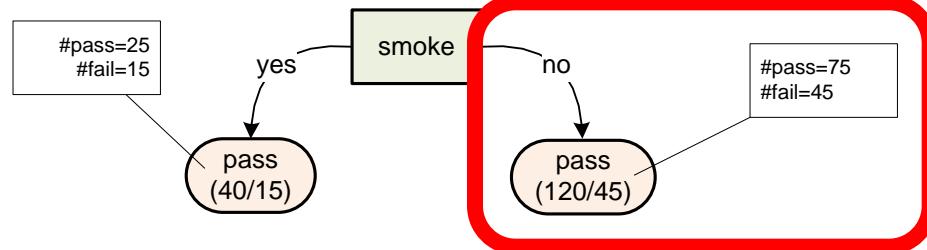
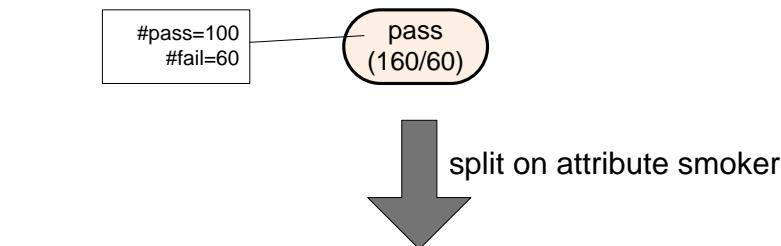
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\ &= 0.9544 \end{aligned}$$

# Answer: Entropy of smokers



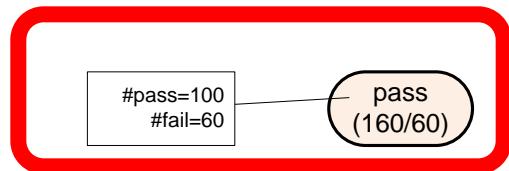
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{25}{40} \log_2\left(\frac{25}{40}\right) + \frac{15}{40} \log_2\left(\frac{15}{40}\right)\right) \\ &= 0.9544 \end{aligned}$$

# Answer: Entropy of non-smokers



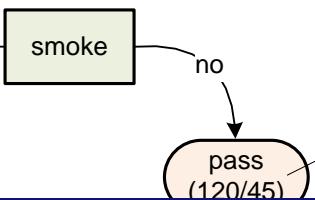
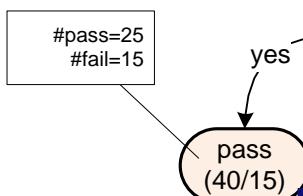
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{75}{120} \log_2\left(\frac{75}{120}\right)\right) + \frac{45}{120} \log_2\left(\frac{45}{120}\right) \\ &= 0.9544 \end{aligned}$$

# Answer: No information gain



$$E = \frac{160}{160} \times 0.9544 = 0.9544$$

split on attribute smoker



could be seen without computation

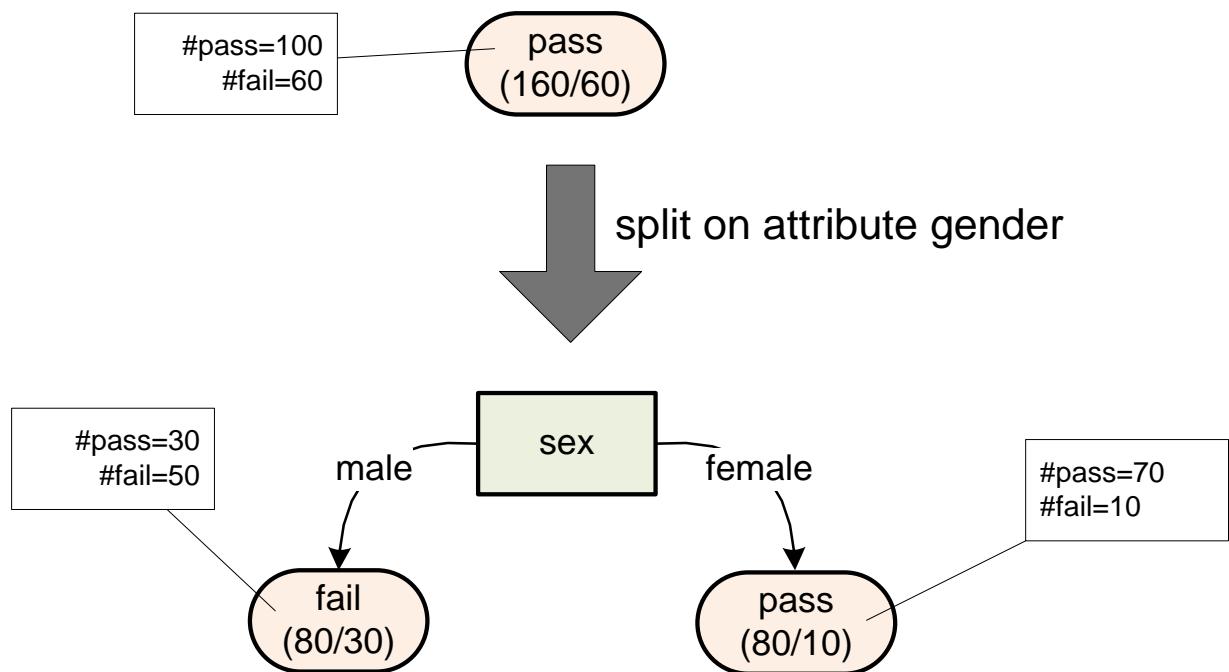
information gain = 0

#fail=45

$$E = \frac{40}{160} \times 0.9544 + \frac{120}{160} \times 0.9544 = 0.9544$$

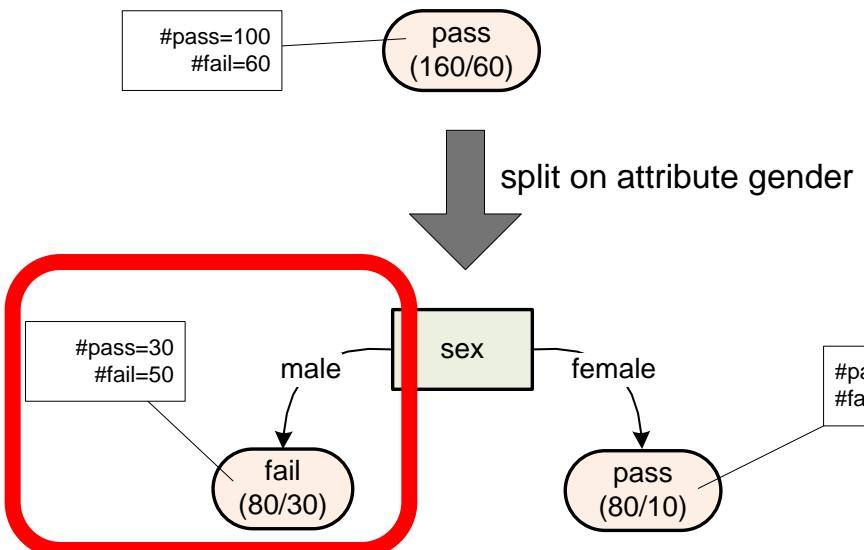


# Question: What is the information gain?



# Answer: Entropy of male students

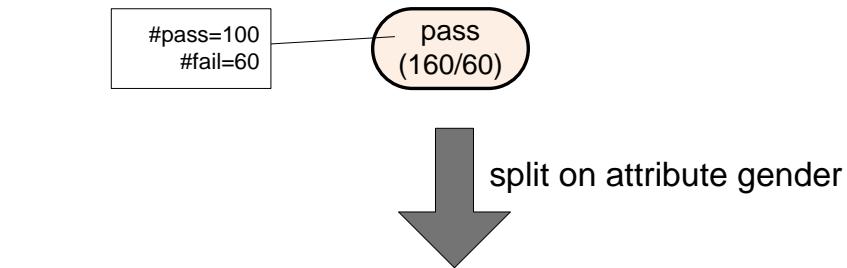
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\ &= 0.9544 \end{aligned}$$



