

# Using natural language processing on clinical notes to predict hospital readmission



# Hospital Readmission Reduction

- targeted as a key metric of patient care
- 2012: Affordable Care Act initiated the Hospital Readmission Reduction Program
- Incentivize improved patient outcomes by financially penalizing hospitals with excessive readmission rates
- \$1.9 billion in penalties in first 5 years (American Hospital Association)

The image features a decorative border made of teal-colored squares, each containing a white triangle pointing towards the center. This border frames a central white rectangular area.

# **About The Data**

---

# Data Collection

---

## MIMIC-III version 1.4

- over 58,000 hospital admissions from critical care units of the Beth Israel Deaconess Medical Center
- 38,645 adults and 7,875 neonates
- data spans June 2001 - October 2012
- Collected as 26 CSV files (6.2GB) and loaded in a PostgreSQL database
- Pulled nursing notes and discharge summary

# Data Preprocessing

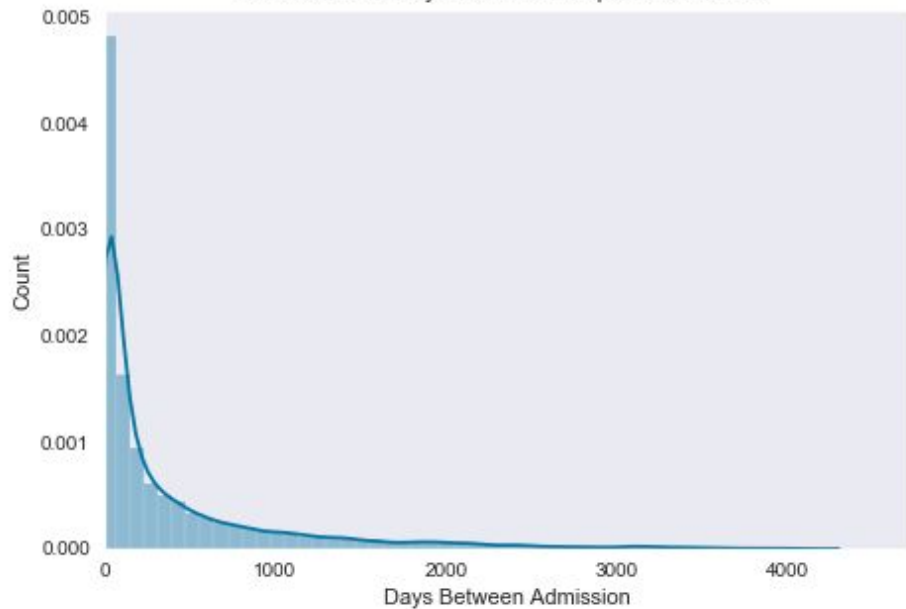
- Defined next admission date for each subject and admission
- Defined admission type (elective, emergency, urgent, newborn)
- Compute number of days between admissions
- Mark elective visits as no readmission
- Combined all notes for each subject and admission into a single string
- Dropped all duplicate and newborn admissions
- Compute target variable using days between admissions
- Split 70% of data into training set, 30% into test set

