A Nonhierarchical Neural Network Approach For Analyzing Textual Data

Brenda L. Battleson, Hao Chen, Carolyn Evans, Joseph Woelfel Department of Communication, University at Buffalo (SUNY), Buffalo, New York 14260 USA

Introduction

Artificial neural networks excel at recognizing patterns in textual data. Its pattern recognition capability allows a neural network engine to assign weights representing the multiple connections among concepts. These weights can then be used to create *dendograms* or otherwise categorize concepts in a hierarchical manner. However this approach has its limitations.

Words often have different meanings depending on the context in which they occur. Yet a hierarchical clustering method is unable to fully describe multiple relationships because it is only able to show concepts connected in one way. Each concept is assigned to only one "best" cluster in the output suggesting

way har toncept is assigned to only one test critical in the optimal suggesting that there is only one meaning of that concept in the data analyzed. The use of a nonhierarchical approach can address this limitation since it allows the researcher to interact with the neural network to explore all possible meanings of a concept. Thus, in the resulting output a concept may appear in as many clusters as are appropriate. In this study a large dataset containing multiple newspaper articles, is examined. Hierarchical and nonhierarchical procedures are compared.

Method

Opinions about the terrorist attacks of September 11, 2001 were of particular opinion and the theorem and the second opportunity of the event in September 2006. To gauge opinion, editorials, opinion pieces and letters to the editors of all U.S. newspapers indexed in the FACTIVA™ database were retrieved for the month of September 2006

The 3.2006 text file was analyzed using the CATPAC[™] text analysis program. Output consisted of scalar products matrix used to generate an artificial neural network (ANN) with output consisting of a weighted input network (WIN) file and the hierarchical clusters represented in both dendogram and 3-dimensional coordinate files. The ORESME ™ software was then used for nonhierarchical analysis of the CATPAC results.

Infinite activities of the CATPAC results. Unlike the traditional forward feed-back propagation neural networks, ORESME¹ is an interactive activation and competition network, and any neuron can be an input, hidden or output neuron. According to Woelfel, the most appropriate representation of ORESME's schematic is Patrick Lemmens' nterpretation of a neural network (Fig.1A.)



In this study, the researcher assigned an input (activation) value to one or more terms and the resulting clusters of "activated concepts" in the ANN were compared. The nodes of an ANN are connected to each other by weights which represent their relative "closeness" in the network. They communicate with each other by a simple linear threshold rule:

The signal sent from any node i to any node j equals the product of the activation value of i and strength of the connection between i and j. Thus the total signal received by any node j will be the sum of the signals received from all the other nodes, or

anet: - $\sum_{i=1}^{N} W_{ij}a_j$

1 CATPAC™ and ORESME™ are PC versions of KAMANDU™ and LISTIAC™ (respectively) both of which are running on mainframe computers at the University at Buffalo's Center for Computational Research.

Results

HIERARCHICAL CLUSTERING The dendogram and 3-dimensional map resulting from the CATPAC analysis show definite clusters:



Figure 2. Dendogram of 2006 op-ed pieces in U.S newspapers for the month of Sentember 2008

Figure 3. ThoughtView 3D map of hierarchical clusters identified using CATPAC.

NONHIERARCHICAL CLUSTERING: ORESME allows the researcher to input or "activate" the neuron representing a concept of concepts. Multiple cycles allow for "learning" and weight adjustment of associated neurons. A threshold value determined by the researcher determines what associated neurons are activated. The results are new and often very different clusters of concepts.

Cure DCP1	e lab	1.00	el a	when done					
ant	these	value	s cl	unge d?					
100									
				TORTURE			detivation.	level	1.888
				TUSH			Activation	10001	005
	8ha11	I thi	ek i	VITHOUT over one			Activation	level -	.832
				LIE			Activation	level -	.083
				TORIURE			Activation	level	1.000
				VITEOUT			Activation	Terre 1	841
	8ha11	I thi	nk i	over one		time?			
				LIE			Activation	level	.083
				TORIURE			Ectivation	level	1.000
				VITEOUT			Activation	lens 1	
	8ha11	I thi	nk i						
				1.110			Activation	level -	.083
				TORIURE			Activation	level	1.000
				BUSH			HOULWALLON	Tenel	885
	Shall	I thi	nk i	over one		time?	85 C 10 X C 101	10401	
				LIE			Activation	level -	.083
				TOPTURE				Terra T.	- T 000
				TOHLORE			RCCIUATION	10001	1.000
	ncey ant lea t. t.	ncept lobe turn ncept lobe ant these leave? t. Shall t. Shall t. Shall	neeps label.GC meeps label.GC ant these value these value c. Shall thi c. Shall thi c. Shall thi	neege label.(CErl a merge label.(CErl a merge label.(CErl a mat there walnes clu- lawr? t. Eball thisk if t. Eball thisk if t. Eball thisk if	eeges label. (Gr.) a chee door menyt label. (Gr.) a chee door set these seales closes? lases? t. Eball (thick I virian? t. Eball (thick I virian?)	Anger 1341-(197) 2 when dawn) meryd 1341-(197) 2 when diwn) anger 1341-(197) 2 when diwn anger 1341-(197) 2 when diwn ange	energy label.(Cr) i when down) energy label.(Cr) i when down) energy label.(Cr) i when down) isource isource isource to final i thick it of the second second to the second second second second second to the second second second second second second to the second sec	anger khol.(Cri 2 when daw) meret halt.(Cri 2 when daw) meret halt.(Cri 2 when daw) term of the when anged term of the when anged term of the when anged term of the week anged term of term of the week anged term of term of	Anger Label. (Cri 2 when dams) Marger Label. (Cri 2 when dams) Marger Label. (Cri 2 when dams) Marger Label. (Cri 2 when dams) A when a state of the state of the state of the state to the state of the state of the state of the state to the state of the state of the state of the state to the state of the state of the state of the state to the state of the state of the state of the state to the state of the state of the state of the state of the state to the state of the st

