

March 7, 2011

Digital Urdu: New Software Improves Data Analysis of Pakistan's National Language

By AUDREY WATTERS of

The extent to which social media sites like Twitter and Facebook play a role in the recent political uprisings in Tunisia, Egypt, Bahrain and so on continues to be a source of [debate](#). What is more clear, however, is that the major languages of these regions are not well-served by electronic resources that make text analysis of these documents and data possible.

But now computer scientists have developed the first software system that will allow for the processing of documents in Urdu, the national language of Pakistan and one of the five most-spoken languages in the world.

The software will help lay the foundation for data mining in Urdu and provide for more accurate transliteration, as well as open the door for projects in similar languages. "This is the first comprehensive, natural language processing system for Urdu," says Rohini Srihari, University of Buffalo associate professor of computer science and engineering.

The work is a joint project between her department and Janya, an Amherst-based company she founded that provides information extraction technology in multiple languages, including Chinese, Arabic, Pashto, and Russian.

Natural Language Processing and Urdu

The problem with data mining and sentiment analysis in other languages is that they don't often have the same sort of "established electronic infrastructures" that we have for English and European languages. "If you are trying to do sentiment analysis - to find out what are the main topics people are talking about in a country, is there intensity building up over something and who is swaying opinion - then you must have an information extraction system," Srihari says.

That is what she has been working on with researchers, something that can perform word segmentation (tagging parts of speech, for example) and entity-tagging (recognizing people, place, and organization names) in a raw, untranslated Urdu document.

"Voice of the Citizen" Through Social Media Data Mining

Srihari says she's focused on the "voice of the citizen" in this project. "Some of the information is political and some of it is not," she says, and despite the turmoil in the region, a lot of the social media chatter in Pakistan is about cricket.

Srihari presented her findings at a recent conference - "Blogs & Bullets: Social Media and the Struggle for Political Change" - at Stanford University. She says she became interested in Urdu because they were looking at blogs from different cultures. Noting that the advent of the Web has caused an explosion in online content in a variety of languages, Srihari [says](#), "When you start looking at blogs in different cultures, you can really start to understand public sentiment and opinions."

Tags: [janya](#), [natural language processing](#), [translation](#), [urdu](#)

Copyright 2011 ReadWriteWeb. All Rights Reserved.

ReadWriteWeb is an independent technology blog. [Read More »](#)