

Natural Language Processing Methods to Parse Birth and Death Records

Matthew Jin¹, Justina Zurauskiene², and Zeynep Madak-Erdogan³

¹Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, ²Institute of Cancer and Genomic Sciences, University of Birmingham, ³Department of Food Science and Human Nutrition, University of Illinois at Urbana-Champaign

INTRODUCTION & AIMS

The goal of this project is to investigate health disparities arising from birth and death outcomes that exist across the residents of Champaign County through the parsing of electronic birth and death records provided by the Champaign Urbana Public Health District (CUPHD) from the last decade. This project is currently focused on identifying maternal health disparities and how these disparities are associated with various birth outcomes. To accomplish this task, we develop a pipeline for digitizing birth and death records by using natural language processing (NLP), which allows our program to read and understand the language in the birth and death records.

This research will provide an autonomous tool for future projects, which require the parsing of not only electronic birth and death records but also other significant records in medicine.

A sample birth certificate file, which is fake filled for training purposes, is shown below. Note the complexity of the document. The data for the fields vary in type, shape, and length. For example, data may come in the form of checks crossed off in small boxes.

SAMPLE BIRTH CERTIFICATE FILE

[illegible]

METHOD

The pipeline will receive the path to a folder of scanned birth certificate files in PDF format. These files are two pages each. First, we convert these PDF files into image files (.jpeg).

Then, we remove the horizontal and vertical lines in the image files by finding all horizontal and vertical contours (`cv2.findContours()`) in the image and “erasing” them by drawing the same contours (`cv2.drawContours()`) with the color white, effectively blurring the horizontal and vertical lines.

With the processed image files, the pipeline calculates the size of the document of each page and zooms in on different parts of the birth certificate image to allow for more accurate character recognition.

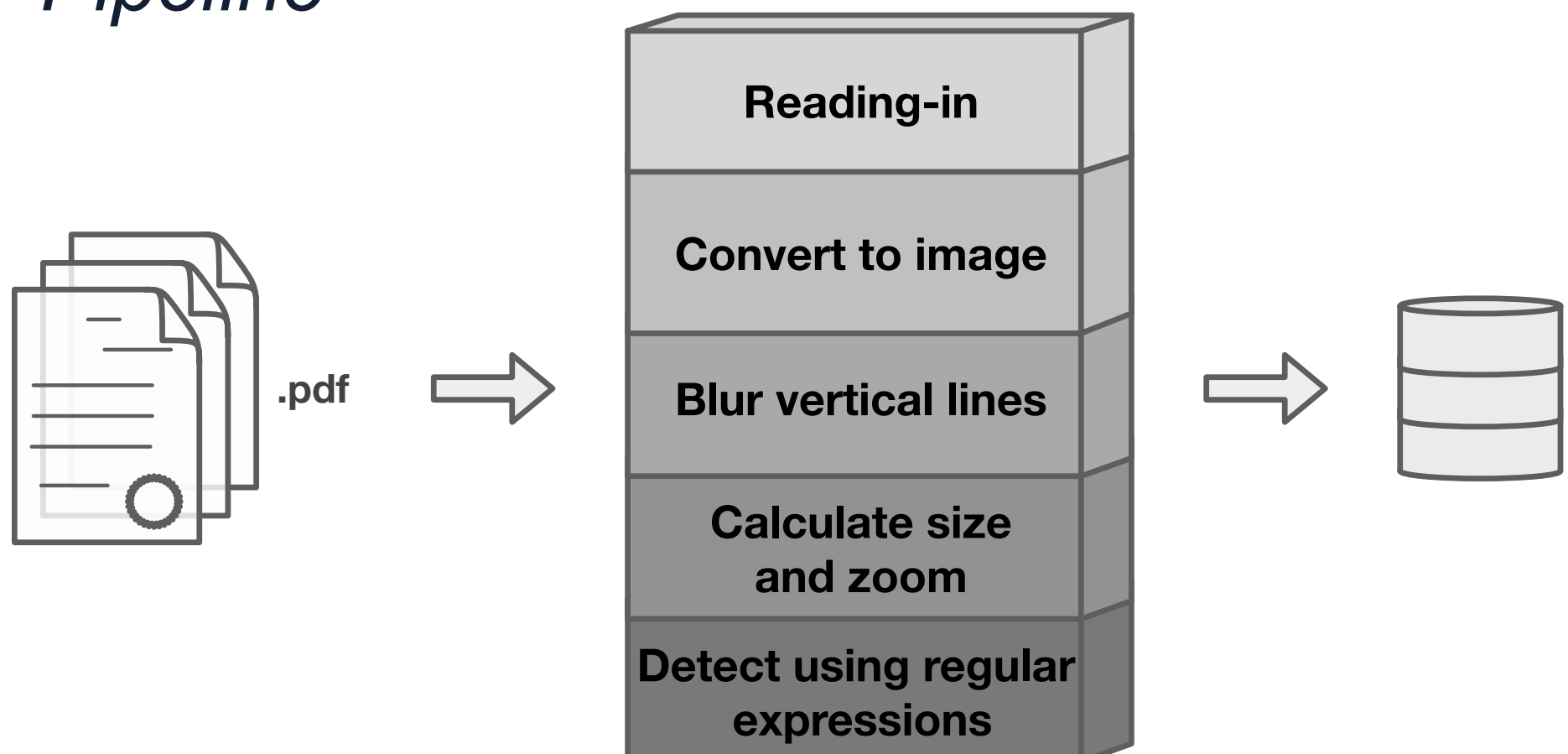
We use Tesseract, an Optical Character Recognition (OCR) software, to read in the zoomed-in portion of the birth certificate. At this step, the pipeline has a string of text describing this zoomed in portion.

