



Review

Utilizing Support Vector Machines for Efficient and Intelligent Data Retrieval

Monika Arora; Uma Kanjilal; Dinesh Varshney

Department of Computer Science, Victoria University of Wellington, P. O. Box 600, Wellington, New Zealand

*Corresponding author

Monika Arora

Department of Computer Science, Victoria University of Wellington, P. O. Box 600, Wellington, New Zealand

Article information

Received: July 30th, 2024; Revised: September 23rd, 2024; Accepted: October 20th, 2024; Published: November 2nd, 2024

Cite this article

Arora M, Kanjilal U, Varshney D. Utilizing support vector machines for efficient and intelligent data retrieval. 2024; 2(2).

doi: <https://doi.org/10.70705/ppp.ir.2024.v02.i02.pp94-97>

ABSTRACT

The information access is the wealth of material that can be retrieved from it; it has developed to provide the main methods or tactics for looking things up online. A search engine is now the de facto standard for locating data on the Internet. Numerous opportunities emerge at various stages where methods may be explored for constructing the effective information retrieval. Support vector machine (SVM) level identification and online document classification are the focus of the current research. The goal of this research is to provide a framework for smart and efficient retrieval. A paradigm for intelligent and efficient retrieval is proposed in this study. The model's stated goal was to identify the critical success criteria for intelligent, efficient retrieval. The goal of the suggested model is to bring together all the different perspectives on information retrieval in order to build a comprehensive theory that can be seen as the foundation of any system. Mapped onto a higher dimensional space using functions, this study examines the use of Support Vector Machine for.

Keywords

Information retrieval; Web information retrieval; Support vector machine.

INTRODUCTION

Appropriate for a big collection of attribute values and polynomial kernels, where the values of the kernels might reach infinite. In order to make vectors that are linearly dependent on n dimensions linearly independent on those dimensions, one may apply the kernel function of order n to them. Upon entering the dimension space, they exhibit linear separability. Regarding the third, the RBF: One method for classifying data is SVMs, or Support Vector Machines [1]. Many believe that SVM is more practical to use in

space occupied by neural networks as well. To improve the retrieval process's accuracy, this method is used. It also handles regular problems with ease. Here is how the SVM works: Data and documents must be initially partitioned into training and testing sets in order to complete this assignment. There are a number of features (observed variables) and a target value in every training set instance. Support vector machines (SVMs) aim to forecast the target value of test data by using a model or objective that is based on training data. The dataset can now sort the label pairs into their respective categories thanks to this training set. To get the best answer for the training data set, support vector machines (SVMs) are used [2, 3]. The data

set is completely linear separable, meaning that the training vectors provide results that are superior to those of the linear kernels. When selecting kernel types for Support Vector Machines, the RBF is by far the most used. The fact that it is limited in its response space and confined throughout the whole actual x -axis range is the primary reason for this. If your data has a class-conditional probability distribution function that is very close to the Gaussian distribution, this kernel will work well with it. It transforms this kind of data into a new space where it can be easily separated into linear components [4].

To see this in action, it's helpful to note that the exponential kernel can be expanded into an infinite series, producing an infinite-dimension polynomial kernel; each of these kernels can transform some dimensions to make them linearly separable. Compared to the Linear and Polynomial kernels, the RBF kernel should, of course, perform much better. Nevertheless, finding the optimal σ and selecting the associated C that solves a particular issue is challenging when designing this kernel. Part of the reason for the RBF-based SVM's error rate is that the SVM becomes very sensitive to training data when certain combinations of σ and C are used. You don't have to provide the weights, number of support vectors N_s , or support vectors s_i when using the RBF kernel; these parameters are automatically produced as part of the training process. This is one benefit of the RBF kernel.



When it comes to classification, the RBF-based kernel works better [5]. Compared to the other three, the fourth, the sigmoid kernel $K(x, y) = \frac{1}{1 + \exp(-\gamma(x^T y + r))}$, is less effective for classification. Actually, a positive definite kernel is required to fulfill Mercer's theorem, which is one of the basic conditions for a valid kernel. Having correctly selected γ and r , the Sigmoid kernel need not be positive definite. It is possible for the SVM to perform worse than chance when the kernel is not positive definite since the outputs will be very inaccurate.

Importantly, given a certain set of γ values, for certain range of values of the same parameters, the kernel assumes the shape of an RBF kernel [4], while for other ranges of values, it acts as a linear kernel. The classification issue may be handled using an SVM scenario. So, it's like an upgraded version of the linear kernel; it provides the necessary adjustment to "enable independence" among the training sets. Since both the linear and this kernel are based on the same idea and are only transformed to point to different spaces, we may anticipate that their performance will be quite similar. The degree to which the data becomes separable is dependent on the polynomial's order p , which in turn affects the performance.