Figure 3. ORESME run showing results of analysis of the concept "TORTURE



Figure 4. Output of ORESME run showing the activation levels of all concepts when an input activation level of 1.0 is assigned to the concept TORTURE. Those concepts exceeding the threshold 0.000 are activated. Terms clustering with TORTURE are LIE. BUSH and WITHOUT







Conclusions

The human brain is the most sophisticated example of a parallel distributed processing machine. The language used to express human ideas, attitudes and emotions is evidence of this sophistication. Yet we try to analyze language and human communication with bounded linear methods. In this study, ideas generated with regard to a single commonly

experienced event produced some predictable results. We see in the hierarchical Cluster analysis that many concepts were grouped as an income SEPTELEVENTH clustered with concepts like ATTACK, BUSH, UNITEDSTATES, IRAQ and most importantly US. Additionally, there was a TORTURE cluster, a NEWS cluster and even a BILL CLINTON cluster However, this obviously represented only part of peoples' thinking with regard to 9/11 and this event's fifth anniversary.

The disadvantage of hierarchical cluster analysis is that we see only part of the picture. Concepts are placed in one "overall best fit" cluster when in reality, they can be in one or many clusters depending on such variables as context, time, place, etc. On the other hand, a nonhierarchical approach allows context, time, prace; etc. On the other hand, a nonineratence approach allows us to see some of those relationships that may not have been statistical "best firs," but are nonetheless important in finding meaning in the text. Consider the concept SEPTELEVENTH which has a clearly defined cluster illustrated in Fig. 3. When it is paired with another term, NEWS, the concepts with which it was originally clustered are not as important. Furthermore, terms with which it was not seemingly related are now much closer. Indeed, SEPTELEVENTH has multidimensional meaning. This study is evidence that software like ORESMETM can be used to

analyze text in a more meaningful way. There are quirks in the software and the output could be more user friendly. But this is of minor concern given the overall results of this research which sunnort the need for a nonhierarchical approach. There are so many conforming variables when it comes to studying human communication. That it is impossible to control these variables only reinforces the importance of using nonhierarchical analysis to discover meaning

References

- Barnett, G.A. & Woelfel, J. (1979). On the dimensionality of psychological processes. *Quality and Quantity*, 13, 215-232. Jain, A.K., Mary, M.N. & Flynn, Y.J. (1999). Data clustering: A review. ACM Computing Surveys, 31(3), 264-323. Mangiameli, P. Chen, S.K. & West, D. (1996). A comparison of SOM neural
- network and hierarchical clustering methods. European Journal of Operational Research, 93, 402-417. McClelland, J.L & Rumelhart, D.E. (1988). Explorations in Parallel Distributed Processing : A Handbook of Models, Programs, and Exercises. Cambridge: MIT Press.
- Mead, G.H. (1934). Mind, Self and Society From the Standpoint of a Social
- Maad, G.H. (1934). Mind. Self and Society From the Standpoint of a Social Behaviorari (2). More, Ed. J. Change, University of Change Press. Reg of the Royal Stantistical Society B. 56(3), 409–45. Wolfeld, J. (1995). Mindea as nonhierarihical clusters in neural networks. In G. A. Barnett & F. J. Boster (Eds.), Progress in Communication Sciences, Vol. 13, Greenwich, CT. Nake Pub. Corp., 13:227.
- Woelfel, J. (1993). Galileo ORESME: User Manual. Royal Oak, MI: Terra Research and Computing. Woelfel, J. & Stoyanoff, N.J. (1993). CATPAC: A Neural Network for Qualitative
- Analysis of Text. Paper presented at the Australian Marketing Asso meeting, Melbourne, Australia. ation annual

More information

tact bl/@buffalo.edu. More information on this and related projects can tained at http://www.informatics.buffalo.edu/faculty/woelfel. A PDF-on of the poster will soon be available in the literature section of the above

University at Buffalo The State University of New York