

Amazon Apparel Recommendations

In this project our objective is to recommend similar products/items in e-commerce website Amazon. So basically what we are doing here is context based recommendation.

WHY IS IT IMPORTANT:-

- Around 35% of the total revenue at Amazon.com is generated using product recommendation which is close to 40 billion dollars, which is a whopping amount. That's why it becomes important to understand how we will use previous and existing data for recommending the similar products.

[4.2] Data and Code:

<https://drive.google.com/open?id=0BwNkduBnePt2VWhCYXhMV3p4dTg>
(<https://drive.google.com/open?id=0BwNkduBnePt2VWhCYXhMV3p4dTg>)

[4.3] Overview of the data

```
In [1]: 1 #import all the necessary packages.
2
3 from PIL import Image
4 import requests
5 from io import BytesIO
6 import matplotlib.pyplot as plt
7 import numpy as np
8 import pandas as pd
9 import warnings
10 from bs4 import BeautifulSoup
11 from nltk.corpus import stopwords
12 from nltk.tokenize import word_tokenize
13 import nltk
14 import math
15 import time
16 import re
17 import os
18 import seaborn as sns
19 from collections import Counter
20 from sklearn.feature_extraction.text import CountVectorizer
21 from sklearn.feature_extraction.text import TfidfVectorizer
22 from sklearn.metrics.pairwise import cosine_similarity
23 from sklearn.metrics import pairwise_distances
24 from matplotlib import gridspec
25 from scipy.sparse import hstack
26 import plotly
27 import plotly.figure_factory as ff
28 from plotly.graph_objs import Scatter, Layout
29
30 plotly.offline.init_notebook_mode(connected=True)
31 warnings.filterwarnings("ignore")
```

```
In [0]: 1 # we have give a json file which consists of all information about
2 # the products
3 # Loading the data using pandas' read_json file.
4 data = pd.read_json('tops_fashion.json')
5
```

```
In [0]: 1 print ('Number of data points : ', data.shape[0], \
2         'Number of features/variables:', data.shape[1])
```

Number of data points : 183138 Number of features/variables: 19

Terminology:

What is a dataset?

Rows and columns

Data-point

Feature/variable

```
In [0]: 1 # each product/item has 19 features in the raw dataset.  
2 data.columns # prints column-names or feature-names.
```

```
Out[35]: Index(['asin', 'author', 'availability', 'availability_type', 'brand', 'color',  
    'editorial_review', 'editorial_review', 'formatted_price',  
    'large_image_url', 'manufacturer', 'medium_image_url', 'model',  
    'product_type_name', 'publisher', 'reviews', 'sku', 'small_image_url',  
    'title'],  
    dtype='object')
```

Of these 19 features, we will be using only 6 features in this workshop.

1. asin (Amazon standard identification number)
2. brand (brand to which the product belongs to)
3. color (Color information of apparel, it can contain many colors as a value ex: red and black stripes)
4. product_type_name (type of the apparel, ex: SHIRT/TSHIRT)
5. medium_image_url (url of the image)
6. title (title of the product.)
7. formatted_price (price of the product)

```
In [0]: 1 data = data[['asin', 'brand', 'color', 'medium_image_url', 'product_type_name']]
```

```
In [0]: 1 print ('Number of data points : ', data.shape[0], \
2       'Number of features:', data.shape[1])
3 data.head() # prints the top rows in the table.
```

Number of data points : 183138 Number of features: 7

Out[37]:

	asin	brand	color	medium_image_url	product_type_name	title	formatte...
0	B016I2TS4W	FNC7C	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	Minions Como Superheroes Ironman Long Sleeve R...	
1	B01N49AI08	FIG Clothing	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	FIG Clothing Womens Izo Tunic	
2	B01JDPCOHO	FIG Clothing	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	FIG Clothing Womens Won Top	
3	B01N19U5H5	Focal18	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	Focal18 Sailor Collar Bubble Sleeve Blouse Shi...	
4	B004GSI2OS	FeatherLite	Onyx Black/ Stone	https://images-na.ssl-images-amazon.com/images...	SHIRT	Featherlite Ladies' Long Sleeve Stain Resistan...	

[5.1] Missing data for various features.

Basic stats for the feature: product_type_name

```
In [0]: 1 # We have total 72 unique type of product_type_names
2 print(data['product_type_name'].describe())
3
4 # 91.62% (167794/183138) of the products are shirts,
5
```

count	183138
unique	72
top	SHIRT
freq	167794
Name:	product_type_name, dtype: object

In [0]:

```

1 # names of different product types
2 print(data['product_type_name'].unique())

```

['SHIRT' 'SWEATER' 'APPAREL' 'OUTDOOR_RECREATION_PRODUCT'
 'BOOKS_1973_AND_LATER' 'PANTS' 'HAT' 'SPORTING_GOODS' 'DRESS' 'UNDERWEAR'
 'SKIRT' 'OUTERWEAR' 'BRA' 'ACCESSORY' 'ART_SUPPLIES' 'SLEEPWEAR'
 'ORCA_SHIRT' 'HANDBAG' 'PET_SUPPLIES' 'SHOES' 'KITCHEN' 'ADULT_COSTUME'
 'HOME_BED_AND_BATH' 'MISC_OTHER' 'BLAZER' 'HEALTH_PERSONAL_CARE'
 'TOYS_AND_GAMES' 'SWIMWEAR' 'CONSUMER_ELECTRONICS' 'SHORTS' 'HOME'
 'AUTO_PART' 'OFFICE_PRODUCTS' 'ETHNIC_WEAR' 'BEAUTY'
 'INSTRUMENT_PARTS_AND_ACCESSORIES' 'POWERSPORTS_PROTECTIVE_GEAR' 'SHIRTS'
 'ABIS_APPAREL' 'AUTO_ACCESSORY' 'NONAPPARELMISC' 'TOOLS' 'BABY_PRODUCT'
 'SOCKSHOSIERY' 'POWERSPORTS RIDING SHIRT' 'EYEWEAR' 'SUIT'
 'OUTDOOR_LIVING' 'POWERSPORTS RIDING JACKET' 'HARDWARE' 'SAFETY_SUPPLY'
 'ABIS_DVD' 'VIDEO_DVD' 'GOLF_CLUB' 'MUSIC_POPULAR_VINYL'
 'HOME_FURNITURE_AND_DECOR' 'TABLET_COMPUTER' 'GUILD_ACCESSORIES'
 'ABIS_SPORTS' 'ART_AND_CRAFT_SUPPLY' 'BAG' 'MECHANICAL_COMPONENTS'
 'SOUND_AND_RECORDING_EQUIPMENT' 'COMPUTER_COMPONENT' 'JEWELRY'
 'BUILDING_MATERIAL' 'LUGGAGE' 'BABY_COSTUME' 'POWERSPORTS_VEHICLE_PART'
 'PROFESSIONAL_HEALTHCARE' 'SEEDS_AND_PLANTS' 'WIRELESS_ACCESSORY']

In [0]:

```

1 # find the 10 most frequent product_type_names.
2 product_type_count = Counter(list(data['product_type_name']))
3 product_type_count.most_common(10)

```

Out[40]: [('SHIRT', 167794),
 ('APPAREL', 3549),
 ('BOOKS_1973_AND_LATER', 3336),
 ('DRESS', 1584),
 ('SPORTING_GOODS', 1281),
 ('SWEATER', 837),
 ('OUTERWEAR', 796),
 ('OUTDOOR_RECREATION_PRODUCT', 729),
 ('ACCESSORY', 636),
 ('UNDERWEAR', 425)]

Basic stats for the feature: brand

In [0]:

```

1 # there are 10577 unique brands
2 print(data['brand'].describe())
3
4 # 183138 - 182987 = 151 missing values.

