



Recipe Recommender

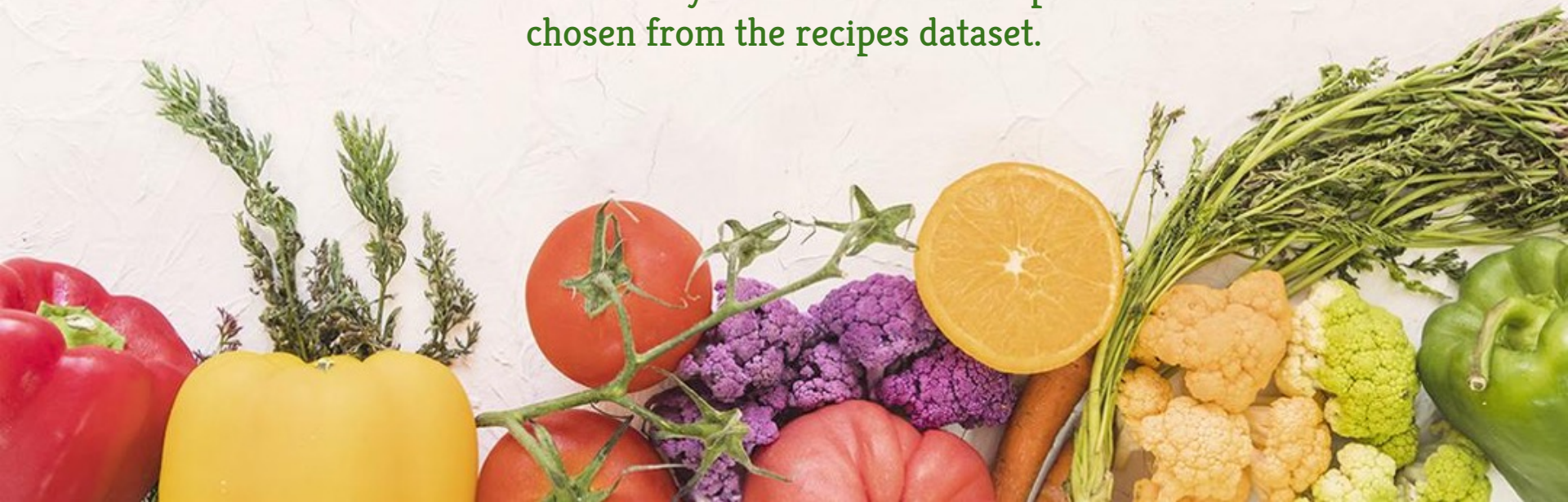
Praveena Suresh
NLP and Unsupervised Learning





Project Goal

The goal of this project is to build a recipe recommendation system based on a recipe that is chosen from the recipes dataset.





Data

The dataset is obtained from Kaggle¹ which includes 2 csv files. One is the interactions file containing recipe IDs rated by user IDs and the other is raw recipes csv file containing 230,000 recipes with names, ID, ingredients, description, steps etc.

Preprocessing



Custom Tokenization

Convert plural to singular and use underscore to keep the full ingredient names to get a better topic model

ground beef, yellow onions, diced tomatoes, tomato paste, tomato soup, kidney beans, chili powder, ground cumin