---

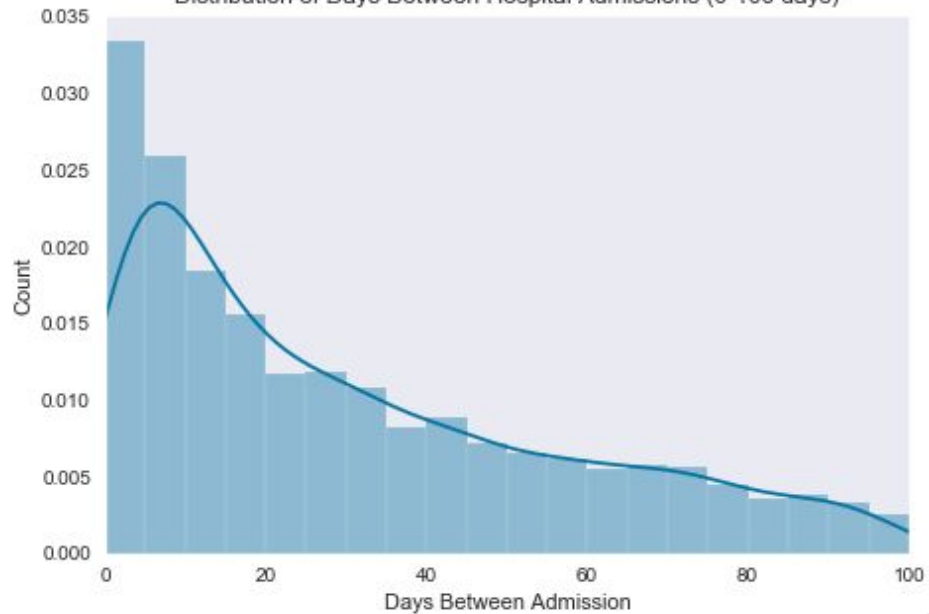
# Exploratory Data Analysis



Distribution of Days Between Hospital Admissions



Distribution of Days Between Hospital Admissions (0-100 days)

