# Predicting the need for medical imaging on patients' information

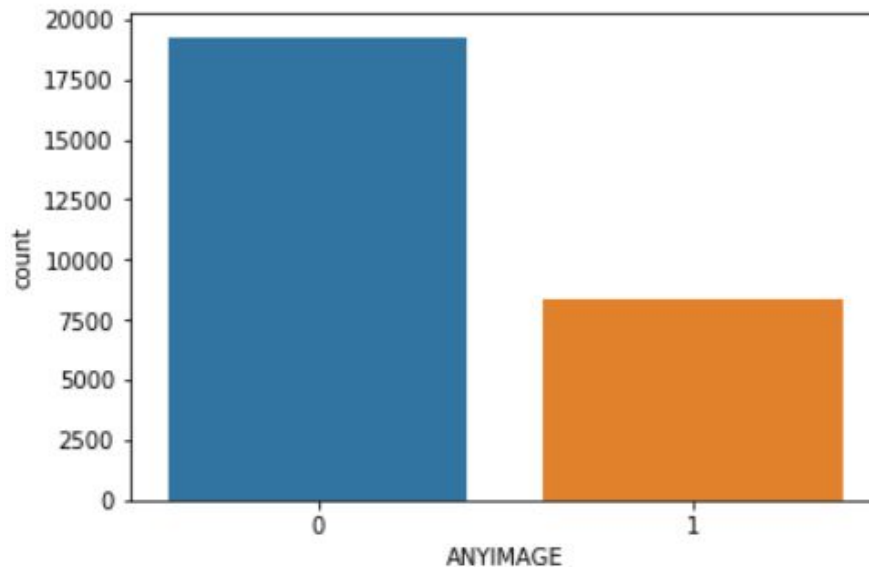Group 2: Keyu Chen, Yinuo Chen, Rongqian Zhang

# Dataset

- Ambulatory Health Care Data

- The National Hospital Ambulatory Medical Care Survey (NHAMCS)

- **Patient** characteristics : age, sex, race, and ethnicity, and so on

- **Visiting** characteristics : temperature, heart rate, systolic blood pressure, pulse oximetry, and so on

- **Text** characteristics : reasons for visit and possible causes of disease

- Response variable : **ANYIMAGE** (X-ray and CT scan)

# Question

**Build a model to predict whether patients need any medical imaging**
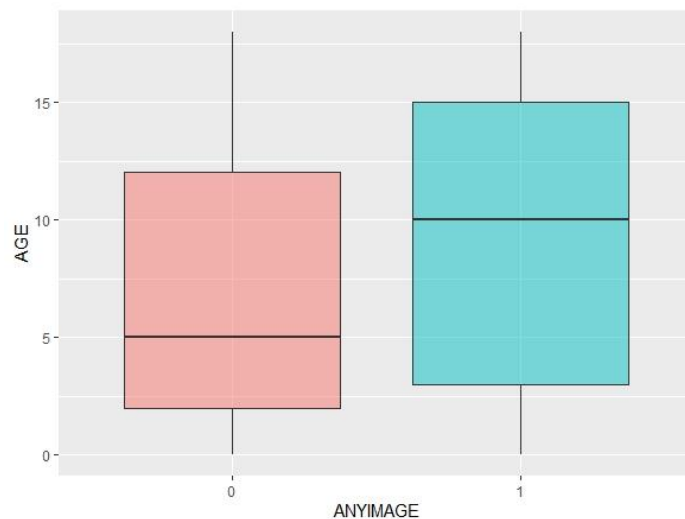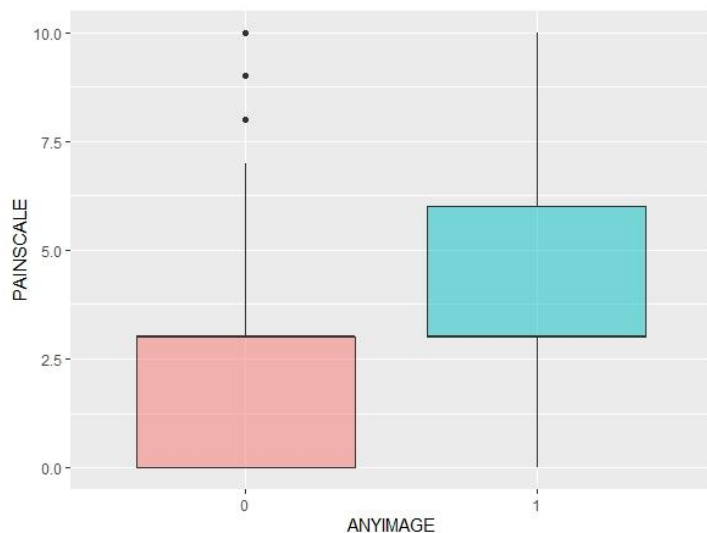
# EDA - Response

