

On the optimum choice of the K Parameter in Hand-Written Digit Recognition by kNN in comparison to SVM

IVAYLO PENEV, MILENA KAROVA, MARIANA TODOROVA

Abstract: - The paper concerns the application of two machine learning algorithms – k-nearest neighbor (kNN) and support vector machines (SVM) for solving the problem of hand-written digit recognition. The main goal of the work is to derive recommendations for the choice of the K parameter in kNN (number of the nearest neighbors) so that the performance of kNN to be the near (or even better than) the performance of SVM – one of the most power machine learning known algorithms. The kNN distance function as well as the method for choosing a class of the recognized digit are explained. The presented experimental results show comparison of the kNN performance to SVM, regarding two criteria – percent of the correctly recognized digit images and run time for recognition. As a final result recommendations for the choice of the K value are summarized.

Key-Words: - machine learning, nearest neighbors, kNN, SVM, hand-written digits, recognition

I. INTRODUCTION

The hand-written digit recognition problem is of significant importance for practice. An example for its application are the mail services of some countries, where the packets are scheduled automatically. Furthermore the problem is a good example for image recognition and is a proper base for research of various algorithms.

The problem is a classification problem. The input data (an image of a hand-written digit) are classified to one of a set of groups (called classes), defined in advance. The classes for this problem are digits from 0 to 9. Classification problems are solved by machine learning algorithms. These algorithms are trained by proper data, for which the pairs input – output data are known. So called

classifier is built, which is then tested with another test data set. Finally the trained algorithm is able to recognize with a high level of probability new data, unknown to the algorithm.

Different machine learning algorithms and methods for solving classification problems are known – for example neural networks, supported vector machines [3]. Each one has advantages and disadvantages. As a result of significant research groups of classification problems and their instances are described, for which some machine learning algorithms are recommended.

One of the most often machine learning algorithms used is the k-nearest neighbor algorithm, notated as kNN. The algorithm is especially suitable for solving the hand-written digit problem. In general the algorithm performs the following main steps:

Calculation of distances between data sets;

Finding the nearest neighbors on the basis of the calculated distances;

Choosing a class for the new data set corresponding to the class of the nearest neighbors.

The main advantages of kNN in comparison to other machine learning algorithms are easy implementation, no necessity of explicit training, easy interpretation of the achieved results [4, 5]. The algorithm needs to be supplied with a data sets of known classes.

The key problem in kNN implementation is the choice of the number of the nearest neighbors (this is the K parameter). The K parameter has significant influence on the performance of the algorithm. If the K value is high, the classifier is precise and more new data sets are correctly classified, but the recognition takes long time. In case of low K value the algorithm completes fast, but produces great recognition error.

The impact of the K parameter on the kNN performance is studied for various problems (e.g. [1, 2]). The common conclusions are, that the choice of K depends on the specific problem and its optimal value is determined experimentally.

The authors are with the Department of Computer Science and Engineering and Department of Automation, Technical University of Varna, 9010 Varna, Studentska Str. 1, BULGARIA. Emails: ivailo.penev@tu-varna.bg, mkarova@ieee.org, mgtodorova@tu-varna.bg <http://cs.tu-varna.bg>

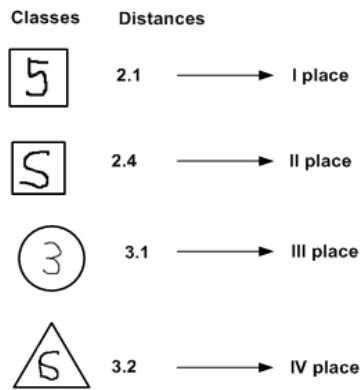


Fig. 3. Nearest neighbors in the case of $K=4$

After the nearest neighbors are found a class for the digit is determined. The class “winner” is the one with most participants in the list of the nearest neighbors. This is the class to which the digit for recognition belongs to. For the rating from fig. 3 the final choice of a class is presented on fig. 4.

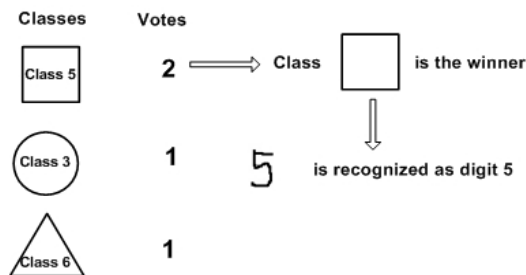


Fig. 4. Choosing a class to which the digit for recognition belongs to

A special case with two or more classes with equal number of participants in the list of nearest neighbors is possible. In this case additional weighted vote is performed. A coefficient is calculated:

$$\frac{1}{D^2(\text{digit_new}, \text{digit_known})} \quad (2)$$

, where digit_new – digit for recognition, digit_known – digit from the candidate class, for which the vote is calculated and D – the calculated distance between digit_new and digit_known .

