

SVM Classifier for Offline Handwritten and Printed Mathematical Expression Recognition

Manisha Bharambe¹^(⊠), Kavita Kobragade², and Poonam Ponde³

¹ MES Abasaheb Garware College, Karve Road, Pune, India mgb.agc@mespune.in ² Fergusson College, F.C. Road, Pune, India ³ N. Wadia College, Pune, India

Abstract. The recognition of Mathematical Expressions (ME) constitutes a challenging problem in character recognition research. A very few studies of offline Mathematical expressions have been so far reported in the literature. This paper focuses on offline handwritten and printed mathematical logical expressions recognition using Support Vector Machine classifier (SVM). In the work of expression recognition, the expressions were segmented into individual characters. The feature extraction method with combination of Normalized chain code and zone based density was used to get the features of a character. The present work considers logical expressions with subscripts for recognition. The experimental results for recognition rates of handwritten and printed expressions are reported. The result shows that the recognition rate of handwritten expression is 84.1% and that for printed expression is 90.3%.

Keywords: Offline handwritten and printed · Mathematical expression recognition · logical expressions · SVM · recognition rate

1 Introduction

Research into recognizing handwritten Mathematical Expressions (ME) began in the late 1960's. The rate of progress was limited in this problem, with absence of benchmark datasets and standard evaluation tools. The input of ME into computers is more difficult than plain text. MEs are more complex in structures. The complexity occurs in the expression because math symbols can have different fonts such as regular roman, italic, bold and calligraphic, and can have different sizes. Hence, recognition problem becomes more challenging. In the sudden outbreak of COVID-19 pandemic, the application of handwritten interaction in education, remote work, and distance learning is the requirement [19].

In the offline printed mathematical expression recognition, some mathematical symbols might be incorrectly recognized or missed, which leads to the insufficient context information. So it is necessary to make a deeper study of the handwritten and printed mathematical expression recognition. Surendra Ramteke et al. [7] had proposed offline handwritten ME recognition system using ANN classifier. The feature set was obtained by using area and centroid feature of bounding box. The recognition result was 90% of the simple mathematical equation. Sanjay S. Gharde et al. [10] have discussed the various steps of recognition process for simple off-line mathematical equations. The feature extraction methods such as zoning, skeleton based direction, projection histogram, boundary values from four directions (top, bottom, left and right) of symbols, and structural features like crossing points, end points and loops were used. ANN and SVM classifiers were used for recognition. Kazuki Ashida et al. [16] developed the performance evaluation of the mathematical formulae recognition system applied to printed formula images. Dong-Yu Zhang et al. [11] had proposed the method for segmentation of touching symbols in printed Mathematical Expressions. Francisco Álvaro and Richard Zanibbi [14] have proposed problem of classifying special relationship between symbols and expressions in online handwritten ME'S. Hans-Jurgen Winkler et al. [15] have proposed research on online segmentation and recognition of mathematical expressions. Dipak Bage et al. [8] had proposed the offline handwritten mathematical symbol recognition system, and adopted the number of feature extraction techniques to recognize the mathematical symbols.

From the literature it was observed that only few studies of ME recognition were reported. Due to large database of math symbols, the researchers used the subset of math symbols in their works. It was observed that some researchers worked on expression recognition only on simple expressions [7, 10] while others have not specified the class of expression under consideration. The segmentation methods used for MEs were not clearly described in the literatures. The expression recognition rate computed by Sanjay S. Gharde et al., [10] and Surendra Ramteke et al [7] was percentage of the number of symbols recognized from the expression.

2 Stages in the Recognition of ME

The Fig. 1 shows the phases in recognition of ME. The scanned expression was preprocessed by using binarization, filtering, and morphological operations. The preprocessed expression was segmented into isolated characters. The feature set was obtained for each segmented character using statistical, topological and moment features. There are various classification methods such as Artificial Neural Network, SVM, Naive Bayes, decision tree. The SVM classifier is commonly used for the classification problems such as handwriting recognition, face detection, email classification, gene classification, and in web pages. The paper uses SVM classifier for ME recognition.



