



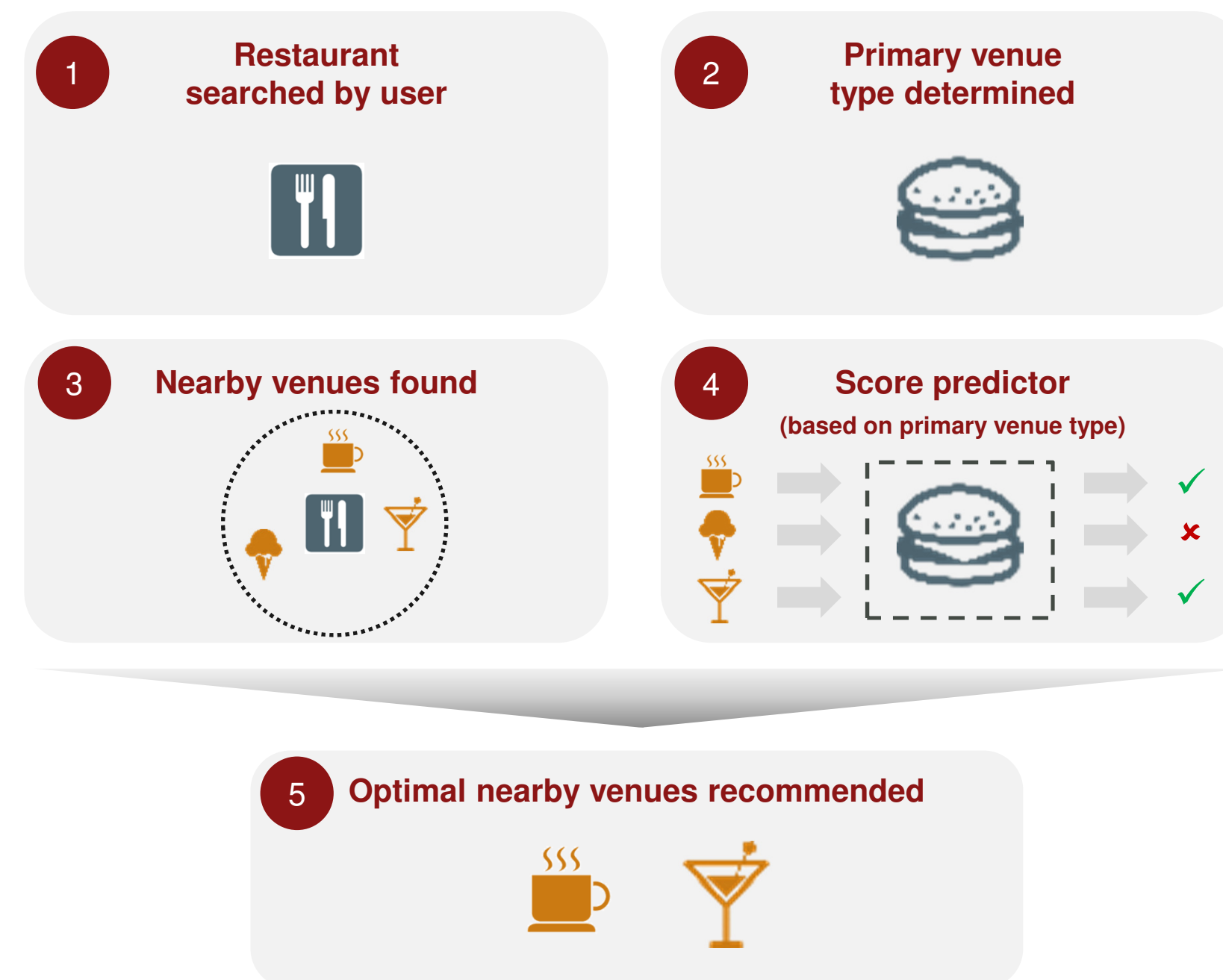
# Complementary Venue Recommendation Model for Yelp

Hyun Sik Kim (hsik@stanford.edu), Ryan Wong (rawong@stanford.edu)