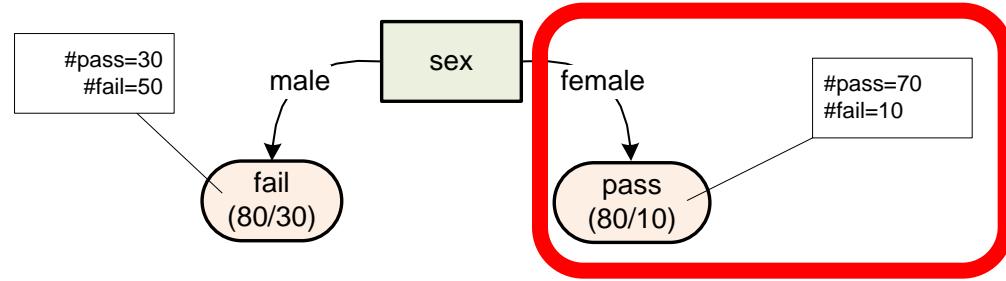
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{30}{80} \log_2\left(\frac{30}{80}\right) + \frac{50}{80} \log_2\left(\frac{50}{80}\right)\right) \\ &= 0.9544 \end{aligned}$$

# Answer: Entropy of female students

$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\ &= 0.9544 \end{aligned}$$



$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{70}{80} \log_2\left(\frac{70}{80}\right) + \frac{10}{80} \log_2\left(\frac{10}{80}\right)\right) \\ &= 0.5436 \end{aligned}$$



# Answer: Information gain

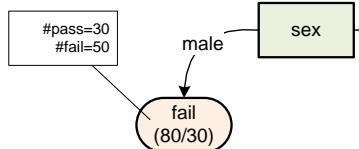
$$\begin{aligned}
 E &= - \sum_{i=1}^k p_i \log_2(p_i) \\
 &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\
 &= 0.9544
 \end{aligned}$$



$$E = \frac{160}{160} \times 0.9544 = 0.9544$$

split on attribute gender

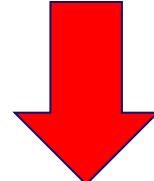
**information gain = 0.2054**



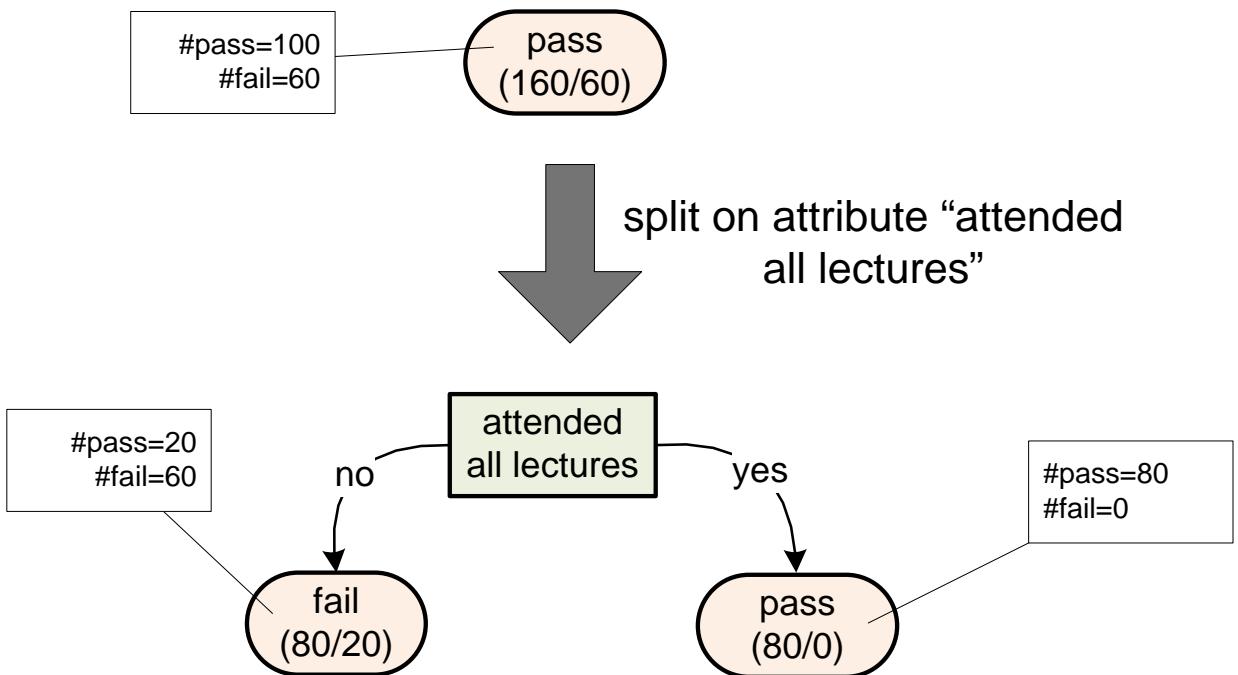
$$E = \frac{80}{160} \times 0.9544 + \frac{80}{160} \times 0.5436 = 0.7490$$

$$\begin{aligned}
 E &= - \sum_{i=1}^k p_i \log_2(p_i) \\
 &= -\left(\frac{30}{80} \log_2\left(\frac{30}{80}\right) + \frac{50}{80} \log_2\left(\frac{50}{80}\right)\right) \\
 &= 0.9544
 \end{aligned}$$

$$\begin{aligned}
 E &= - \sum_{i=1}^k p_i \log_2(p_i) \\
 &= -\left(\frac{70}{80} \log_2\left(\frac{70}{80}\right) + \frac{10}{80} \log_2\left(\frac{10}{80}\right)\right) \\
 &= 0.5436
 \end{aligned}$$

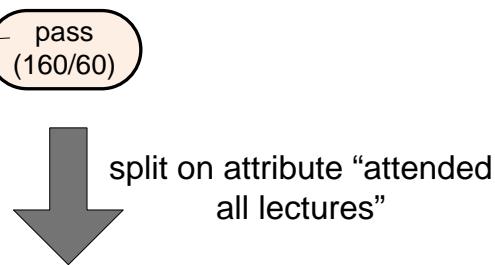


# Question: What is the information gain?

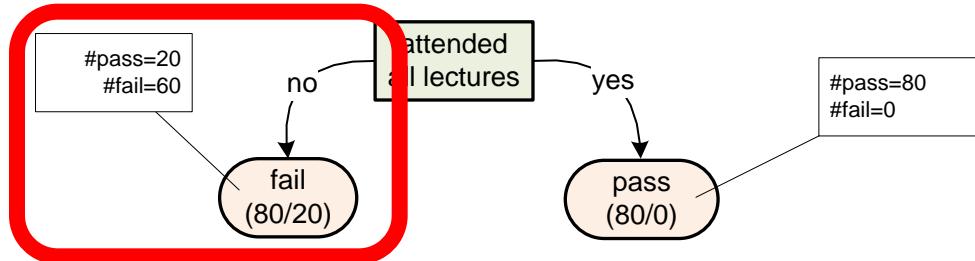


# Answer: Entropy of missing students

$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\ &= 0.9544 \end{aligned}$$