Fig. 1. Phases of Mathenmatical Expression (ME) recognition system



Fig. 2. Dataset of 12 characters (10 logical symbols and 2 digits)

The recognition was performed on 48 handwritten expressions. Each handwritten expression was written by 20 different writers, resulting in 960 expressions. The total number of characters of handwritten expressions after segmentation was 8640 (80% used for training and 20% used for testing). An experiment was carried out on 35 printed expressions. The number of occurrences of each expression varies from 5 to 15 times with different fonts, results in total 288 printed expressions. A number of 288 printed expressions were segmented into 3559 images (80% used for training and 20% used for testing). The dataset of 26 alphabets (a-z), 10 logical symbols (Fig. 2) and numerals 0 and 1 were used in the expressions. Hence there were 38 class labels.

The handwritten or printed training expressions were preprocessed by using filtering, binarization morphological operations and cropping. The preprocessed expression was segmented which results in the isolated characters. Segmentation is a difficult task in the handwritten ME due to the following reasons:

- ME structure is a 2-Dimentional structure, containing overlapping characters, i.e. the neighboring characters are written such that they share the area of others in a symbol block.
- Touching characters reduces the rate of segmentation.
- The characters in the ME have different spatial relationship like subscripts.



Fig. 3. (a) and (c) Original handwritten images, (b) and (d) Segmented image with BB and Centroid plot.

Two methods of segmentation are used to segment the characters from the image, namely, projection profile cutting and connected component labeling. Projection profile cutting works quite well for one dimensional structures. However, ME is a two dimensional structure that contains symbols having special relationship or one inside other symbol. Hence, Projection profile cutting is not well suited for the ME. The connected component labeling method was used for segmentation of ME [3]. Connected components of the image, Labeled L, were obtained for all the characters in the expression. The values of Bounding Box (BB) were obtained by using properties of image region. Centroid of each BB is found to identify the subscripts which is less than the half of the height of the expression image. Each subscript denoted with 'sb' and other characters with 'up'. Each segmented character had label structure contained three fields: class label, sb (subscript)/up (no subscript), and position. The proposed method is used to classify the characters into 38 different classes of characters (26 alphabets, 10 math symbols, 2 numerals (subscripts). A class label from the set C(yj) \in {1, 2, 3,..., 38} was assigned to each character to recognize the characters in the expression (Fig. 3).

Once the character is segmented, the feature extraction methods: Normalized chain code (16 features), zone based density (26 features), diagonal and intersection points (128 features), Moment invariant features (7 features), Projection Histogram (127 features) were used to get the features of the characters of training expressions [4, 5]. The proposed system was carried out the classification with combinations of these features to get better recognition rate. It was observed that the performance of classifier was high using combination of zone based density features and normalized chain code features with feature vector of size 42. The SVM classifier was used to classify a character of the expression. The expression is reconstructed by using label structure. All the correctly recognized characters in the expression were used for training. The recognition process of test expressions using SVM classifier is shown in the Fig. 4.



Fig. 4. Expression recognition system using SVM

3 Testing Process

The result of classification of each character was a class label from the set of class $C(yj) \in \{1, 2, 3, ..., 38\}$. The true label (label of the character with respect to its class) implies the correct recognition of the character. However, just the features and class labels do not give sufficient information to recognize the expressions. The information about the location of the character and grouping of the characters are also needed in expression recognition. The expression was reconstructed using label and relative position of subscripts and main character. An algorithm for expression recognition is discussed below.

3.1 Algorithm for Expression Recognition

Input: Set of Handwritten or Printed Expressions.

Output: Class labels of the characters in the expressions.

Method:

Step 1. Read scanned handwritten or printed expression.

Step 2. Preprocess the expression image.

Step 3. Segment the expression to separate characters using the segmentation algorithm and obtain the isolated characters and its position information.