## Overview

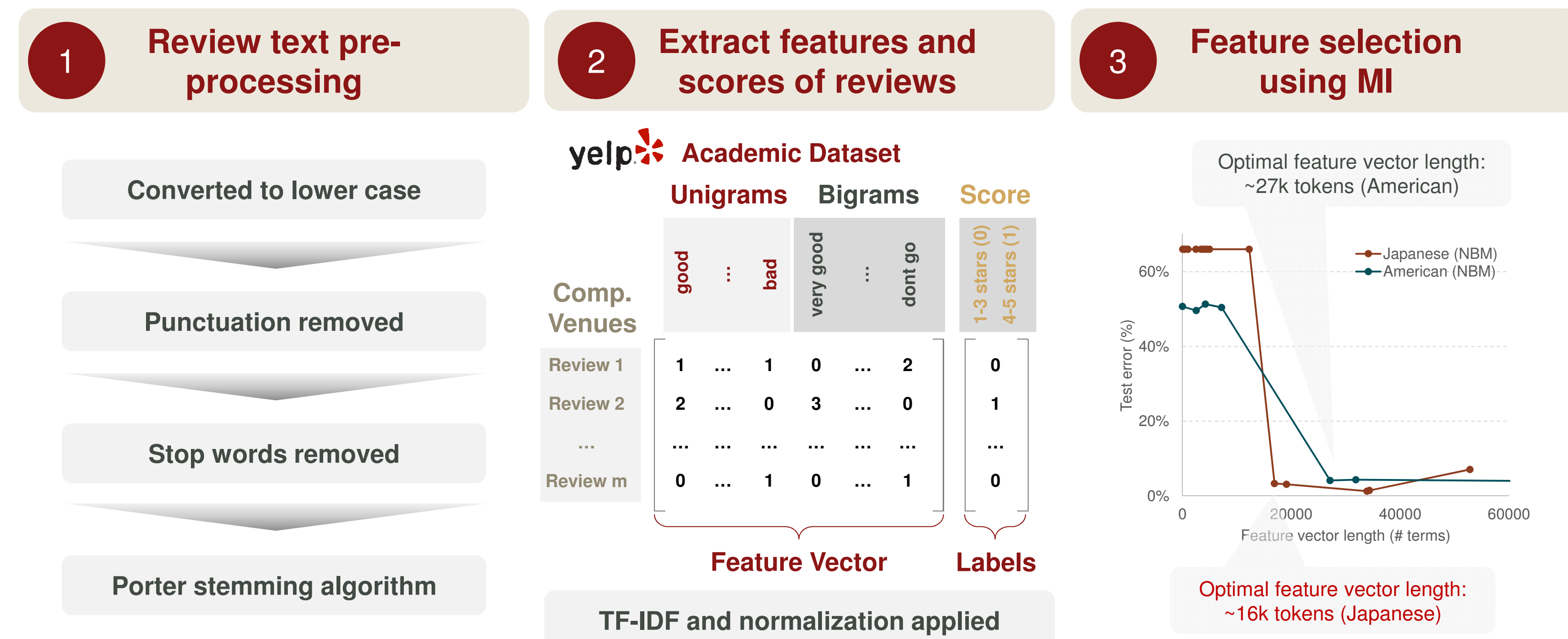
We developed a **novel machine learning model** that **recommends nearby complementary venues** (e.g. café) that the user is likely to enjoy based on a restaurant searched using Yelp.

- To simplify the search process for selecting multiple venues for a single outing using Yelp
- Our model **recommends complementary venues** using binary classifiers that predict the venue score using review text:
  - > **Great venue (1)**: 4-5 star rating
  - > **Mediocre / poor venue (0)**: 1-3 star rating
- Optimal recommendation/s based on:
  - > **Highest predicted score**
  - > **Proximity** (within 1 mile of restaurant)
- Model trained based on review text written by a **common reviewer** of the restaurant and complementary venue



## Data and features

Feature vector comprised the count of **unigrams and bi-grams** in **Yelp reviews** of **complementary venues** from the **Yelp academic dataset**.



## Model performance and evaluation

Our **Naive Bayes with TF-IDF and normalization** implementation exhibited **~3-4% test error performance**. This was comparable to SVM performance.