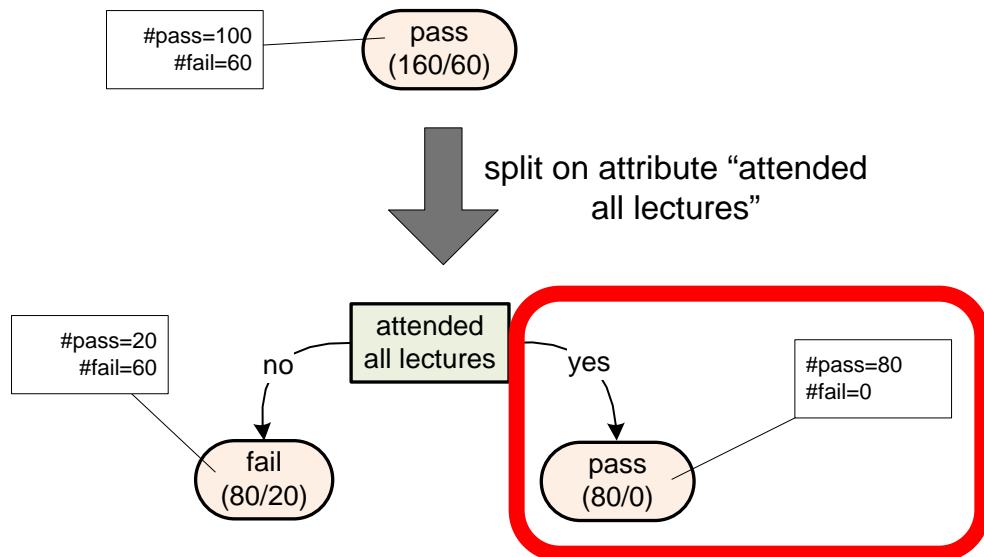


$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{20}{80} \log_2\left(\frac{20}{80}\right) + \frac{60}{80} \log_2\left(\frac{60}{80}\right)\right) \\ &= 0.8113 \end{aligned}$$



# Answer: Entropy of attending students

$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\ &= 0.9544 \end{aligned}$$



$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{80}{80} \log_2\left(\frac{80}{80}\right)\right) \\ &= 0 \end{aligned}$$

# Answer

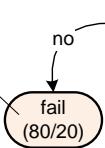
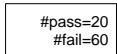
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{100}{160} \log_2\left(\frac{100}{160}\right) + \frac{60}{160} \log_2\left(\frac{60}{160}\right)\right) \\ &= 0.9544 \end{aligned}$$



$$E = \frac{160}{160} \times 0.9544 = 0.9544$$

information gain = 0.5488

split on attribute "attended all lectures"



attended  
all lectures

$$E = \frac{80}{160} \times 0.8113 + \frac{80}{160} \times 0 = 0.4056$$

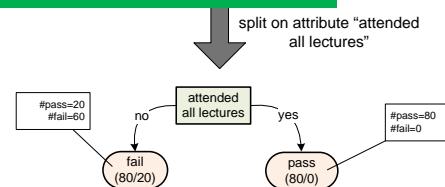
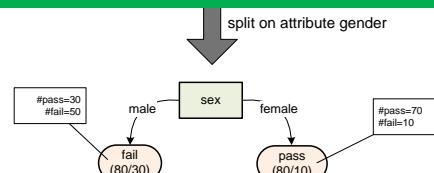
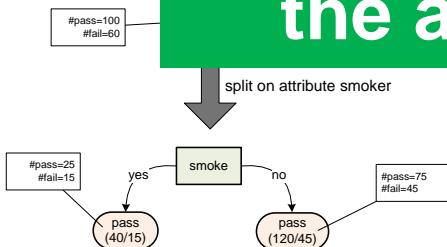
$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{20}{80} \log_2\left(\frac{20}{80}\right) + \frac{60}{80} \log_2\left(\frac{60}{80}\right)\right) \\ &= 0.8113 \end{aligned}$$

$$\begin{aligned} E &= - \sum_{i=1}^k p_i \log_2(p_i) \\ &= -\left(\frac{80}{80} \log_2\left(\frac{80}{80}\right)\right) \\ &= 0 \end{aligned}$$

# Comparing information gains



So we should split the root node on the attribute "attend all lectures"!