Step 4. Apply the feature extraction algorithm to each image to get feature set of each character, and obtain the feature set of an expression.

Step 5: Input the feature set of an expression to the SVM classifier to obtain the recognized characters of an expression.

Step 6: Reconstruct the expression with the help of position numbers of the characters. **Step 7**: Display the recognized expression.

4 Experimental Results

To predict the class of the test characters in the expressions, the unknown pattern vector of the test expressions and the feature set representing training images were fed to classifiers SVM.

Normalized chain code and zone based density method was used to extract the features of segmented characters from the test expression. A recognition accuracy of the SVM classifier was calculated. The expression recognition rate (ERR) was computed as percentage of correctly recognized expressions from the set of same expressions.

Define:

CR = Number of correctly recognized expressions, out of correctly segmented

TE = Total number of occurrences of a single expression

Then, for each expression the recognition rate was computed as,

$$\text{ERR} = \frac{CR}{TE} * 100$$

The average recognition rate of handwritten expressions and printed expressions was 84.1% and 90.3% respectively. It was observed that the recognition rate of printed expressions was higher than that of handwritten expressions. The experiment was performed using Matlab 8.1. The recognition of some expressions using SVM is shown in the Fig. 5 and Fig. 6. The numeral values in the output shows the class labels of the characters. The class labels for letters 'a' to 'z', 10 math symbols, numerals '0' and '1' were assigned as 1 to 26, 27 to 36, 37, and 38 respectively. The error rate in output defines the percentage of misclassified characters. The resultant recognized expression is shown as the output in the command window after reconstructing the expression using position of the characters (Tables 1, 2, 3).

#Exp	Handwritten Expressions	HERR
		(%)
1	a, vb,	90
2	(va, nb)vc,	85
3	$(\sim P, \vee \sim q)$	90
4	$\neg a, \Rightarrow b_1$	75
5	R, VY,	90
6	a.v.b.	85
7	~a,vb,	80
8	$P_{n} \land P_{n} \Rightarrow 9$	75
9	$m_r \Rightarrow m_r \land m_s$	80
10	aVa.⇒o	70
11	$f(g) \Rightarrow a_0 \wedge a_1$	75
12	(a, b, b)	85
13	X. V Y.	70
14	~ OVL >C	75
15	QuAbi	80
16	$(P \land q) \rightarrow (\sim P)$	90
17		100
18	Tarb	95
19	NNavab	90
20	$P \land (r \rightarrow s)$	100
21	NPA(QVC)	90
22	$P \rightarrow (P \lor q)$	85
23	(PAt)V(CANA)	85
24	$S \rightarrow (r v t)$	90
25	(NPAR)A 19 - PR	85
26	$(PVS) \rightarrow (q, AV)$	85
27	((PAr)V(NQANY))V(NPANY)-31	70
28	$(P \rightarrow q) \land (q \rightarrow p) \rightarrow r$	95
29	$((P \rightarrow Q) \land Q) \rightarrow P$	90
30	N(PANY) V(NQ.VS)	80
31	(PAO) ~ (~PV~q)	75

Table 1. Handwritten Expression Recognition Rate (HERR) of handwritten test expressions

(continued)

32	$P \rightarrow (P \vee (P \vee \neg 9))$	80
33	(lavb)v (marc)	75
34	$\alpha(\alpha P) \Rightarrow P$	85
35	2レ1=>1	85
36	an (bvc) - an 1	85
37	(avb)AC	90
38	andra	95
39	$\sim \alpha \Rightarrow b$	85
40	~~aV~b^~c	95
41	$i \vee j$	85
42	i @ j	85
43	$P \oplus q \rightarrow (P \vee q) \land (\neg q \vee P)$	85
44	$P \oplus 2 \rightarrow \neg ((P \land 2) \lor (\neg P \land \neg 2))$	75
45	a (b (c () d () e	80
46	$(z \oplus y) \rightarrow (z \oplus y)$	85
47	$f(a \oplus b) \Rightarrow f(a) \lor f(b)$	75
48	$!(x \vee Y) \rightarrow \neg x \vee \neg Y$	80