Finally, the pipeline uses regular expressions pattern matching to extract desired patterns such as dates, city names, and more.

Regular Expression Pattern

```
re.compile(r'[A-Za-z]+\d{2},.\d{4}')
```

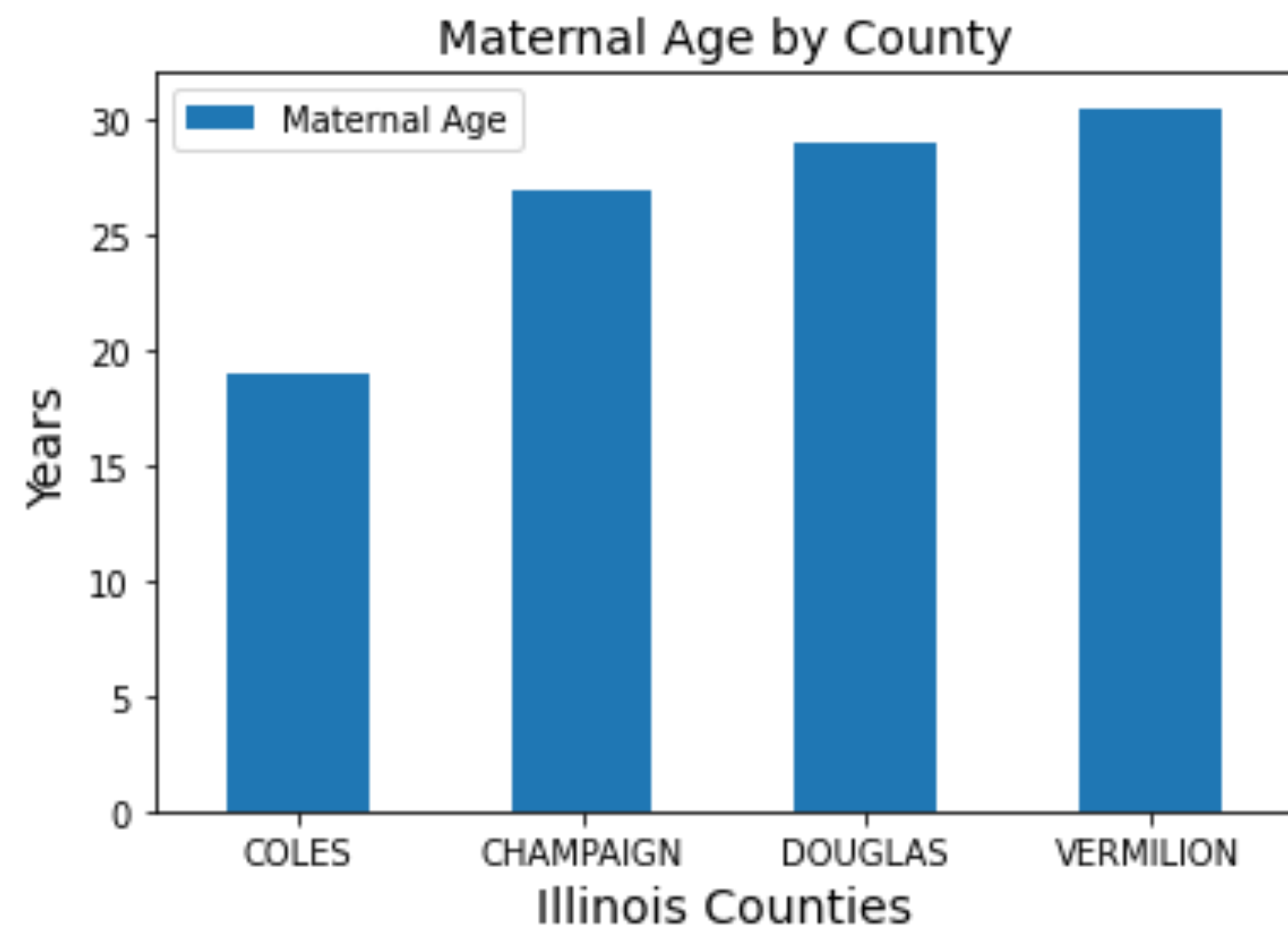
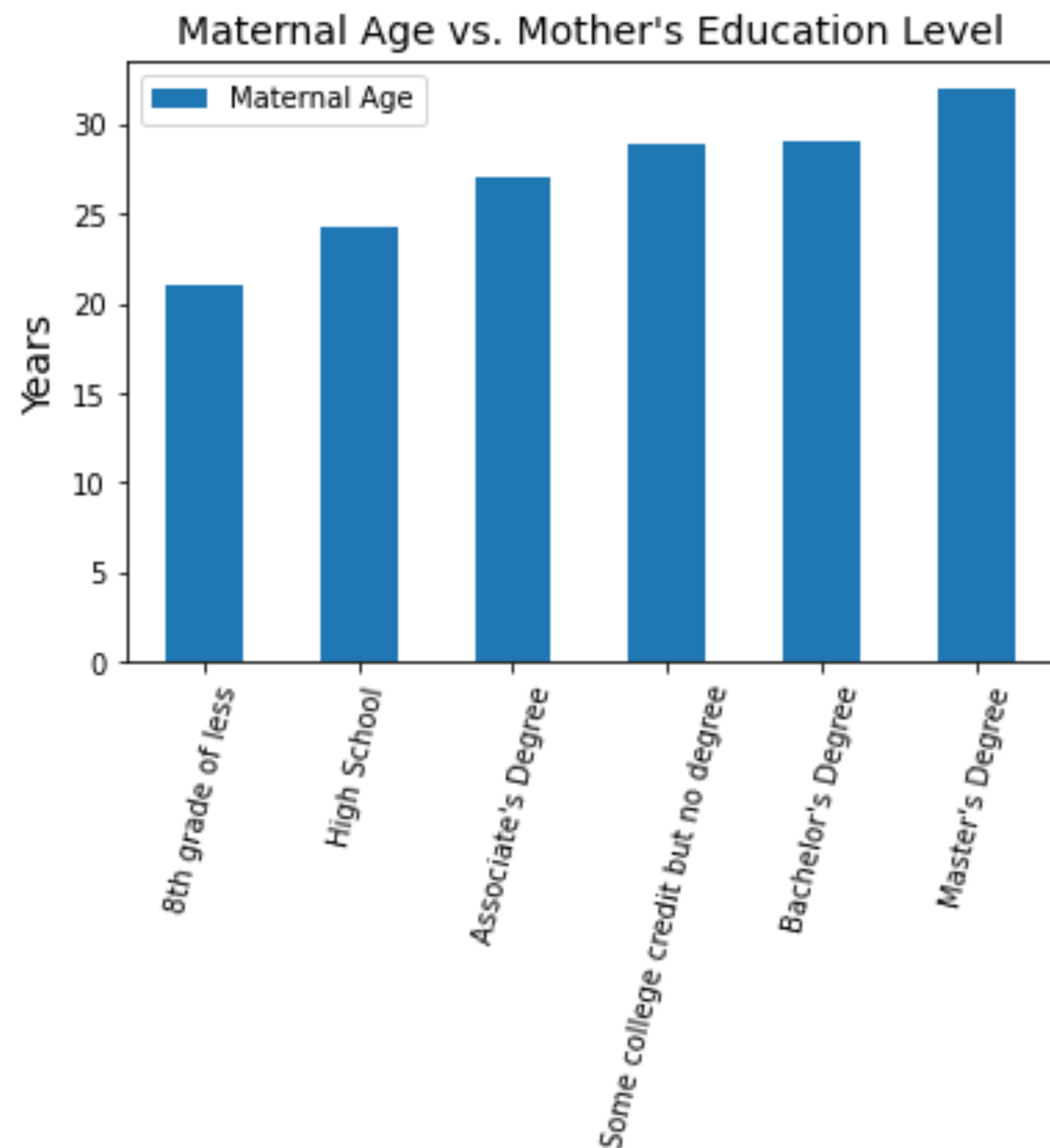
Diagram of Pipeline