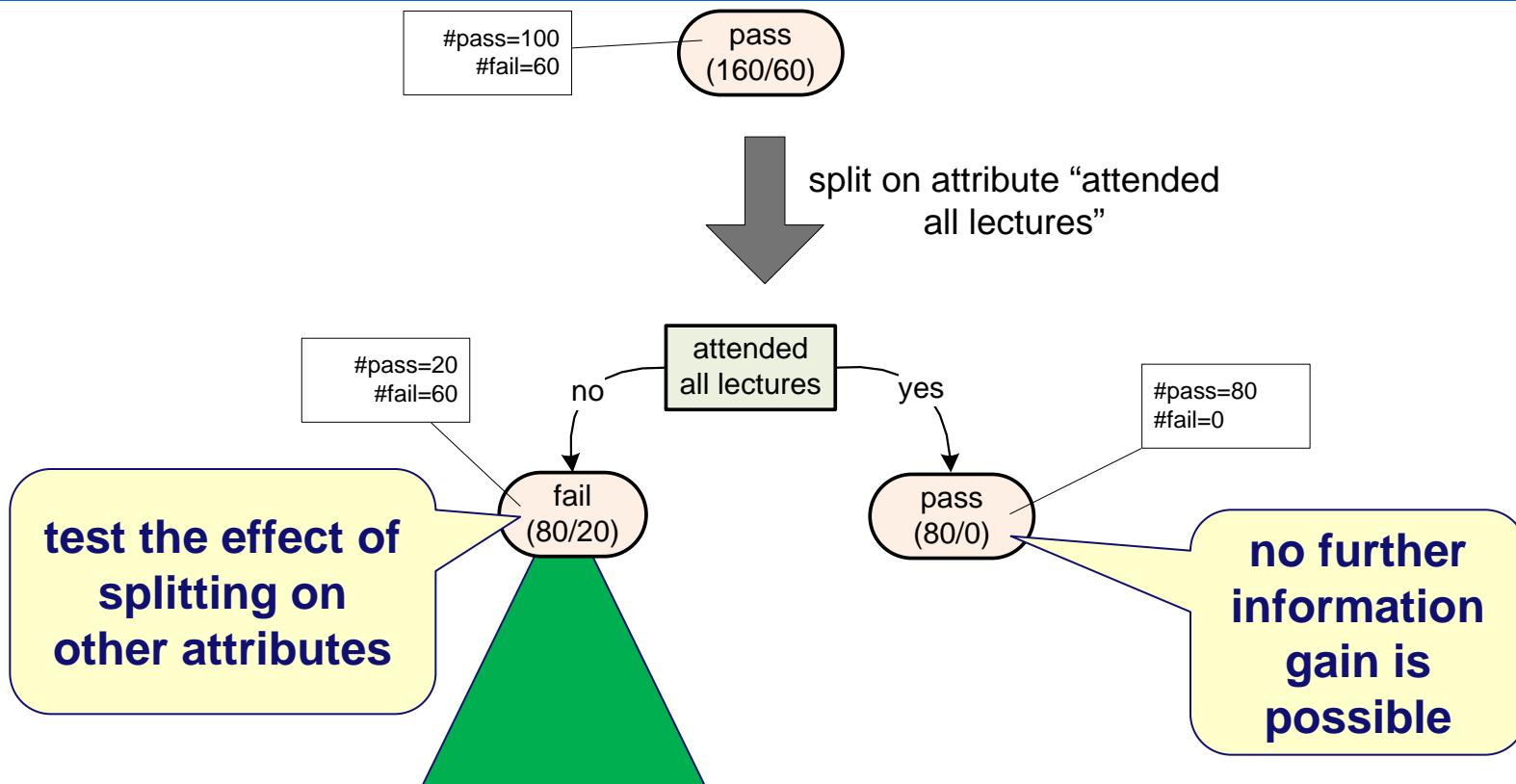


information gain = 0

information gain = 0.2054

information gain = 0.5488

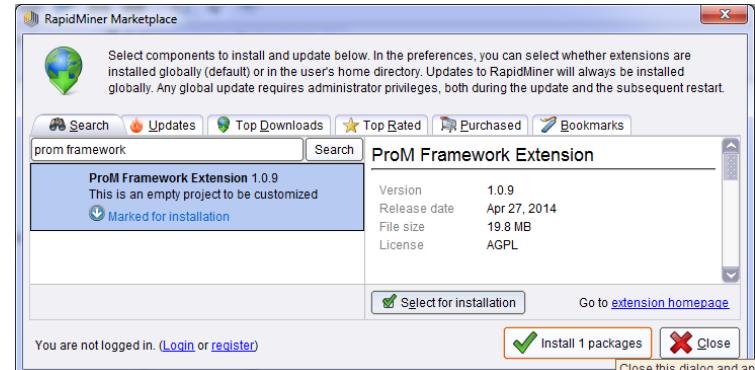
# Iterate until no significant gain is possible



# RapidMiner

(installation is optional)

- An integrated **extendible** environment for **machine learning, data mining, text mining, and predictive analytics.**
- RapidMiner Marketplace also provides a **ProM extension** for process mining.
- Commercial and open-source versions of the software.



# Decision trees in RapidMiner



gender	age	smoker	car brand	claim
female	47	yes	Volvo	no
male	31	no	Alfa Romeo	yes
ma	<b>CSV file contains information about 999</b>			
ma	<b>customers of an insurance company.</b>			
male	44	no	BMW	no
fema	<b>The company wants to know which</b>			
ma	<b>customers claim insurance.</b>			
...	...	...	...	...

# Decision trees in RapidMiner



gender	age	smoker	car brand	claim
female	47	yes	Volvo	no
male	31	no	Alfa Romeo	yes
male	59	no	Alfa Romeo	yes
male	28	no	Fiat	no
male	44	no	BMW	yes

**Response variable (dependent variable): claim.**

**Predictor variables (independent variables):  
gender, age, smoker, car brand.**

# Data in RapidMiner

The screenshot shows the RapidMiner Data import wizard in Step 5 of 5. The interface is titled "Data import wizard - Step 5 of 5". The main content area displays a file tree under "Local Repository (wil)". The tree structure is as follows:

- Local Repository (wil)
  - data (wil)
    - MOOC (wil)
      - Golf (wil - v1, 1/16/14 11:19 AM - 507 bytes)
      - decision-tree (wil - v1, 1/16/14 11:38 AM - 2 kB)
      - food-poisoning (wil - v1, 1/24/14 1:38 AM - 239 kB)
      - food-poisoning-simple (wil - v1, 1/24/14 4:17 AM - 122 kB)
      - insurance-claims (wil - v1, 1/18/14 12:06 PM - 8 kB)
      - pampers (wil - v1, 1/23/14 1:49 PM - 1 kB)
    - processes (wil)

Below the tree, there is a status message: "0 errors." and a "Name" input field containing "insurance-data-decision-tree". At the bottom, there is a "Location" field with the value "/Local Repository/data/MOOC/insurance-data-decision-tree". Navigation buttons include "Previous", "Next", "Finish", and "Cancel".

//Local Repository/processes/insurance-claims\* – RapidMiner 5.3.015 @ nbwin1027

File Edit Process Tools View Help

Result Overview ExampleSet (/Local Repository/data/MOOC/insurance-data-decision-tree)

Data View Meta Data View Plot View Advanced Charts Annotations

ExampleSet (999 examples, 1 special attribute, 4 regular attributes) View Filter (999 / 999): all

Row No.	claim	gender	age	smoker	car brand
1	no	female	47	yes	Volvo
2	yes	male	31	no	Alfa Romeo
3	yes	male	59	no	Alfa Romeo
4	no	male	28	no	Fiat
5	no	male	44	no	BMW
6	no	female	27	no	Fiat
7	no	male	29	no	Subaru
8	yes	male	44	yes	Subaru
9	no	male	39	no	BMW
10	yes	male	35	no	Subaru
11	no	male	43	no	Subaru
12	yes	male	25	no	BMW
13	no	male	39	no	Volkswagen
14	yes	male	37	no	Alfa Romeo
15	no	female	30	no	Fiat
16	no	female	24	no	Fiat
17	yes	male	26	no	Alfa Romeo
18	no	male	43	no	BMW
19	no	male	46	no	BMW
20	no	female	25	no	Fiat
21	no	female	27	no	Nissan
22	no	female	31	no	Nissan
23	no	male	29	yes	Volkswagen
24	yes	male	42	no	BMW
25	no	male	26	no	Fiat
26	yes	male	27	no	Alfa Romeo

Data is stored in repository.  
Now we can apply an analysis workflow to it.

Repository Browser

Select a repository location.

Samples (none)

DB

Local Repository (wil)

data (wil)

MOOC (wil)

insurance-data-decision-tree (wil - v1, 7/16/14 10:13 PM - 8 kB)

Golf (wil - v1, 1/16/14 11:19 AM - 507 bytes)

decision-tree (wil - v1, 1/16/14 11:38 AM - 2 kB)

food-poisoning (wil - v1, 1/24/14 1:38 AM - 239 kB)

food-poisoning-simple (wil - v1, 1/24/14 4:17 AM - 1 kB)

insurance-claims (wil - v1, 1/18/14 12:06 PM - 8 kB)

pampers (wil - v1, 1/23/14 1:49 PM - 1 kB)

processes (wil)

**insurance-data-decision-tree**

Data Table  
Number of examples = 999  
5 attributes:

Role	Name	Type	Range	Missing	Comments
	gender	binominal	=[female, ...]	= 0	
	age	integer	= [20 - 73]	= 0	
	smoker	binominal	= [no, yes]	= 0	
	car brand	polynomial	= [Alfa Ro...]	= 0	
label	claim	binominal	= [no, yes]	= 0	

