



UNIVERSITÀ DI PISA

# Intelligent Systems For Pattern Recognition

Midterm 2 | Assignment 4

This assignment involves implementing a Bayesian network from scratch.

The network was implemented using 3 classes:

- 1) **BayesianNetwork** - to handle operations on the entire underlying graph.
- 2) **BayesianNode** - to manage the single node.
- 3) **CPT\_Item** - an auxiliary class to manage Conditional Probabilities Tables.

Gennaro Daniele **Acciaro**

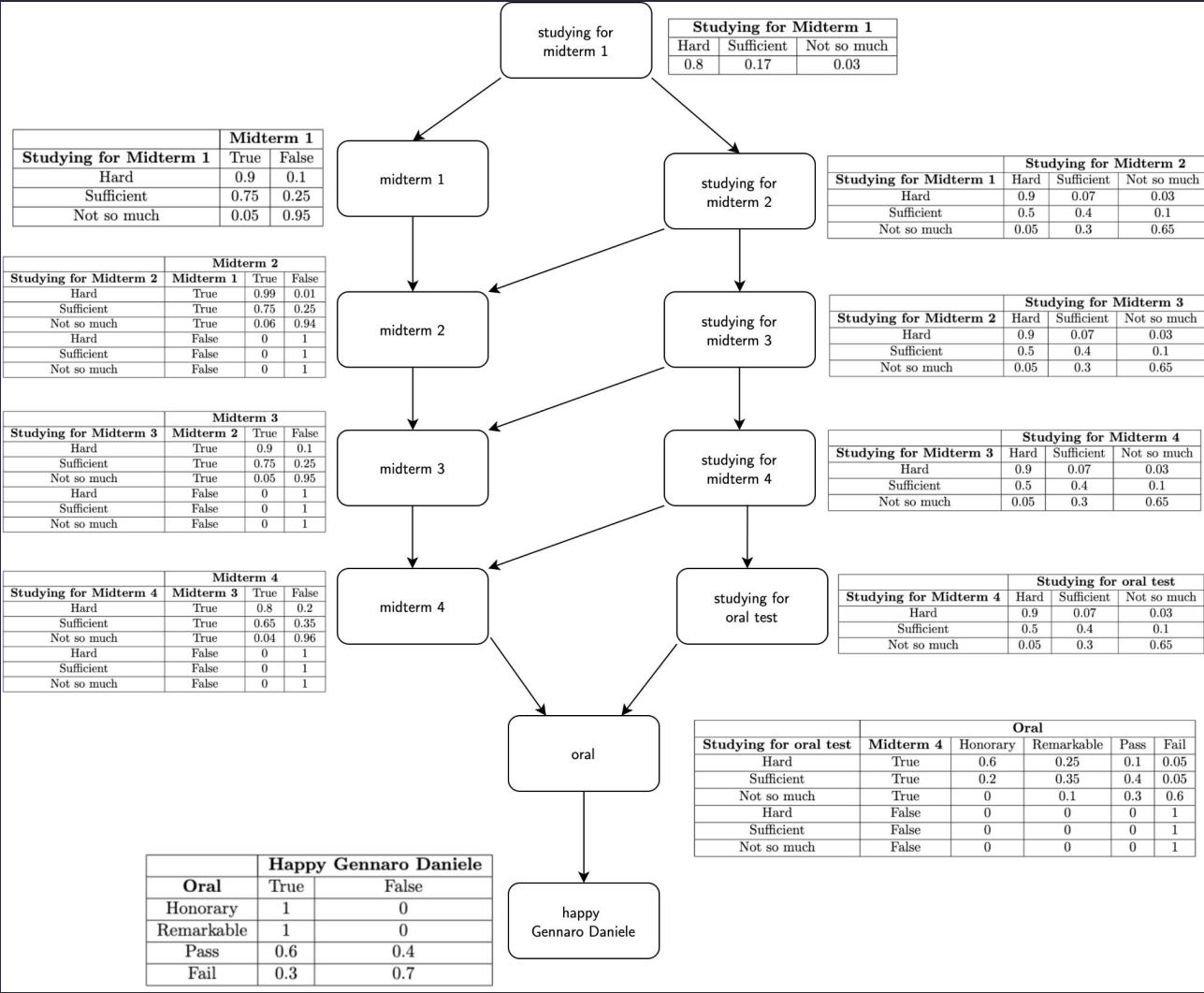
MAT. **635009**

# My process

My problem is the path to passing **this exam through the midterms.**

The random variable are:  $\{studying\_m1, studying\_m2, studying\_m3, studying\_m4, studying\_oral, midterm1, midterm2, midterm3, midterm4, oral, happy\_gennaro\_daniele\}$

