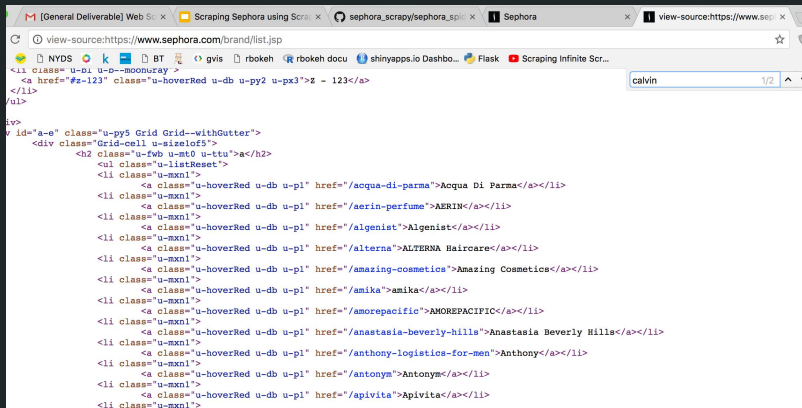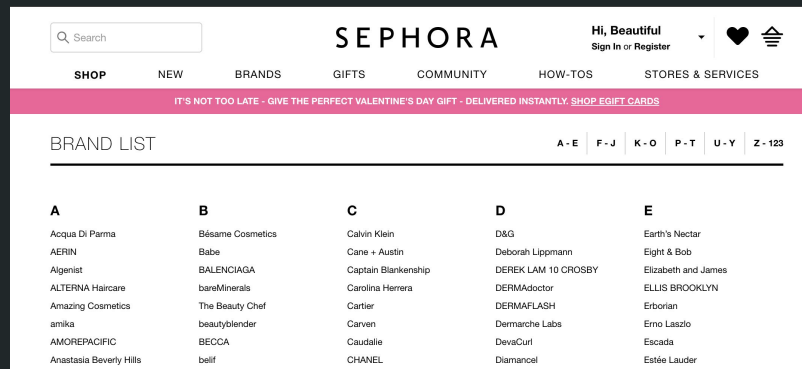# What is my goal for this project?

1. Scrape Sephora
2. Learning how to use Flask and creating an Application which can do simple queries of Sephora reviews data using Flask and sqlite3

# How to Scrape Sephora's website

1. Getting the brand links
2. Getting the product links
3. Getting product details
4. Constructing the API link using information from (#2)
5. Getting data from the API - Reviews and Reviewer details
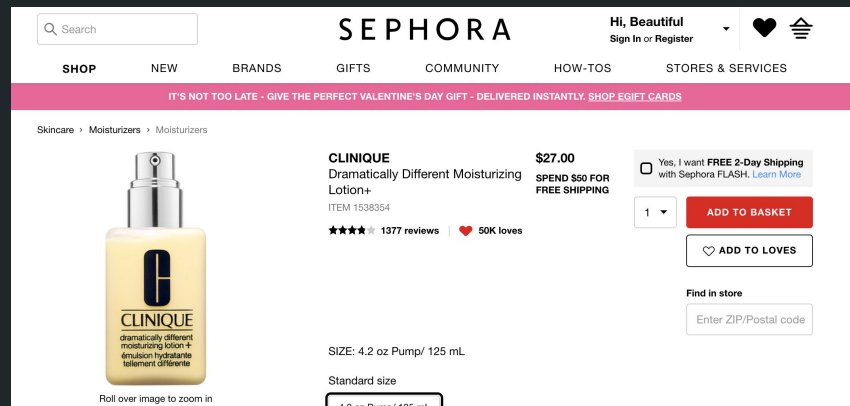
# Getting the brand links