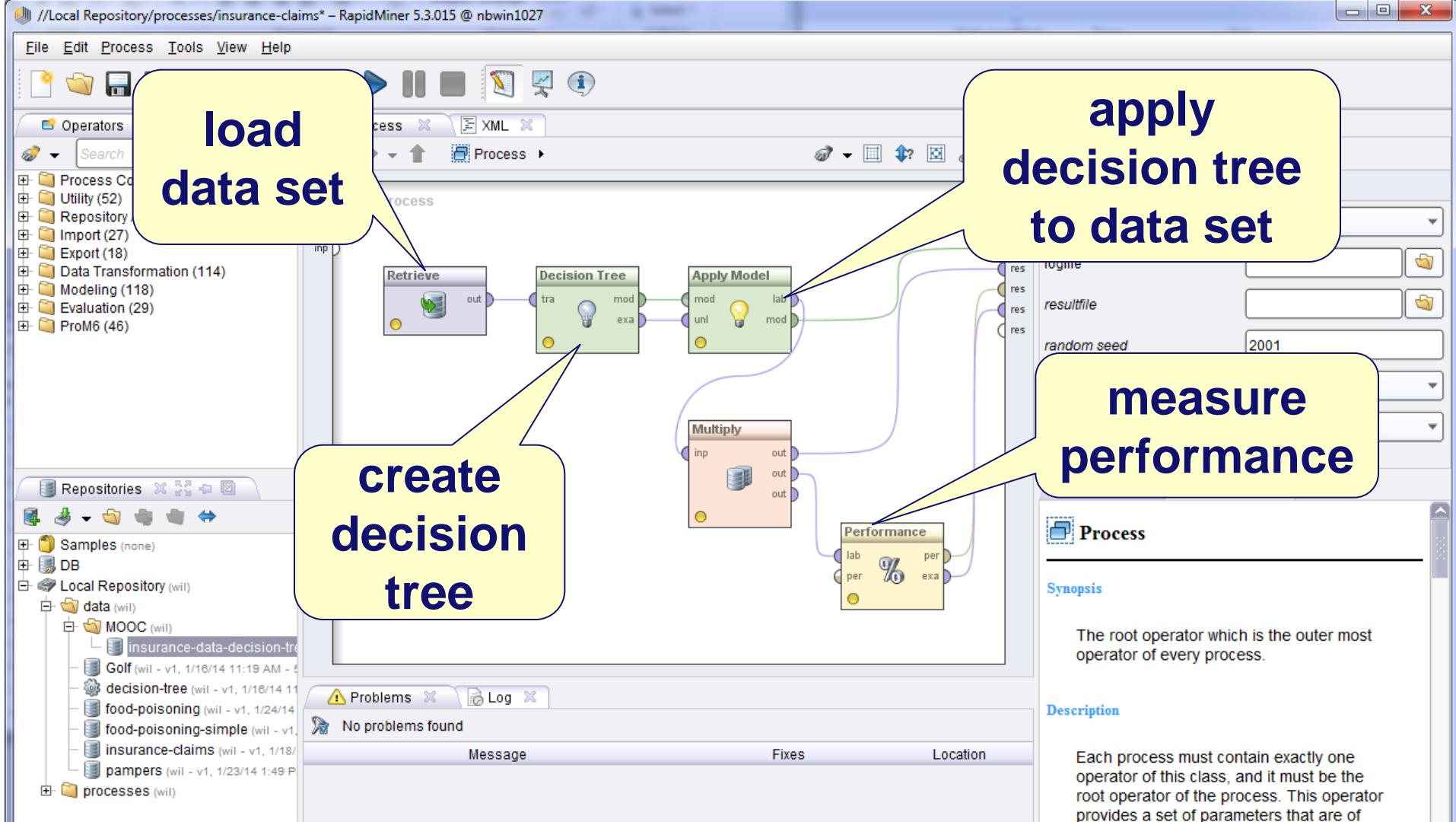
Name insurance-data-decision-tree

Location ./data/MOOC/insurance-data-decision-tree

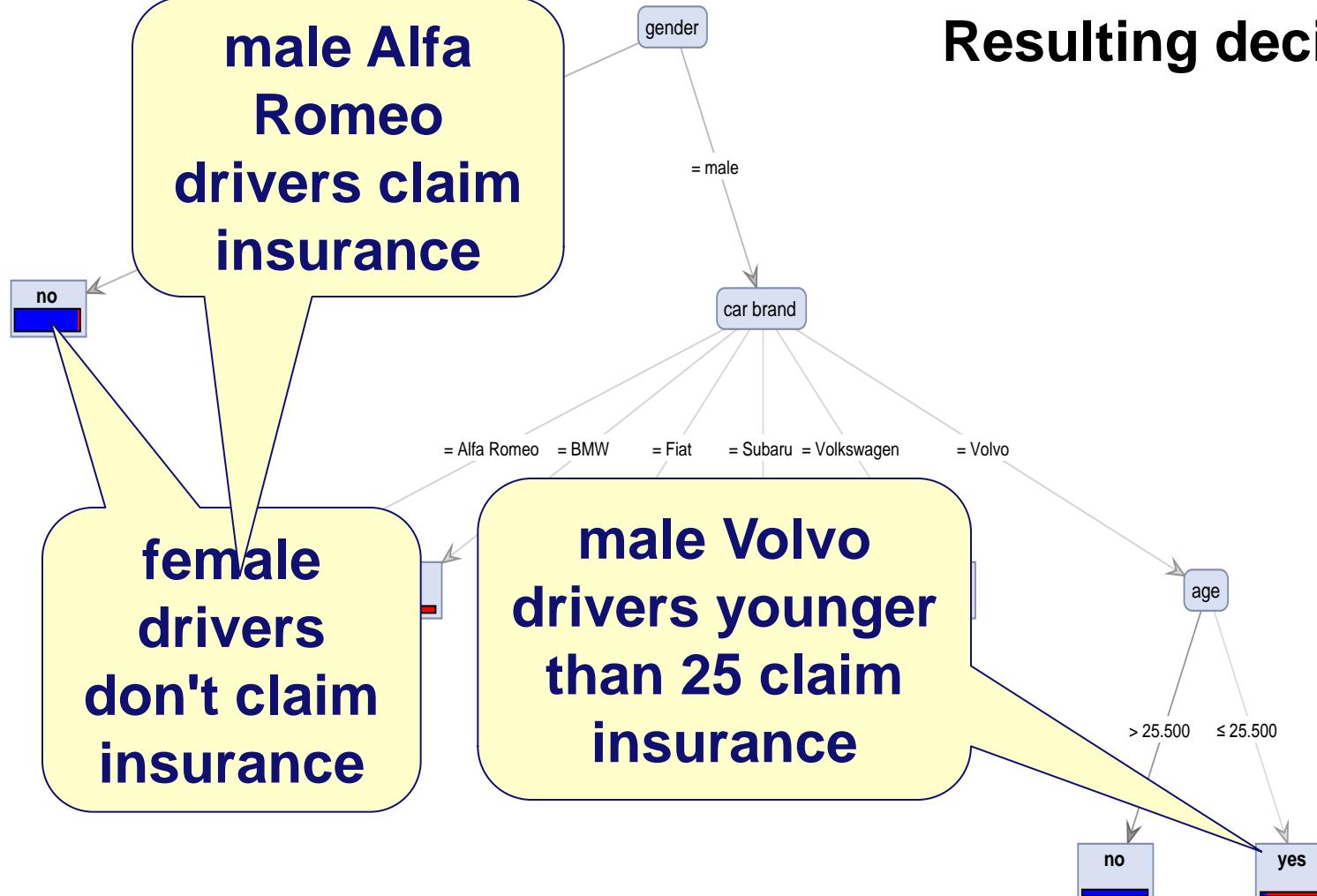
Resolve relative to //Local Repository/processes

Press "F3" for focus.

OK Cancel



# Resulting decision tree



513 females  
498 did not claim

83 male BMW drivers  
64 did not claim  
19 claimed  
22.9% wrong !!

