

Talking Turkish: Using N-Grams for Syndromic Surveillance in a Turkish Emergency Department without the Need for English Translation

Sylvia Halasz PhD¹, Philip Brown¹, Cem Oktay MD², Arif Alper Cevik MD³,

Isa Kilicaslan MD², Colin Goodall PhD¹, Dennis G Cochrane MD^{4,5},

John R Allegra MD, PhD^{4,5}, Guy Jacobson PhD¹, Simon Tse PhD¹

AT&T Labs – Research¹, Akdeniz Üniversitesi, Turkey², Eskisehir Osmangazi Üniversitesi, Turkey³,

Emergency Medical Associates of NJ Research Foundation⁴,

Morristown Memorial Hospital Residency in Emergency Medicine⁵

Introduction: Previously we used an “N-Gram” classifier for syndromic surveillance of emergency department (ED) chief complaints (CC) in English for bioterrorism. The classifier is trained on a set of ED visits for which both the ICD diagnosis code and CC are available by measuring the associations of text fragments within the CC (e.g. 3 characters for a “3-gram”) with a syndromic group of ICD codes. Because the ICD system is language independent, the technique has the potential advantage of rapid automated deployment in multiple languages. Our objective was to apply the N-Gram method to a training set of Turkish ED data to create a Turkish CC classifier for the respiratory syndrome (RESP) and determine its performance in a test set.

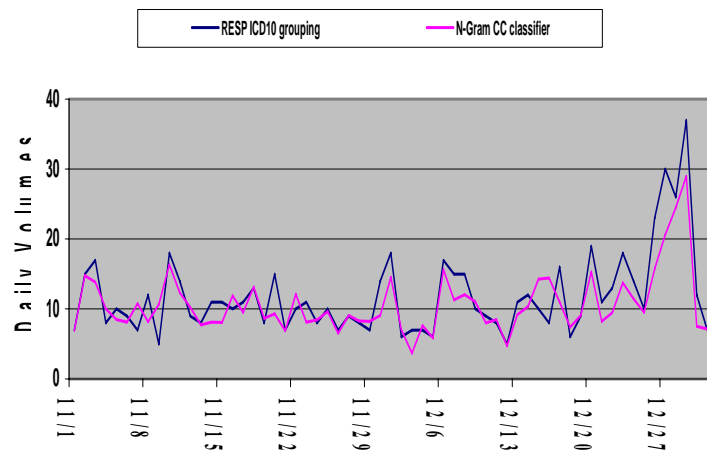
Objectives: To determine how closely the performance of an ngram CC classifier for the Respiratory (RESP) syndrome matched the performance of the ICD9 classifier.

Methods: Design: Retrospective cohort. Setting: University hospital ED in Turkey. Participants: All ED visits in the database for 2002. Protocol: Two of the authors created by consensus a respiratory grouping of ICD10 codes chosen to be similar to a standard respiratory grouping of ICD9 codes created by the ESSENCE-CDC project. We then used an N-Gram method adapted from AT&T Labs' technologies applied to the first 10 months of data as a training set to create a Turkish CC RESP classifier. We next applied the classifier to the test set of visits for the last two months and determined the correlation for daily volumes measured by the CC classifier versus the RESP ICD10 grouping.

Results: The Turkish ED database contained 30,157 visits. The correlation coefficient was $R = 0.88$.

Conclusions: The N-Gram method successfully created a CC RESP classifier in Turkish that performed similarly to the ICD10 RESP grouping. This was accomplished without translating the Turkish CC to English or understanding the Turkish language. This approach has promise in that it may offer a complementary method to using manual and natural-language techniques and has the advantages of systematic, consistent and rapid deployment as well as language independence

Time Series of Daily Volumes by the Resp ICD10 Groupings and the N-Gram CC Classifier (Test Set 11-1-02 to 12-31-02)



Daily Volumes by N-Gram CC classifier versus RESP ICD10 Grouping (Test Set)