 Table 1. (continued)

From the experiments, it was observed that the recognition errors were due to misclassification characters of similar shape in handwritten expression recognition. The confusing pairs of letter 'v' and symbol 'V'; 'c' and '(' gives the incorrect recognition. The correct classification of handwritten expressions using SVM classifier was shown in the Fig. 5. The recognition rate of the expressions was poor due to the misclassified character 'v' in the Fig. 6. The resultant output of classification is consists of class label 22 represent character 'v' instead of symbol 'or'.

#Exp	Printed Expressions	PERR
	78	(%)
1	$a \vee (b \wedge c)$	100
2	$(a \vee b) \wedge a$	100
3	$\mathbf{x}_0 \vee \mathbf{y}_0$	90
4	$(\mathbf{p} \rightarrow \mathbf{q}) \land \mathbf{q} \rightarrow \mathbf{p}$	90.9
5	$(a \lor b) \land (a \lor 0)$	83.33
6	$\sim (\mathbf{p} \rightarrow \sim \mathbf{q})$	80
7	$a_1 \wedge b_1 \rightarrow 1$	90
8	$p \wedge (r \rightarrow s)$	100
9	$(a \lor 1) \land (a \lor c)$	87.5
10	\sim (p \land \sim r) \lor (\sim q \lor s)	90
11	$(p \lor s) \rightarrow (q \land r)$	100
12	$(p \lor \sim q) \rightarrow (r \land p)$	86.6
13	$(\mathbf{q} \lor \sim \mathbf{r}) \land (\mathbf{p} \lor \mathbf{q})$	90
14	$\mathbf{p} \rightarrow \mathbf{q} \rightarrow (\mathbf{p} \land \mathbf{q}) \lor (\mathbf{p} \land \mathbf{q})$	70
15	$(\mathbf{p} \wedge \mathbf{q} \wedge \mathbf{r}) \vee (\sim \mathbf{p} \vee (\mathbf{q} \wedge \sim \mathbf{r}))$	80
16	$(\mathbf{p} \wedge \mathbf{r}) \vee (\sim \mathbf{q} \wedge \sim \mathbf{r}) \vee (\sim \mathbf{p} \wedge \sim \mathbf{r})$	75
17	$f(x_0) \Rightarrow g(f(x_1))$	80
18	$p \oplus q \stackrel{\textbf{\textbf{-}}}{\textbf{\textbf{-}}} ((p \lor \neg p) \land (\neg q \lor \neg p)) \land ((p \lor q) \land (\neg q \lor q))$	q ^{`70}
19	$p \Rightarrow q \land q \Rightarrow p$	100
20	$\sim p \implies \sim q$	100
21	$f(a \oplus b) \Rightarrow f(a) \lor f(b)$	90
22	$a \oplus b$	100
23	$(\mathbf{p} \rightarrow \mathbf{q}) \land \neg (\mathbf{p} \rightarrow \neg \mathbf{q})$	85.7
24	$(p \land t) \lor (c \land \neg q)$	100
25	!a∧b	100

 Table 2. Printed expression Recognition rate (PERR) of printed test expressions.

(continued)

26	$(p \to r) \lor (\neg s \to \neg t) \lor (\neg u \to v)$	80
27	$((p \to (q \to r)) \to ((p \to q) \to (p \to r))$	80
28	$(\mathbf{p} \wedge \mathbf{q}) \rightarrow \sim (\mathbf{p} \rightarrow \sim \mathbf{q})$	85.7
29	$(p \lor \neg q \lor \neg r) \land p \lor (q \land r)$	87.5
30	$\sim \sim a v \sim b$	100
31	$p \oplus q \rightarrow \neg((p \land q) \lor (\neg p \land \neg q))$	100
32	$((\underline{p} \rightarrow q) \land (q \rightarrow r)) \rightarrow (p \rightarrow r)$	100
33	$(\neg q \land (p \rightarrow q)) \rightarrow \neg p$	100
34	$x \oplus (y \oplus z) \Rightarrow (x \oplus y) \oplus z$	100
35	$(p \oplus q) \to (p \oplus \neg q)$	100