90 m  
4

86 claimed  
4.6% wrong

no

yes



no



no



yes



car brand

= Alfa Romeo

= BMW

= Fiat

= Subaru = Volks

//Local Repository/processes/insurance-claims\* – RapidMiner 5.3.015 @ nbwin1027

File Edit Process Tools View Help

real class predicted class

ExampleSet (Multiply) ExampleSet (/Local Repository/data/MOOC... (Multiply)

Data View Met... Data View Plot View Advanced Charts Annotations

ExampleSet (999 examples, 4 special attributes, 4 regular attributes) View Filter (999 / 999): all

Row No.	claim	confidence(...)	confidence(...)	prediction(c...)	gender	age	smoker	car brand
1	no	0.971	0.029	no	female	47	yes	Volvo
2	yes	0.044	0.956	yes	male	31	no	Alfa Romeo
3	yes	0.044	0.956	yes	male	59	no	Alfa Romeo
4	no	0.827	0.173	no	male	28	no	Fiat
5	no	0.771	0.229	no	male	44	no	BMW
6	no	0.971	0.029	no	female	27	no	Fiat
7	no	0.275	0.725	yes	male	29	no	Subaru
8	yes	0.275	0.725	yes	male	44	yes	Subaru
9	no	0.771	0.229	no	male	39	no	BMW

Row No.	claim	confidence(...	confidence(...	prediction(c...	gender	age	smoker	car brand
9	no	0.771	0.229	no	male	39	no	BMW
10	yes	0.275	0.725	yes	male	35	no	Subaru
11	no	0.275	0.725	yes	male	43	no	Subaru
12	yes	0.771	0.229	no	male	25	no	BMW
13	no	0.740	0.260	no	male	39	no	Volkswagen
14	yes	0.044	0.956	yes	male	37	no	Alfa Romeo

## Which instances are classified incorrectly?

**11:** A male 43-year old non-smoking Subaru driver was predicted to claim but did not.

**12:** A male 25-year old non-smoking BMW driver was predicted to not claim, but actually did claim insurance.

//Local Repository/processes/insurance-claims\* – RapidMiner 5.3.015 @ nbwin1027

File Edit Process Tools View Help

Tree (Decision Tree) ExampleSet (/Local Repository/data/MOOC/insurance-data-decision-tree)  
Result Overview ExampleSet (Multiply) PerformanceVector (Performance) ExampleSet (Multiply)

Table / Plot View Text View Annotations

Criterion Selector  
accuracy  
f\_measure  
false\_positive  
false\_negative  
true\_positive  
true\_negative

Multiclass Classification Performance  
Annotations  
Table View Plot View

accuracy: 90.79%

	true no	true yes	class precision
pred. no	761	68	91.80%
pred. yes	24	146	85.88%
class recall	96.94%	68.22%	

	true no	true yes	class precision
pred. no	761	68	91.80%
pred. yes	24	146	85.88%
class recall	96.94%	68.22%	

**5000 parties ate at an Italian restaurant.**

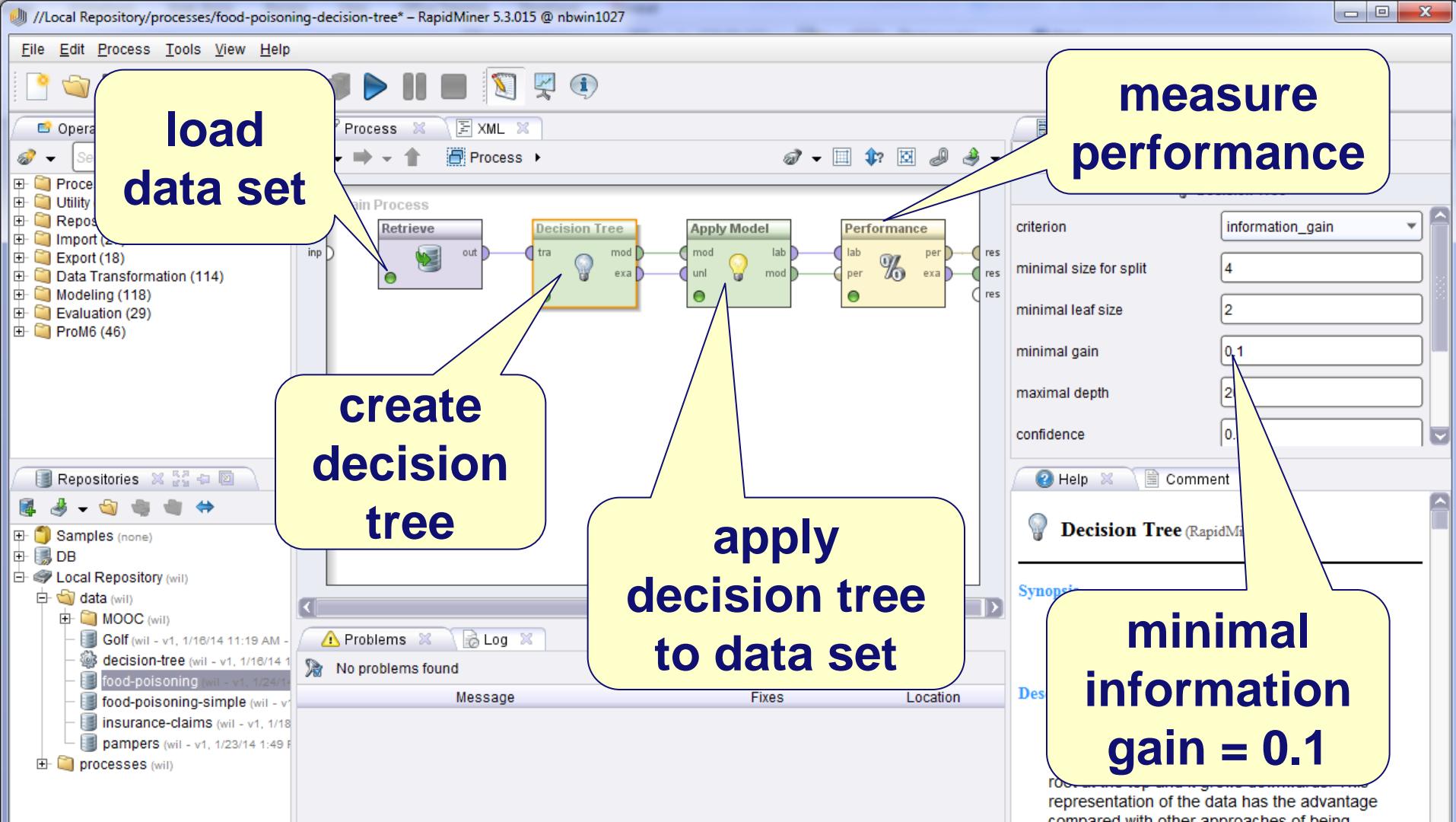


**Menu includes: pizza margherita, pizza romana, pizza marinara, pizza capricciosa, pizza siciliana, lasagna, spaghetti carbonara, spaghetti alla diavola, vino rosso, vino bianco, birra, and espresso.**



csv file loaded into the repository

ExampleSet (5000 examples, 1 special attribute, 12 regular attributes)														View Filter (5000 / 5000): all	
Row No.	class	pizza margh...	pizza romana	pizza marin...	pizza capric...	pizza sicilia...	lasagna	spaghetti c...	spaghetti al...	vino rosso	vino bianco	birra	espresso		
1	not sick	2	0	0	0	1	1	0	0	0	0	3	1		
2	not sick	1	3	1	0	4	0	1	1	2	1	0	2		
3	not sick	0	3	0	0	4	1	0	3	0	1	1	1		
4	not sick	0	0	0	0	0	0	0	0	0	0	3	0		
5	not sick	2	0	1	1	0	3	0	0	0	0	3	0		
6	not sick	0	2	0	1	3	0	0	1	0	2	1	4		
7	not sick	0	0	0	0	0	2	1	2	1	1	0	0		
8	not sick	2	0	1	0	0	3	0	0	0	0	2	0		
9	nauseous	0	1	0	0	4	1	0	1	1	1	1	2		
10	not sick	0	0	0	0	0	0	0	0	0	0	1	0		
11	not sick	0	1	0	0	0	3	3	1	1	0	3	0		
12	not sick	0	0	0	0	0	2	2	2	1	1	1	0		
13	very sick	3	0	3	0	0	1	0	0	0	0	3	1		
14	not sick	1	2	0	0	3	1	0	2	2	1	0	0		
15	nauseous	1	3	0	0	4	0	0	4	3	0	0	2		
16	not sick	0	2	0	0	4	0	1	1	2	2	0	2		
17	not sick	0	0	0	0	0	3	3	2	0	0	1	1		
18	not sick	0	1	0	2	4	1	1	2	1	2	1	2		
19	not sick	3	0	0	0	0	1	1	0	0	0	2	0		
20	not sick	0	0	0	0	0	3	3	2	0	0	2	2		