*https://www.sephora.com/brand/list.jsp*

# Getting the product links



view-source:https://www.sephora.com/clinique?products=all

# Getting the product details

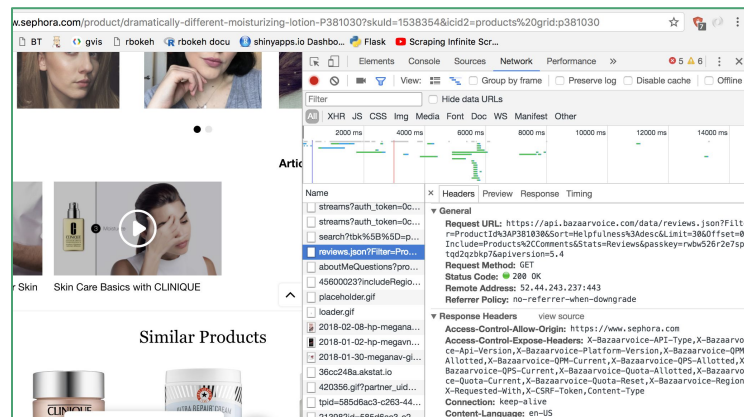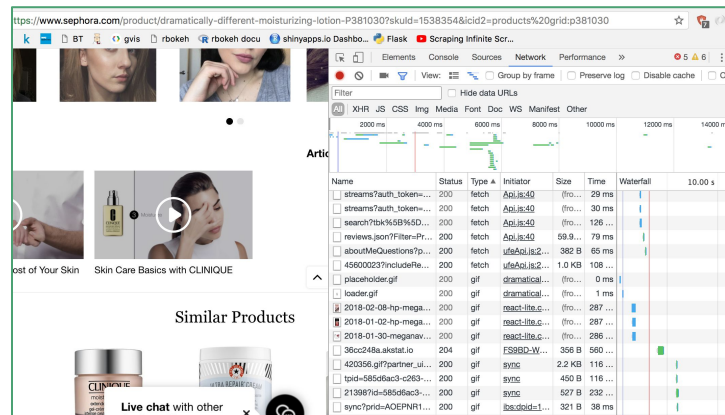<html lang="en" class="css-1mok2x0" data-comp="Index"><head data-comp="Head"><title>Dramatically Different Moisturizing Lotion+ – CLINIQUE | Sephora</title><meta name="viewport" content="width=1024"/><meta name="description" content="Shop CLINIQUE's Dramatically Different Moisturizing Lotion+ at Sephora. This fast-absorbing, lightweight daily moisturizer is for dry to combination skin types."/><link rel="canonical" href="https://www.sephora.com/product/dramatically-different-moisturizing-lotion-P381030"/><link rel="alternate" media="only screen and (max-width: 640px)" href="https://m.sephora.com/product/dramatically-different-moisturizing-lotion-P381030"/><meta name="og:description" content="Shop CLINIQUE's Dramatically Different Moisturizing Lotion+ at Sephora. This fast-absorbing, lightweight daily moisturizer is for dry to combination skin types."/><meta name="og:title" content="Dramatically Different Moisturizing Lotion+ – CLINIQUE | Sephora"/><meta name="og:type" content="website"/><meta name="og:url" content="https://www.sephora.com/product/dramatically-different-moisturizing-lotion-P381030"/><meta name="og:image" content="/productimages/sku/s1538354-main-hero-300.jpg"/><meta name="apple-mobile-web-app-capable" content="yes"/><meta name="format-detection" content="telephone=no"/><script>if (typeof global === "undefined") window.global = window;global.Sephora = global.Sephora || {};Sephora.targetersToInclude = "?";Sephora.template = "Product/ProductPage";Sephora.renderedData = {"rendered":"2018-02-13 04:01:39,850","template":"Product/ProductPage","channelProp":"FS","renderHost":"ph626201.bwi40g","pageRenderTime":167.293};Sephora.renderQueryParams = {"hash":"e65809e2a080be4ade592b80131166101083fc08","channel":"FS","country":"US","urlPath":"%2Fproduct%2Fdramatically-different-moisturizing-lotion-P381030"};</script><script>Sephora.productPage = { defaultSkuId: 1538354 }</script><script>"use strict";try{var ce=new window.CustomEvent("test");if(ce.preventDefault(),!0!==ce.defaultPrevented)throw new Error("Could not prevent default")}catch(e){var CustomEvent=function(e,t){var n,r;return t=t||{bubbles:!1,cancelable:!1,detail:void 0},n=document.createEvent("CustomEvent"),n.initCustomEvent(e,t.bubbles,t.cancelable,t.detail),r=n.preventDefault,n.preventDefault=function()

https://www.sephora.com/product/dramatically-different-moisturizing-lotion-P381030?skuId=1538354&icid2=products%20grid:p381030

# Getting the Reviews - Infinite Scroll Problem

How to solve this?

A. Inspect Element

B. Go to the Network, reload the page

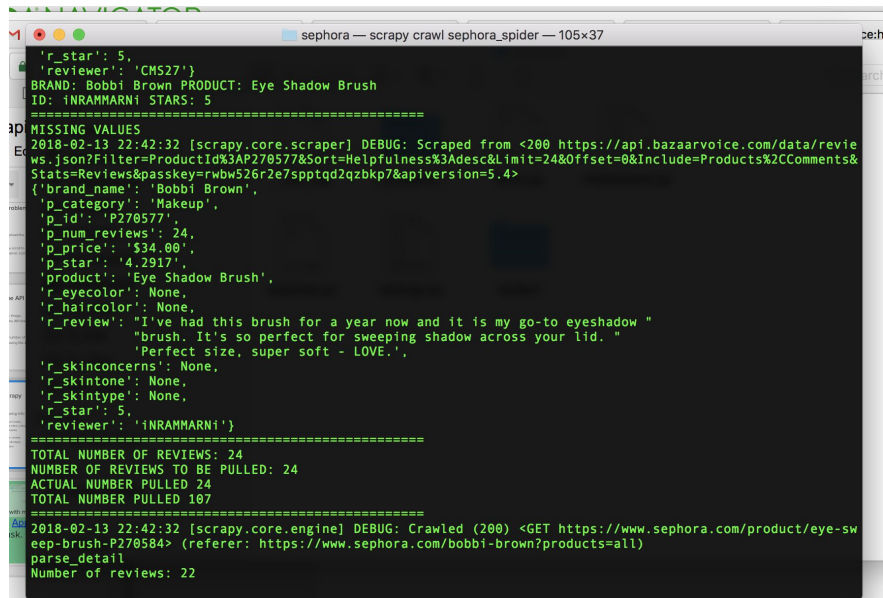C. Notice that when you scroll to make the reviews appear, a json fetch item appears

# Constructing the API link

- As we can see in the image, there's a pattern to the API link

- Product ID

- And we can set the number of reviews we can call using the API

https://api.bazaarvoice.com/data/reviews.json?Filter=ProductId%3A**P381030**&Sort=Helpfulness%3Adesc&**Limit=30**&Offset=0&Include=Products%2CComments&Stats=Reviews&passkey=rwbw526r2e7spptqd2qzbkp7&apiversion=5.4

# Running the Scrapy program

- I collected the following info:

    ○ Brand Name, product name, product ID, average stars, category, price, number of reviews

    ○ Reviewer nickname, review, reviewer skintone, skintype, haircolor and eyecolor

    ○ I collected around **80k+** rows of data

What did I do with my data? I made an [Application](#) using Flask.

# Conclusion

To sum things up, I was able to show that Sephora's website can be scraped relatively easily and I learned how to create an Application using Flask and sqlite3, achieving my goals for the project.