RESULTS

The pipeline is currently able to extract the first page of birth certificate files with over 98% accuracy. The first page contains fields with mostly socioeconomic data. Currently, our CSV database reveals patient data using a subset of 21 scanned birth certificate records provided by the CUPHD.

On the first page of the birth certificate file, there is information about the date of birth of the child and the mother. With this information, we can plot how the maternal age varies depending on different fields of interest.



DATABASE

The data pipeline will ultimately output a CSV database of patient information. This table is stored as a Pandas DataFrame during the execution of the pipeline.

The row names of the DataFrame are filenames with unique numbers for patient identification. The column names of the DataFrame are the field names within the birth certificate and death certificate files. A snippet of the birth data table is show below.

21. BIRTHING PARENT'S EDUCATION	
Birth cert 12B.pdf	Some college credit but no degree
Birth cert 3B.pdf	Master's Degree
Birth cert 2B.pdf	High School
Birth cert 13B.pdf	Some college credit but no degree
Birth cert 5B.pdf	Bachelor's Degree

FUTURE DIRECTION

- Continue to extract data from birth certificate files, including the second page of these files which contains the medical history of the mother and infant
- Analyze and relate different fields of interest with natural language processing and deep learning
- Start parsing data from death records to create a new CSV database

ACKNOWLEDGEMENTS

NCSA SPIN Program

- Dr. Zurauskiene
- Dr. Madak-Erdogan

Champaign Urbana Public Health District