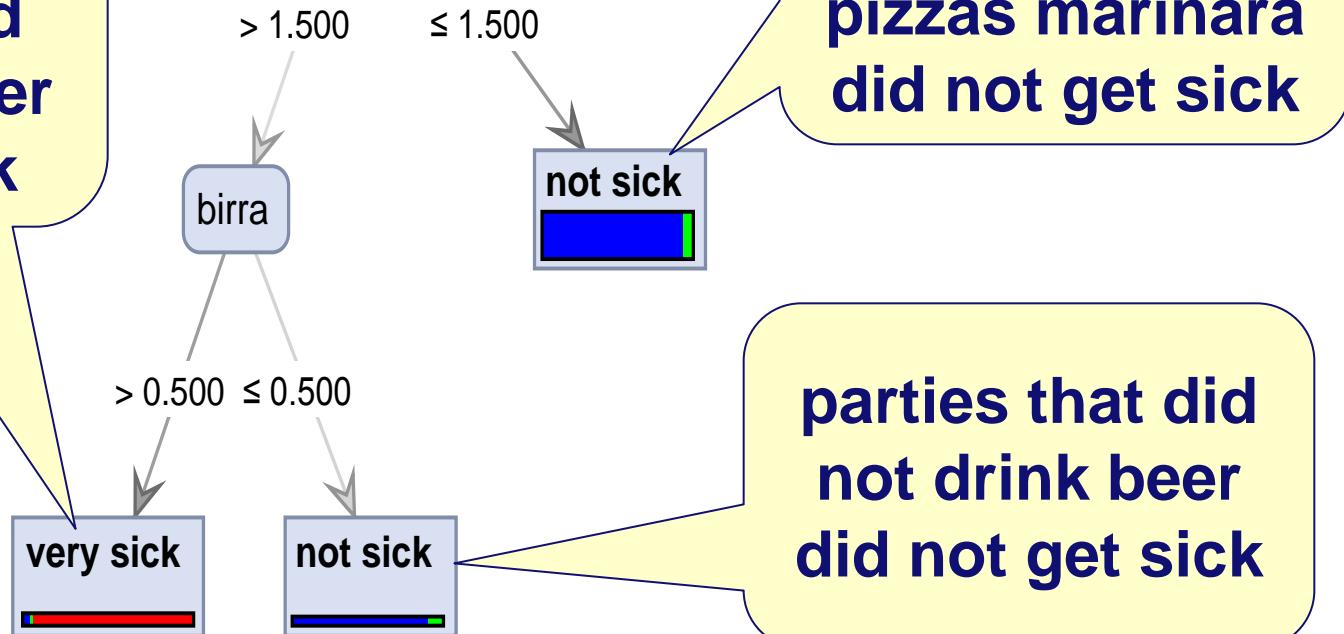


# Decision tree indicates that a combination of pizzas marinara and beer caused sickness.

parties  
multipl

marinara and  
that drank beer  
got very sick

blue = not sick  
red = very sick  
green = nauseous



total information gain = 0.1

307 of the 313 parties that  
were nauseous were  
classified as "not sick"

accuracy: 93.26%

	true not sick	true nauseous	true very sick	class precision
pred. not sick	4193	307	0	
pred. nauseous	0	0	6	
pred. very sick	24	6	0	
class recall	99.43%	0.00%		

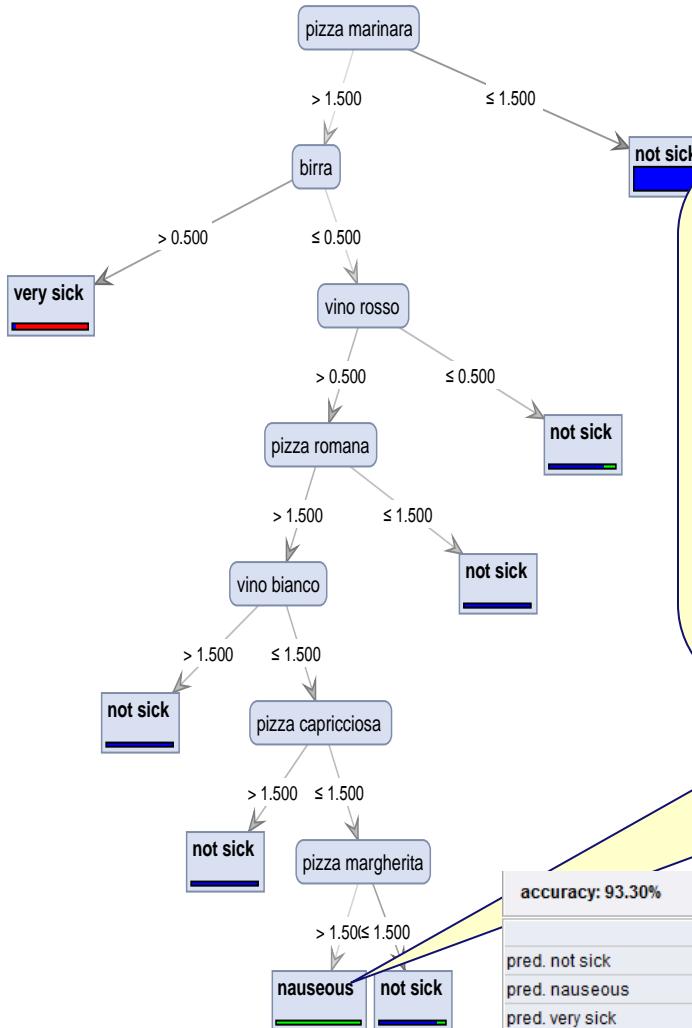
6 of the 313 parties that  
were nauseous were  
classified as "very sick"

$> 0.500 \leq 0.500$

The decision tree does not explain why  
some parties were **nauseous**.

blue = not sick  
red = not nauseous  
green = nauseous

minimal information gain = 0.05



people that ate multiple pizzas marinara, pizzas romana, pizzas margherita, but at most one pizza capricciosa, and drank red wine but not multiple glasses of white wine and did not drink any beer got nauseous.