- Each intermediate exam requires the student to have passed the previous midterm and to have studied for that particular midterm. The domain of each midterm exam is boolean.
- Each "studying" node indicates how much a student has studied for a midterm or final oral. The domain of these nodes is [ "Hard", "Sufficient", "Not so much" ]
- The domain for the oral variable indicates 4 possible grades : [ "Honorary", "Remarkable", "Pass", "Fail" ]
- The random variable happy\_gennaro\_daniele indicates my own happiness based on the grade on this exam.



# How I represented CPTs

Each node has a Conditional Probability Table that is implemented as a set of *CPT\_Items*, a pair (cause, probability).

Causes are defined by means of a list, and each element in this list contains the cause object and a domain value of this object.

Passing the object directly may not be the most efficient way, computationally speaking, to create CPTs, but I decided to do it this way to make more user-friendly the insertion of values in CPTs. In fact libraries like [pomegranate](#), allow to create CPTs considering using only the domain values. Personally I think this can create confusion and errors if this activity is done manually, even more so if we have to describe a very large CPT.

On the next slide there is a comparison between my version and pomegranate.

Finally, probabilities are handled via a simple list of values. This list follows the order of the domain values of the node itself.

# How I represented CPTs

```
1  bn.add_probabilities(oral,
2      {
3      CPT_Item(causes=[[studying_oral, "Hard"], [midterm4, True]],
4                probabilities=[0.6, 0.25, 0.1, 0.05]),
5      CPT_Item(causes=[[studying_oral, "Sufficient"], [midterm4, True]],
6                probabilities=[0.2, 0.35, 0.4, 0.05]),
7      CPT_Item(causes=[[studying_oral, "Not so much"], [midterm4, True]],
8                probabilities=[0., 0.1, 0.3, 0.6]),
9
10     CPT_Item(causes=[[studying_oral, "Hard"], [midterm4, False]],
11               probabilities=[0, 0, 0, 1]),
12     CPT_Item(causes=[[studying_oral, "Sufficient"], [midterm4, False]],
13               probabilities=[0, 0, 0, 1]),
14     CPT_Item(causes=[[studying_oral, "Not so much"], [midterm4, False]],
15               probabilities=[0, 0, 0, 1]),
16     }
17 )
```

My CPT

```
1 oral = ConditionalProbabilityTable(
2     [
3         ["Hard", True, "Honorary", 0.6],
4         ["Hard", True, "Remarkable", 0.25],
5         ["Hard", True, "Pass", 0.1],
6         ["Hard", True, "Fail", 0.05],
7
8         ["Sufficient", True, "Honorary", 0.2],
9         ["Sufficient", True, "Remarkable", 0.35],
10        ["Sufficient", True, "Pass", 0.4],
11        ["Sufficient", True, "Fail", 0.05],
12
13        ["Not so much", True, "Honorary", 0.],
14        ["Not so much", True, "Remarkable", 0.1],
15        ["Not so much", True, "Pass", 0.3],
16        ["Not so much", True, "Fail", 0.6],
17
18        ["Hard", False, "Honorary", 0.],
19        ["Hard", False, "Remarkable", 0.],
20        ["Hard", False, "Pass", 0.],
21        ["Hard", False, "Fail", 1.],
22
23        ["Sufficient", False, "Honorary", 0.],
24        ["Sufficient", False, "Remarkable", 0.],
25        ["Sufficient", False, "Pass", 0.],
26        ["Sufficient", False, "Fail", 1.],
27
28        ["Not so much", False, "Honorary", 0.],
29        ["Not so much", False, "Remarkable", 0.],
30        ["Not so much", False, "Pass", 0.],
31        ["Not so much", False, "Fail", 1.],
32    ], [studying_oral, midterm4]
33 )
34
```

Pomegranate's CPT

# How I calculate probabilities

Given a sequence of random variables, we can calculate the probability of the joint distribution of those variables using BN.

In this project, this computation was performed by visiting the graph through the **BFS** graph visitation algorithm. The BFS starts with the independent nodes and adds to the queue all the children of the current node that have not yet been visited.

For each node then, one probability must be selected from its CPT.

