

Assignment 1

COMP9418 – Advanced Topics in Statistical Machine Learning

Lecturer: Gustavo Batista

Last revision: Sunday 25th September, 2022 at 21:44

Instructions

Submission deadline: Sunday, 16th October 2022, at 18:00:00 AEDT.

Late Submission Policy: The penalty is set at 5% per late day for a maximum of 5 days. This is the UNSW standard late penalty. For example, if an assignment receives an on-time mark of 70/100 and is submitted three days late, it will receive a mark reduction of $70/100 * 15\%$. After five days, the assignment will receive a mark reduction of 100%.

Form of Submission: This is an **individual** or group of **two students** assignment. Write the name(s) and zID(s) in this Jupyter notebook. **If submitted in a group, only one member should submit the assignment. Also, create a group on WebCMS by clicking on Groups and Create and include both group members.**

You can reuse any piece of source code developed in the tutorials.

You can submit your solution via WebCMS.

Alternatively, you can submit your solution using give. On a CSE Linux machine, type the following on the command line:

```
$ give cs9418 ass1 solution.zip
```

Recall the guidance regarding plagiarism in the course introduction: this applies to this assignment. If evidence of plagiarism is detected, it will result in penalties ranging from loss of marks to suspension.

The dataset and breast cancer domain description in the Background section are from the assignment developed by Peter Lucas, Institute for Computing and Information Sciences, Radboud Universiteit.

Introduction

In this assignment, you will develop some sub-routines in Python to implement operations on Graphical Models. You will code an efficient independence test, learn parameters from complete data, and classify examples. Our classifiers will be based on Bayesian networks, Naïve Bayes and Tree-augmented Naïve Bayes classifiers.

We will apply these classifiers to the diagnosis of breast cancer. We start with some background information about the problem.

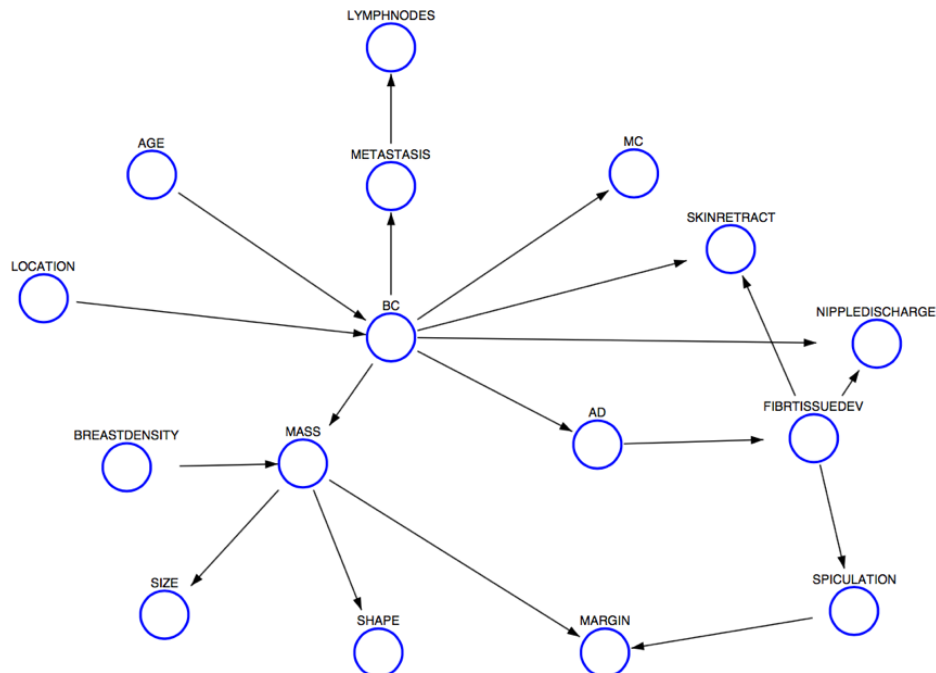
Background

Breast cancer is the most common form of cancer and the second leading cause of cancer death in women. Every 1 out of 9 women will develop breast cancer in their lifetime. Although it is impossible to say what