beef, yellow onion, tomato, tomato paste, kidney beans, chili powder, cumin



beef, yellow_onion, tomato, tomato_paste, kidney_beans, chili_powder, cumin

Workflow

Preprocessed Data

Pandas
Numpy
NLTK

Vectorize

Custom
Tokenization
and Vectorize

Topic Modeling

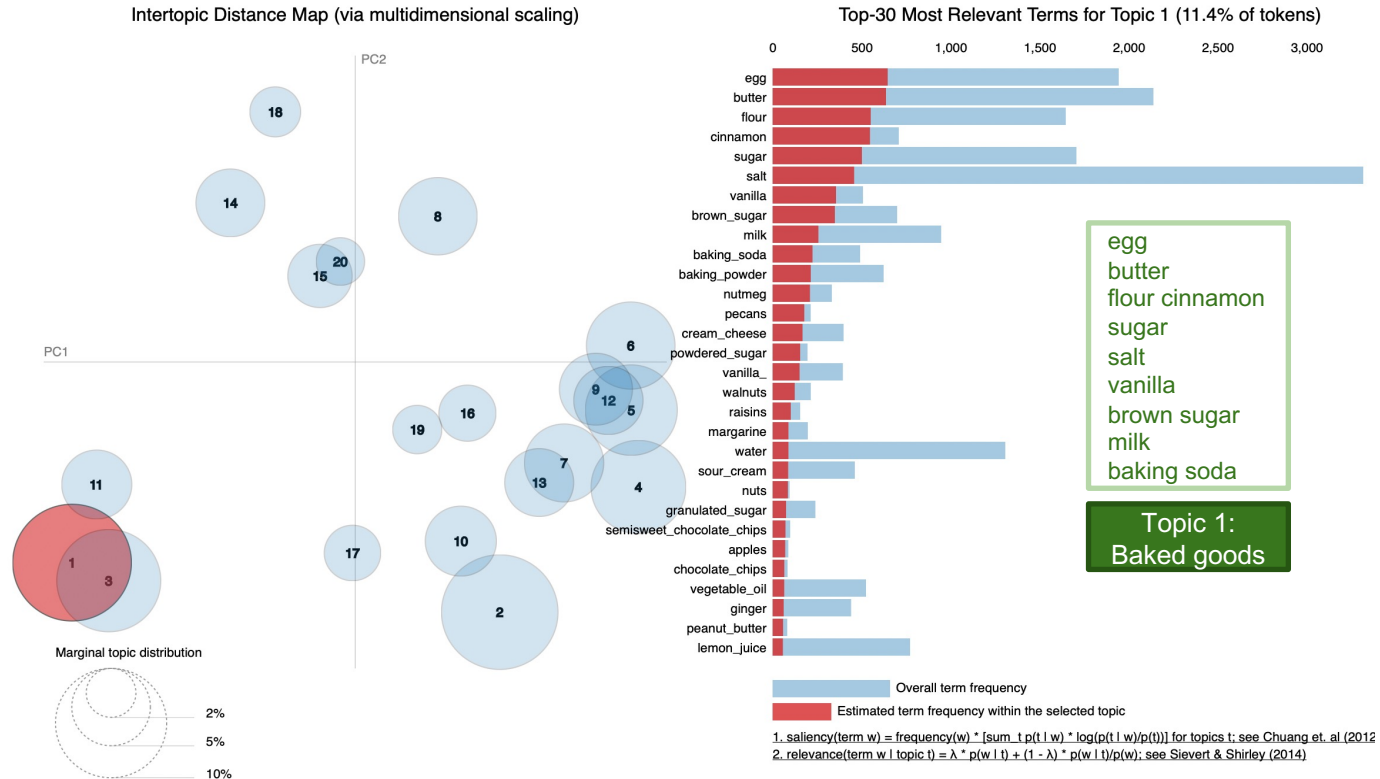
NMF
LDA

Recommender

Cosine
Similarity



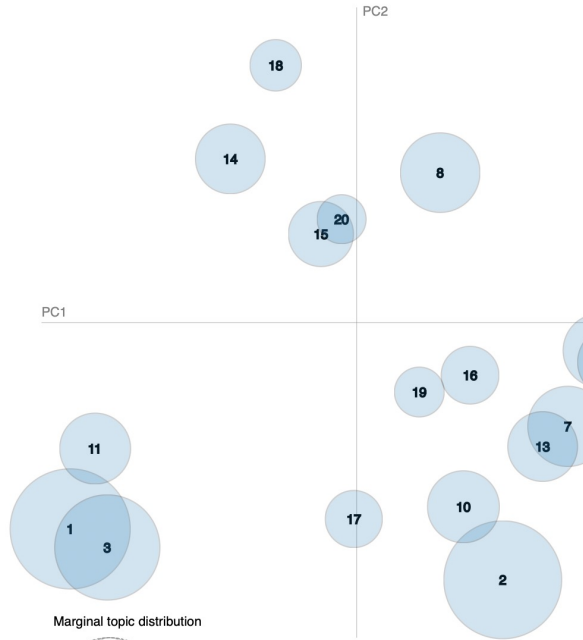
Interactive Topic Visualization with pyLDavis



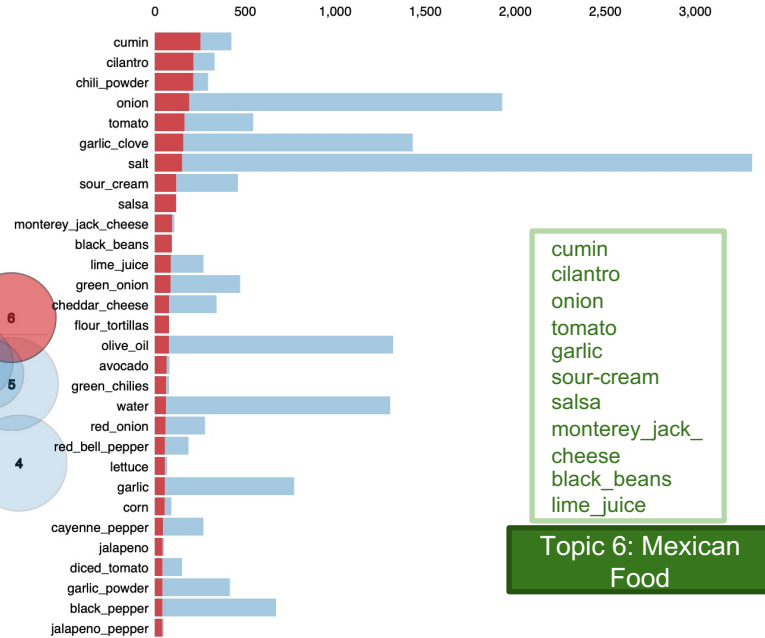
Interactive Topic Visualization with pyLDavis



Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 6 (6.3% of tokens)



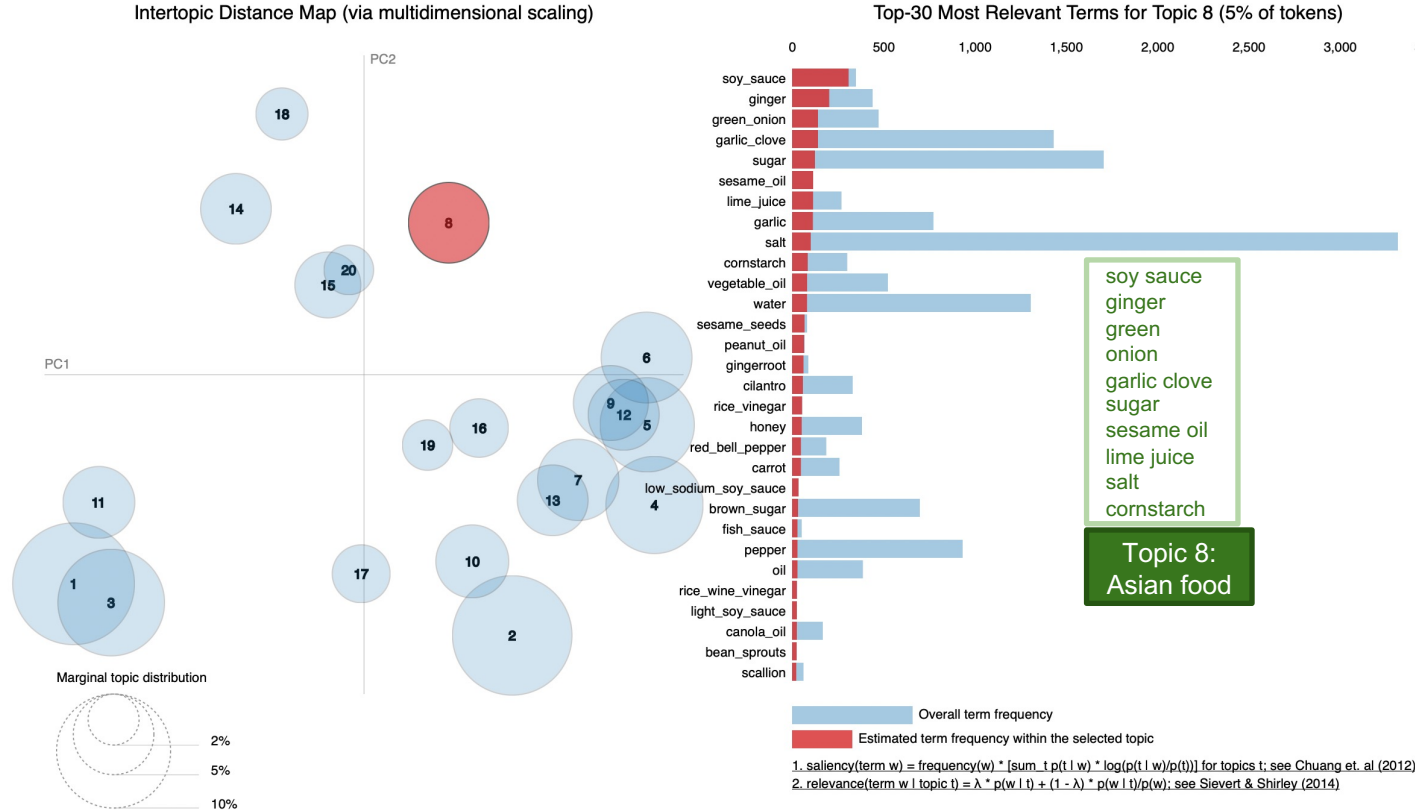
cumin
cilantro
onion
tomato
garlic
sour-cream
salsa
monterey_jack_
cheese
black_beans
lime_juice

