

Comparative Study of Data Classifiers Using Rapidminer

Abhishek Kori

Assistant Professor, IT Department, SVVV Indore, India

Abstract--Data mining, the extraction of hidden predictive information from large databases, is a powerful new technology with great potential to help to focus on the most important information in data warehouses. Classification is the process of organizing data into categories for its most effective and efficient use. A well-planned data classification system makes essential data easy to find and retrieve. In this paper our focus is text classification by Naïve Bayes Classification and KNN and determine accuracy of the classifier using rapidminer tool.

Keywords: Classification, Naïve Bayes, KNN

I. INTRODUCTION

RapidMiner is a data science software platform developed by the company of the same name that provides an integrated environment for machine learning, deep learning, text mining, and predictive analytics. It is used for business and commercial applications as well as for research, education, training, rapid prototyping, and application development and supports all steps of the machine learning process including data preparation,

results visualization, validation and optimization.^[1] RapidMiner is developed on an open core model. The RapidMiner (free) Basic Edition, which is limited to 1 logical processor and 10,000 data rows, is available under the AGPL license.

II. LITERATURE SURVEY

Classification is a learning function that maps a given data item into one of several predefined classes. It is a data analysis technique to extract models describing important data classes and predict future values. Data mining uses classification technique uses with machine learning, image processing, natural language processing, statistical and visualization techniques to discover and present knowledge in an understandable format. Most of the classification algorithms in literature are memory resident, typically assuming a small data size.

III. CLASSIFICATION ALGORITHM

3.1 Naïve Bayes Classifier

Naive Bayes is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. It is not a single algorithm for training such classifiers, but a family of

algorithms based on a common principle: all naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable.

$$p(c/d) = \frac{p(c)p(d/c)}{p(d)}$$

Where P(d) plays no role in selecting c.

3.2 KNN Classifier

k-NN is a type of instance-based learning, or lazy learning, where the function is only approximated locally and all computation is deferred until classification. The k-NN algorithm is among the simplest of all machine learning algorithms.

IV. EXPERIMENTATION & RESULTS

In this paper, experiment is carried out using tool Rapidminer 7.4 we are taking into consideration the data set Iris, given as sample data inside the repository panel of the tool. We apply validation tool on the data set which in turn contains training and testing operations. In training column, we took naïve bayes operator and in testing column we took apply model and performance tools respectively. After connecting all operators we execute the tool which in turn shows the accuracy as result. In the paper, we have consider this process for both naïve bayes and KNN and compare the results given in the figure shows below.



Figure 4.2 Applying Validation Operator



Figure 4.3 Applying KNN Operator

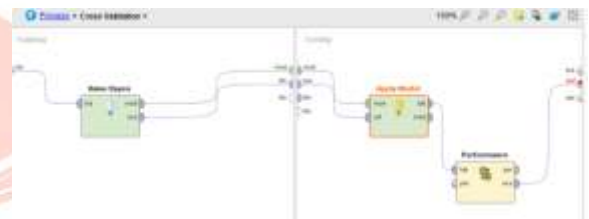


Figure 4.4 Applying Naïve Bayes Operator

	class iris-setosa	class iris-versicol	class iris-virginica	class precision
pred iris-setosa	30	0	0	100.00%
pred iris-versicol	3	47	2	94.00%
pred iris-virginica	3	3	47	94.00%
class recall	100.00%	94.00%	94.00%	

Figure 4.5 Performance Evaluation KNN

	class iris-setosa	class iris-versicol	class iris-virginica	class precision
pred iris-setosa	30	3	3	100.00%
pred iris-versicol	0	47	4	92.18%
pred iris-virginica	3	3	46	93.50%
class recall	100.00%	94.00%	92.00%	