exactly causes breast cancer, some factors may increase or change the risk for the development of breast cancer. These include age, genetic predisposition, history of breast cancer, breast density and lifestyle factors. Age, for example, is the most significant risk factor for non-hereditary breast cancer: women aged 50 or older have a higher chance of developing breast cancer than younger women. The presence of BRCA1/2 genes leads to an increased risk of developing breast cancer, irrespective of other risk factors. Furthermore, breast characteristics, such as high breast density, are determining factors for breast cancer.

The primary technique used currently for the detection of breast cancer is mammography, an X-ray image of the breast. It is based on the differential absorption of X-rays between the various tissue components of the breast, such as fat, connective tissue, tumour tissue and calcifications. On a mammogram, radiologists can recognise breast cancer by the presence of a focal mass, architectural distortion or microcalcifications. Masses are localised findings, generally asymmetrical to the other breast, distinct from the surrounding tissues. Masses on a mammogram are characterised by several features which help distinguish between malignant and benign (non-cancerous) masses, such as size, margin, and shape. For example, a mass with an irregular shape and ill-defined margin is highly suspicious for cancer, whereas a mass with a round shape and well-defined margin is likely to be benign. Architectural distortion is focal disruption of the normal breast tissue pattern, which appears on a mammogram as a distortion in which surrounding breast tissues appear to be “pulled inward” into a focal point, often leading to spiculation (star-like structures). Microcalcifications are tiny bits of calcium, which may show up in clusters, or in patterns (like circles or lines) and are associated with extra cell activity in breast tissue. They can also be benign or malignant. It is also known that most cancers are located in the upper outer quadrant of the breast. Finally, several physical symptoms of breast cancer are nipple discharge, skin retraction, and palpable lump.

Breast cancer develops in stages. The early stage is in situ (“in place”), meaning that cancer remains confined to its original location. When it has invaded the surrounding fatty tissue and possibly has spread to other organs or the lymph, so-called metastasis is referred to as invasive cancer. It is known that early detection of breast cancer can help improve survival rates.



[10 Marks] Task 1 – Efficient d-separation test

In this part of the assignment, you will implement an efficient version of the d-separation algorithm. Let us start with a definition for d-separation:

Definition. Let \mathbf{X} , \mathbf{Y} and \mathbf{Z} be disjoint sets of nodes in a DAG G . We will say that \mathbf{X} and \mathbf{Y} are d-separated by \mathbf{Z} , written $\text{dsep}(\mathbf{X}, \mathbf{Z}, \mathbf{Y})$, iff every path between a node in \mathbf{X} and a node in \mathbf{Y} is blocked by \mathbf{Z} where a path is blocked by \mathbf{Z} iff there is at least one inactive triple on the path.

This definition of d-separation considers all paths connecting a node in \mathbf{X} with a node in \mathbf{Y} . The number of such paths can be exponential. The following algorithm provides a more efficient test implementation that does not require enumerating all paths.

Algorithm. Testing whether \mathbf{X} and \mathbf{Y} are d-separated by \mathbf{Z} in a DAG G is equivalent to testing whether \mathbf{X} and \mathbf{Y} are disconnected (ignoring edge direction) in a new DAG G' , which is obtained by pruning DAG G as follows:

1. We delete any leaf node W from DAG G as long as W does not belong to $\mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z}$. This process is repeated until no more nodes can be deleted.
2. We delete all edges outgoing from nodes in \mathbf{Z} .

More detail on this algorithm can be found on page 66 (and surrounding pages) in the textbook (Darwiche).

Implement the efficient version of the d-separation algorithm in a function `d_separation(G, X, Z, Y)` that return a boolean: true if \mathbf{X} is d-separated from \mathbf{Y} given \mathbf{Z} and false otherwise.

