

**IST 707 HOMEWORK 8**

Joy (Qiaoyi) Liu

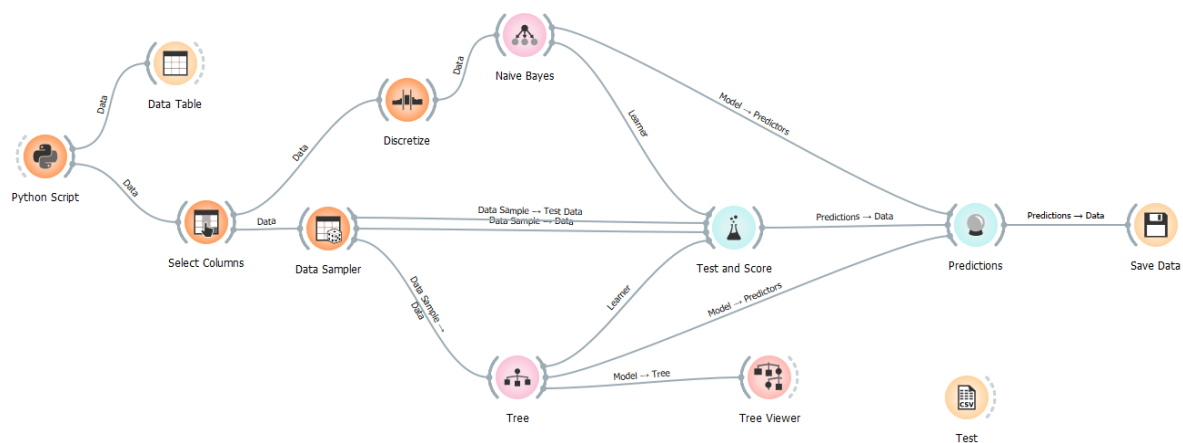
IST 707 Applied Machine Learning

Professor Joshua Introne

October 23<sup>rd</sup>, 2022

## Introduction

The Modified National Institute of Standards and Technology (MNIST) dataset consists of 784 pixels/columns indicating each image with 28 pixels in height and 28 pixels in width. By showing an integer value between 0 to 255 for each pixel, it represents the lightness or darkness of that pixel to deduce the possible number of that image. With a train dataset of 42,000 rows and a label indicating the number, the Decision Tree and Naïve Bayes model could predict the possible number/label of the test dataset (28,000 rows).



## Decision Tree

I first imported the train dataset into Orange via the Python Script widget. Then I selected the label column as the target for the predictions afterward. I compared the performance of the decision tree by tuning the sampling proportion from 70% to 90% in 3-fold validation. The predictions from the decision tree did not vary (AUC = .995, CA = .948, F1 = .948, Precision = .948, Recall = .948).

## Naïve Bayes

I added the discretize widget before Naïve Bayes and changed the setting from Equal-frequency discretization to Equal-width discretization. There was only a slight difference in its

prediction accuracy. The Equal-frequency discretization has  $F1 = .841$  and the Equal-width discretization has  $F1 = .842$ .

