

ARABIC CHARACTER RECOGNITION USING LEARNING VECTOR QUANTIZATION

Keumala Anggraini¹, Lestari Handayani²

Informatics Engineering Department
Faculty of Science and Technology
Universitas Islam Negeri Sultan Syarif Kasim Riau
Pekanbaru - Indonesia

Email :¹keumalaanggraini93@gmail.com, ²lestari.handayani@uin-suska.ac.id

ABSTRAK – Arabic character recognition should be research again. Arabic character have 28 characters with 4 different positions in sentence, then Arabic character has 28-100 characters. The method is used for Arabic character recognition is learning vector quantization neural network. It is because, the learning vector quantization could classify input in category defined on training network. The objective of this study is to testing LVQ method in Arabic characters recognition. The experiment conducted using all types' position of character in sentence, there are isolated, begin, middle, and end. The testing data of Arabic character passed preprocessing phase to get vector number that was the size of matrix is used as input for learning vector quantization. The size of matrix was 8x10 for isolated, middle; end and 7x12 for begin. The success accuracy rate for isolated was 76, 43%, begin was 65, 45%, middle was 62, 73%, and end was 80%. The success accuracy percentage for all Arabic character was 72%.

Keywords: Arabic Character, Character Recognition, Learning Vector Quantization, Size of Matrix, Success Accuracy.

I. INTRODUCTION

Pattern recognition is a technique used to classify an object by main feature or main properties [1]. Pattern recognition done by identifying pattern in an object on one of group. Pattern recognition can be done on characters of letters.

A letter is smallest unit or smallest information from a sentence that needs to be defined so that information on a sentence could understand. A letter have different form between one another, to distinguish it adapted characteristic one of the letters [2]. The letters have differences pronunciation, although the letters have same characteristic. One of the letters has difference characteristic and pronunciation is Arabic character.

Arabic characters have 28 difference characters between one another. Arabic characters have difference characters in accordance with position on a sentence. Parts of Arabic characters have 4 form differences position on sentences, another part of Arabic characters have 2 form differences position on sentences. Therefore, Arabic characters have 28-100 characters, with addition of changes

in the form characters. Difference characters will change if the characters on isolated, begin, middle and end [3].

Learning vector quantization (LVQ) was introduced by Tuevo Kohonen, who also introduced Kohonen method. LVQ is one of artificial neural network and supervised network. LVQ classify input in category that has been defined through supervised network (Putra, 2010). On research which title Difference between Kohonen Neural Network and Learning Vector Quantization on The Real Time Handwritten System have purpose that recognition with LVQ better than Kohonen in terms of accuracy.

Another research entitle Analysis and Implementation of the Kohonen Neural Network for Arabic Character Recognition [5], featuring error result of experiment kohonen neural network. On that research, the error result was 43, 64% on begin of Arabic characters. In this study, used LVQ method decrease the error of Arabic character recognition. Based on [5] and this study same in Arabic characters, there were begin, middle, end, and isolated. But different in site of matrix image of Arabic characters. In this study found best practice for matrix image have various size. It was 8x10 for isolated, middle, and end, 7x12 for begin of Arabic characters.

II. CHARACTER RECOGNITION

Almost all type of character recognition use the same phases to recognition the characters. Generally, procedure to character recognition is used by many researches [5].

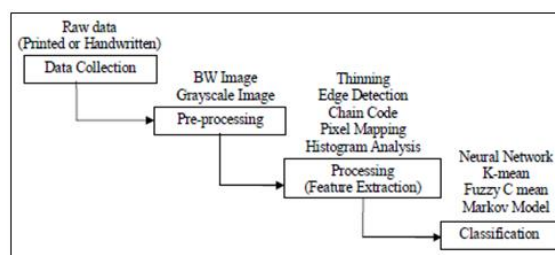


Fig 1. Character Recognition Phases [5]

On this research, Arabic character recognition using learning vector quantization phases in figure 2:

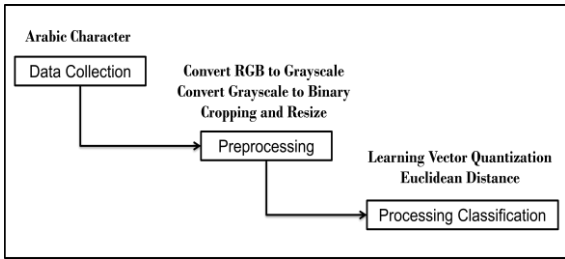


Fig 2. Arabic Character Recognition using Learning Vector Quantization Phases

The necessary data collection is handwritten Arabic character on canvas area on system and the character stored in image. Next, the processing phase, the step is cropping, conversion image to gray scale type, conversion into binary and resize. Next, the processing classification phase is process used learning vector quantization and Euclidean distance.