 Table 2. (continued)

Table 3. Average Expression Recognition Rate

Expression Type	#exp	Wrongly Segmented	Wrongly Recognized	Correctly recognized	ERR	
Handwritten	960	105	48	807	84.1	
Printed	288	20	8	260	90.3	

In case of printed expression recognition, normally, the expressions were not correctly segmented due to touching characters. The correctly segmented expressions results in the high accuracy in recognition step. Most of the errors in the recognition occurs due to misclassified characters 'a' with ')'. The correct classifications of printed expression is shown in the Fig. 7.



Fig. 5. Recognition of handwritten expression (a, v b)

Co	ommand	Windo	w													\odot
	ypred	-														*
	3	D	34	16	28	30	18	35	22	34	30	17	29	19	35	
	erro :	-														н
	0	.071	4													
6	error	rate	e= 0.	071429)											
JX.	~	(p	^	~ r)	V (~	q V	S)>>					*

Fig. 6. Recognition of handwritten expression $\sim (P \land \sim \gamma) \lor (\sim q \lor s)$

Co	omman	d Winc	low													۲
	ypre	d =														^
		6	34	24	37	35	33	7	34	6	34	24	38	35	35	
	erro	-														н
	erro	0 r rat	e= 0.	00000	0											
fx.		f	(x	0) =	> g	(f	(x 1))>>					-

Fig. 7. Recognition of printed expression $f(x_0) \Rightarrow g(f(x_1))$

5 Conclusion

SVM classifier was used to classify the characters in the expressions. The accuracy of the expression was computed. The average recognition rate of handwritten was 84.1%

and printed expression was 90.3%. The recognition rate of printed expression was higher than that of handwritten expression.