Topic 6: Mexican Food

Overall term frequency
Estimated term frequency within the selected topic

1. saliency(term, w) = $\text{frequency}(w) \cdot [\sum_t p(t|w) \cdot \log(p(t|w)/p(t))]$ for topics t : see Chuang et. al (2012)
2. relevance($\text{term } w | \text{topic } t$) = $\lambda \cdot p(w|t) + (1 - \lambda) \cdot p(w|t)/p(w)$: see Sievert & Shirley (2014)

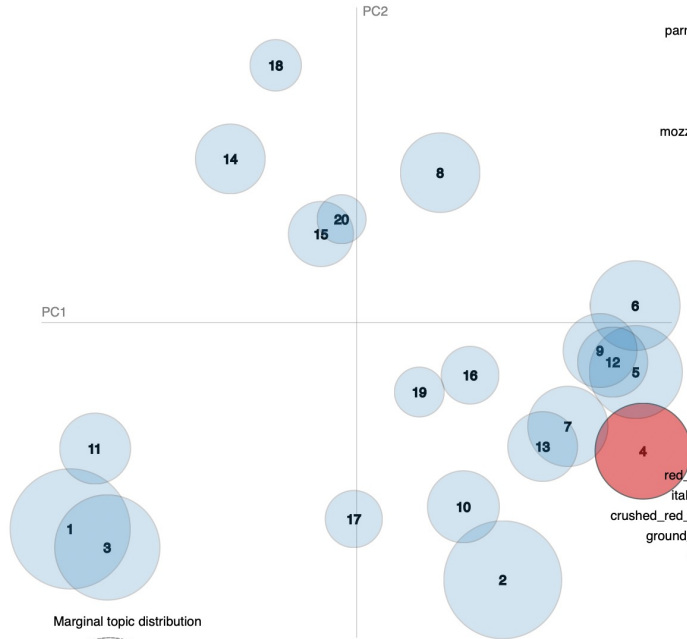
Interactive Topic Visualization with pyLDavis



Interactive Topic Visualization with pyLDavis



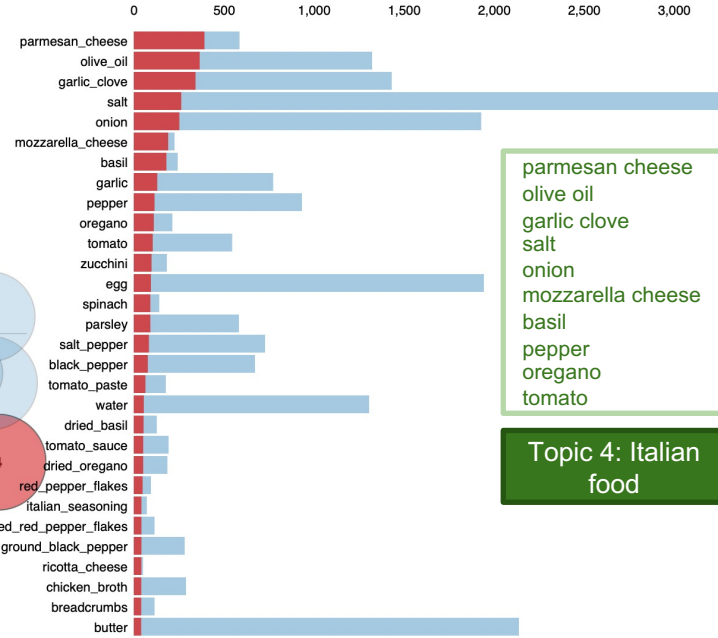
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 4 (7.3% of tokens)



parmesan cheese
olive oil
garlic clove
salt
onion
mozzarella cheese
basil
pepper
oregano
tomato

Topic 4: Italian food

Overall term frequency
Estimated term frequency within the selected topic

1. $s_{\text{allency}}(\text{term } w) = \text{frequency}(w) \cdot [\sum_t p(t|w) \cdot \log(p(t|w)/p(t))]$ for topics t ; see Chuang et. al (2012)
2. $\text{relevance}(\text{term } w | \text{topic } t) = \lambda \cdot p(w | t) + (1 - \lambda) \cdot p(w | t)/p(w)$; see Sievert & Shirley (2014)



Recommender

baby green salad



asian pear salad
blackened steak salad
roasted beet salad with
raspberry balsamic vinaigrette





Recommender

classic hummus



lemon dill hummus light
detoxifying hummus creamy
roasted garlic hummus



Future Work



1. Improve tokenization by more NLP pre-processing.
2. Better model selection by tuning hyperparameters
3. Build a recommender system using collaborative filtering