The vote for each candidate class is calculated as follows:

$$\text{vote} = \sum_{i=1}^n \frac{1}{D^2(\text{digit_new}, \text{digit_known})} \quad (3)$$

, where n – number of known digits from the candidate class, for which the vote is calculated.

Equation 3 shows, that less distance between the digit for recognition and the known digits has greater impact on the calculated vote. The winner is the class with most votes.

III EXPERIMENTAL RESULTS

III.a EXPERIMENTAL ENVIRONMENT

The kNN and SVM algorithms are implemented in C#. The machine learning framework Accord.NET is used. It provides powerful tools for audio, image processing, statistics as well as other features, implemented in C# [6].

The application is installed and tested in the following computer configuration:

- Processor Intel Xeon E5450 3.0 GHz;
- 4 GB RAM;
- Windows 7 Ultimate Service Pack 1 64-bit (x64);
- .NET Framework 4.5.

III.b EXPERIMENTS

The experiments aim to measure the performance of the kNN algorithm in comparison to the SVM (Support Vector Machines) algorithm for the hand-written recognition problem. The experiments are carried out in two stages. On the first stage the two algorithms are tested separately. The three best variations of the algorithms from the first stage are then compared on the second stage. The tests are carried out with 946 examples of images of digits from 0 to 9. The performance criteria are percent of correctly recognized digits and run time of the algorithms.

III.b.1 FIRST STAGE OF THE EXPERIMENTS

III.b.1.1 PERFORMANCE OF THE KNN ALGORITHM

Before running the tests the algorithm is provided by 1934 images of digits with known classes. The algorithm is run with a set of 946 digits (from 0 to 9) for recognition. The tests are performed with different values of the K parameter. Table 1 presents the results from the tests.

Table 1. Results from the tests of kNN for various *K*

k	Number of correct recognized digits	Number of incorrect recognized digits	% of the correct recognized digits	Time for kNN execution (ms)
1	935	12	98,73	7,591
2	933	13	98,63	7,532
3	934	12	98,73	7,582
4	936	10	98,94	7,483
5	929	17	98,2	7,445
6	929	17	98,2	7,47
7	926	20	97,89	7,571
8	927	19	97,99	7,601
9	926	20	97,89	7,574
10	927	19	97,99	7,504
15	922	24	97,46	7,51
20	921	25	97,36	7,524
25	916	30	96,83	7,544
30	912	34	96,41	7,577
35	908	38	95,98	7,504
40	903	43	95,45	7,535
45	901	45	95,24	7,504
Average time for kNN execution				7.532

Fig. 5 presents the part of the correct recognized digits from all the tests for various *K* values.

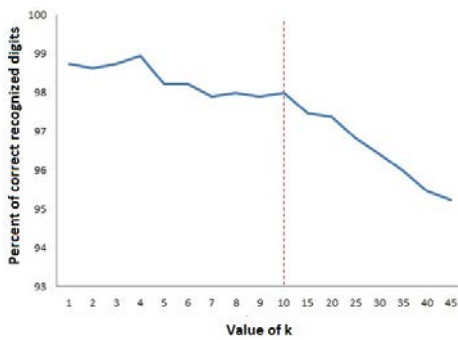


Fig. 5. Part of the correct recognized digits by kNN

The impact of the *K* value on the time for algorithm completion is also important (fig. 6).

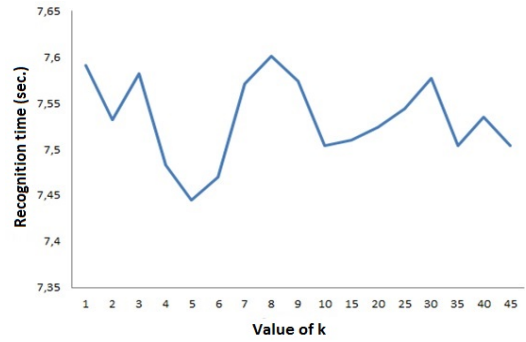


Fig. 6. Recognition time for different *K*

The experimental results show, that the part of the correct recognized digits is greatest for *K* = 4 (98.94%). The recognition time is least for *K* = 5 (7.445 sec.).