References

- Ahmad M., Harold Mouchere, Christian Viard Gaudin, (2009), "Towards Handwritten Math ematicalExpression recognition", International Conference on Document Analysis and Recognition (ICDAR), IEEE 978-0-7095, pp:1046–50.
- Ahmad Montaser Awal, Harold Mouchere, Christian Viard Gaudin, The problem of handwritten mathematical expression recognition, International Conference on Frontiers in Handwriting Recognition (ICOFHR), ISBN, 978-0-7695-4221-8.
- Manisha Bharambe, "Segmentation of Offline Printed and Handwritten Mathematical Expressions" International Journal of Computer Applications (0975–8887) National Conference on Digital Image and Signal Processing 2016
- Manisha Bharambe, "Recognition of Offline Handwritten Mathematical Expressions" International Journal of Computer Applications (0975–8887) National conference on Digital Image and Signal Processing, DISP 2016
- Manisha Bharambe "Logical Symbol Recognition using Normalized Chain code and Density Features", International Journal of Engineering Research & Technology (IJERT), ISSN: 2278–0181, Vol. 3, Issue-12, December 2014, pp: 619–623
- 6. Kang kim, Taik Rhee, Jae LEE, (2009), "Utilizing consistency context for handwritten mathematical expression recognition", ICDAR, 978-0-7695-3725-2, IEEE, pp: 1051–1055.
- 7. Surendra Ramteke, Dhanashree Patil, Nilima Patil, (2012), "Neural network Approach to Mathematical Expression Recognition System", International Journal of Engineering Research & Technology (IJERT), Vol. 1 Issue 10, December-2012, ISSN: 2278–0181.
- Dipak D.Bage, K.P. Aditya, Sanjay Gharde, (2013), "A new approach for recognizing offline handwritten mathematical symbols using character geometry", International Journal of Innovative research in science, engineering and Technology, vol. 2, Issue 7, ISSN: 2319–8753, pp: 2823–2830.
- H. Mouchere, C. Viard-Gaudin, D. H. Kim, J. H. Kim, U Garain, (2012), "ICFHR 2012-Competition on Recognition of online mathematical Expressions", (CROHME 2012), 078-0-7695-4774-9/12, IEEE.
- Sanjay S. Gharde, Baviskar Pallavi V, K. P. Adhiya, (2013), "Identification of Handwritten Simple Mathematical Equation Based on SVM and Projection Histogram", International Journal of Soft Computing and Engineering (IJSCE) ISSN: 2231–2307, Volume-3, Issue-2, pp: 425–430.
- Dong-Yu Zhang, Xue-Dong Tian, and Xin-Fu Li, (2010), "An Improved Method for Segmentation of Touching Symbols in Printed Mathematical Expressions", International Conference on advance Computer Control, IEEE, ISBN: 978-1-4244-5848-6, pp: 251–253.
- Ernesto Tapia and Raul Rojas, (2004), "Recognition of On-line Handwritten Mathematical Expressions Using a Minimum Spanning Tree Construction and Symbol Dominance", Llados and Y.B.Kwon(Eds.): GREC 2003, LNCS 3088, pp: 329–340, 2004, Springer.
- Ernesto Tapia and Raul Rojas, (2005), "Recognition of On-Line Handwritten Mathematical Expressions in the E-Chalk System - An Extension", Proceedings of the 2005 Eight International Conference on Document Analysis and Recognition (ICDAR'05), 1520–5263/05, IEEE
- Francisco Álvaro, Richard Zanibbi, (2013), "A Shape-Based Layout Descriptor for Classifying Spatial Relationships in Handwritten Math", DocEng'13, ACM 978-1-4503-1789-4/13/09

- Hans Jurgen Winkler and Manfred Lang, (1997), "On-Line symbol segmentation and recognition in handwritten mathematical expressions", International conference on Acoustics, Speech and Signal Processing, ICASSP, 0-8186-7919-0/97, IEEE.
- Kazuki Ashida, Masayuki Okamoto, Hiroki Imai, (2006), "Performance Evaluation of a Mathematical Formula Recognition System with a large scale of printed formula images", Proceedings of the Second International Conference on Document Image Analysis for Libraries (DIAL'06), 0-7695-2531-8/06, IEEE
- Lei Gao, Shulin Pan, Shen Jiao, (2013), "An Analytic Hierarchy Process Based Method to Process Mathematical Expressions", Journal Of Theoretical And Applied Information Technology, Vol. 47 No.3, ISSN: 1992–8645.
- Fanglin Wang , Fukeng He , Ning Bia, Ching Y Suen , Jun T., Offline Handwritten Mathematical Expression Recognition Based on Square Loss Sequence-to-Sequence Neural Network, https://doi.org/10.2139/ssrn.4212220
- Dmytro Z., Viktor Z., And Olga R, Online Handwritten Mathematical Expression Recognition and Applications: A Survey, IEEE, Volume 9, 2021, https://doi.org/10.1109/ACCESS2021. 3063413.
- Riddhi A, Shila P, Anilkumar T, Gaurav H, Survey of Mathematical Expression Recognition for Printed and Handwritten Documents, https://doi.org/10.1080/02564602.2021.2008277
- 21. Yuging W, Zhengu W, Zhaokun Z, Shuajian Z., Yuesheng Z., Dual Branch Network Toswards Accurate Printed Mathematical Expression Recognition, ICANN (2022), pp 594–606
- Lyzandra D'souza, Maruska ascarenhas, Offline Handwritten Mathematical Expression Recognition using Convolutional Neural Network, IEEE *Xplore*: 15 November 2018, https:// doi.org/10.1109/ICICET.2018.8533789

Open Access This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (http://creativecommons.org/licenses/by-nc/4.0/), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