- 2012 - 2016: 27,665 observations
- Positive : Negative = 1 : 2.3

# EDA - Structured Data

- Some variables have significantly different distribution in two groups
- Ex: Pain scale, age

# EDA - Text Data

- Data cleaning
  - Remove urls, punctuations, single letter
  - Remove stop-words
  - Lowercase
- Reasons for visiting
  - Most common words
  - Word cloud

| word | n |
| --- | --- |
| pain | 9315 |
| fever | 5046 |
| sore | 4279 |
| unspecifi | 4132 |
| cough | 3979 |
| ach | 2941 |

# EDA - Text Data

- Possible causes of disease
  - Most common words
  - Word cloud

| word | n |
| --- | --- |
| NA | 19574 |
| fall | 3444 |
| activ | 1879 |
| oth | 1713 |
| involv | 1562 |
| occurr | 1535 |

# Model - structured variables

- Multilayer perceptron

# Model - text variable

- Word embedding - GloVe
  - Global Vectors for Word Representation
  - Based on word-word co-occurrence matrix
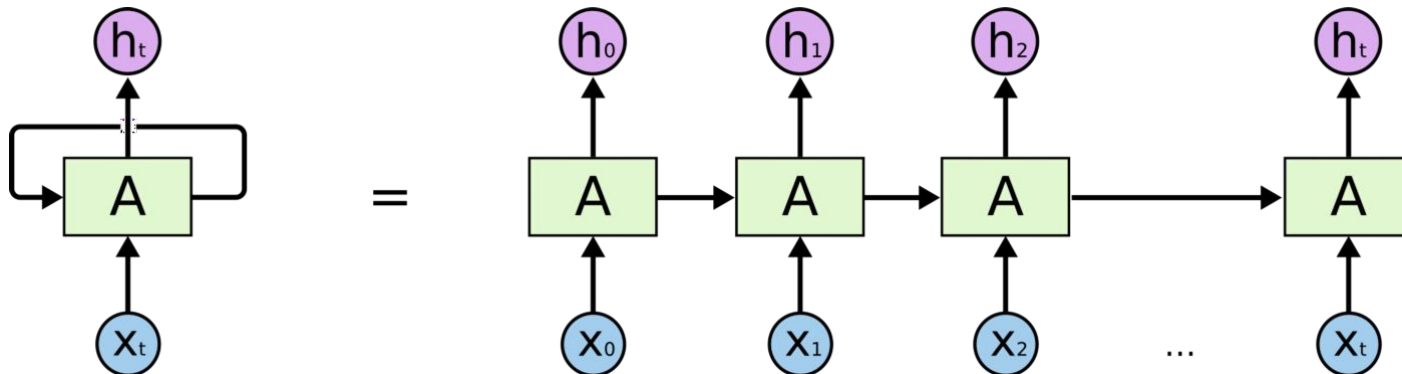  - Similar words will have similar representation

The cat sat on the mat.
The dog sat on the mat.

co-occurrence matrix
with a window size of 1.