II. ARABIC CHARACTER

Arabic character's features [6]:

1. Arabic character have 28 characters and written right to left.
2. Arabic character didn't use different between uppercase and lowercase.
3. The width of each other is different.
4. Arabic characters have 4 difference forms according to isolated, begin, middle and end.

Isolated	End	Middle	Beginning	FEAF	FEB0	FED5	FED6	FED8	FED7
FE83 أ	FE84 آ			ز	ز	ق	ق	ق	ق
FE8F ب	FE90 ب	FE92 ب	FE91 ب	FE81 س	FE82 س	FE84 س	FE83 س	FE8A ك	FE8C ك
FE95 ت	FE96 ت	FE98 ت	FE97 ت	FE85 ش	FE86 ش	FE89 ش	FE87 ش	FE8D ل	FE8E ل
FE99 ث	FE8A ث	FE8C ث	FE8B ث	FE89 ص	FE8A ص	FE8C ص	FE8B ص	FE81 م	FE82 م
FE8D ج	FE8E ج	FE89 ج	FE8F ج	FE8D ض	FE8E ض	FE89 ض	FE8F ض	FE85 ن	FE86 ن
FEA1 ح	FEA2 ح	FEA4 ح	FEA3 ح	FE81 ط	FE82 ط	FE84 ط	FE83 ط	FE89 ه	FE8A ه
FEA5 خ	FEA6 خ	FEA8 خ	FEA7 خ	FE85 ظ	FE86 ظ	FE89 ظ	FE87 ظ	FE8D و	FE8E و
FEA8 د	FEAA د			FE89 ع	FE8A ع	FE8C ع	FE8B ع	FE81 ي	FE82 ي
FEAD ذ	FEAC ذ			FE8D غ	FE8E غ	FE89 غ	FE8F غ	FE85 ف	FE86 ف
FEAD ر	FEAE ر			FE81 ف	FE82 ف	FE84 ف	FE83 ف		

Fig 3. Arabic Character [5]

III. ANALYSIS

The analysis is the process analysis. The process analysis is analysis steps and process recognition using learning vector quantization, there are data collection, pre-processing and processing classification.

3.1 Data Collection

Data collection is the input which will be incorporated to do recognition. Input stored in image. The input is Arabic character.

3.2 Pre-processing

Pre-processing is the process change colors in image and separation between object and character. Pre-processing phase is conversion to gray scale type, conversion into binary, cropping and resize.

Conversion to Grayscale

Calculation to obtain the value for a pixel gray scale is as:

$$\text{Grayscale} = (0,2989 * R) + (0,5870 * G) + (0,1141 * B) \dots\dots\dots (1)$$

Conversion into Binary

Calculation to obtain the value for a pixel gray scale is as:

$$g(x,y) = \begin{cases} 1, & \text{jika } f(x,y) \geq T \\ 0, & \text{jika } f(x,y) < T \end{cases} \dots\dots\dots (2)$$

With $g(x,y)$ is the binary from grayscale type and T is middle value.

Cropping and Resize

Cropping step is separation process object and background in image. Resize step is conversion process size of image matrix into $m*n$ size of matrix. On resize step, the used matrix is $8x10$ for isolated; middle and end, $7x12$ for begin.



Fig 4. Character Sheen (ش)

Sheen (ش) has passed cropping step, conversion to grayscale and conversion into binary, sheen (ش) must pass resize step. Size of matrix sheen (ش) is $8x10$, it's because, and sheen (ش) is isolated.

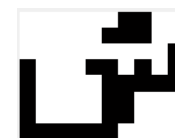


Fig 5. $8x10$ of Sheen (ش)

Binary representation from resize for sheen is:

Table 1. Binary of 8x10 Sheen (ش)

1	1	1	1	1	1	0	0	1	1
[1,1]	[1,2]	[1,3]	[1,4]	[1,5]	[1,6]	[1,7]	[1,8]	[1,9]	[1,10]
1	1	1	1	1	0	0	0	1	1
[2,1]	[2,2]	[2,3]	[2,4]	[2,5]	[2,6]	[2,7]	[2,8]	[2,9]	[2,10]
1	1	1	1	1	1	1	1	1	0
[3,1]	[3,2]	[3,3]	[3,4]	[3,5]	[3,6]	[3,7]	[3,8]	[3,9]	[3,10]
0	1	1	1	0	0	1	0	1	0
[4,1]	[4,2]	[4,3]	[4,4]	[4,5]	[4,6]	[4,7]	[4,8]	[4,9]	[4,10]
0	1	1	1	1	0	0	0	0	0
[5,1]	[5,2]	[5,3]	[5,4]	[5,5]	[5,6]	[5,7]	[5,8]	[5,9]	[5,10]
0	1	1	1	1	0	0	1	1	1
[6,1]	[6,2]	[6,3]	[6,4]	[6,5]	[6,6]	[6,7]	[6,8]	[6,9]	[6,10]
0	1	1	1	1	0	0	1	1	1
[7,1]	[7,2]	[7,3]	[7,4]	[7,5]	[7,6]	[7,7]	[7,8]	[7,9]	[7,10]
0	0	0	0	0	0	1	1	1	1
[8,1]	[8,2]	[8,3]	[8,4]	[8,5]	[8,6]	[8,7]	[8,8]	[8,9]	[8,10]

