

Bayesian Network Classifiers

ALESSIO FALAI

University of Bologna
alessio.falai@studio.unibo.it

July 7, 2020

Abstract

This report explains the work done and the results achieved in the Knowledge Representation class project, specifically in the module about uncertainty and probabilistic reasoning, held by professor Paolo Torroni. In this work, we used the Adult dataset to classify world citizens that perceive a high income, based on different features, that were accurately pre-processed and discretized. The classification task was performed using various Bayesian network structures, based on the Naive Bayes model. Probabilistic inference results on a test set were then compared with ground-truth data to evaluate the accuracy of the dataset, along with other classification-related measures.

I. INTRODUCTION

IN this work, we tested the capabilities of various Bayesian network structures in a classification task, over the standard Adult dataset [UCI, 2010], which aims at separating people whose income is greater than 50 thousands dollars per year from the rest.

The first operation that needed to be done was data cleaning:

- Useless features, like `fnlwgt`, were removed
- Redundant features, like `education-num`, were removed too
- Rows containing null values were removed, since there were only a few

The second operation that needed to be done was data discretization, to simplify the following construction of the Bayesian network structures:

- The age variable was divided into 4 bins (child, young adult, adult and senior)
- The hours-per-week variable was divided into 4 bins (part time, full time, over time and too much time)
- The capital-gain and capital-loss variables were binned according to different quantiles distributions

The output of this pre-processing step will be used as input for the following steps.

II. BAYESIAN NETWORKS

Different Bayesian network structures were used to compare classification capabilities over the same dataset, so as to assess which structure could be more suitable to solve the presented problem [Cheng et al., 1999].

Every tested structure is actually based on the Naive Bayes model and compared against it.

In each of the tested Bayesian networks, the various CPDs were estimated from the training dataset using MLE (Maximum Likelihood Estimation).

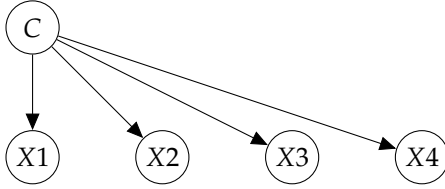
i. Naive Bayes (NB)

The Naive Bayes model, as shown in figure 1, has been extensively used in classification tasks, with good accuracy results, because of its simplicity.

It does not require any structural learning, since it has a fixed structure, where the classification variable is the parent node of every other feature variable.

The Naive Bayes model works by assuming full independence between each pair of variables, given the classification node.

Figure 1: Example of a NB model

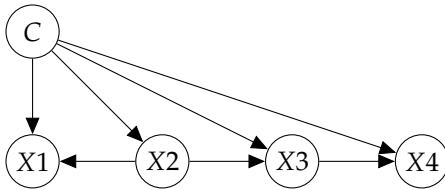


ii. Tree-Augmented Naive Bayes (TAN)

The TAN model [Friedman et al., 1997], as shown in figure 2, is just like a Naive Bayes model, so there is a connection from the classification node to every other feature node, with the exception that, without the classification node and all its related edges, the Bayesian network becomes a tree.

The TAN model needs to be learned from the training data, by using a modified version of the Chow-Liu algorithm [Chow, Liu, 1968].

Figure 2: Example of a TAN model



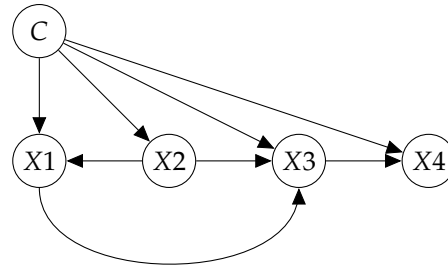
iii. BN-Augmented Naive Bayes (BAN)

The BAN model [Friedman et al., 1997], as shown in figure 3, is just like a Naive Bayes model, so there is a connection from the classification node to every other feature node, with the exception that, without the classification node and all its related edges, the Bayesian network becomes a DAG (Directed Acyclic Graph).

The BAN model needs to be learned from the training data, by using a modified CBL2 algorithm [Cheng et al., 1997b].

The actual implementation of the CBL2 algorithm resulted in much higher running times and similar results w.r.t the simpler and faster TAN model.

Figure 3: Example of a BAN model



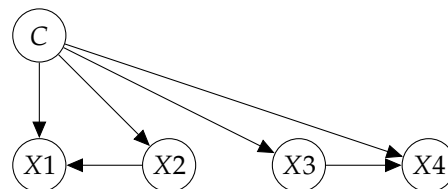
iv. Forest-Augmented Naive Bayes (FAN)

The FAN model [Jiang et al., 2005], as shown in figure 4, is just like a Naive Bayes model, so there is a connection from the classification node to every other feature node, with the exception that, without the classification node and all its related edges, the Bayesian network becomes a forest.

The FAN model needs to be learned from the training data, by using a similar reasoning as what has been done with the TAN model.

This model gives results similar to the ones achieved by the TAN model, but with a slightly higher running time.

Figure 4: Example of a FAN model



III. RESULTS

Table 1 shows the results obtained with the variable elimination inference algorithm, over the various Bayesian network structures, based on different binary classification measures, i.e. accuracy, precision, recall and F-score.

The reported results are averages over different runs.

Table 1: Results summary

	NB	TAN	BAN	FAN
Accuracy	80%	85%	84%	84%
Precision	60%	73%	69%	72%
Recall	78%	63%	67%	63%
F-score	68%	68%	68%	67%

IV. CONCLUSIONS

We developed different variants of the famous Naive-Bayes model, all of which augmented the model with a specific structure over the non-classification variables.

Results show that augmenting the Naive-Bayes model yields better accuracy over a classification task on the Adult dataset. Further investigation would be required to check if this boost in accuracy does not depend on the type of inference algorithm used.

Additional studies would also be required to actually assess which augmentation method should be preferred in specific situations.

Moreover, practical programming limitations didn't allow us to create hybrid models with both discrete and continuous variables. These limitations led to the discretization of certain features, which would have been better represented as continuous ones.

In conclusion, this work shows that augmenting the simple Naive Bayes model can boost classification-related measures, like accuracy, but different augmentations didn't provide much gain w.r.t. each other in terms of those same measurements.

REFERENCES

- [UCI, 2010] A. Frank, A. Asuncion. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>.
- [Chow, Liu, 1968] Chow, C.K. and Liu, C.N.. Approximating discrete probability distributions with dependence trees. *IEEE Trans. on Information Theory*, 14 (pp. 462-467).
- [Cheng et al., 1997a] Cheng, J., Bell, D.A. and Liu, W.. An Algorithm for Bayesian Belief Network Construction from Data. *Proceedings of AI & STAT'97* (pp. 83-90), Florida.
- [Cheng et al., 1997b] Cheng, J., Bell, D.A. and Liu, W.. Learning Belief Networks from Data: An Information Theory Based Approach. *Proceedings of ACM CIKM'97*.
- [Cheng et al., 1999] Jie Cheng, Russell Greiner. Comparing Bayesian Network Classifiers. *UAI'99: Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence* (pp. 101-108).
- [Friedman et al., 1997] Jie Cheng, Russell Greiner. Bayesian Network Classifiers. *Machine Language*, Volume 29, Issue 2-3.
- [Jiang et al., 2005] Jiang, Zhang, Cai, Su. Learning Tree Augmented Naive Bayes for Ranking. *DASFAA'05: Proceedings of the 10th international conference on Database Systems for Advanced Applications* (pp. 688-698).