[5 Marks] Task 2 – Markov blanket

The Markov blanket for a variable X is a set of variables that, when observed, will render every other variable irrelevant to X . If the distribution is induced by DAG G , then a Markov blanket for variable X can be constructed using X 's parents, children, and spouses in G . A variable Y is a spouse of X if the two variables have a common child in G .

You will implement a function `Markov(G, X)` that returns a python set with the Markov blanket of X in G .

After implementing this function, you can use the d-separation test to assess the correctness of your Markov blanket implementation and vice-versa. You can use the definition of Markov blankets below to understand the connection between these two concepts.

A Markov blanket for a variable $X \in \mathbf{X}$ is the set of variables $\mathbf{B} \subseteq \mathbf{X}$ such that $X \notin \mathbf{B}$ and $X \perp \mathbf{X} \setminus (\mathbf{B} \cup \{X\}) | \mathbf{B}$

[5 Marks] Task 3 - Learning the outcome space from data

In this and the following tasks, we will implement three classifiers: Bayesian network, Naïve Bayes and a Tree-augmented Naïve Bayes classifier. We will need to create an *outcome space* for all these implementations.

We defined the outcome space as a python dictionary in the tutorials, and it maps each variable to a tuple with its possible outcomes. In the tutorials, we manually defined the outcome spaces for our examples. However, as we work with more extensive problems, learning the outcome space from data becomes more efficient.

Thus, in this task, you will implement a function `learn_outcome_space(dataframe)` that learns the outcome space from the pandas `dataframe`. Refer to our tutorials for more details about how we defined the outcome space dictionary.

[5 Marks] Task 4 – Estimate Bayesian network parameters from data

Estimating the parameters of a Bayesian Network is a relatively simple task if we have complete data. The file `bc.csv` has 20,000 complete instances, i.e., without missing values. This task will estimate and store the conditional probability tables for each graph node. As we will see in more detail in the Naïve Bayes and

Bayesian Network learning lectures, the Maximum Likelihood Estimate (MLE) for those probabilities are simply the empirical probabilities (counts) obtained from data.

In this task, you will implement a method `model.learn_parameters(data, alpha)` that learns the parameters of the Bayesian network `model`. This method should work the same as the function built in tutorials, except also implement additive smoothing with parameter α .

[5 Marks] Task 5 – Bayesian network classification

The Bayesian Network previously described has a variable that plays a central role in the analysis. The variable `BC` (Breast Cancer) can assume the values `No`, `Invasive` and `InSitu`. Accurately identifying its correct value would lead to an automatic system that could help early breast cancer diagnosis.

First, the variables `metastasis` and `lymphnodes` must be removed since these two variables can be understood as pieces of information derived from `BC` and they may not be available at the point when `BC` is classified. This is done for you at the beginning of the notebook.

Design a new method `model.predict(class_var, evidence)` that implements the classification with complete data. This function should return the MPE value for the attribute `class_var` given the `evidence`. As we are working with complete data, `evidence` is an instantiation for all variables in the Bayesian network `model` but `class_var`. **Implement the efficient classification procedure discussed in the lectures. Assure that you only join the necessary factors.**

[5 Marks] Task 6 - Bayesian network accuracy estimation

In this task, you will implement a function to measure the accuracy of the Bayesian network classifier. Design a new function `assess_bayes_net(model, dataframe, class_var)` that uses the test cases in `dataframe` to assess the performance of the Bayesian network. Such a function should return the classifier accuracy (a value between 0 and 1) defined as:

$$Acc = \frac{TP + TN}{TP + TN + FP + FN}$$

where TP are the true positives, TN are the true negatives, FP are the false positives, and FN are the false negatives.

[5 Marks] Task 7 - Bayesian network assessment with cross-validation

We will use k -fold cross-validation to generate training and test set pairs to assess our classifiers. In this task, you will implement a function called `cross_validation_bayes_net(G, dataframe, class_var, k)` to compute and report the average accuracy of the Bayesian network classifier specified by the graph `G` over $k = 10$ cross-validation runs as well as the standard deviation.