3.3 Processing Classification

Processing classification is learning vector quantization and Euclidean distance.

Learning Vector Quantization

Learning Vector Quantization (LVQ) is supervised network which training of competitive layer. Competitive layer will learn automatically to classify input [7].

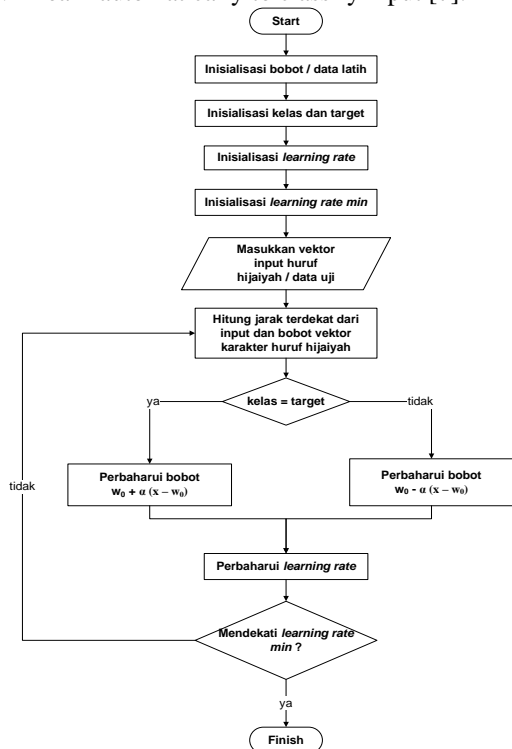


Fig 6. Flowchart of Learning Vector Quantization for Arabic Character Recognition

LVQ algorithm [1]:

LVQ Algorithm LVQ is [1] :

1. The first is determined each output, width and learning rate.
2. Compare each input and each width b measuring the distance between width (w_0) and input (x_p).
3. Distance measurement is calculated using Euclidean distance.
4. Minimum value from comparison result will determine class from vector input and change in weight from that class vector input.
5. Changes in new weight calculated.

Calculation for input and width has same class.

$$w_0' = w_0 + \alpha(x - w_0) \dots\dots\dots (3)$$

Calculation for input and width has same class.

$$w_0' = w_0 - \alpha(x - w_0) \dots\dots\dots (4)$$

Euclidean Distance

Euclidean distance is done to get smallest distance from new width used LVQ. Smallest distance from result will get pattern resemble input. Calculation Euclidean distance is:

$$j(v_1, v_2) = \sqrt{\sum_{k=1}^N (v_1(k) - v_2(k))^2} \dots\dots\dots (5)$$

IV. IMPLEMENTATION AND RESULT

Learning Vector Quantization for Arabic character recognition implemented by processor Intel core i3 2, 4 GHz, RAM 2GB, hard disk 320 GB, operation system Windows 7, and programming Matlab 2010a, which used with not network.

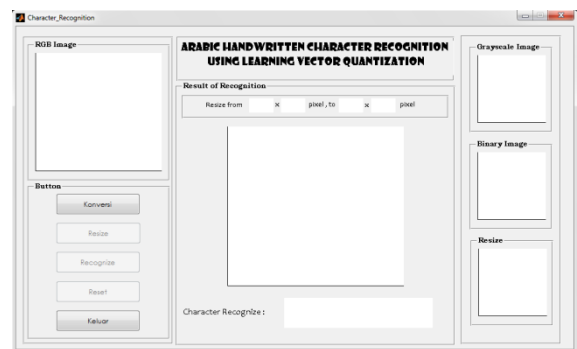


Fig 6. Character Recognition Interface

Character will be drawing on canvas interface. Character which drawing on canvas is Sheen (ش). Sheen (ش) will pass pre-processing and processing classification.