Figure 4.6 Performance Evaluation Naïve Bayes

id	sepal length	sepal width	petal length	petal width	class
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	4.4	1.5	0.4	setosa
7	4.8	3.4	1.6	0.2	setosa
8	5.2	3.7	1.4	0.3	setosa
9	5.2	3.4	1.6	0.4	setosa
10	4.7	3.4	1.3	0.2	setosa
11	4.7	3.0	1.6	0.3	setosa
12	4.8	3.4	1.4	0.3	setosa
13	4.9	3.6	1.4	0.3	setosa
14	4.8	3.6	1.0	0.2	setosa
15	5.1	3.5	1.4	0.3	setosa
16	4.8	3.0	1.4	0.3	setosa
17	4.7	3.2	1.3	0.2	setosa
18	4.9	3.1	1.5	0.2	setosa
19	5.4	4.4	1.5	0.4	setosa
20	5.2	3.7	1.4	0.3	setosa
21	5.2	3.4	1.6	0.4	setosa
22	4.7	3.4	1.3	0.2	setosa
23	4.7	3.0	1.6	0.3	setosa
24	4.8	3.4	1.4	0.3	setosa
25	4.9	3.6	1.4	0.3	setosa
26	4.8	3.6	1.0	0.2	setosa
27	5.1	3.5	1.4	0.3	setosa
28	4.8	3.0	1.4	0.3	setosa
29	4.7	3.2	1.3	0.2	setosa
30	4.9	3.1	1.5	0.2	setosa
31	5.4	4.4	1.5	0.4	setosa
32	5.2	3.7	1.4	0.3	setosa
33	5.2	3.4	1.6	0.4	setosa
34	4.7	3.4	1.3	0.2	setosa
35	4.7	3.0	1.6	0.3	setosa
36	4.8	3.4	1.4	0.3	setosa
37	4.9	3.6	1.4	0.3	setosa
38	4.8	3.6	1.0	0.2	setosa
39	5.1	3.5	1.4	0.3	setosa
40	4.8	3.0	1.4	0.3	setosa
41	4.7	3.2	1.3	0.2	setosa
42	4.9	3.1	1.5	0.2	setosa
43	5.4	4.4	1.5	0.4	setosa
44	5.2	3.7	1.4	0.3	setosa
45	5.2	3.4	1.6	0.4	setosa
46	4.7	3.4	1.3	0.2	setosa
47	4.7	3.0	1.6	0.3	setosa
48	4.8	3.4	1.4	0.3	setosa
49	4.9	3.6	1.4	0.3	setosa
50	4.8	3.6	1.0	0.2	setosa
51	5.1	3.5	1.4	0.3	setosa
52	4.8	3.0	1.4	0.3	setosa
53	4.7	3.2	1.3	0.2	setosa
54	4.9	3.1	1.5	0.2	setosa
55	5.4	4.4	1.5	0.4	setosa
56	5.2	3.7	1.4	0.3	setosa
57	5.2	3.4	1.6	0.4	setosa
58	4.7	3.4	1.3	0.2	setosa
59	4.7	3.0	1.6	0.3	setosa
60	4.8	3.4	1.4	0.3	setosa
61	4.9	3.6	1.4	0.3	setosa
62	4.8	3.6	1.0	0.2	setosa
63	5.1	3.5	1.4	0.3	setosa
64	4.8	3.0	1.4	0.3	setosa
65	4.7	3.2	1.3	0.2	setosa
66	4.9	3.1	1.5	0.2	setosa
67	5.4	4.4	1.5	0.4	setosa
68	5.2	3.7	1.4	0.3	setosa
69	5.2	3.4	1.6	0.4	setosa
70	4.7	3.4	1.3	0.2	setosa
71	4.7	3.0	1.6	0.3	setosa
72	4.8	3.4	1.4	0.3	setosa
73	4.9	3.6	1.4	0.3	setosa
74	4.8	3.6	1.0	0.2	setosa
75	5.1	3.5	1.4	0.3	setosa
76	4.8	3.0	1.4	0.3	setosa
77	4.7	3.2	1.3	0.2	setosa
78	4.9	3.1	1.5	0.2	setosa
79	5.4	4.4	1.5	0.4	setosa
80	5.2	3.7	1.4	0.3	setosa
81	5.2	3.4	1.6	0.4	setosa
82	4.7	3.4	1.3	0.2	setosa
83	4.7	3.0	1.6	0.3	setosa
84	4.8	3.4	1.4	0.3	setosa
85	4.9	3.6	1.4	0.3	setosa
86	4.8	3.6	1.0	0.2	setosa
87	5.1	3.5	1.4	0.3	setosa
88	4.8	3.0	1.4	0.3	setosa
89	4.7	3.2	1.3	0.2	setosa
90	4.9	3.1	1.5	0.2	setosa
91	5.4	4.4	1.5	0.4	setosa
92	5.2	3.7	1.4	0.3	setosa
93	5.2	3.4	1.6	0.4	setosa
94	4.7	3.4	1.3	0.2	setosa
95	4.7	3.0	1.6	0.3	setosa
96	4.8	3.4	1.4	0.3	setosa
97	4.9	3.6	1.4	0.3	setosa
98	4.8	3.6	1.0	0.2	setosa
99	5.1	3.5	1.4	0.3	setosa
100	4.8	3.0	1.4	0.3	setosa

Figure 4.1 Sample Data Set Iris

V.CONCLUSION

This is because KNN is non-parametric, i.e. it makes no assumption about the data distribution. Contrast this to NB, which assumes that attributes are conditionally independent to each other given the class, and are normally distributed (for real-valued attributes). The experiment carried out shows that accuracy of KNN is greater than Naïve Bayes but this situation is not similar always, as it may vary with different data set.

REFERENCES

- [1] Ting, S. L., W. H. Ip, and Albert HC Tsang. "Is Naive Bayes a good classifier for document classification?." *International Journal of Computer Applications* 28.2 (2011): 37-40.
- [2] Williamson, Eric R., and Saurabh Chakravarty. "CS5604 Fall 2016 Classification Team Final Report." (2016).
- [3] Ikonomakis, M., S. Kotsiantis, and V. Tampakas. "Text classification using machine learning techniques." *WSEAS transactions on computers* 4.8 (2005): 966-974.
- [4] Kamruzzaman, S. M., Farhana Haider, and Ahmed Ryadh Hasan. "Text classification using data mining." *arXiv preprint arXiv:1009.4987* (2010).
- [5] Joachims, Thorsten. "Text categorization with support vector machines: Learning with many relevant features." *European conference on machine learning*. Springer Berlin Heidelberg, 1998.
- [6] Menaka, S., and N. Radha. "Text classification using keyword extraction technique." *International Journal of Advanced Research in Computer Science and Software Engineering* 3.12 (2013).
- [7] Williamson, Eric R., and Saurabh Chakravarty. "CS5604 Fall 2016 Classification Team Final Report." (2016).
- [8] Dalal, Mita K., and Mukesh A. Zaveri. "Automatic text classification: a technical review." *International Journal of Computer Applications* 28.2 (2011): 37-40.
- [9] Ting, S. L., W. H. Ip, and Albert HC Tsang. "Is Naive Bayes a good classifier for document classification?." *International Journal of Software Engineering and Its Applications* 5.3 (2011): 37-46.
- [10] Mahesh Kini M , Saroja Devi H , Prashant G Desai, Niranjana Chiplunkar." Text Mining Approach to Classify Technical Research Documents using Naïve Bayes" *International Journal of Advanced Research in Computer and Communication Engineering Vol. 4, Issue 7, July 2015*
- [11]. Gongde Guo, Hui Wang, David Bell, Yaxin Bi and Kieran Greer, "KNN Model-Based Approach in Classification", Proc. ODBASE pp- 986– 996, 2003