- If the node is independent, it will have only one probability in the CPT and then I will use that.
- If the node is dependent, I choose the probability whose causes are present in the causes passed as input to the algorithm.

# How I calculate probabilities

```
1 def query(self, param):
2     """ This function performs a BFS on the Bayesian network to find the probabilities to be multiplied.
3         It starts by visiting independent nodes, and then adds their children to the queue.
4         As it visits the queue, it adds the probability given the parents
5         (which have already been visited) to the result.
6     """
7
8     """ Preconditions """
9     if len(param) != len(self.nodes):
10         raise Exception("A query needs every and only the variables in the BN")
11
12     """ The graph visit needs to start from independent nodes """
13     independent_nodes = self.get_independent_nodes()
14
15     """ Once we found independent nodes, we put them to the BFS queue """
16     queue = independent_nodes
17
18     """ Mark all the vertices as not visited """
19     visited = {}
20     for node in self.nodes:
21         visited[node] = False
22
23     prob_found = list()
```

# How I

```
24 while queue:
25     item = queue.pop(0) #Remove from head
26
27     """ I execute this for cycle because I don't want
28         to put constraints in the order of appearance of the random variables in the query """
29     input_value = None
30     for p in param:
31         input_item = p[0]
32         value = p[1]
33         if item == input_item:
34             input_value = value
35             break
36
37     value_index = item.domain_values.index(input_value)
38
39     """ safety first """
40     if not hasattr(item, "prob"):
41         raise Exception("No CPT added to the node <"+str(item)+">")
42
43     for p in item.prob:
44         selected_probability = p.probabilities[value_index]
45         if not p.causes:
46             """If the node is independent"""
47             prob_found.append(selected_probability)
48         else:
49             """If the node is dependent, I choose the probability
50                 whose causes are present in the causes passed as input to the algorithm."""
51             match_counter = 0
52             for single_cause in p.causes:
53                 if single_cause in param:
54                     match_counter+=1
55
56             if match_counter == len(p.causes):
57                 prob_found.append(selected_probability)
58                 break
59
60     """ Adding my unvisited children to the queue """
61     for child in item.children:
62         if not visited[child]:
63             queue.append(child)
64             visited[child] = True
65
66     result = 1
67     for prob in prob_found:
68         result *= prob
69
70     return result
```



# Ancestral Sampling

*Ancestral Sampling* is a sampling process that involves sampling nodes without parents, and then sampling their children, conditioned by the probability of the parents, and so on.

I implemented this sampling technique through the BFS visitation algorithm.

For each node in the graph I call the **.sample()** function which, given a cause, generates a sample for that node given the probability described in the node's CPT for that particular cause. (For independent nodes, the cause is left blank).

Once the probability is found, I use Numpy's **np.random.choice()** method to generate the node's sample (which is a random value from its domain given the probability).

# Ancestral Sampling

```
1  def ancestral_sampling(self):
2      """ The graph visit needs to start from independent nodes """
3      independent_nodes = self.get_independent_nodes()
4
5      """ Once we found independent nodes, we put them to the BFS queue """
6      queue = independent_nodes
7
8      """ Mark all the vertices as not visited """
9      visited = {}
10     for node in self.nodes:
11         visited[node] = False
12
13     sample_causes = []
14     while queue:
15         item = queue.pop(0) # Remove from head
16         sample = item.sample(sample_causes)
17
18         sample_causes.append([item, sample])
19
20         """ Adding my unvisited children to the queue """
21         for child in item.children:
22             if not visited[child]:
23                 queue.append(child)
24                 visited[child] = True
25
26     return self.query(sample_causes), sample_causes
```

Bayesian **Network**'s method

# Ancestral Sampling

```
1 def sample(self, selected_causes = None):
2     #Precondition for dependent nodes:
3     for p in self.prob:
4         if p.causes and selected_causes is None:
5             raise Exception("You can't sample the node "+str(self)+" without pass the causes")
6
7     selected_prob = None
8
9     for p in self.prob:
10        if p.causes: # If the causes list is not empty
11
12            #The selected_causes input list may contain more causes than are needed for this node
13            #So, I check if all my causes are contained into the input list, ignoring all the others
14
15            all_causes_are_selected = 0
16            for cause in p.causes: # For each my cause ..
17
18                for single_selected_cause in selected_causes:
19                    if cause == single_selected_cause: # .. checks if it's contained
20                        all_causes_are_selected+=1
21
22            #If all the causes are contained, I found my probabilities
23            if len(p.causes) == all_causes_are_selected:
24                selected_prob = p.probabilities
25                break
26        else:
27            """ independent nodes """
28            selected_prob = p.probabilities
29            break
30
31    picked_value = np.random.choice(self.domain_values, 1, p=selected_prob)
32    return picked_value[0]
```