II. ANALYZING DATA SET PREPARATION FOR SVM TESTING

In the experiment, the test data evaluates the retrieval performance of various relevance feedback methods on document based IR. A category and subcategories are assumed for the classification of the document of the individuals. The details of the test data are described in Figure

6.4.1. The dataset parameters are first picked from the database randomly, and this category is assumed to be the user's query target. The individual can be male or female based on the category IT for 1, HR for 2, Finance 3 and Marketing for 4.

The training dataset is primarily focused on IT1, IT2, and IT3 but may also include information from any other topic area, such as HR1. For each individual's specific area of interest, a dataset is kept. After then, the system takes it a step further by useful comments. Documents are selected from the database and classified as relevant or non-relevant according to the database's ground truth in each iteration of the relevance feedback process. In the first round, we use the identical starting data points to run all techniques, and we randomly choose two relevant papers and three irrelevant ones. An example from the actual world is shown in Table 1. Both the data format and the parameter sequence are detailed in the table. Following is a breakdown of the dataset's 14 parameters.

Information Technology (IT), Human Resources (HR), Finance (Finance), and Marketing (Marketing) make up the bulk of the dataset. The second parameter determines whether the data is male or female. The set of criteria for 3–14 determines the Yes/No category, which is equal to 1/0. The IT category is assigned the area 3-5, the HR category 6-8, the marketing category 9-11, and the finance category 12-14. In order to get the model ready, we use the SVM-Train and SVM-Predict functions with a 100% accuracy assumption. After then, in subsequent rounds, different methods use different display set selection algorithms to choose which documents to show. The model's accuracy is 100% when tested with the provided training

data.

We use the algorithms for the SVM-based methods in the experiment by changing the codes in the libsvm library [6]. The importance of the experimental circumstances in influencing the assessment outcomes is highlighted. It selects the same kernel and settings for all SVM-based algorithms to offer an objective and bias-free performance measure. We conducted an experiment to compare the performance of several kernels in order to choose the optimal one for the given dataset. Table 1 lists the kernel functions used in the experiment. There are four groups that the datasets fall under (4-Cat). There are four texts in each category that all belong to the same semantic class. By calculating the average accuracy on the top 16 retrieval results, we can assess the performance of various kernel functions. The training dataset model may then be constructed with this assistance.

Because the data is dispersed and subject to error due to its vast distribution, the optimal parameter may be altered by the size of the data collection. On the other hand, in reality, it comes from It is important to ensure that the data set is appropriate for the whole training set when using cross-validation. Using the scaling values, you can prevent characteristics with larger numerical ranges from overpowering those with smaller ranges. Additionally, it is used to sidestep numerical issues that may arise along the computation. The parameter's use in inner products determines the kernel values.

Example: while working with feature vectors like linear and polynomial kernels, it's possible to run into numerical issues when dealing with attributes with very big values. It is recommended that the training and testing datasets utilize the same set of methods with scalability. As an example, you may scale the first characteristic from both the training and testing sets of data.

III. MODEL SELECTION

There are four common kernels available to use. The RBF kernel is more popular in usage because it uses the penalty parameter C for the parameter of the kernel chosen. Also the RBF kernel is rationally the first choice because it uses the nonlinearly maps samples and that to a higher dimensional space. But linear kernel provides up with the relationship between the class labels and attributes. The linear kernel is one of special case of RBF, which uses with a penalty parameter C . The performance as the RBF kernel with some parameters (C ; γ) are same [7]. Also, the sigmoid kernel behaves like RBF for some of the parameters [4]. It is a polynomial kernel which considers more hyper-parameters than the RBF kernel.

Cross validation and Grid search are considered as two different parameters for an RBF kernel: C and γ for any problem. The model creation uses a parameter for selection for developing the best model. The main objective is to make out the best suited values for (C ; γ) so that the classifier can accurately predict on the unknown data as testing data. The common approach applied is to divide the data set into two parts. The prediction accuracy is obtained from that unknown set to more precisely that reflect the performance on categorize an independent set of data. The cross-validation process called the improved version also uses a that is based on v -fold cross-validation. The data set in cross-validation process divide the training set into equal size v subsets. Also the one subset is tested by using the



classifier qualified on the left over $(v - 1)$ subsets. The whole training set is predicted or tested once so the cross-validation accuracy in terms of the percentage of data which are correctly classified. The circles and triangles which are filled are the training data set while hollow circles and triangles from the testing data. The testing accuracy of the classifier in Figures 2(a) and (b) is not good since it over fits the training data. The training and testing data used on the training and confirmation sets in cross-validation is not appropriate. It gives better testing accuracy as well as cross-validation using Grid search. It recommend a grid-search on C and using cross-validation to various pairs of $(C; \gamma)$ values are applied and achieve the best cross-validation accuracy is picked. It found that trying exponentially rising sequences of C and γ in a practical method.s accuracy. Also, there are several advanced methods used similar to the cross-validation rate. Also, there are two motivations for the application (a) Training data and an over fitting classifier (b) Applying an over fitting classifier on testing data