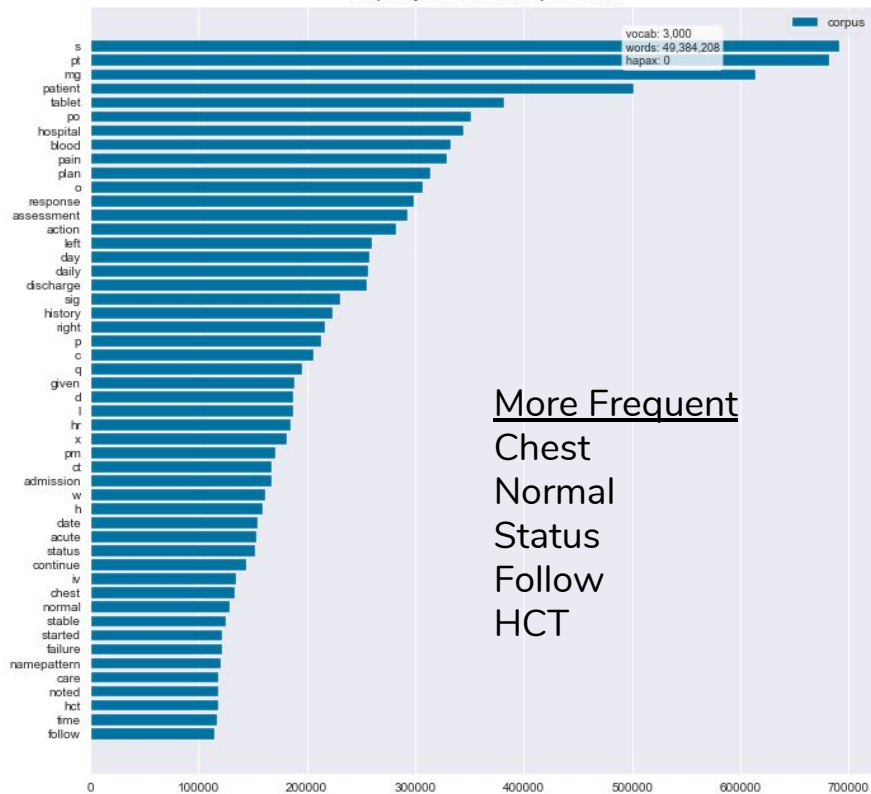


Readmissions peak early



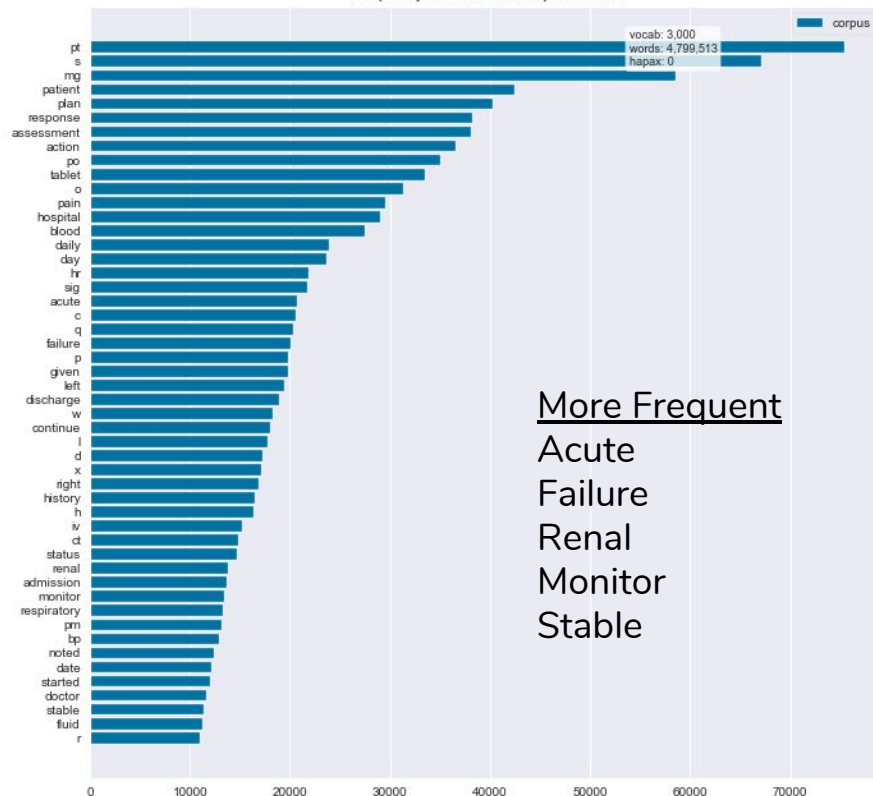
# NOT READMITTED

Frequency Distribution of Top 50 tokens



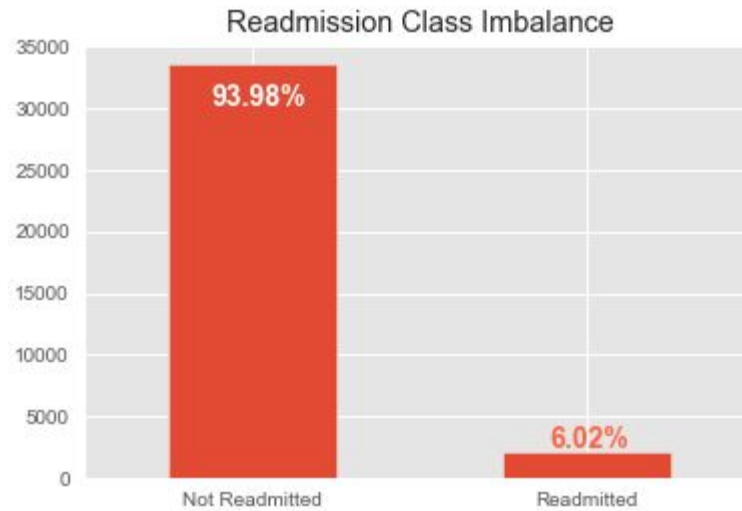
# READMITTED

Frequency Distribution of Top 50 tokens






# Imbalanced Dataset





T-SNE MAP



# **Feature Engineering**

# Bag-of-words

Removed  
punctuation and  
numbers

Lowercase

Tokenized

# Word Embeddings



Cleaned and tokenized

Window size = 6

200-dimension word vectors

Stemmed

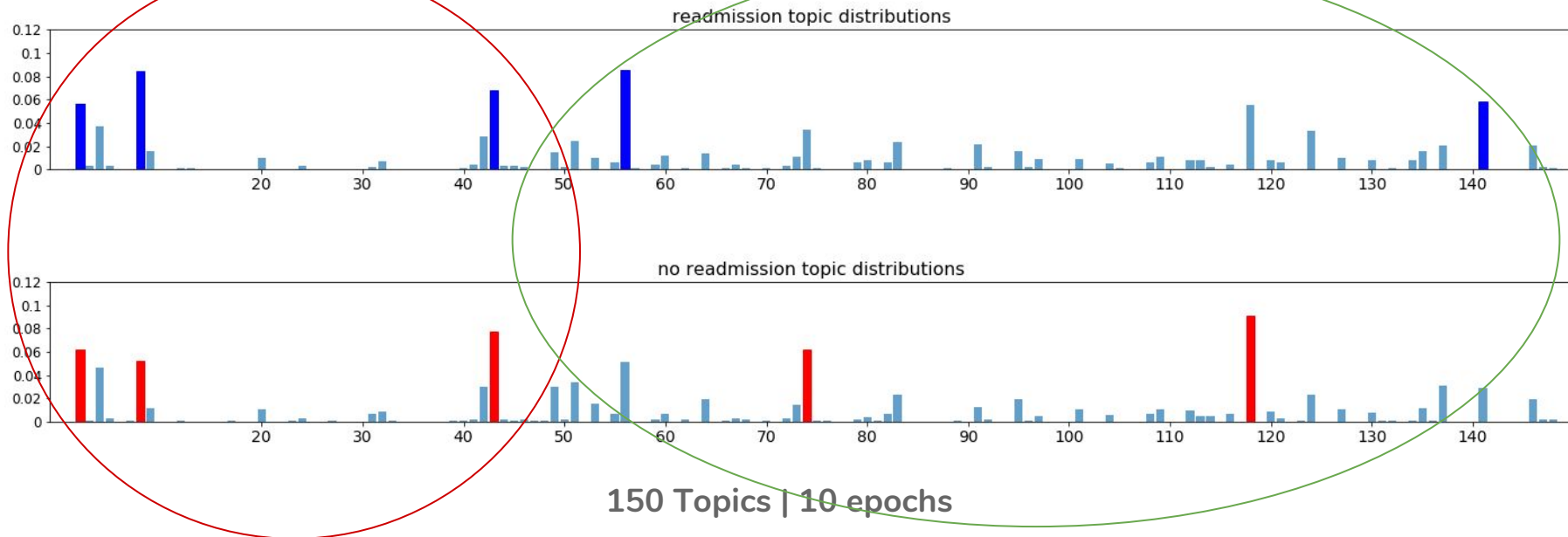
Epochs = 6

Lemmatized



# Latent Dirichlet Allocation

Average Probability of Topic



Top Topics for Readmitted Patients	
Topic	Top Words
2	tablet, daili, sig, cardiac, ventricular
8	tablet, daili, sig, hospital1, pt
43	statu, unit, show, number, also
56	tablet, sig, daili, need, capsul
141	daili, tablet, cultur, sig, neg

Top Topics for Not Readmitted Patients	
Topic	Top Words
2	tablet, daili, sig, cardiac, ventricular
8	tablet, daili, sig, hospital1, pt
43	statu, unit, show, number, also
74	tablet, daili, sig, disp, refil
118	arteri, coronari, qd, postop, statu