```

count	182987
unique	10577
top	Zago
freq	223
Name:	brand, dtype: object

```
In [0]: 1 brand_count = Counter(list(data['brand']))
2 brand_count.most_common(10)
```

```
Out[42]: [('Zago', 223),
('XQS', 222),
('Yayun', 215),
('YUNY', 198),
('XiaoTianXin-women clothes', 193),
('Generic', 192),
('Boohoo', 190),
('Alion', 188),
('Abetteric', 187),
('TheMogan', 187)]
```

Basic stats for the feature: color

```
In [0]: 1
2 print(data['color'].describe())
3
4
5 # we have 7380 unique colors
6 # 7.2% of products are black in color
7 # 64956 of 183138 products have brand information. That's approx 35.4%.
```

```
count      64956
unique     7380
top        Black
freq       13207
Name: color, dtype: object
```

```
In [0]: 1 color_count = Counter(list(data['color']))
2 color_count.most_common(10)
```

```
Out[44]: [(None, 118182),
('Black', 13207),
('White', 8616),
('Blue', 3570),
('Red', 2289),
('Pink', 1842),
('Grey', 1499),
('*', 1388),
('Green', 1258),
('Multi', 1203)]
```

Basic stats for the feature: formatted_price

In [0]:

```

1 print(data['formatted_price'].describe())
2
3
4 # Only 28,395 (15.5% of whole data) products with price information

```

```

count      28395
unique     3135
top        $19.99
freq       945
Name: formatted_price, dtype: object

```

In [0]:

```

1 price_count = Counter(list(data['formatted_price']))
2 price_count.most_common(10)

```

Out[46]:

```

[(None, 154743),
 ('$19.99', 945),
 ('$9.99', 749),
 ('$9.50', 601),
 ('$14.99', 472),
 ('$7.50', 463),
 ('$24.99', 414),
 ('$29.99', 370),
 ('$8.99', 343),
 ('$9.01', 336)]

```

Basic stats for the feature: title

In [0]:

```

1 print(data['title'].describe())
2
3 # All of the products have a title.
4 # Titles are fairly descriptive of what the product is.
5 # We use titles extensively in this workshop
6 # as they are short and informative.
7

```

```

count          183138
unique         175985
top           Nakoda Cotton Self Print Straight Kurti For Women
freq            77
Name: title, dtype: object

```

In [0]:

```

1 data.to_pickle('pickels/180k_apparel_data')

```

We save data files at every major step in our processing in "pickle" files. If you are stuck anywhere (or) if some code takes too long to run on your laptop, you may use the pickle files we give you to speed things up.

```
In [0]: 1 # consider products which have price information
2 # data['formatted_price'].isnull() => gives the information
3 #about the dataframe row's which have null values price == None/Null
4 data = data.loc[~data['formatted_price'].isnull()]
5 print('Number of data points After eliminating price=NULL :', data.shape[0])
```

Number of data points After eliminating price=NULL : 28395

```
In [0]: 1 # consider products which have color information
2 # data['color'].isnull() => gives the information about the dataframe row's w
3 data = data.loc[~data['color'].isnull()]
4 print('Number of data points After eliminating color=NULL :', data.shape[0])
```

Number of data points After eliminating color=NULL : 28385

We brought down the number of data points from 183K to 28K.

We are processing only 28K points so that most of the workshop participants can run this code on their laptops in a reasonable amount of time.

For those of you who have powerful computers and some time to spare, you are recommended to use all of the 183K images.

```
In [0]: 1 data.to_pickle('pickels/28k_apparel_data')
```

```
In [0]: 1 # You can download all these 28k images using this code below.
2 # You do NOT need to run this code and hence it is commented.
3
4
5 ...
6 from PIL import Image
7 import requests
8 from io import BytesIO
9
10 for index, row in images.iterrows():
11     url = row['large_image_url']
12     response = requests.get(url)
13     img = Image.open(BytesIO(response.content))
14     img.save('images/28k_images/'+row['asin']+'.jpeg')
15
16
17 ...
```

```
Out[52]: "\nfrom PIL import Image\nimport requests\nfrom io import BytesIO\n\nfor index,\nrow in images.iterrows():\n    url = row['large_image_url']\n    response = requests.get(url)\n    img = Image.open(BytesIO(response.content))\n    img.save('workshop/images/28k_images/'+row['asin']+'.jpeg')\n\n\n"
```

[5.2] Remove near duplicate items

[5.2.1] Understand about duplicates.

In [0]:

```

1 # read data from pickle file from previous stage
2 data = pd.read_pickle('pickels/28k_apparel_data')
3
4 # find number of products that have duplicate titles.
5 print(sum(data.duplicated('title')))
6 # we have 2325 products which have same title but different color
7

```

2325

These shirts are exactly same except in size (S, M,L,XL)

 :B00AQ4GMCK :B00AQ4GM

 :B00AQ4GMLQ :B00AQ4GN

These shirts exactly same except in color

 :B00G278GZ6 :B00G278W6C

 :B00G278Z2A :B00G2786X8

In our data there are many duplicate products like the above examples, we need to de-dupe them for better results.

[5.2.2] Remove duplicates : Part 1

In [0]:

```

1 # read data from pickle file from previous stage
2 data = pd.read_pickle('pickels/28k_apparel_data')

```

In [0]: 1 data.head()

Out[103]:

	asin	brand	color	medium_image_url	product_type_name	title	format
4	B004GSI2OS	FeatherLite	Onyx Black/ Stone	https://images-na.ssl-images-amazon.com/images...	SHIRT	Featherlite Ladies' Long Sleeve Stain Resistan...	
6	B012YX2ZPI	HX-Kingdom Fashion T-shirts	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	Women's Unique 100% Cotton T - Special Olympic...	
11	B001LOUGE4	Fitness Etc.	Black	https://images-na.ssl-images-amazon.com/images...	SHIRT	Ladies Cotton Tank 2x1 Ribbed Tank Top	
15	B003BSRPB0	FeatherLite	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	FeatherLite Ladies' Moisture Free Mesh Sport S...	
21	B014ICEDNA	FNC7C	Purple	https://images-na.ssl-images-amazon.com/images...	SHIRT	Supernatural Chibis Sam Dean And Castiel Short...	

In [0]:

```
1 # Remove ALL products with very few words in title
2 data_sorted = data[data['title'].apply(lambda x: len(x.split())>4)]
3 print("After removal of products with short description:", data_sorted.shape[
```

After removal of products with short description: 27949

```
In [0]: 1 # Sort the whole data based on title (alphabetical order of title)
2 data_sorted.sort_values('title',inplace=True, ascending=False)
3 data_sorted.head()
```

Out[105]:

		asin	brand	color	medium_image_url	product_type_name	title	fc
61973	B06Y1KZ2WB	Éclair	Black/Pink		https://images-na.ssl-images-amazon.com/images...	SHIRT	Éclair Women's Printed Thin Strap Blouse Black...	
133820	B010RV33VE	xiaoming		Pink	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Womens Sleeveless Loose Long T-shirts...	
81461	B01DDSDLNS	xiaoming		White	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Women's White Long Sleeve Single Brea...	
75995	B00X5LYO9Y	xiaoming	Red Anchors		https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Stripes Tank Patch/Bear Sleeve Anchor...	
151570	B00WPJG35K	xiaoming		White	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Sleeve Sheer Loose Tassel Kimono Woma...	



Some examples of duplicate titles that differ only in the last few words.

Titles 1:

- 16. woman's place is in the house and the senate shirts for Womens XXL White
- 17. woman's place is in the house and the senate shirts for Womens M Grey

Title 2:

25. tokidoki The Queen of Diamonds Women's Shirt X-Large
26. tokidoki The Queen of Diamonds Women's Shirt Small
27. tokidoki The Queen of Diamonds Women's Shirt Large

Title 3:

61. psychedelic colorful Howling Galaxy Wolf T-shirt/Colorful Rainbow Animal Print Head Shirt for woman Neon Wolf t-shirt
62. psychedelic colorful Howling Galaxy Wolf T-shirt/Colorful Rainbow Animal Print Head Shirt for woman Neon Wolf t-shirt
63. psychedelic colorful Howling Galaxy Wolf T-shirt/Colorful Rainbow Animal Print Head Shirt for woman Neon Wolf t-shirt
64. psychedelic colorful Howling Galaxy Wolf T-shirt/Colorful Rainbow Animal Print Head Shirt for woman Neon Wolf t-shirt