This result is explained with the specifics of the kNN algorithm. In the case of *K*=4 (four nearest neighbors) there is high probability for additional vote between two classes with equal number of nearest neighbors, which increases the recognition time. In the case of *K*=5 (five nearest neighbors) the probability for additional vote is less. Consequently the recommended value of *K* in the hand-written digit recognition by kNN is $K \approx \sqrt{N}$, where *N* is the number of classes. The conclusion is, that in order to achieve better recognition time, the *K* value should be odd (table 2).

Table 2. Recommendations for the *K* value in hand-written digit recognition by kNN algorithm

Target criteria	Value of <i>K</i>	Recommended value of k
High part of the correct recognized digits	$K \approx \sqrt{N}$, where <i>N</i> – number of classes (in the given problem <i>N</i> = 10)	<i>K</i> = 4
Less recognition time		<i>K</i> = 5 – odd value

III.b.1.2 PERFORMANCE OF THE SVM ALGORITHM

The SVM algorithm is run with 946 images of digits from 0 to 9. The algorithm is tested with the following kernel functions:

- Linear;
- Polynomial;
- Quadratic;
- Sigmoid;
- Tstudent;
- Wave;
- Gaussian.

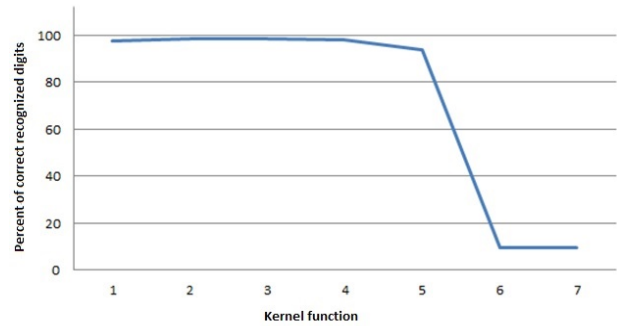
Before the tests the SVM algorithm is learnt by 1934 images of digits with known classes. Table 3 presents the results from the tests.

Table 3. Results from the tests of SVM for kernel functions

Kernel function	№	Number of correct recognized digits	Number of incorrect recognized digits	%	Time for SVM execution (ms)
Linear	1	924	22	97,67	2,525
Polynomial	2	932	14	98,52	2,497
Quadratic	3	932	14	98,52	2,477
Sigmoid	4	929	17	98,2	2,337
Tstudent	5	885	61	93,55	6,037
Wave	6	89	857	9,41	0,925
Gaussian	7	89	857	9,41	5,922
Average time for SVM execution					3,246

The results show the significant impact of the kernel function choice on the SVM performance. The overall percent of the correct recognized digit images is high. The two exclusions are Wave and Gaussian kernel functions. The reason is, that these functions are not suitable for input data in numeric format (the digit images are presented by 0s and 1s). In all other kernel functions the percent of the correct recognized digit images is significant. Its lowest value is 93.55%. The best results vary from 97.67% to 98.52%. The winners are the polynomial and the quadratic kernel functions.

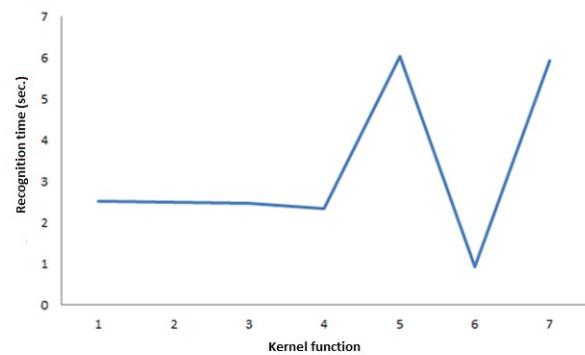
Fig. 7 presents the part of the correct recognized digits from all the tests of SVM for different kernel functions.



1 - Linear 2 - Polynomial 3 - Quadratic 4 - Sigmoid 5 - Tstudent 6 - Wave 7 - Gaussian

Fig. 7. Part of the correct recognized digits by SVM

The time for the algorithm completion for different kernel functions is presented on fig. 8.



1 - Linear 2 - Polynomial 3 - Quadratic 4 - Sigmoid 5 - Tstudent 6 - Wave 7 - Gaussian

Fig. 8. Recognition time for different kernel functions

The kernel functions with the least completion time are linear, polynomial, quadratic and sigmoid.

