

Development of a Natural Language Processing Machine to Enhance Identification of Opioid Use Disorder in Electronic Health Records

Avery Chadd¹, Rebecca Silvola², Yana Vorontsova², Andrea Broyles³, Jonathan Cummins³, Daniel Hood³, Robert Stratford², Sara Quinney^{2,4}

¹Indiana University School of Medicine; ²Indiana University School of Medicine, Division of Clinical Pharmacology; ³Regenstrief Institute; ⁴Indiana University School of Medicine, Department of Obstetrics and Gynecology

Background/Objective: Real-world data, including electronic health records (EHRs), has shown tremendous utility in research relating to opioid use disorder (OUD). Traditional analysis of EHR data relies on explicit diagnostic codes and results in incomplete capture of cases and therefore underrepresentation of OUD rates. Machine learning can rectify this by surveying free clinical notes in addition to structured codes. This study aimed to address disparities between true OUD rates and cases identified using traditional ICD codes by developing a natural language processing (NLP) machine for identifying affected patients from EHRs.

Methods: Patients (≥ 12 years old) who had received an opioid prescription from IU Health or Eskenazi Health between 1/1/2009 and 12/31/2015 were identified by the Regenstrief Institute. Exclusion criteria included any cancer, sickle cell anemia, or palliative care diagnoses. Cases of OUD were identified through ICD codes and NLP. The NLP machine was developed using a dictionary of key OUD terms and a training corpus of 300 patient notes. A testing corpus of 148 patient notes was constructed and validated by manual review. The NLP machine and ICD 9/10 codes were independently tested against this corpus.

Results: Although ICD codes identified OUD cases with high specificity (98.08%), this method demonstrated moderate sensitivity (53.13%), accuracy (68.92%), and F1 score (68.92%). Testing using the NLP method demonstrated increased sensitivity (93.75%), increased accuracy (89.19%), and increased F1 score (91.84%); specificity mildly decreased (80.77%).

Conclusion: Our revised NLP machine was more effective at capturing OUD cases in EHRs than traditional identification using ICD codes. This illustrates NLP's enhanced capability of identifying OUD cases compared to structured data.

Potential Impacts: These findings establish a role for NLP in OUD research involving large datasets. Ultimately, this is intended to improve identification of risk factors for OUD, which is of significant clinical importance during a public health crisis.