In [0]:

```
1 indices = []
2 for i, row in data_sorted.iterrows():
3     indices.append(i)
```

```
In [0]: 1 import itertools
2 stage1_dedupe_asins = []
3 i = 0
4 j = 0
5 num_data_points = data_sorted.shape[0]
6 while i < num_data_points and j < num_data_points:
7
8     previous_i = i
9
10    # store the list of words of ith string in a, ex: a = ['tokidoki', 'The',
11    a = data['title'].loc[indices[i]].split()
12
13    # search for the similar products sequentially
14    j = i+1
15    while j < num_data_points:
16
17        # store the list of words of jth string in b, ex: b = ['tokidoki', 'T
18        b = data['title'].loc[indices[j]].split()
19
20        # store the maximum length of two strings
21        length = max(len(a), len(b))
22
23        # count is used to store the number of words that are matched in both
24        count = 0
25
26        # itertools.zip_longest(a,b): will map the corresponding words in bot
27        # example: a =['a', 'b', 'c', 'd']
28        # b = ['a', 'b', 'd']
29        # itertools.zip_longest(a,b): will give [('a', 'a'), ('b', 'b'), ('c', '
30        for k in itertools.zip_longest(a,b):
31            if (k[0] == k[1]):
32                count += 1
33
34            # if the number of words in which both strings differ are > 2 , we ar
35            # if the number of words in which both strings differ are < 2 , we ar
36            if (length - count) > 2: # number of words in which both sensences di
37                # if both strings are differ by more than 2 words we include the
38                stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[i]])
39
40
41            # start searching for similar apperals corresponds 2nd string
42            i = j
43            break
44        else:
45            j += 1
46    if previous_i == i:
47        break
```

```
In [0]: 1 data = data.loc[data['asin'].isin(stage1_dedupe_asins)]
```

We removed the duplicates which differ only at the end.

```
In [0]: 1 print('Number of data points : ', data.shape[0])
```

Number of data points : 17593

```
In [0]: 1 data.to_pickle('pickels/17k_apperal_data')
```

[5.2.3] Remove duplicates : Part 2

In the previous cell, we sorted whole data in alphabetical order of titles. Then, we removed titles which are adjacent and very similar title

But there are some products whose titles are not adjacent but very similar.

Examples:

Titles-1

86261. UltraClub Women's Classic Wrinkle-Free Long Sleeve Oxford Shirt, Pink, XX-Large

115042. UltraClub Ladies Classic Wrinkle-Free Long-Sleeve Oxford Light Blue XXL

Titles-2

75004. EVALY Women's Cool University Of UTAH 3/4 Sleeve Raglan Tee

109225. EVALY Women's Unique University Of UTAH 3/4 Sleeve Raglan Tees

120832. EVALY Women's New University Of UTAH 3/4-Sleeve Raglan Tshirt

```
In [0]: 1 data = pd.read_pickle('pickels/17k_apperal_data')
```

one thing we must keep in mind here is that we are recommending a product of other brand and category, so color, size and price of two shirts of same

```
In [0]: 1 # This code snippet takes significant amount of time.
2 # O(n^2) time.
3 # Takes about an hour to run on a decent computer.
4
5 indices = []
6 for i, row in data.iterrows():
7     indices.append(i)
8
9 stage2_dedupe_asins = []
10 while len(indices)!=0:
11     i = indices.pop()
12     stage2_dedupe_asins.append(data['asin'].loc[i])
13     # consider the first apperal's title
14     a = data['title'].loc[i].split()
15     # store the list of words of ith string in a, ex: a = ['tokidoki', 'The',
16     for j in indices:
17
18         b = data['title'].loc[j].split()
19         # store the list of words of jth string in b, ex: b = ['tokidoki', 'T'
20
21         length = max(len(a),len(b))
22
23         # count is used to store the number of words that are matched in both
24         count = 0
25
26         # itertools.zip_longest(a,b): will map the corresponding words in bot
27         # example: a =['a', 'b', 'c', 'd']
28         # b = ['a', 'b', 'd']
29         # itertools.zip_longest(a,b): will give [('a', 'a'), ('b', 'b'), ('c', '
30         for k in itertools.zip_longest(a,b):
31             if (k[0]==k[1]):
32                 count += 1
33
34             # if the number of words in which both strings differ are < 3 , we ar
35             if (length - count) < 3:
36                 indices.remove(j)
```

```
In [0]: 1 # from whole previous products we will consider only
2 # the products that are found in previous cell
3 data = data.loc[data['asin'].isin(stage2_dedupe_asins)]
```

```
In [0]: 1 print('Number of data points after stage two of dedupe: ', data.shape[0])
2 # from 17k apperals we reduced to 16k apperals
```

Number of data points after stage two of dedupe: 16042

```
In [0]: 1 data.to_pickle('pickels/16k_apperal_data')
2 # Storing these products in a pickle file
3 # candidates who wants to download these files instead
4 # of 180K they can download and use them from the Google Drive folder.
```

6. Text pre-processing

```
In [0]: 1 data = pd.read_pickle('pickels/16k_appral_data')
2
3 # NLTK download stop words. [RUN ONLY ONCE]
4 # goto Terminal (Linux/Mac) or Command-Prompt (Window)
5 # In the temrinal, type these commands
6 # $python3
7 # $import nltk
8 # $nltk.download()
```

```
In [0]: 1 # we use the list of stop words that are downloaded from nltk lib.
2 stop_words = set(stopwords.words('english'))
3 print ('list of stop words:', stop_words)
4
5 def nlp_preprocessing(total_text, index, column):
6     if type(total_text) is not int:
7         string = ""
8         for words in total_text.split():
9             # remove the special chars in review like '#$@!%^&*()_+-~?>< etc
10            word = ("").join(e for e in words if e.isalnum())
11            # Conver all letters to lower-case
12            word = word.lower()
13            # stop-word removal
14            if not word in stop_words:
15                string += word + " "
16        data[column][index] = string
```

list of stop words: {'such', 'and', 'hers', 'up', 'she', 'd', 'further', 'all', 'than', 'under', 'is', 'off', 'both', 'most', 'few', 'should', 're', 'very', 'just', 'then', 'didn', 'myself', 'in', 'too', 's', 'shouldn', 'herself', 'because', 'how', 'itself', 'what', 'shan', 'weren', 'doing', 'them', 'couldn', 'their', 'so', 'ain', 'haven', 'yourself', 'now', 'll', 'isn', 'about', 'over', 'into', 'before', 'during', 'on', 'as', 'aren', 'against', 'above', 'down', 'they', 'below', 'me', 'again', 'for', 'why', 'been', 'yourselves', 'more', 'her', 'that', 'can', 'am', 'was', 'themselves', 'mightn', 'does', 'those', 'only', 'hasn', 'any', 'ma', 'are', 'nor', 'out', 'you', 'ourselves', 'the', 'an', 'has', 'where', 'i', 'while', 'ours', 'its', 'your', 'had', 'were', 'being', 'no', 'or', 'needn', 've', 'y', 'a', 'each', 'have', 'through', 'when', 'mustn', 'by', 'won', 'from', 'own', 'will', 'there', 't', 'him', 'these', 'doesn', 'theirs', 'my', 'did', 'of', 'who', 'until', 'wouldn', 'we', 'do', 'having', 'yours', 'other', 'wasn', 'it', 'with', 'once', 'here', 'don', 'o', 'whom', 'this', 'if', 'but', 'hadn', 'our', 'some', 'm', 'not', 'between', 'himself', 'same', 'at', 'be', 'he', 'after', 'which', 'to', 'his'}

```
In [0]: 1 start_time = time.clock()
2 # we take each title and we text-preprocess it.
3 for index, row in data.iterrows():
4     nlp_preprocessing(row['title'], index, 'title')
5 # we print the time it took to preprocess whole titles
6 print(time.clock() - start_time, "seconds")
```

3.572722000000006 seconds

In [0]: 1 data.head()

Out[6]:

		asin	brand	color	medium_image_url	product_type_name	title	format
4	B004GSI2OS	FeatherLite	Onyx Black/ Stone		https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies long sleeve stain resistant...	
6	B012YX2ZPI	HX-Kingdom Fashion T-shirts	White		https://images-na.ssl-images-amazon.com/images...	SHIRT	womens unique 100 cotton special olympics wor...	
15	B003BSRPB0	FeatherLite	White		https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies moisture free mesh sport sh...	
27	B014ICEJ1Q	FNC7C	Purple		https://images-na.ssl-images-amazon.com/images...	SHIRT	supernatural chibis sam dean castiel neck tshi...	
46	B01NACPBG2	Fifth Degree	Black		https://images-na.ssl-images-amazon.com/images...	SHIRT	fifth degree womens gold foil graphic tees jun...	