We provide a scaffold for this function that uses the Scikit-learn class `KFold`.

[5 Marks] Task 8 - Naïve Bayes classifier structure

Let's work now on the Naïve Bayes classifier. This classifier is a Bayesian network with a pre-defined graph structure. This pre-defined structure can be directly learned from the dataset definition, i.e., the class attribute and the features.

You will start creating a function `learn_naive_bayes_structure(dataframe, class_var)` that learns the Naïve Bayes graph structure from a pandas `dataframe` using `class_var` as the class variable. This function should return a graph object with the learned graph. The graph class was implemented in the first tutorial.

[5 Marks] Task 9 – Naïve Bayes classification

As the Naïve Bayes classifier is a Bayesian network, we can use the existing `BayesNet` class to create a new class `NaiveBayes`. Thus, the Naïve Bayes class will inherit all the methods we implemented for the Bayesian networks.

We will also implement a new method for the Naïve Bayes class in this task. The method `model.predict_log(class_var, evidence)` implements the classification with complete data. It differs from the implementation in Task 5 by **using the log probability trick discussed in the lectures**.

The method should return the MPE value for the attribute `class_var` given the `evidence`. As we work with complete data, `evidence` is an instantiation for all variables but `class_var`.

[5 Marks] Task 10 - Naïve Bayes accuracy estimation

This task is similar to Task 6. You should implement a function `assess_naive_bayes(model, dataframe, class_var)` that uses the test cases in `dataframe` to assess the performance of the Naïve Bayes classifier `model` for the class variable `class_var`. This function will return the accuracy of the classifier according to the test examples in `dataframe`.

[5 Marks] Task 11 - Naïve Bayes assessment with cross-validation

This task is similar to Task 7. You will implement a function called `cross_validation_naive_bayes(dataframe, class_var, k)`, compute and report the average accuracy of a Naïve Bayes classifier over $k = 10$ -fold cross-validation runs as well as the standard deviation. We provide a scaffold for this function that uses the Scikit-learn class `KFold`.

[15 Marks] Task 12 - Tree-augmented naïve Bayes structure

Let's work now with the Tree-augmented Naïve Bayes classifier (TAN). The TAN classifier is a mid-term between a Naïve Bayes and a Bayesian network, and it allows a richer graph structure learned directly from data using mutual information.

We will start by creating a new function, `learn_tan_structure(dataframe, class_var)`, that learns the Tree-augmented graph structure from a pandas dataframe. Refer to Lecture 6, slide 24 for the algorithm that describes the steps to learn the graph structure from data. Remind that mutual information (MI) was defined in Lecture 3, slide 29.

Also, observe that the TAN classifier is a Bayesian network. Therefore, we can use the existing `BayesNet` class to help implementing this classifier.

[5 Marks] Task 13 - Tree-augmented naïve Bayes assessment with cross-validation

In our final task, we will implement a function called `cross_validation_tan(dataframe, class_var, k)`, compute and report the average accuracy over $k = 10$ -fold cross-validation runs as well as the standard deviation. This task is similar to Tasks 7 and 11. Notice we are not asking you to implement an `assess_tan` function since it would be the same implemented for the Bayesian network classifier.

[20 Marks] Task 14 – Report

Write a report (**with less than 500 words**) summarising your findings in this assignment. Your report should address the following:

- a. Make a summary and discussion of the experimental results. You can analyse your results from different aspects such as accuracy, runtime, coding complexity and independence assumptions. You can use plots to illustrate your results.
- b. Discuss the time and memory complexity of the implemented algorithms.

Use Markdown and Latex to write your report in the Jupyter notebook. If you want, develop some plots using Matplotlib to illustrate your results. Be mindful of the maximum number of words. Please, be concise and objective.