Extremely small improvement at the cost of overfitting.

	true not sick	true nauseous	true very sick	class precision
pred. not sick	4193	305	0	93.22%
pred. nauseous	0	2	0	100.00%
pred. very sick	24	6	470	94.00%
class recall	99.43%	0.64%	100.00%	

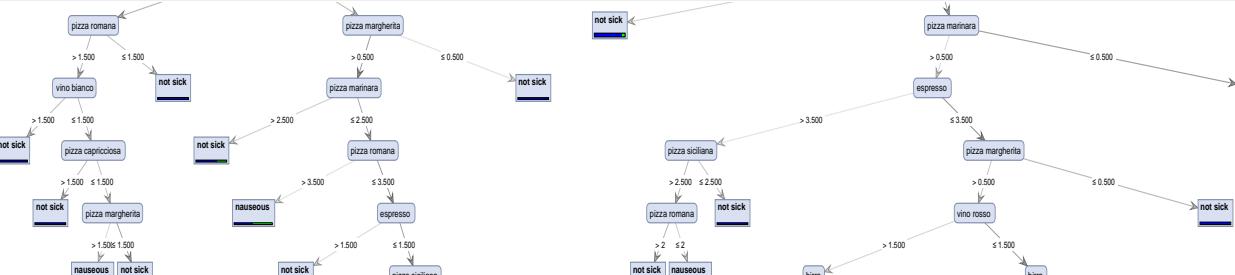
not sick



# underfitting

accuracy: 84.34%

	true not sick	true nauseous	true very sick	class precision
pred. not sick	4217	313	470	84.34%
pred. nauseous	0	0	0	0.00%
pred. very sick	0	0	0	0.00%
class recall	100.00%	0.00%	0.00%	

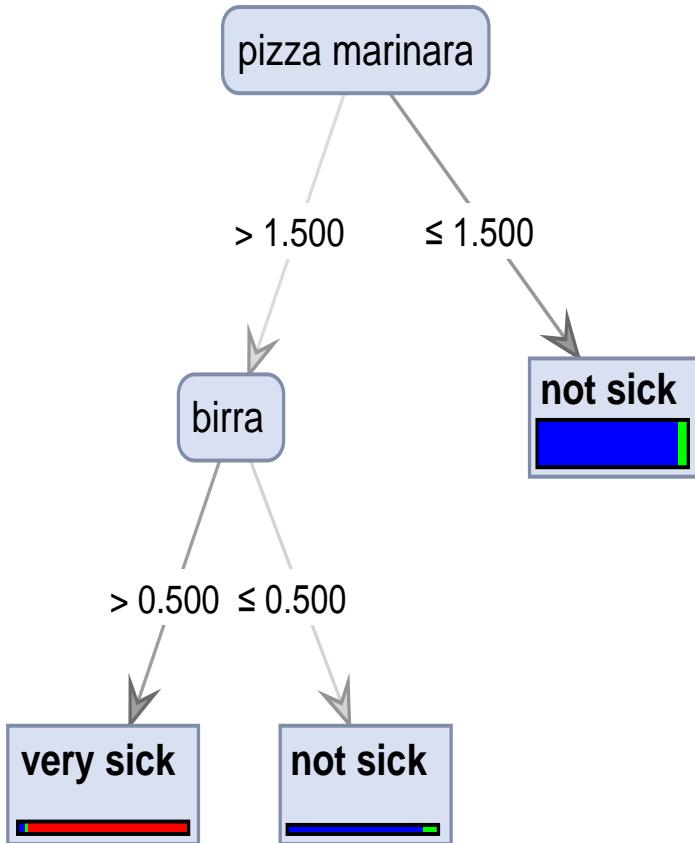


# overfitting

accuracy: 93.48%

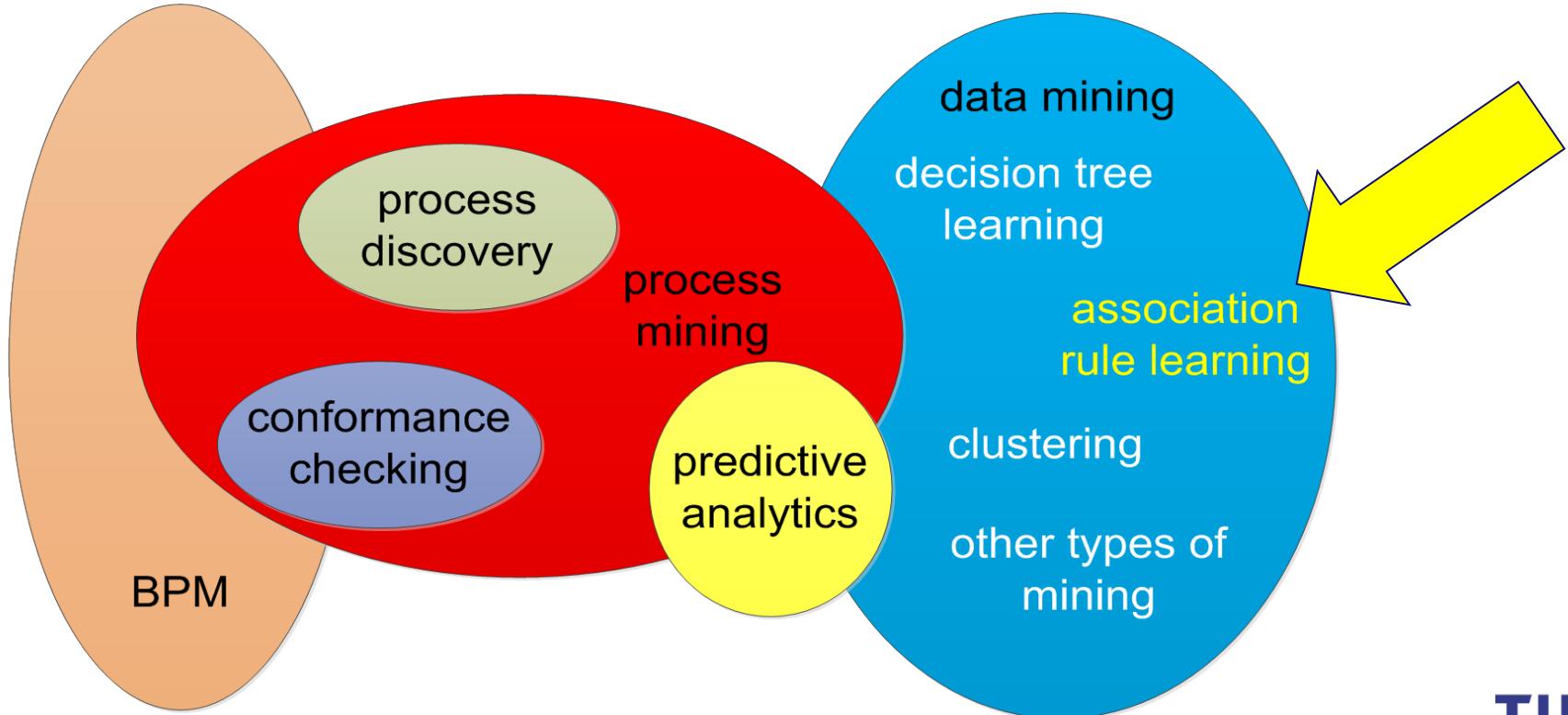
	true not sick	true nauseous	true very sick	class precision
pred. not sick	4190	293	0	93.46%
pred. nauseous	3	14	0	82.35%
pred. very sick	24	6	470	94.00%
class recall	99.36%	4.47%	100.00%	





- Reasonable balance between underfitting and overfitting.
- Can be used to understand what is happening.
- Can be used for predictions and recommendations.

# Next



## *Part I: Preliminaries*

**Chapter 1**

Introduction

**Chapter 2**

Process Modeling and Analysis

**Chapter 3**

Data Mining

## *Part III: Beyond Process Discovery*

**Chapter 7**

Conformance Checking

**Chapter 8**

Mining Additional Perspectives

**Chapter 9**

Operational Support

## *Part II: From Event Logs to Process Models*

**Chapter 4**

Getting the Data

**Chapter 5**

Process Discovery: An Introduction

**Chapter 6**

Advanced Process Discovery Techniques

**Chapter 10**

Tool Support

**Chapter 11**

Analyzing “Lasagna Processes”

**Chapter 12**

Analyzing “Spaghetti Processes”

## *Part IV: Putting Process Mining to Work*

**Chapter 13**

Cartography and Navigation

**Chapter 14**

Epilogue