In [0]: 1 data.to_pickle('pickels/16k_apperal_data_preprocessed')

Stemming

In [0]:

```

1 from nltk.stem.porter import *
2 stemmer = PorterStemmer()
3 print(stemmer.stem('arguing'))
4 print(stemmer.stem('fishing'))
5
6
7 # We tried using stemming on our titles and it didnot work very well.
8

```

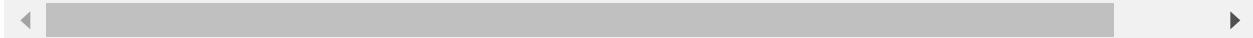
argu
fish

[8] Text based product similarity

```
In [0]: 1 data = pd.read_pickle('pickels/16k_apperal_data_preprocessed')
2 data.head()
```

Out[10]:

		asin	brand	color	medium_image_url	product_type_name	title	format
4	B004GSI2OS	FeatherLite	Onyx Black/ Stone		https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies long sleeve stain resistant...	
6	B012YX2ZPI	HX-Kingdom Fashion T-shirts	White		https://images-na.ssl-images-amazon.com/images...	SHIRT	womens unique 100 cotton special olympics wor...	
15	B003BSRPB0	FeatherLite	White		https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies moisture free mesh sport sh...	
27	B014ICEJ1Q	FNC7C	Purple		https://images-na.ssl-images-amazon.com/images...	SHIRT	supernatural chibis sam dean castiel neck tshi...	
46	B01NACPBG2	Fifth Degree	Black		https://images-na.ssl-images-amazon.com/images...	SHIRT	fifth degree womens gold foil graphic tees jun...	



In [33]:

```

1 # Utility Functions which we will use through the rest of the workshop.
2
3
4 #Display an image
5 def display_img(url,ax,fig):
6     # we get the url of the apparel and download it
7     response = requests.get(url)
8     img = Image.open(BytesIO(response.content))
9     # we will display it in notebook
10    plt.imshow(img)
11    #plotting code to understand the algorithm's decision.
12    def plot_heatmap(keys, values, labels, url, text):
13        # keys: list of words of recommended title
14        # values: len(values) == len(keys), values(i) represents the occurer
15        # Labels: len(labels) == len(keys), the values of labels depends on i
16            # if model == 'bag of words': labels(i) = values(i)
17            # if model == 'tfidf weighted bag of words': labels(i) = tfidf
18            # if model == 'idf weighted bag of words': labels(i) = idf(key)
19        # url : apparel's url
20
21        # we will devide the whole figure into two parts
22        gs = gridspec.GridSpec(2, 2, width_ratios=[4,1], height_ratios=[4,1])
23        fig = plt.figure(figsize=(25,3))
24
25        # 1st, plotting heat map that represents the count of commonly occurred
26        ax = plt.subplot(gs[0])
27        # it displays a cell in white color if the word is intersection(lists)
28        ax = sns.heatmap(np.array([values]), annot=np.array([labels]))
29        ax.set_xticklabels(keys) # set that axis labels as the words of title
30        ax.set_title(text) # apparel title
31
32        # 2nd, plotting image of the the apparel
33        ax = plt.subplot(gs[1])
34        # we don't want any grid lines for image and no labels on x-axis and
35        ax.grid(False)
36        ax.set_xticks([])
37        ax.set_yticks([])
38
39        # we call dispaly_img based with paramete url
40        display_img(url, ax, fig)
41
42        # displays combine figure ( heat map and image together)
43        plt.show()
44
45    def plot_heatmap_image(doc_id, vec1, vec2, url, text, model):
46
47        # doc_id : index of the title1
48        # vec1 : input apparels's vector, it is of a dict type {word:count}
49        # vec2 : recommended apparels's vector, it is of a dict type {word:count}
50        # url : apparels image url
51        # text: title of recomonded apparel (used to keep title of image)
52        # model, it can be any of the models,
53            # 1. bag_of_words
54            # 2. tfidf
55            # 3. idf
56

```

```

57     # we find the common words in both titles, because these only words contr
58     intersection = set(vec1.keys()) & set(vec2.keys())
59
60     # we set the values of non intersecting words to zero, this is just to sh
61     for i in vec2:
62         if i not in intersection:
63             vec2[i]=0
64
65     # for labeling heatmap, keys contains list of all words in title2
66     keys = list(vec2.keys())
67     # if ith word in intersection(list of words of title1 and list of words o
68     values = [vec2[x] for x in vec2.keys()]
69
70     # Labels: len(labels) == len(keys), the values of labels depends on the m
71     # if model == 'bag of words': labels(i) = values(i)
72     # if model == 'tfidf weighted bag of words': labels(i) = tfidf(keys(i))
73     # if model == 'idf weighted bag of words': labels(i) = idf(keys(i))
74
75     if model == 'bag_of_words':
76         labels = values
77     elif model == 'tfidf':
78         labels = []
79         for x in vec2.keys():
80             # tfidf_title_vectorizer.vocabulary_ it contains all the words in t
81             # tfidf_title_features[doc_id, index_of_word_in_corpus] will give i
82             if x in tfidf_title_vectorizer.vocabulary_:
83                 labels.append(tfidf_title_features[doc_id, tfidf_title_vectoriz
84             else:
85                 labels.append(0)
86     elif model == 'idf':
87         labels = []
88         for x in vec2.keys():
89             # idf_title_vectorizer.vocabulary_ it contains all the words in t
90             # idf_title_features[doc_id, index_of_word_in_corpus] will give i
91             if x in idf_title_vectorizer.vocabulary_:
92                 labels.append(idf_title_features[doc_id, idf_title_vectorizer.
93             else:
94                 labels.append(0)
95
96     plot_heatmap(keys, values, labels, url, text)
97
98
99     # this function gets a list of words along with the frequency of each
100    # word given "text"
101    def text_to_vector(text):
102        word = re.compile(r'\w+')
103        words = word.findall(text)
104        # words stores list of all words in given string, you can try 'words = te
105        return Counter(words) # Counter counts the occurrence of each word in list
106
107
108
109    def get_result(doc_id, content_a, content_b, url, model):
110        text1 = content_a
111        text2 = content_b
112
113        # vector1 = dict{word11:#count, word12:#count, etc.}

```

```
114     vector1 = text_to_vector(text1)
115
116     # vector1 = dict{word21:#count, word22:#count, etc.}
117     vector2 = text_to_vector(text2)
118
119     plot_heatmap_image(doc_id, vector1, vector2, url, text2, model)
```

[8.2] Bag of Words (BoW) on product titles.

In [0]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 title_vectorizer = CountVectorizer()
3 title_features = title_vectorizer.fit_transform(data['title'])
4 title_features.get_shape() # get number of rows and columns in feature matrix
5 # title_features.shape = #data_points * #words_in_corpus
6 # CountVectorizer().fit_transform(corpus) returns
7 # the a sparse matrix of dimensions #data_points * #words_in_corpus
8
9 # What is a sparse vector?
10
11 # title_features[doc_id, index_of_word_in_corpus] = number of times the word
12
13
```

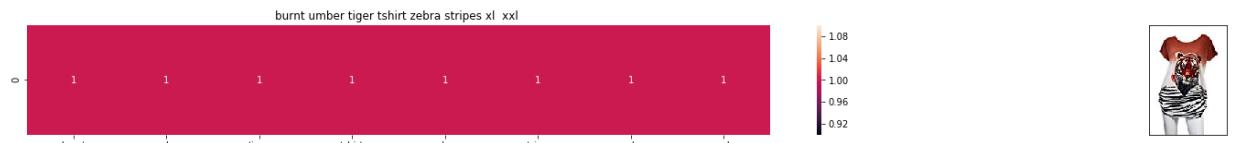
Out[17]: (16042, 12609)

In [0]:

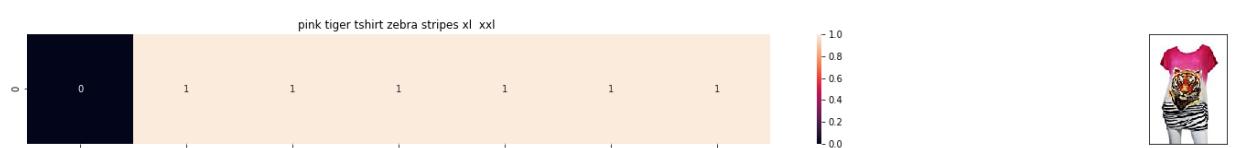
```