## Algorithm Performance Comparison

Predictions - Orange

Show probabilities for: Classes in data Restore Original Order

	Naive Bayes	Tree	label	Tree	Naive Bayes	Tree (0)	Tree (1)
1	0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.97:0.02:0.00:0.00:0.00:0.01:0.00:0.00...	2	2	2	0	0
2	0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00...	0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00...	6	4	6	0	0
3	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.95:0.01:0.01:0.00:0.01:0.01:0.00:0.01:0.00...	1	1	1	0	0.955752
4	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.97:0.01:0.00:0.00:0.00:0.00:0.00:0.01:0.00...	1	1	1	0	0.955752
5	0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00...	0.00:0.00:0.00:0.33:0.00:0.00:0.67:0.00:0.00:0.00...	6	1	6	0	0.956853
6	1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.96:0.00:0.01:0.00:0.00:0.00:0.02:0.00:0.00:0.01...	0	0	0	0.976119	0
7	0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00:0.00...	5	5	5	0	0
8	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00...	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00...	9	9	9	0	0
9	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00...	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00...	9	3	9	0	0
10	0.00:0.00:0.00:0.00:0.91:0.00:0.00:0.00:0.00:0.09...	0.00:0.00:0.00:0.00:0.96:0.00:0.00:0.00:0.00:0.02...	4	4	4	0	0
11	0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.03:0.95:0.03:0.00:0.00:0.00:0.00:0.00:0.00...	2	2	2	0	0
12	0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00...	0.01:0.00:0.01:0.01:0.00:0.00:0.98:0.00:0.00:0.00...	6	6	6	0	0
13	0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	2	5	4	0	0
14	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00...	0.00:0.04:0.00:0.00:0.00:0.00:0.00:0.00:0.96:0.00...	8	8	8	0	0.00943396
15	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.95:0.01:0.01:0.00:0.01:0.01:0.00:0.01:0.00...	1	1	1	0	0.956853
16	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.95:0.01:0.01:0.00:0.01:0.01:0.00:0.01:0.00...	1	1	1	0	0.956853
17	0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00...	0.01:0.00:0.01:0.01:0.00:0.00:0.98:0.00:0.00:0.00...	6	6	6	0	0
18	0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00...	0.00:0.01:0.01:0.01:0.00:0.00:0.00:0.96:0.00:0.01...	7	7	7	0	0.00478469
19	0.00:0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.00:0.99:0.00:0.00:0.00:0.00:0.01:0.00...	3	3	3	0	0
20	1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.96:0.00:0.01:0.00:0.00:0.00:0.02:0.00:0.00:0.01...	0	0	0	0.976119	0
21	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00...	0.00:0.00:0.00:0.00:0.03:0.00:0.00:0.00:0.01:0.96...	9	9	9	0	0
22	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.01:0.00:0.01:0.01:0.00:0.00:0.98:0.00:0.00:0.00...	6	6	1	0	0
23	0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00...	1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0	4	6	0	0
24	0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.25:0.00:0.00:0.00:0.00:0.25:0.00:0.50...	7	9	4	0	0
25	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.95:0.01:0.01:0.00:0.01:0.01:0.00:0.01:0.00...	1	1	1	0	0.956853
26	0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00...	0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00:0.00...	7	4	7	0	0
27	0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.03:0.95:0.03:0.00:0.00:0.00:0.00:0.00:0.00...	2	2	2	0	0
28	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00...	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:1.00:0.00...	8	8	8	0	0
29	0.00:0.00:0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.00:0.00:0.96:0.00:0.00:0.00:0.00:0.02...	4	4	4	0	0
30	0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.00:0.05:0.00:0.95:0.00:0.00:0.00:0.00...	3	5	5	0.00952381	0
31	0.00:0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.00:0.97:0.02:0.00:0.00:0.00:0.01:0.00:0.00...	2	2	2	0	0
32	0.00:1.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00:0.00...	0.00:0.95:0.01:0.01:0.00:0.01:0.01:0.00:0.01:0.00...	1	1	1	0	0.956853

Show performance scores Target class: (Average over classes)

Model	AUC	CA	F1	Precision	Recall
Naive Bayes	0.984	0.843	0.844	0.847	0.843
Tree	0.996	0.951	0.951	0.951	0.951

7350 | 7350 | 2x7350

In general, the Decision Tree performs better than Naïve Bayes. The decision tree has a higher accuracy in predicting the labels and visualizing the results clearly. However, it has disadvantages in overfitting. The Decision Tree has the F-measure of .948, which is highly possible the model is overfitting to the dataset. However, to be certain whether the accuracy indicates overfitting would require a validation curve. I only found the solution to do so in Python, not in

Orange. Naïve Bayes is less likely to overfit due to the fact they “ignore” irrelevant features. Therefore, they have less accuracy in prediction. I personally did not recognize the difference in speed in generating the Decision Tree and Naïve Bayes. So I looked up and found that it is commonly accepted the Naïve Bayes is computationally faster in making decisions because it simplifies the computation by assuming there are no dependencies amongst attributes (Ashari et al., 2013).

### Reference

Ashari, A., Paryudi, I., & Tjoa, A. M. (2013). Performance comparison between Naïve Bayes, decision tree and k-nearest neighbor in searching alternative design in an energy simulation tool. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 4(11).