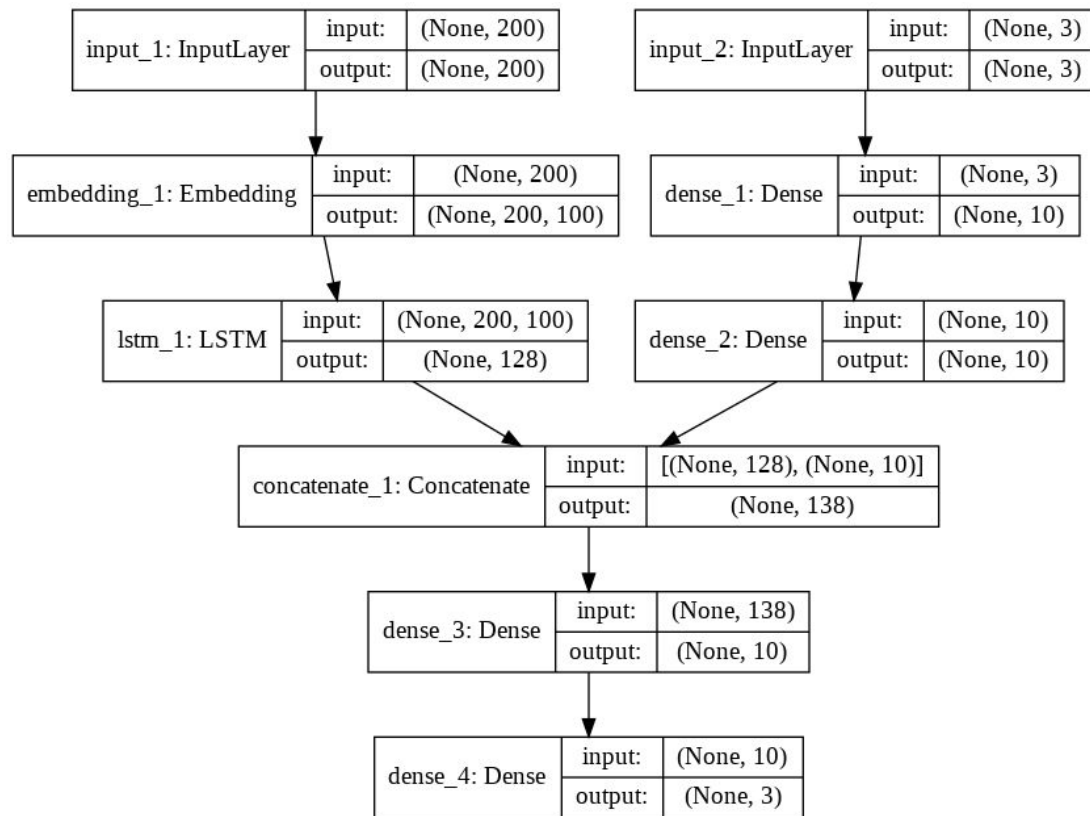|  | the | cat | sat | on | mat | dog |
|---|---|---|---|---|---|---|
| the | 0 | 1 | 0 | 2 | 2 | 1 |
| cat | 1 | 0 | 1 | 0 | 0 | 0 |
| sat | 0 | 1 | 0 | 2 | 0 | 1 |
| on | 2 | 0 | 2 | 0 | 0 | 0 |
| mat | 2 | 0 | 0 | 0 | 0 | 0 |
| dog | 1 | 0 | 1 | 0 | 0 | 0 |

# Model - text variable

- LSTM (Long short-term memory)
  - Similar to the way human read sentence
  - LSTM read words sequentially
  - At each time stamp, it uses current word and the activation from last timestamp as input

# Model - mixed variables

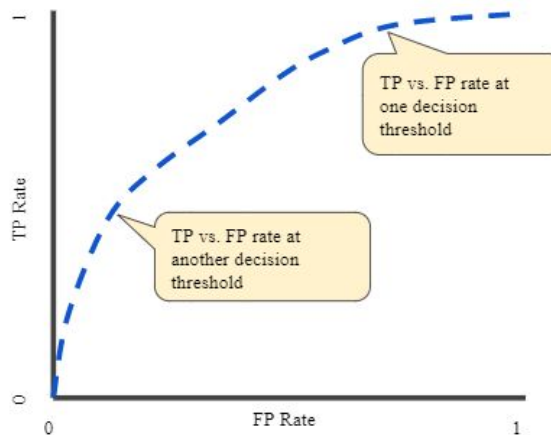- Hybrid neural network
  - Concatenating MLP and LSTM

# Results

- Evaluation
- AUC (Area Under The Curve)

$$TPR = \frac{TP}{TP + FN} \qquad FPR = \frac{FP}{FP + TN}$$

# Thank you!