Appendix



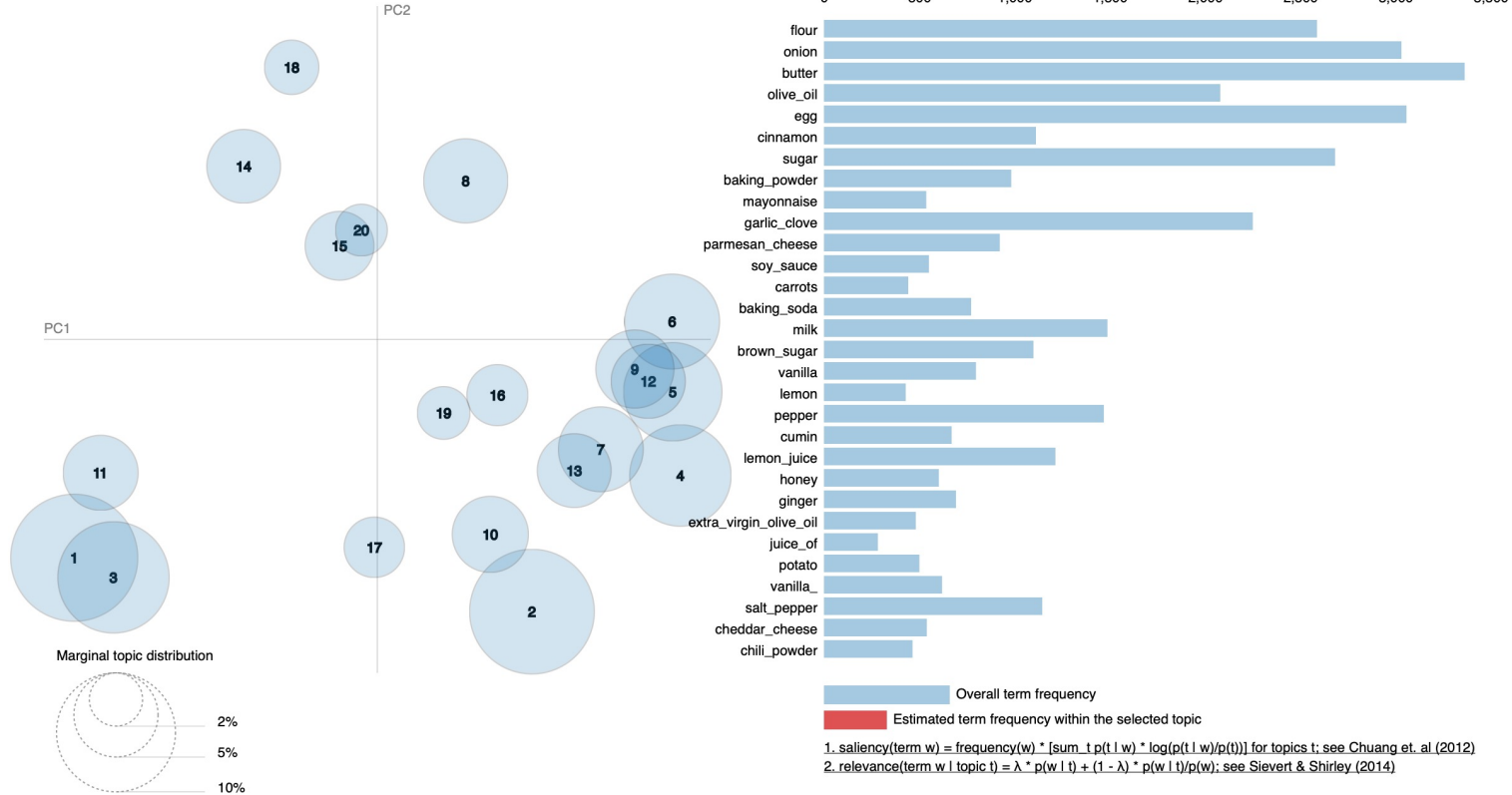
1. Food.com Recipes and Interactions Kaggle dataset

<https://www.kaggle.com/shuyangli94/food-com-recipes-and-user-interactions>

2. Evaluate Topic Models: Latent Dirichlet Allocation (LDA)

<https://towardsdatascience.com/evaluate-topic-model-in-python-latent-dirichlet-allocation-lda-7d57484bb5d0>

Appendix





Thanks!

Do you have any questions?

CREDITS: This presentation template was created by [Slidesgo](#), including icons by [Flaticon](#), and infographics & images by [Freepik](#) and illustrations by [Storyset](#)