Bayesian **Node's** method

# Sampling Results

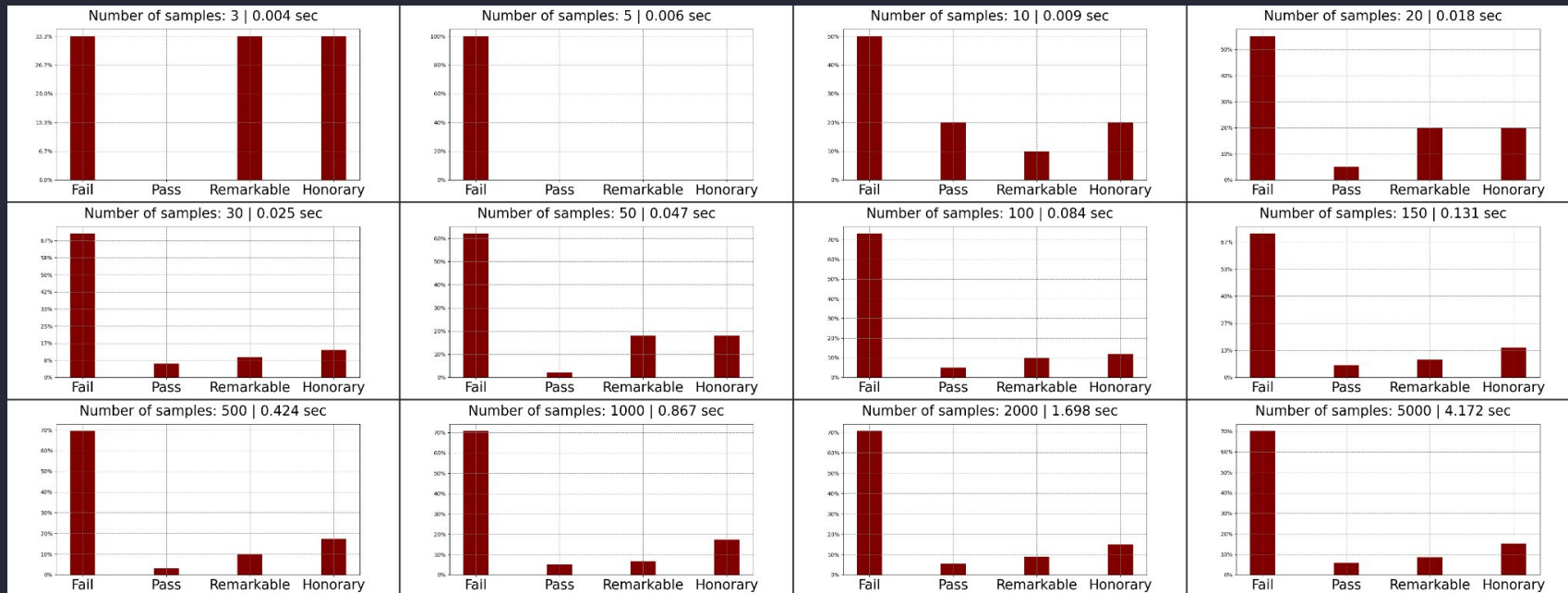
As we can see from the graphs on the next slide, a fair trade-off between the number of samples (and execution time) and the quality of the distribution is about **100** samples.

In fact the distribution of results with 100 samples is very similar to that with 5000 samples, although the time required for the first run is 0.084 secs while the execution of 2000 samples takes 4.172 secs.

Another interesting observation might be that the distribution follows the idea that, unfortunately, it is much easier to fail this exam (about 70% of total cases) than to pass it with an honors grade (about 15% of total cases).

This was an expected result.

# Sampling Results



# Finals observations

Building a Bayesian network from scratch was quite a challenge.

Computationally, performing a BFS for each query is expensive, it may be possible to improve the complexity of the project by performing a single visit to the graph after adding all the edges and save this topological order in a data structure and reuse it, whenever it is necessary to visit the graph.

Lessons learned:

- Experience with Bayesian Network
- Experience with Pomegranate (to compare my BN with this library)
- First experience with sampling approaches