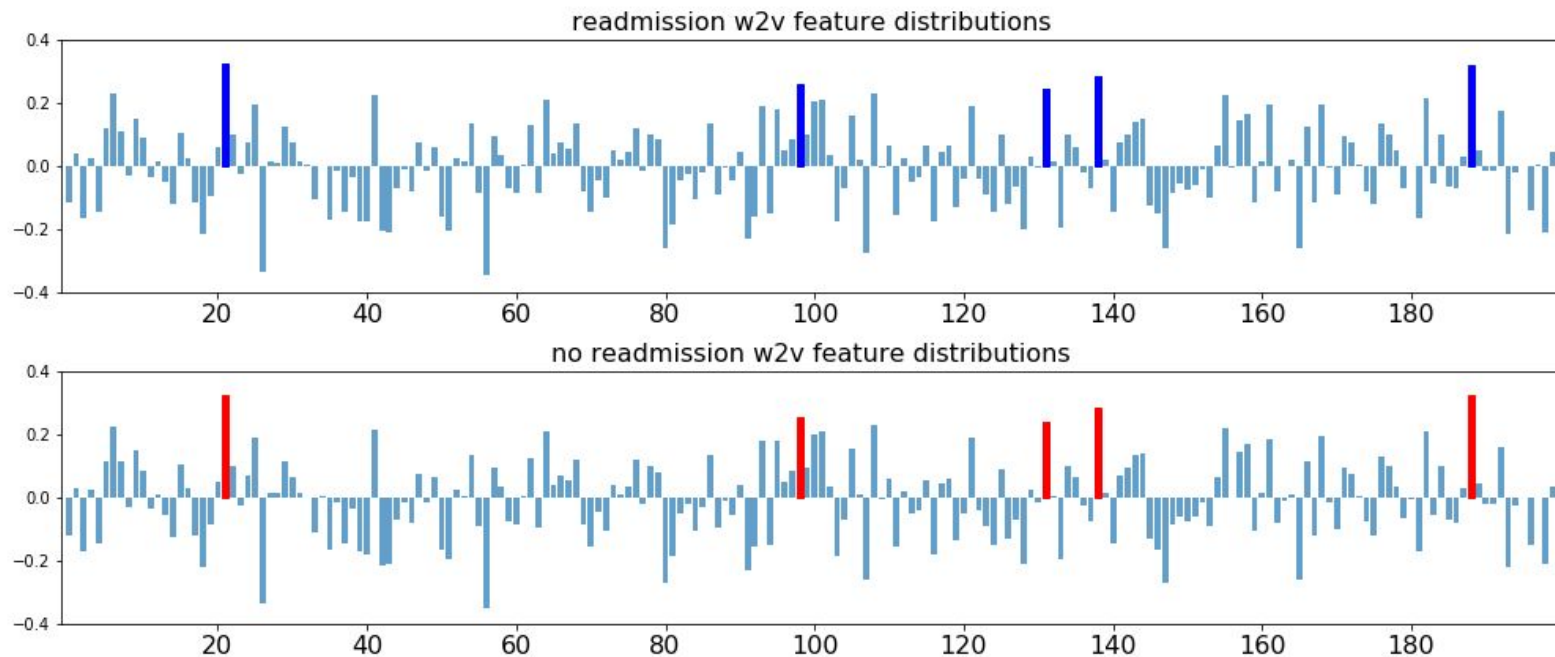


Defining LDA Topics By Examining the Top Words





# Word2Vec Feature Distribution





# Predictive Modeling



	ROC-AUC
Random Forest (under-sampling)	0.7076
Word2Vec & LDA w/ Logistic Regression	0.6796

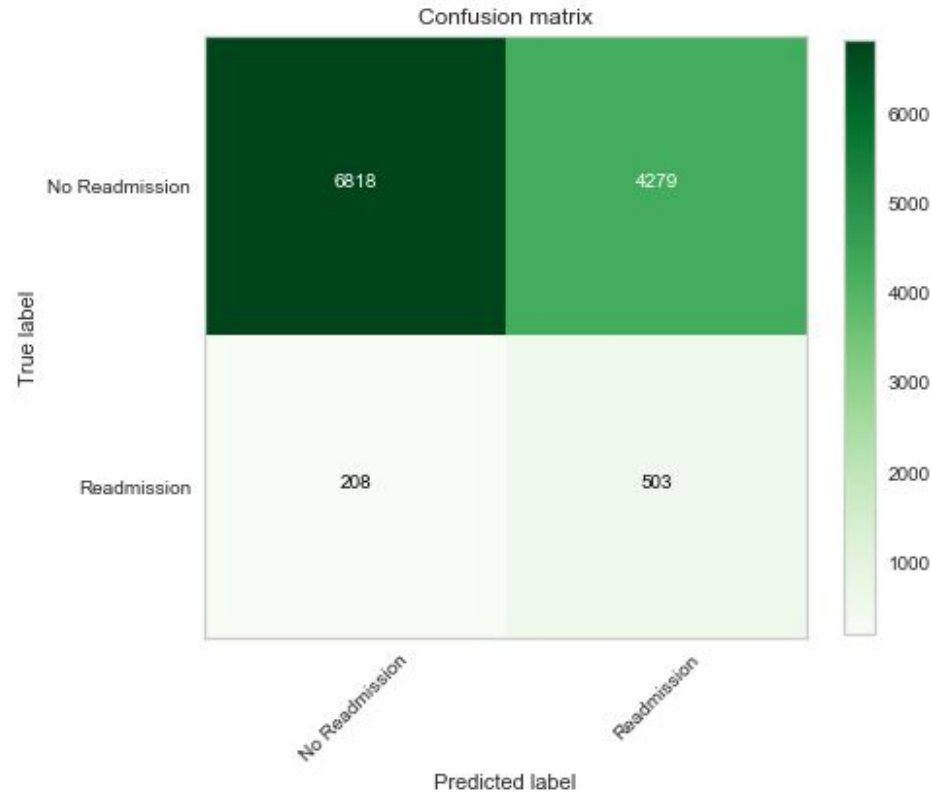
# Random Forest

---

Random undersampling

500 trees

Maximum depth = 25



 Bag-of-words with Random Forest and Random Undersampling

**Moving  
Forward**



# Next steps in model improvement

## Word2Vec

Use random forest or other classifier

Grid search:

- window size
- learning rate
- number of epochs
- downsampling threshold

Alter threshold (or class weight)

## Random Forest

More extensive grid search:

- Number of estimators
- Tree depth
- Leaf and node parameters
- Etc.

Balance class weight

Additional feature engineering:

- LDA
- Lab work, pharmacy, etc



# Next steps in production

- EHR flag for readmission risk
- Research key risk factors
- Test effectiveness with RCT

---

# Any Questions?

E. Chris Lynch

[echrislynch@gmail.com](mailto:echrislynch@gmail.com)

[github.com/TheeChris](https://github.com/TheeChris)

[linkedin.com/in/echrislynch](https://linkedin.com/in/echrislynch)