

Chief Complaint Preprocessing Evaluated on Statistical and Non-statistical Classifiers

Jagan Dara, MS¹, John N. Dowling, MD, MS¹, Debbie Travers, PhD, RN², Gregory F. Cooper, MD, PhD¹, Wendy W. Chapman, PhD¹

Department of Biomedical Informatics, University of Pittsburgh¹ and School of Nursing, University of North Carolina²

OBJECTIVE

To determine whether preprocessing chief complaints before automatically classifying them into syndromic categories improves classification performance.

METHODS

Two preprocessors— Chief Complaint Processor (CCP) and Emergency Medical Text Processor (EMT-P)¹—were applied to two chief complaint classifiers: CoCo², a naïve Bayesian classifier; and the NYC Department of Health and Mental Hygiene coder (NYC), a keyword-based classifier. We measured sensitivity and specificity of classification into seven syndromes before and after preprocessing for 10,161 chief complaints.

RESULTS

Syndrome	Sensitivity			Specificity		
	Before	CCP	EMTP	Before	CCP	EMTP
Botulinic	55.3	50.6	76.5	99.9	99.9	99.9
Constitutional	50.0	53.1	84.3	98.9	98.8	96.9
Gastrointestinal	67.2	67.1	94.9	99.0	99.0	98.7
Hemorrhagic	60.9	66.2	77.2	99.7	99.7	99.3
Neurological	53.9	52.4	61.7	98.8	98.6	98.2
Rash	75.4	78.5	90.0	99.8	99.8	99.7
Respiratory	75.5	77.9	91.8	99.0	99.2	98.4

Table 1- CoCo: Sensitivity of classification before and after preprocessing.

Syndrome	Sensitivity			Specificity		
	Before	CCP	EMTP	Before	CCP	EMTP
Botulinic	41.1	43.5	41.2	99.9	99.9	99.9
Constitutional	80.0	77.0	85.3	96.5	96.5	96.2
Gastrointestinal	86.0	88.7	88.4	100.0	99.5	99.4
Hemorrhagic	86.5	88.9	89.2	99.1	99.1	99.1
Neurological	58.5	58.2	60.1	94.8	95.1	94.8
Rash	78.5	78.5	80.0	99.9	99.9	99.9
Respiratory	92.7	93.3	93.6	94.2	94.7	94.7

Table 2 - NYC keyword search: Sensivity and Specificity of classification before and after preprocessing.

Tables 1 and 2 show classification performance before preprocessing, after applying CCP, and after applying EMT-P. CCP exhibited high accuracy (85%) in preprocessing chief complaints. However, CCP did not exhibit an overall improvement in classification performance for CoCo (Table 1) or for NYC (Table 2), showing only a small increase in sensitivity for a few syndromes. Preprocessing with EMT-P showed a large and significant improvement

(bold in table) in CoCo's classification performance, boosting sensitivity between eight (Neurological) and 34 points (Constitutional). NYC's classification performance slightly increased with EMT-P.

CONCLUSIONS

Preprocessing steps such as spelling correction, synonym replacement, and truncation expansion, which both CCP and EMT-P perform, only slightly improved classification performance of a statistical and a keyword-based classifier. CoCo's sensitivity increased dramatically with EMT-P, because EMT-P splits chief complaints into multiple problems before classifying them into syndromic categories, whereas CoCo only assigns one syndromic category to a chief complaint regardless of the number of clinical problems in the chief complaint. NYC allows multiple classifications and was, therefore, not affected by EMT-P's splitting module. After pre-processing with EMT-P, sensitivity for CoCo exceeded that for NYC for Botulinic, GI, Hemorrhagic, and Rash and was similar for the other three syndromes.

Evaluation of preprocessing systems should not be limited to technical accuracy of the preprocessor but should include the effect of preprocessing on syndromic classification. Our results suggest that splitting chief complaints into multiple problems is important for CoCo and that other preprocessing steps only slightly improve classification performance for statistical and keyword-based classifiers.

REFERENCES

[1] Travers DA, Haas SW. Evaluation of Emergency Medical Text Processor, a system for cleaning chief complaint data. *Academic Emergency Medicine*. 2004; 11(11): 1170-1176.

[2] Olszewski RT. Bayesian classification of triage diagnoses for the early detection of epidemics. In: *Recent Advances in Artificial Intelligence: Proceedings of the Sixteenth International FLAIRS Conference*; 2003: AAAI Press; 2003. p. 412-416.

This material is based upon work supported by the National Science Foundation under grant number 0325581

Further Information: Jagan Dara,
jdara@cbmi.pitt.edu