# Know Your Disease!

**ASU** IRA A. FULTON SCHOOLS OF
**engineering**
ARIZONA STATE UNIVERSITY

WebMD

Kashyap Bhansali, 1209278261, kashyap.bhansali@asu.edu
Saloni Shah, 1209404115 , saloni.shah@asu.edu
Raj Dalvi, 1209232176, rsdalvi@asu.edu
Shruti Mahajan, 1210431622, smahaja7@asu.edu

## Motivation & Relevant Work:

- Online health discussion forums are popular information seeking sites, which provide reliable health and medical news.

- The health information retrieving and organizing patterns of a user can be recognized through this alternative called Question-Answering Systems and forums. (Dolares M.et al, 2010).

## Research Questions:

1. What kind of questions are asked on the forum? What are the most commonly discussed Topics & Diseases?

2. What are the Symptoms people with particular Disease have? What parts of the human body are affected by this Disease?

3. Prevalence of the Diseases in the population over the years?

## Data Collection:

- Collected data about diseases, symptoms and posted date of question and answers from the WebMD dataset provided. Also extracted diseases and symptoms from Q&A data.

- Scrapped additional diseases and symptoms data from Q&A forum: WebMD (www.webmd.com) and from Mayo Clinic (www.mayoclinic.org)

- Processed the data to map the Diseases and Symptoms to Topics provided in the dataset. Collected data about the Human body parts where the diseases affect.

## Findings:

- There is a lot of discussion about Women health problems, accompanied by Q&A's for General Symptoms and Sensory Organs of Human body.

- Pain and Infection are the most common complications associated with a gamut of diseases.

- Complicated issues like Cancer, Blood problems, Surgery has been less discussed and talked about on the forum, since they're critical enough to directly consult a doctor in-person.
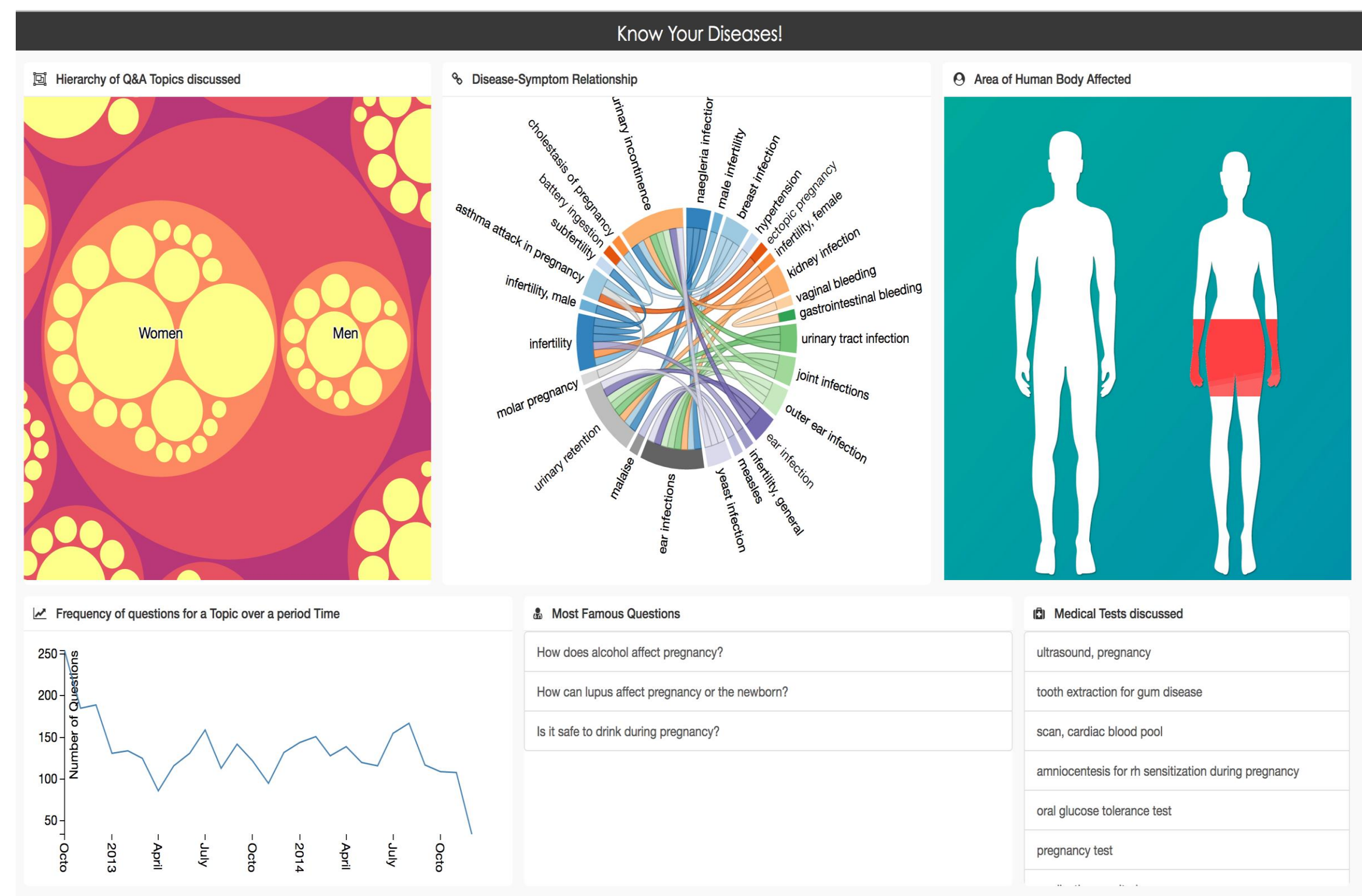


## Methodology:

- Used the scraped data of Diseases and Symptoms to extract the relevant close matches from the Q&As in the dataset using *Sequence Matching*.

- *Mapped the relationships* between Diseases and Symptoms that occur in the same Q&As and represented using a *Chord diagram*.

- Based on the body parts affected by the disease, relevant area is highlighted on a 2D human body figure.

- Majority of the body parts were mapped by extracting data from the dataset while a few manual inputs were necessary for improving the accuracy.

- Sequence Matching was preferred over TF-IDF because it reduced the occurrence of meaningless data and also handles typos better. Usually, the complex names that may be misspelled would be handled and identified accurately.

- Using an interactive line chart we demonstrate trend discussion over a particular Disease topic.

- With a *bag-of-words* approach we perform some clustering with manual inputs for correctness, to build a hierarchical structure and represent using *Zoomable-Bubble* chart.