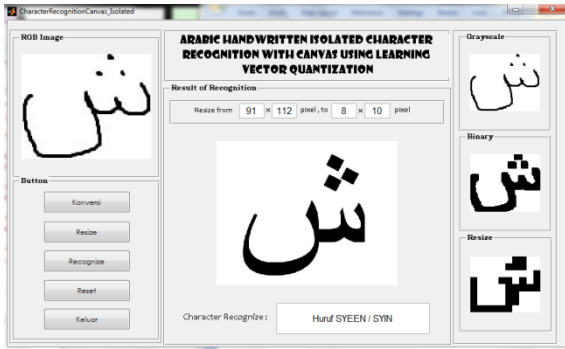


Fig 7. Interface Character Recognition Sheen (ش)

Each Arabic character will be tested. The tested character is isolated, begin, middle and end. The experiment is Arabic character will be tested by 5 characters each Arabic character.

Table 2. The Experiment and Result Arabic Character Recognition

Character Type	Character	Tested Character	Result	Percentage (%)
Isolated	ا	ا	ا	100
		ا	ا	
		ا	ا	
		ا	ا	
		ا	ا	

	ي	ي	ي	60
		ي	ي	
		ي	ي	
		ي	ي	
ي		ي		
Begin	ب	ب	ب	100
		ب	ب	
		ب	ب	

Character Type	Character	Tested Character	Result	Percentage (%)	
Middle	ب	ب	ب	80	
		ب	ب		
			
		ب	ب		
		ب	ب		
	...	ب	ب	ب	100
			ب	ب	
			ب	ب	
			ب	ب	
			ب	ب	
...	ب	ب	ب	20	
		ب	ب		
		ب	ب		
		ب	ب		
		ب	ب		
End	ا	ا	ا	100	
		ا	ا		
		ا	ا		

Character Type	Character	Tested Character	Result	Percentage (%)
		أ	ا	
		ا	ا	
	ي	ي	ي	40
		ي	ج	
		ي	ي	
		ي	ف	
		ي	ف	
	⋮	⋮	⋮	⋮

The experiment used size of matrix tested, and then percentage calculation of success and error from the experiment is calculated. The calculation for success percentage:

$$\frac{\text{Number of Success}}{\text{Number of Tested}} \times 100\% \dots\dots\dots (5)$$

And the calculation for error percentage:

$$\text{Error} = 100 \% - \text{Success Percentage} \dots\dots\dots (6)$$

Table 3. Percentage of Result Arabic Character Recognition

Character Type	Number of Tested Character	Number of recognize Character	Percentage (%)	Error (%)
Isolated	140	107	76,43	23,57
Begin	110	72	65,45	34,55
Middle	110	69	62,73	37,27
End	140	112	80	20
Total	500	360		
Average			72	28

V. CONCLUSION

1. Size of matrix is 8x10 for isolated, middle, end and 7x12 for begin.
2. Learning Vector Quantization can recognize with percentage rate of success is 76, 43 for isolated, 65,

45 for begin, 62, 73 for middle, 80% for end. An average success is 72% for 500 tested Arabic characters.

3. LVQ method for classification in Arabic characters more less error than [5]. Shown from table 3 with average error 28%.
4. The high percentage from every tested percentage is end character. It's because, end character have clear differences each other.

REFERENCES

- [1] Putra, Darma. *Pengolahan Citra Digital*. Yogyakarta: Penerbit Andi. 2010.
- [2] P, Tacbir Hendro, Agus Komaruddin, dan Dila Fadhilah. "Pengenalan Pola Huruf Arab menggunakan Jaringan Saraf Tiruan dengan Metode Backpropagation.", *Seminar Nasional Teknologi Informasi dan Komunikasi (SNASTIKOM 2012)*. 2012.
- [3] Batawi, Yusof A, dan Osama A Abulnaja. "Accuracy Evaluation of Arabic Optical Character Recognition Voting Technique : Experimental Study ." *International Journal of Electrical & Computer Sciences IJECS-IJENS Vol: 12 No: 01*. 2012.
- [4] Asworo. "Perbandingan antara Metode Kohonen Neural Network dan Learning Vector Quantization pada Sistem Pengenalan Tulisan Tangan secara Real Time.". 2010.
- [5] Handayani, Lestari, Iwan Iskandar, dan Weli Andrian.. "Analysis and Implementation of the Kohonen Neural Network for Arabic Character Recognition ." *ICoSTechS | http://www.icostechs.org*. 2014.
- [6] AbdelRaouf, Ashraf, Colin A Higgins, dan Mahmoud Khalil. "A Database for Arabic Printed Character Recognition." *A. Campilho and M. Kamel (Eds.): ICIAR 2008*. 2008.
- [7] Kusumadewi, Sri. *Membangun Jaringan Syaraf Tiruan Menggunakan Matlab & Excel Link*. Yogyakarta: Graha Ilmu. 2004