### Naive Bayes comparable to SVM – linear was optimal SVM kernel

American Restaurants	Training	Test <sup>(1)</sup>
Naive Bayes (TF-IDF <sup>(2)</sup> )	2.2%	4.1%
SVM with linear kernel	2.9%	5.5%
SVM with linear kernel (TF-IDF <sup>(2)</sup> )	1.4%	5.1%
SVM with RBF kernel	1.8%	6.2%

Naive Bayes robust despite conditional independence assumption

Linear outperformance over other kernels suggest text classification problem linearly separable

Both models benefited from TF-IDF and normalization

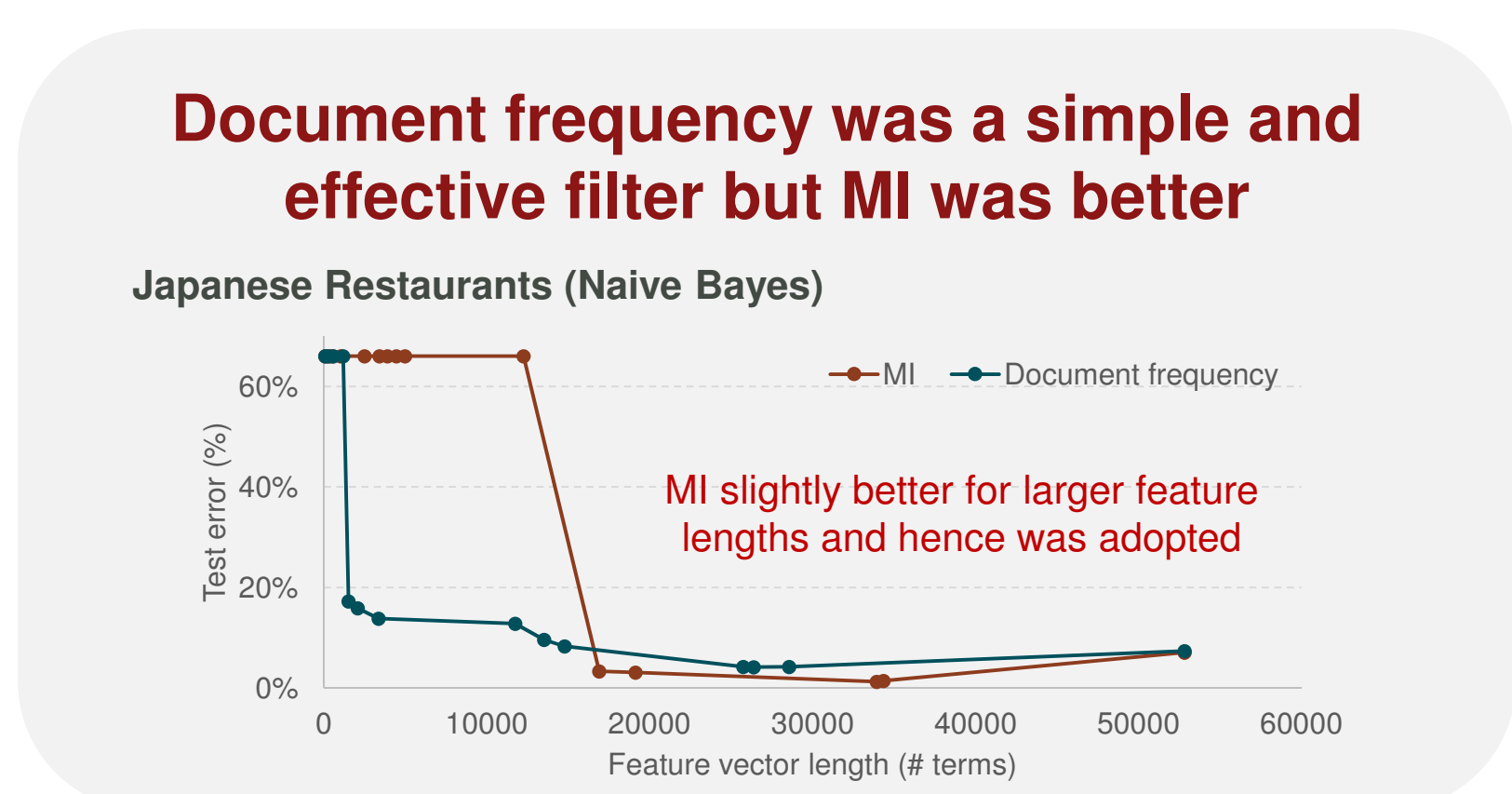
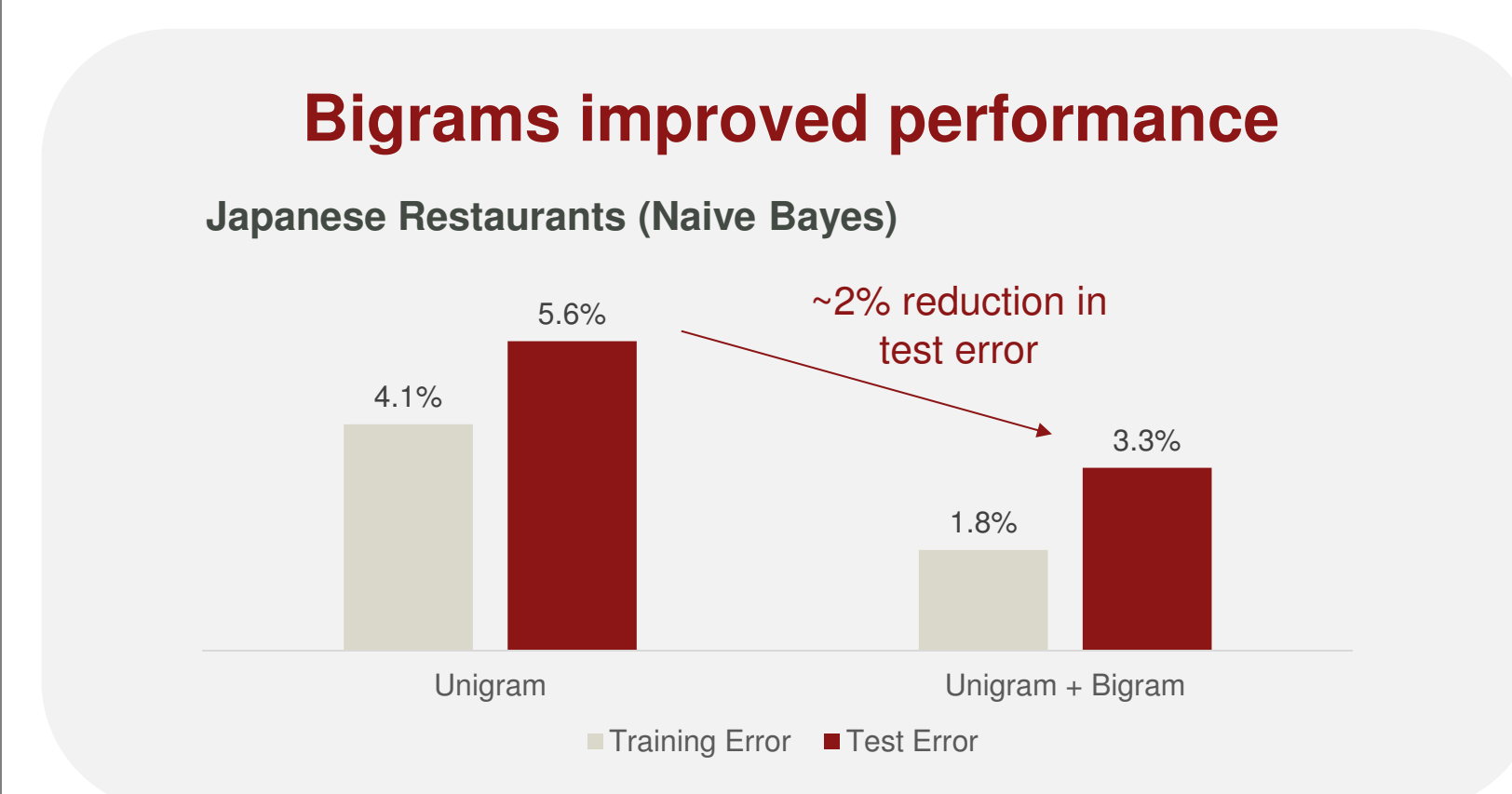
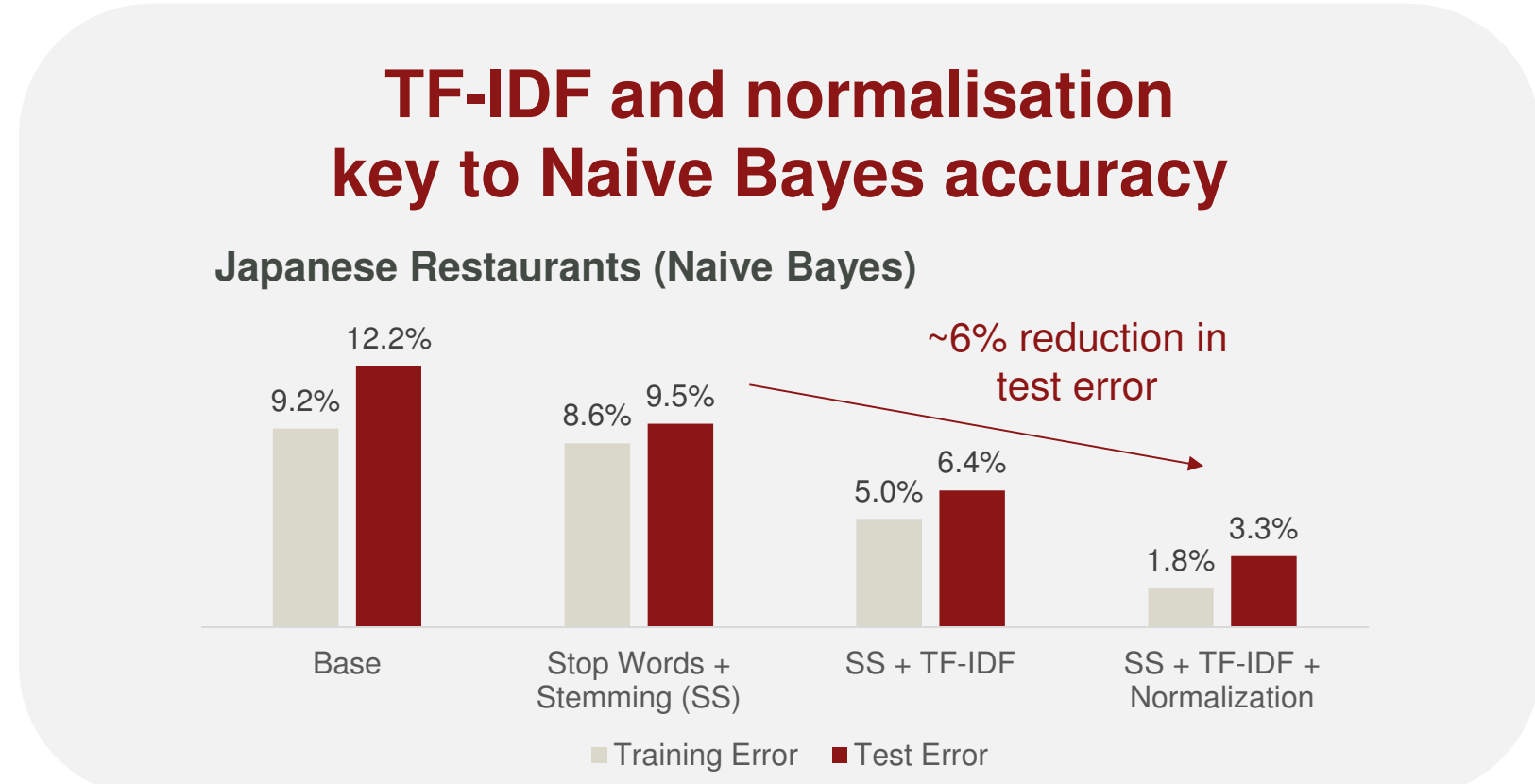
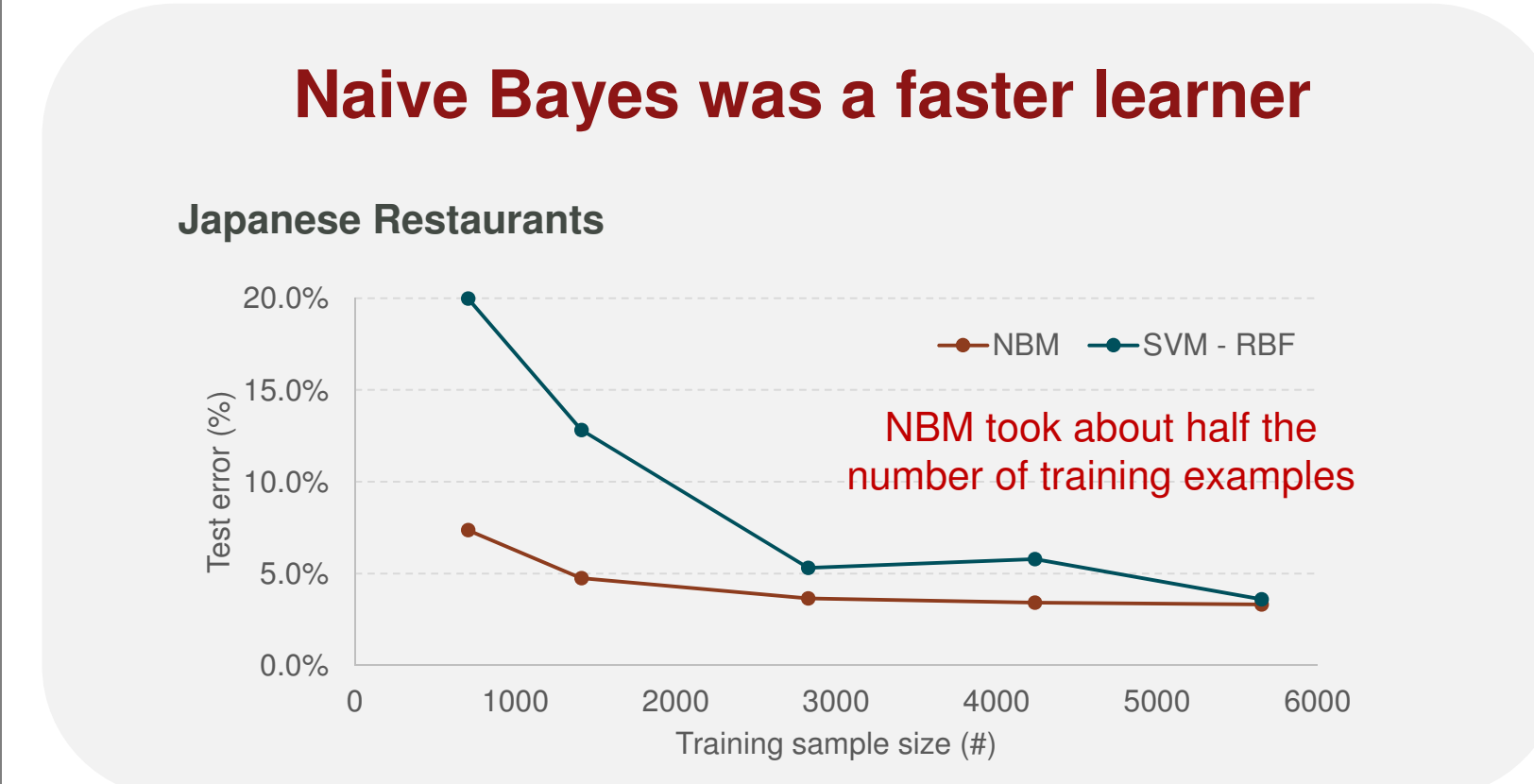
Sample	Training	Test <sup>(1)</sup>
No. of complementary venues (#)	51	51
No. of reviews (#)	5,652	565
Average review length (words)	130	130

Japanese Restaurants	Training	Test <sup>(1)</sup>
Naive Bayes (TF-IDF <sup>(2)</sup> )	1.8%	3.3%
SVM with linear kernel	2.5%	3.8%
SVM with linear kernel (TF-IDF <sup>(2)</sup> )	0.5%	2.4%
SVM with RBF kernel	1.9%	3.6%

Sample	Training	Test <sup>(1)</sup>
No. of complementary venues (#)	20	20
No. of reviews (#)	5,652	565
Average review length (words)	136	136

(1) Test error based on k-fold cross validation with k = 10 over the training sample. Hence, test sample essentially refers to the randomly selected folds, each one-tenth the size of the training sample.  
(2) TF-IDF and normalization.



## Future work

- Second classifier that distinguishes between 4 and 5 stars to provide finer granularity of venue scores
- Vector representations of words (word2vec)
- Exploration of n-grams to explore predictive power of phrases and idioms