III.b.1 SECOND STAGE OF THE EXPERIMENTS

On the second stage of the experiments the three variations of the algorithms with best results from the first stage are performed again. Each digit image from the test set is used as an input of the algorithms. The results are the run times of the algorithms for correct recognized digit images. Table 4 presents the results.

Table 4. Comparison of the run times of kNN and SVM

digit	Run time for kNN			Run time for SVM		
	k=4	k=3	k=2	Polynomial	Quadratic	Linear
0	0,221	0,221	0,222	0,215	0,215	0,215
1	0,237	0,225	0,224	0,217	0,213	0,218
2	0,22	0,228	0,224	0,214	0,214	0,215
3	0,223	0,223	0,221	0,215	0,216	0,213
4	0,223	0,221	0,224	0,215	0,216	0,216
5	0,226	0,22	0,223	0,214	0,216	0,212
6	0,242	0,224	0,224	0,214	0,218	0,219
7	0,224	0,222	0,222	0,218	0,218	0,218
8	0,222	0,222	0,234	0,216	0,217	0,218
9	0,222	0,222	0,222	0,216	0,216	0,217

The run times of the algorithms are presented on fig. 9.

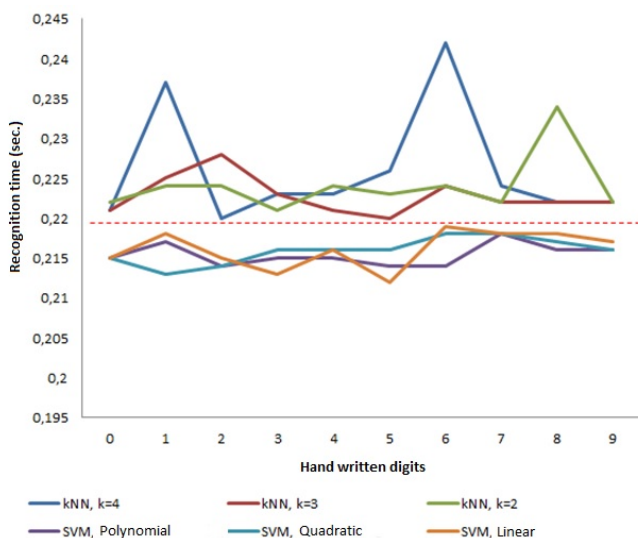


Fig. 9. Run times of kNN and SVM

The upper part of the figure concern kNN algorithm for different values of the K parameter. There are three hops of the run time. All of them are for even values of K ($K=4$, $K=2$). The reason is the necessity of additional vote for classes with equal number of voices to be performed. This result confirms the conclusion that the odd values are recommended for the K parameter in hand-written digit recognition by kNN.

It is obvious, that the run times for the SVM algorithm with polynomial, quadratic and linear kernel function are less than the run time of kNN. Furthermore the curve of the SVM run time is more stable, without sharp increase and decrease. It should be noted, however, that the presented run times do not include the additional training stage, which is an essential part of SVM. The kNN algorithm does not need to be trained in advance.

IV CONCLUSION

The recommendations for the choice of K , summarized in the previous part, combined with the positive features of the kNN algorithm, could increase its performance for solving the hand-written digit recognition. The good recognition time makes the algorithm proper to work in real-time systems, where large data sets should be processed quickly (for example in robots and controllers).

The performance of kNN for the hand-written digit recognition is close to the performance of SVM concerning the correct recognized instances and the algorithm run time. Unlike SVM, however, kNN does not need training in advance.

The future work will focus the following directions:

- Comparison of kNN with other algorithms for hand-written digit recognition;
- Experiments with the kNN in other machine learning problems to confirm or reject the recommendation for K , given in this paper;
- Implementation of kNN for hand-written digit recognition in real-time robot systems.

References:

- [1] A. Ghosh, On optimum choice of k in nearest neighbor classification, *Computational Statistics & Data Analysis*, 2006.
- [2] K. Weinberger, J. Blitzer, L. Saul, Distance metric learning for large margin NN classification, *Advances in neural information processing systems*, 2005, pp.1473-1480.
- [3] P. Harrington, *Machine learning in action*, Manning Publications, ISBN: 9781617290183, 2012.
- [4] Y. Lecun, Comparison of learning algorithms for handwritten digit recognition, *International conference on artificial neural networks*, 1995, pp. 53-60.
- [5] Y. Song, IKNN: Informative k -nearest neighbor pattern classification, *European Conference on Principles of Data Mining and Knowledge Discovery*, Springer Berlin Heidelberg, 2007, pp. 248-264.
- [6] <http://accord-framework.net>