# CS6200 Project Proposal

Yuting Sun

## Background

It is always a challenging task to try out a new cosmetic product. Different people have different skin types and product preference. A product that is good for one person may not be suitable for another. Improper products would waste a lot of money and even cause severe skin problems. Since the COVID-19 pandemic, it's not convinient for people to go to the stores and try the products on themselves. Therefore, it would be helpful if we could build a searching tool for cosmetic products.

Ingredients are key to the property and effectiveness of a cosmetic product. For example, salicylic acid is widely used in treating acnes, titanium dioxide is used for UV absorbance in sunscreen products. Ingredient analysis can be very helpful for people to search for proper products. Some people likes the certain propety of a product, so s/he would propably like a product with similar ingredients. Some people are allergic to a certain ingredient, so it would be necessary to avoid products with the same ingredient. However, most cosmetic products contain more than 30 ingredients, which is difficult for people to read and analysis. It would be a good idea to combine information retrieval techniques with ingredient analysis. This would help customers find out new product that are most suitable for them, as well as getting more insights on cosmetic ingredients.

## Solution

In this project, I would like to propose a searching tool for cosmetic products that contains two parts:

**1) Search for relevant products based on product type, name, and description.**

- The product category (eg. skincare, makeup, etc) can be used as filters for search.
- Product category and brand are indexed as single terms.

  Product name and description are indexed with n-gram overlap (where n=1 to 3) with tf-idf weighting. Stopwords are removed.
- Top-K relevant results can be obtained by BM25 algorithm using python ElasticSearch. The results are by default sorted by relevant scores.

**2) Recommend similar product based on ingredient analysis.**

- Ingredients of each product is indexed and vectorized as a single document. Similarity is calculated via cosine score.
- When a user opens a product page, a list of similar products (sorted by similarity) are shown at the right side.

## Dataset

The Sephora dataset can be obtained from [Kaggle](#). This dataset contains 9000+ products with product information, detailed description and ratings. More than 7500 products have ingredient information.

Product Information Example:



## Future Improvement

The current dataset contains 9k products. If product description and ingredients are considered as different documents, we have 18k documents in total. This might be a bit small for the dataset.

The current dataset contains number of reviews for each product but does not contain review text. The total number of reviews is ~2.7 million. Therefore, we can expand the existing dataset by parsing the review text to conduct better search for relevant product.