

CS412 Project Proposal:

What's Cooking?

Team: AllEatable

1 Team Member

He Huang Ye Liu Zhu Wang Congying Xia Lichao Sun Fan Zhu

2 Dataset

Our project is about a competition named “What’s Cooking?” on Kaggle¹. The dataset is provided by Yummy². Food is an important part of most cultures. Some of our strongest geographic and cultural associations are tied to a region’s local foods. In the dataset, we have a train file and a test file. The train file includes the recipe id, the type of cuisine, and the list of ingredients of each recipe (of variable length). The data is stored in JSON format. In the test file, the format of a recipe is the same as the train file, only the cuisine type is removed, as it is the target variable we are going to predict. In the train file, there are 39774 recipes and 20 kinds of cuisines in total. In the test file, there are 9944 recipes that we need to predict. You can find the competition from the website³. And you can download the dataset from the website⁴.

3 Machine Learning Task

Our goal with this project is to use recipe ingredients to predict the cuisine.

3.1 Task: Multi-class Text Classification

- **Features:** all V different ingredients, and each recipe is represented as a vector of length V .
- **Classes:** all 20 kinds of cuisines.

¹www.kaggle.com

²www.yummly.com

³www.kaggle.com/c/whats-cooking

⁴www.kaggle.com/c/whats-cooking/data

4 Machine Learning Techniques:

1. SVM (multi-class)
2. Random Forest
3. Naive Bayes
4. Class-based Association Rule Mining
5. Neural Network
6. Latent Dirichlet Allocation (LDA)

4.1 Visualization

To reduce V -dimensional vectors to 2- or 3- dimensional ones

4.1.1 Visualization Methods:

- Principle Component Analysis (PCA)
- t-SNE

5 Contributions

1. Using many traditional machine learning techniques.
2. Comparing the performances of each used technique on this task.
3. Visualizing the recipes in 2-D/3-D space to show the distances of different cuisines.