When $s=2$ and $t=50$, the default values for c and γ are taken into account in order to apply the training data. When dealing with issues involving hundreds or more data points, the aforementioned method is effective. A practical method for handling massive data sets is to randomly choose a section of the dataset and run the grid-search on it. The whole dataset was subjected to the better-region-only grid-search ia. The testing data is produced using the serial number instead of the class or category number. In order to test 50 records, the model that was created for 16 records is used. All fifty participants' data is grouped into four groups according to the parameters' values, as per the specification. Classification on a grand scale was also shown by this method. Given the available data, this dataset, which is designed for only four groups, may theoretically be extended to thousands of additional categories. Both both feature-rich and feature-poor data are guided by these in our practice data. Before feeding the data to SVM, assume that there may be data relevant for using a subset of the characteristics, especially if there are hundreds of them.

4. A Result from an Experiment Based on Sample Data

Here we compare the accuracy of the suggested approach. The three issues listed in the table below are addressed by the experiments. To get the precision, we use the LIBSVM program [8]. That precision is a direct result of the data used for training and testing. As a second point, it need to demonstrate the disparity between the accuracy levels with and without scaling, which is not taken into account for this test results. In order to build a model that can recover the other parameters of the training set's qualities, they are evaluated and tested on the testing set. Thirdly, the accuracy of the suggested approach for ideal model selection is also given. Lastly, the LIBSVM utility automates the application of the whole dataset. In order to classify the massive dataset prior to retrieval, the same parameter selection technique will be used.

When it comes to specialization and the number of connected individuals who manage to achieve fame, the most popular class is 2. The model's training dataset specifies which individuals meet specific criteria for classification. An HR professional may have an interest in IT, and vice versa; this is just one example of how different fields

may overlap. Based on the training data, this model is able to interpret and classify. The results will be obtained by applying the model that was trained on the training data to the testing data.

Figure 3 shows that the results are evenly distributed among all four groups. If the data distribution in the figure is accurate, then these are the top retrieval results generated by the same relevance feedback system for all four datasets. The following conclusions are drawn from the experiment: The use of relevant feedback mechanisms has enhanced the retrieval outcomes following iteration. This proves that using relevant feedback strategies may enhance document based IR's retrieval performance. Consequently, these methods can only obtain locally relevant photos; they can't enhance retrieval performance by retrieving more relevant images.

IV. CONCLUSION

Many people think about how the parameters interact with one other. A connection to a class is established via these parameters or documents. This class finds the dataset and puts it in the right category so you can find the information you're looking for. It is the information-seeking criteria that are the cause of the document interactions. Data representation in categories may be a crucial step in assessing relevance feedback systems for real-world information retrieval. Three distinct characteristics are extracted to characterize the data, or its qualities. It is explored that document-based IR relevance feedback issues may be addressed using SVM-based relevance feedback methods. The relevance feedback issue with unbalanced datasets and the suggestion of a new relevance feedback method using Support Vector Machine.

We outline the benefits of our suggested methods and show how they stack up against more conventional methods. Both the training data and the actual dataset may be used for the experiments. The experimental findings show that the relevance feedback technique based on support vector machines is a good candidate for increasing document-based information retrieval performance.

REFERENCES

- (Quek, C. Y., 1997). Classification of Documents on the World Wide Web. Honors thesis submitted to Carnegie Mellon University in 1997.
- [2] A training approach for optimum margin classifiers by Bosser, Guyon, and Vapnik. On pages 144–152, you may find the proceedings of the Fifth Annual Workshop on Computational Learning Theory. Publ. by ACM in 1992.
- [3] Support-vector network by Cortes and Vapnik. Volume 20, pages 273-297, 1995, Machine Learning.
- A research on sigmoid kernels for support vector machines and the training of non-PSD kernels using SMO-type approaches was published in March 2003 by Lin H.-T. and Lin C.-J.
- [4]. The work "Orthogonal Least Squares Learning Algorithm for Radial Basis Function Networks" was published in the IEEE Transactions on Neural Networks in 1991. The authors of the article are Chen, Cowan, and Grant.

Using social network analysis to promote knowledge generation



and sharing: a bird's-eye perspective (IBM Corporation, 2002)

[6] Cross, R. Referenced in

[7] Keerthi and Lin (2003). Behaviors of support vector machines

when they approach infinity with a Gaussian kernel. *Journal of Neural Computing*, 15, 7, 1667-1689.

[8] LIBSVM: a library for support vector machines, 2001, by