1 def bag_of_words_model(doc_id, num_results):
2     # doc_id: apparel's id in given corpus
3
4     # pairwise_dist will store the distance from given input apparel to all r
5     # the metric we used here is cosine, the coside distance is mesured as K(
6     # http://scikit-Learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(title_features,title_features[doc_id])
8
9     # np.argsort will return indices of the smallest distances
10    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
11    #pdists will store the smallest distances
12    pdists = np.sort(pairwise_dist.flatten())[0:num_results]
13
14    #data frame indices of the 9 smallest distace's
15    df_indices = list(data.index[indices])
16
17    for i in range(0,len(indices)):
18        # we will pass 1. doc_id, 2. title1, 3. title2, url, model
19        get_result(indices[i],data['title'].loc[df_indices[0]], data['title']
20        print('ASIN :',data['asin'].loc[df_indices[i]])
21        print ('Brand:', data['brand'].loc[df_indices[i]])
22        print ('Title:', data['title'].loc[df_indices[i]])
23        print ('Euclidean similarity with the query image :', pdists[i])
24        print('='*60)
25
26    #call the bag-of-words model for a product to get similar products.
27    bag_of_words_model(12566, 20) # change the index if you want to.
28    # In the output heat map each value represents the count value
29    # of the label word, the color represents the intersection
30    # with inputs title.
31
32    #try 12566
33    #try 931

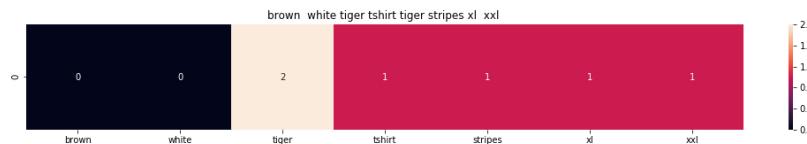
```



ASIN : B00JXQB5FQ
 Brand: Si Row
 Title: burnt umber tiger tshirt zebra stripes xl xxl
 Euclidean similarity with the query image : 0.0
 =====



ASIN : B00JXQASS6
 Brand: Si Row
 Title: pink tiger tshirt zebra stripes xl xxl
 Euclidean similarity with the query image : 1.73205080757
 =====



ASIN : B00JXQCWT0

Brand: Si Row

Title: brown white tiger tshirt tiger stripes xl xxl

Euclidean similarity with the query image : 2.44948974278

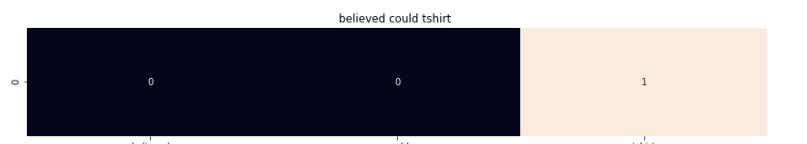


ASIN : B00JXQCUIC

Brand: Si Row

Title: yellow tiger tshirt tiger stripes l

Euclidean similarity with the query image : 2.64575131106

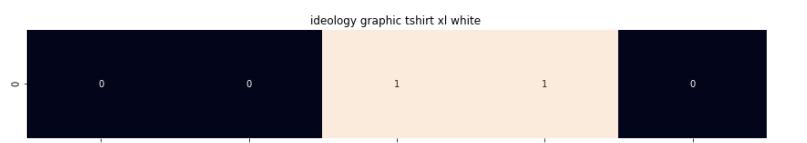


ASIN : B07568NZX4

Brand: Rustic Grace

Title: believed could tshirt

Euclidean similarity with the query image : 3.0

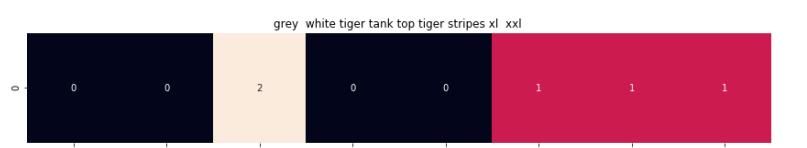


ASIN : B01NB0NKRO

Brand: Ideology

Title: ideology graphic tshirt xl white

Euclidean similarity with the query image : 3.0

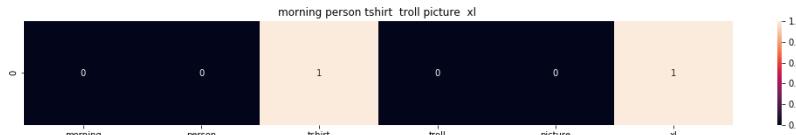


ASIN : B00JXQAFZ2

Brand: Si Row

Title: grey white tiger tank top tiger stripes xl xxl

Euclidean similarity with the query image : 3.0



ASIN : B01CLS8LMW

Brand: Awake

Title: morning person tshirt troll picture xl

Euclidean similarity with the query image : 3.16227766017

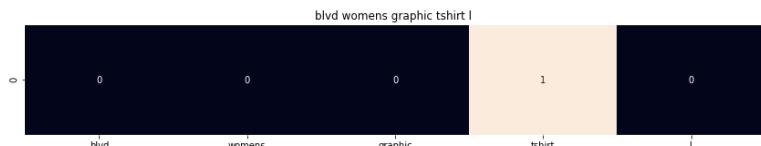


ASIN : B01KVZUB6G

Brand: Merona

Title: merona green gold stripes

Euclidean similarity with the query image : 3.16227766017



ASIN : B0733R2CJK

Brand: BLVD

Title: blvd womens graphic tshirt l

Euclidean similarity with the query image : 3.16227766017

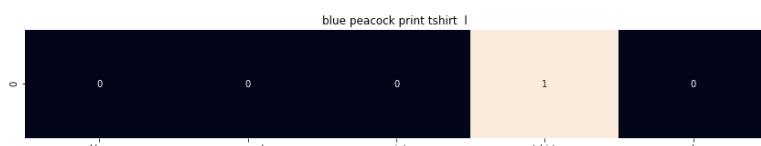


ASIN : B012VQLT6Y

Brand: KM T-shirt

Title: km tiger printed sleeveless vest tshirt

Euclidean similarity with the query image : 3.16227766017

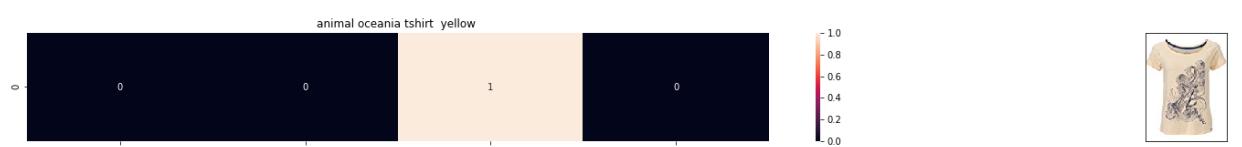
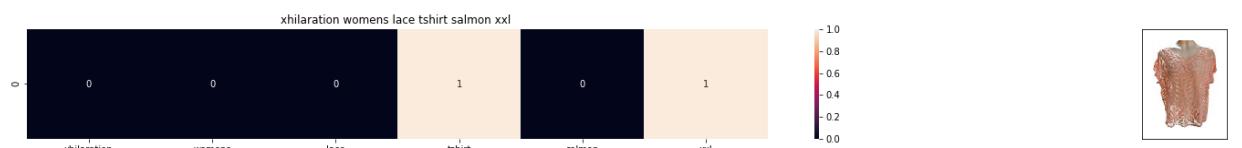
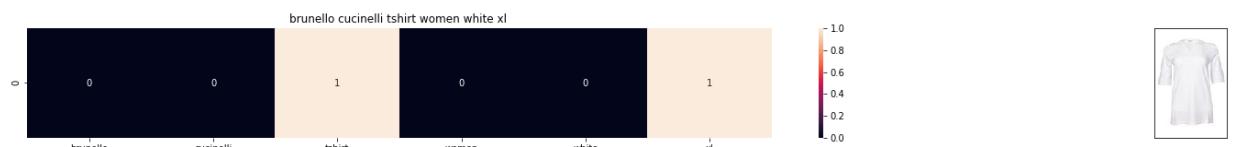


ASIN : B00JXQC8L6

Brand: Si Row

Title: blue peacock print tshirt l

Euclidean similarity with the query image : 3.16227766017



ASIN : B06X6GX6WG

Brand: Animal

Title: animal oceania tshirt yellow

Euclidean similarity with the query image : 3.16227766017



ASIN : B017X8PW9U

Brand: Diesel

Title: diesel tserraf tshirt black

Euclidean similarity with the query image : 3.16227766017



ASIN : B00IAA4JIQ

Brand: I Love Lucy

Title: juniors love lucywaaaaahhhh tshirt size xl

Euclidean similarity with the query image : 3.16227766017

[8.5] TF-IDF based product similarity

In [0]:

```

1 tfidf_title_vectorizer = TfidfVectorizer(min_df = 0)
2 tfidf_title_features = tfidf_title_vectorizer.fit_transform(data['title'])
3 # tfidf_title_features.shape = #data_points * #words_in_corpus
4 # CountVectorizer().fit_transform(courpus) returns the a sparase matrix of di
5 # tfidf_title_features[doc_id, index_of_word_in_corpus] = tfidf values of the

```

In [0]:

```

1 def tfidf_model(doc_id, num_results):
2     # doc_id: apparel's id in given corpus
3
4     # pairwise_dist will store the distance from given input apparel to all r
5     # the metric we used here is cosine, the coside distance is mesured as K(
6     # http://scikit-Learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(tfidf_title_features,tfidf_title_featu
8
9     # np.argsort will return indices of 9 smallest distances
10    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
11    #pdists will store the 9 smallest distances
12    pdists = np.sort(pairwise_dist.flatten())[0:num_results]
13
14    #data frame indices of the 9 smallest distace's
15    df_indices = list(data.index[indices])
16
17    for i in range(0,len(indices)):
18        # we will pass 1. doc_id, 2. title1, 3. title2, url, model
19        get_result(indices[i], data['title'].loc[df_indices[0]], data['title']
20        print('ASIN :',data['asin'].loc[df_indices[i]])
21        print('BRAND :',data['brand'].loc[df_indices[i]])
22        print ('Eucliden distance from the given image :', pdists[i])
23        print('*'*125)
24    tfidf_model(12566, 20)
25    # in the output heat map each value represents the tfidf values of the Label

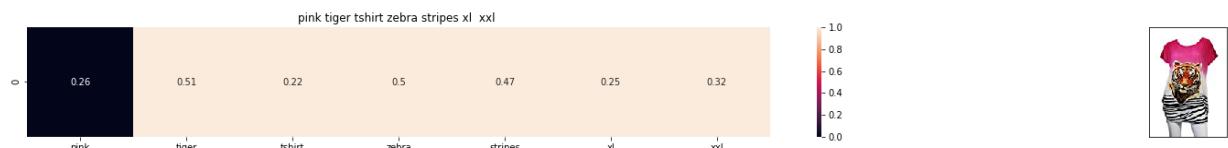
```



ASIN : B00JXQB5FQ

BRAND : Si Row

Eucliden distance from the given image : 0.0



ASIN : B00JXQASS6

BRAND : Si Row

Eucliden distance from the given image : 0.753633191245



ASIN : B00JXQCWT0

BRAND : Si Row

Eucliden distance from the given image : 0.935764394377



ASIN : B00JXQAFZ2

BRAND : Si Row

Euclidean distance from the given image : 0.95861535242



ASIN : B00JXQCUIC

BRAND : Si Row

Euclidean distance from the given image : 1.00007496145



ASIN : B00JXQA094

BRAND : Si Row

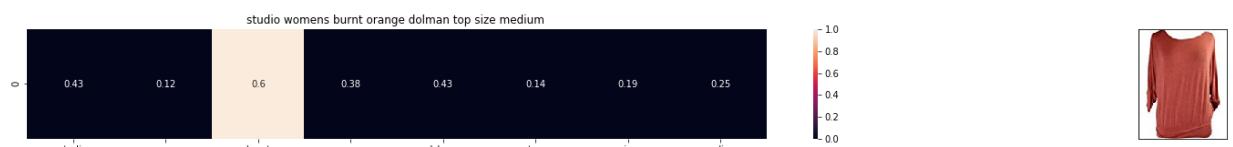
Euclidean distance from the given image : 1.02321555246



ASIN : B00JXQUWA

BRAND : Si Row

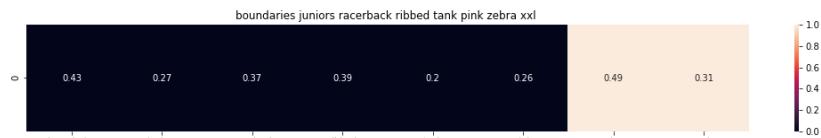
Euclidean distance from the given image : 1.0319918463



ASIN : B06XSCVFT5

BRAND : Studio M

Euclidean distance from the given image : 1.21068436704



ASIN : B06Y2GTYPM

BRAND : No Boundaries

Eucliden distance from the given image : 1.21216838107



ASIN : B012VQLT6Y

BRAND : KM T-shirt

Eucliden distance from the given image : 1.21979064028



ASIN : B06Y1VN8WQ

BRAND : Black Swan

Eucliden distance from the given image : 1.220684966



ASIN : B00Z6HEXWI

BRAND : Black Temptation

Eucliden distance from the given image : 1.22128139212



ASIN : B074TR12BH

BRAND : Ultra Flirt

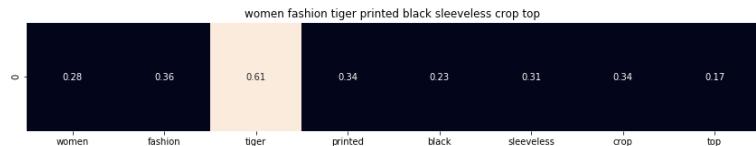
Eucliden distance from the given image : 1.23133640946



ASIN : B072R2JXKW

BRAND : WHAT ON EARTH

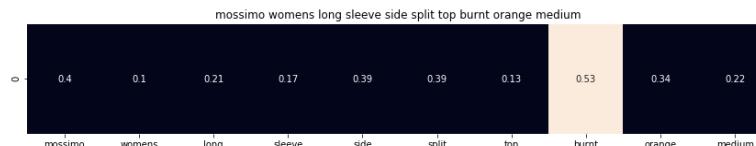
Eucliden distance from the given image : 1.23184519726



ASIN : B074T8ZYGX

BRAND : MKP Crop Top

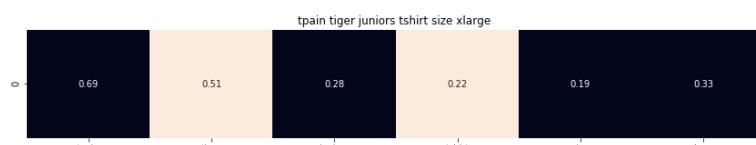
Eucliden distance from the given image : 1.23406074574



ASIN : B071ZDF6T2

BRAND : Mossimo

Eucliden distance from the given image : 1.23527855777



ASIN : B01K0H02OG

BRAND : Tultex

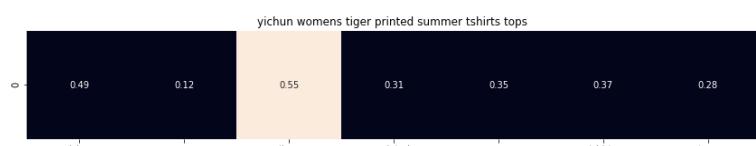
Eucliden distance from the given image : 1.23645729881



ASIN : B00H8A6ZLI

BRAND : Vivian's Fashions

Eucliden distance from the given image : 1.24996155053



ASIN : B010NN9RX0

BRAND : YICHUN

Eucliden distance from the given image : 1.25354614209

=====

=====



ASIN : B06XBY5QXL

BRAND : Liz Claiborne

Eucliden distance from the given image : 1.25388329384

=====

=====

[8.5] IDF based product similarity

In [18]:

```

1 idf_title_vectorizer = CountVectorizer()
2 idf_title_features = idf_title_vectorizer.fit_transform(data['title'])
3
4 # idf_title_features.shape = #data_points * #words_in_corpus
5 # CountVectorizer().fit_transform(courpus) returns the a sparase matrix of di
6 # idf_title_features[doc_id, index_of_word_in_corpus] = number of times the w

```

In [19]:

```

1 def nContaining(word):
2     # return the number of documents which had the given word
3     return sum(1 for blob in data['title'] if word in blob.split())
4
5 def idf(word):
6     # idf = log(#number of docs / #number of docs which had the given word)
7     return math.log(data.shape[0] / (nContaining(word)))

```

In [20]:

```

1 # we need to convert the values into float
2 idf_title_features = idf_title_features.astype(np.float)
3
4 for i in idf_title_vectorizer.vocabulary_.keys():
5     # for every word in whole corpus we will find its idf value
6     idf_val = idf(i)
7
8     # to calculate idf_title_features we need to replace the count values wit
9     # idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0]
10    for j in idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonze
11
12        # we replace the count values of word i in document j with idf_value
13        # idf_title_features[doc_id, index_of_word_in_courpus] = idf value of
14        idf_title_features[j,idf_title_vectorizer.vocabulary_[i]] = idf_val
15

```

In [0]:

```

1 def idf_model(doc_id, num_results):
2     # doc_id: apparel's id in given corpus
3
4     # pairwise_dist will store the distance from given input apparel to all r
5     # the metric we used here is cosine, the coside distance is mesured as K(
6     # http://scikit-Learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(idf_title_features,idf_title_features[
8
9         # np.argsort will return indices of 9 smallest distances
10    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
11    #pdists will store the 9 smallest distances
12    pdists = np.sort(pairwise_dist.flatten())[0:num_results]
13
14    #data frame indices of the 9 smallest distace's
15    df_indices = list(data.index[indices])
16
17    for i in range(0,len(indices)):
18        get_result(indices[i],data['title'].loc[df_indices[0]], data['title'][i])
19        print('ASIN :',data['asin'].loc[df_indices[i]])
20        print('Brand :',data['brand'].loc[df_indices[i]])
21        print ('euclidean distance from the given image :', pdists[i])
22        print('='*125)
23
24
25
26 idf_model(12566,20)
27 # in the output heat map each value represents the idf values of the Label wo

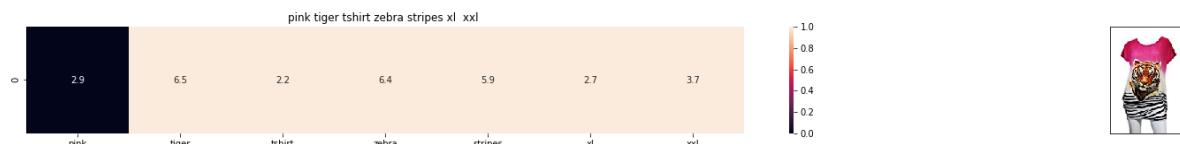
```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from the given image : 0.0



ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from the given image : 12.2050713112

[9] Text Semantics based product similarity

In [0]:

```
1 # credits: https://www.kaggle.com/c/word2vec-nlp-tutorial#part-2-word-vectors
2 # Custom Word2Vec using your own text data.
3 # Do NOT RUN this code.
4 # It is meant as a reference to build your own Word2Vec when you have
5 # lots of data.
6
7
8 ...
9 # Set values for various parameters
10 num_features = 300      # Word vector dimensionality
11 min_word_count = 1       # Minimum word count
12 num_workers = 4          # Number of threads to run in parallel
13 context = 10             # Context window size
14 downsampling = 1e-3     # Downsample setting for frequent words
15
16 # Initialize and train the model (this will take some time)
17 from gensim.models import word2vec
18 print ("Training model...")
19 model = word2vec.Word2Vec(sen_corpus, workers=num_workers, \
20                           size=num_features, min_count = min_word_count, \
21                           window = context)
22
23 ...
```

In [0]:

```
1 from gensim.models import Word2Vec
2 from gensim.models import KeyedVectors
3 import pickle
4
5 # in this project we are using a pretrained model by google
6 # its 3.3G file, once you load this into your memory
7 # it occupies ~9Gb, so please do this step only if you have >12G of ram
8 # we will provide a pickle file which contains a dict ,
9 # and it contains all our corpus words as keys and model[word] as values
10 # To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
11 # from https://drive.google.com/file/d/0B7XkCwpI5KDYNLNUTTLSS21pQmM/edit
12 # it's 1.9GB in size.
13
14 ...
15 model = KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin'
16 ...
17
18 #if you do NOT have RAM >= 12GB, use the code below.
19 with open('word2vec_model', 'rb') as handle:
20     model = pickle.load(handle)
21
```

In [31]:

```

1 # Utility functions
2
3 def get_word_vec(sentence, doc_id, m_name):
4     # sentence : title of the apparel
5     # doc_id: document id in our corpus
6     # m_name: model information it will take two values
7         # if m_name == 'avg', we will append the model[i], w2v representation
8         # if m_name == 'weighted', we will multiply each w2v[word] with the i
9     vec = []
10    for i in sentence.split():
11        if i in vocab:
12            if m_name == 'weighted' and i in idf_title_vectorizer.vocabulary:
13                vec.append(idf_title_features[doc_id, idf_title_vectorizer.voc
14            elif m_name == 'avg':
15                vec.append(model[i])
16            else:
17                # if the word in our corpus is not there in the google word2vec
18                vec.append(np.zeros(shape=(300,)))
19    # we will return a numpy array of shape (#number of words in title * 300
20    # each row represents the word2vec representation of each word (weighted/
21    return np.array(vec)
22
23 def get_distance(vec1, vec2):
24     # vec1 = np.array(#number_of_words_title1 * 300), each row is a vector of
25     # vec2 = np.array(#number_of_words_title2 * 300), each row is a vector of
26
27     final_dist = []
28     # for each vector in vec1 we calculate the distance(euclidean) to all vec
29     for i in vec1:
30         dist = []
31         for j in vec2:
32             # np.linalg.norm(i-j) will result the euclidean distance between
33             dist.append(np.linalg.norm(i-j))
34         final_dist.append(np.array(dist))
35     # final_dist = np.array(#number of words in title1 * #number of words in
36     # final_dist[i,j] = euclidean distance between vectors i, j
37     return np.array(final_dist)
38
39
40 def heat_map_w2v(sentence1, sentence2, url, doc_id1, doc_id2, model):
41     # sentence1 : title1, input apparel
42     # sentence2 : title2, recommended apparel
43     # url: apparel image url
44     # doc_id1: document id of input apparel
45     # doc_id2: document id of recommended apparel
46     # model: it can have two values, 1. avg 2. weighted
47
48     s1_vec = np.array(#number_of_words_title1 * 300), each row is a vector(w
49     s1_vec = get_word_vec(sentence1, doc_id1, model)
50     s2_vec = np.array(#number_of_words_title1 * 300), each row is a vector(w
51     s2_vec = get_word_vec(sentence2, doc_id2, model)
52
53     # s1_s2_dist = np.array(#number of words in title1 * #number of words in
54     # s1_s2_dist[i,j] = euclidean distance between words i, j
55     s1_s2_dist = get_distance(s1_vec, s2_vec)
56

```

```

57
58
59 # devide whole figure into 2 parts 1st part displays heatmap 2nd part dis
60 gs = gridspec.GridSpec(2, 2, width_ratios=[4,1],height_ratios=[2,1])
61 fig = plt.figure(figsize=(15,15))
62
63 ax = plt.subplot(gs[0])
64 # ploting the heap map based on the pairwise distances
65 ax = sns.heatmap(np.round(s1_s2_dist,4), annot=True)
66 # set the x axis labels as recommended apparels title
67 ax.set_xticklabels(sentence2.split())
68 # set the y axis labels as input apparels title
69 ax.set_yticklabels(sentence1.split())
70 # set title as recommended apparels title
71 ax.set_title(sentence2)
72
73 ax = plt.subplot(gs[1])
74 # we remove all grids and axis labels for image
75 ax.grid(False)
76 ax.set_xticks([])
77 ax.set_yticks([])
78 display_img(url, ax, fig)
79
80 plt.show()

```

In [0]:

```

1 # vocab = stores all the words that are there in google w2v model
2 # vocab = model.wv.vocab.keys() # if you are using Google word2Vec
3
4 vocab = model.keys()
5 # this function will add the vectors of each word and returns the avg vector
6 def build_avg_vec(sentence, num_features, doc_id, m_name):
7     # sentace: its title of the apparel
8     # num_features: the lenght of word2vec vector, its values = 300
9     # m_name: model information it will take two values
10    # if m_name == 'avg', we will append the model[i], w2v representation
11    # if m_name == 'weighted', we will multiply each w2v[word] with the i
12
13    featureVec = np.zeros((num_features,), dtype="float32")
14    # we will intialize a vector of size 300 with all zeros
15    # we add each word2vec(wordi) to this festureVec
16    nwords = 0
17
18    for word in sentence.split():
19        nwords += 1
20        if word in vocab:
21            if m_name == 'weighted' and word in idf_title_vectorizer.vocabulary:
22                featureVec = np.add(featureVec, idf_title_features[doc_id, idf_title_vectorizer.vocabulary[word]])
23            elif m_name == 'avg':
24                featureVec = np.add(featureVec, model[word])
25        if(nwords>0):
26            featureVec = np.divide(featureVec, nwords)
27    # returns the avg vector of given sentance, its of shape (1, 300)
28    return featureVec

```

[9.2] Average Word2Vec product similarity.

In [0]:

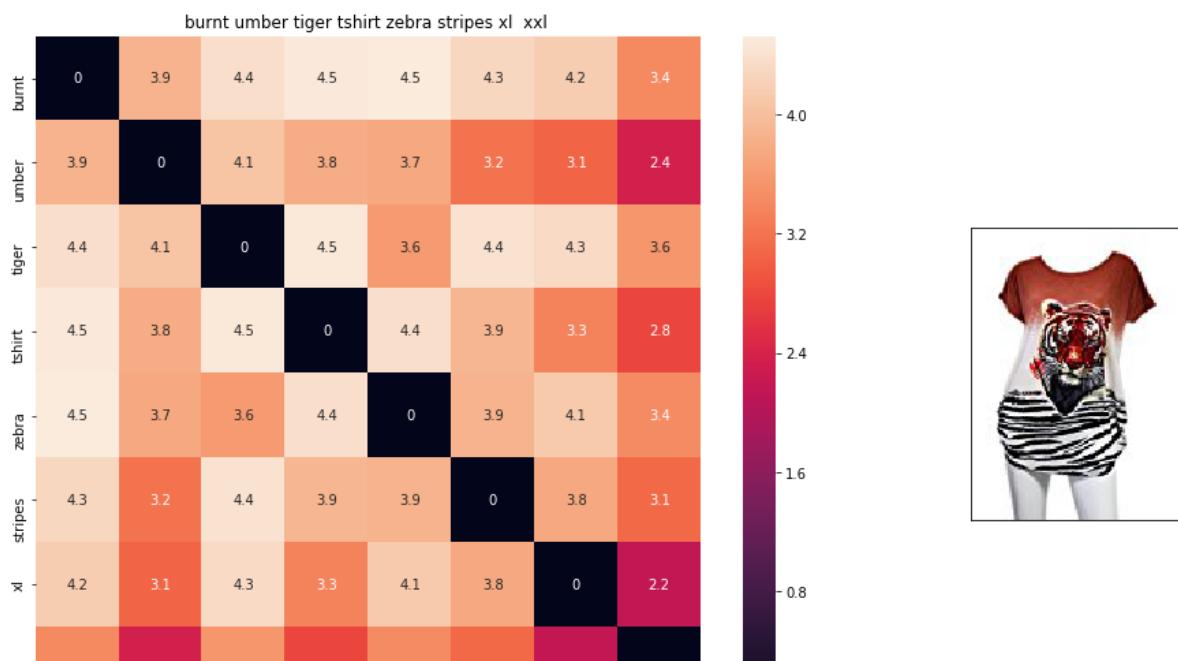
```
1 doc_id = 0
2 w2v_title = []
3 # for every title we build a avg vector representation
4 for i in data['title']:
5     w2v_title.append(build_avg_vec(i, 300, doc_id, 'avg'))
6     doc_id += 1
7
8 # w2v_title = np.array(# number of doc in courpus * 300), each row correspond
9 w2v_title = np.array(w2v_title)
10
```

In [0]:

```

1 def avg_w2v_model(doc_id, num_results):
2     # doc_id: apparel's id in given corpus
3
4     # dist(x, y) = sqrt(dot(x, x) - 2 * dot(x, y) + dot(y, y))
5     pairwise_dist = pairwise_distances(w2v_title, w2v_title[doc_id].reshape(1,
6
7         # np.argsort will return indices of 9 smallest distances
8         indices = np.argsort(pairwise_dist.flatten())[0:num_results]
9         #pdists will store the 9 smallest distances
10        pdists = np.sort(pairwise_dist.flatten())[0:num_results]
11
12    #data frame indices of the 9 smallest distance's
13    df_indices = list(data.index[indices])
14
15    for i in range(0, len(indices)):
16        heat_map_w2v(data['title'].loc[df_indices[0]],data['title'].loc[df_in
17        print('ASIN :',data['asin'].loc[df_indices[i]])
18        print('BRAND :',data['brand'].loc[df_indices[i]])
19        print ('euclidean distance from given input image :', pdists[i])
20        print('*'*125)
21
22
23 avg_w2v_model(12566, 20)
24 # in the give heat map, each cell contains the euclidean distance between wor

```



[9.4] IDF weighted Word2Vec for product similarity

In [0]:

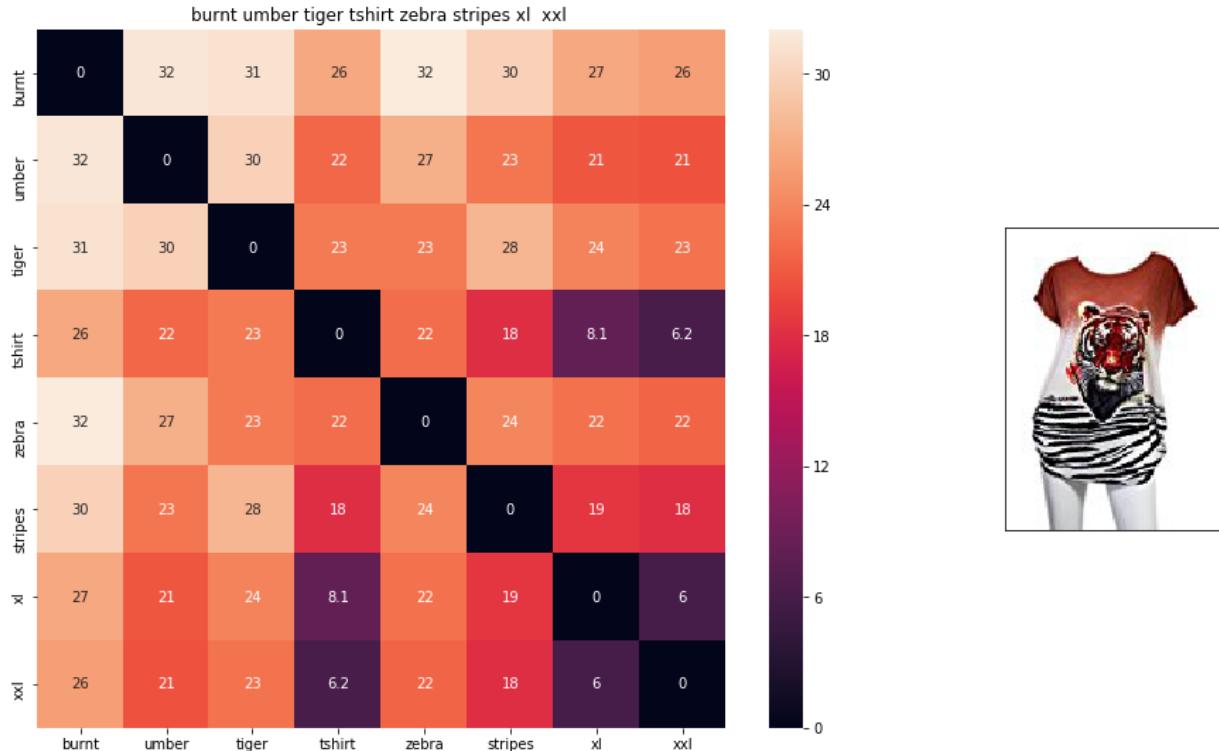
```
1 doc_id = 0
2 w2v_title_weight = []
3 # for every title we build a weighted vector representation
4 for i in data['title']:
5     w2v_title_weight.append(build_avg_vec(i, 300, doc_id, 'weighted'))
6     doc_id += 1
7 # w2v_title = np.array(# number of doc in courpus * 300), each row correspond
8 w2v_title_weight = np.array(w2v_title_weight)
```

In [0]:

```

1 def weighted_w2v_model(doc_id, num_results):
2     # doc_id: apparel's id in given corpus
3
4     # pairwise_dist will store the distance from given input apparel to all r
5     # the metric we used here is cosine, the coside distance is mesured as K(
6     # http://scikit-Learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_
8
9     # np.argsort will return indices of 9 smallest distances
10    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
11    #pdists will store the 9 smallest distances
12    pdists = np.sort(pairwise_dist.flatten())[0:num_results]
13
14    #data frame indices of the 9 smallest distace's
15    df_indices = list(data.index[indices])
16
17    for i in range(0, len(indices)):
18        heat_map_w2v(data['title'].loc[df_indices[0]],data['title'].loc[df_in
19        print('ASIN :',data['asin'].loc[df_indices[i]])
20        print('Brand :',data['brand'].loc[df_indices[i]])
21        print('euclidean distance from input :', pdists[i])
22        print('*125)
23
24 weighted_w2v_model(12566, 20)
25 #931
26 #12566
27 # in the give heat map, each cell contains the euclidean distance between wor

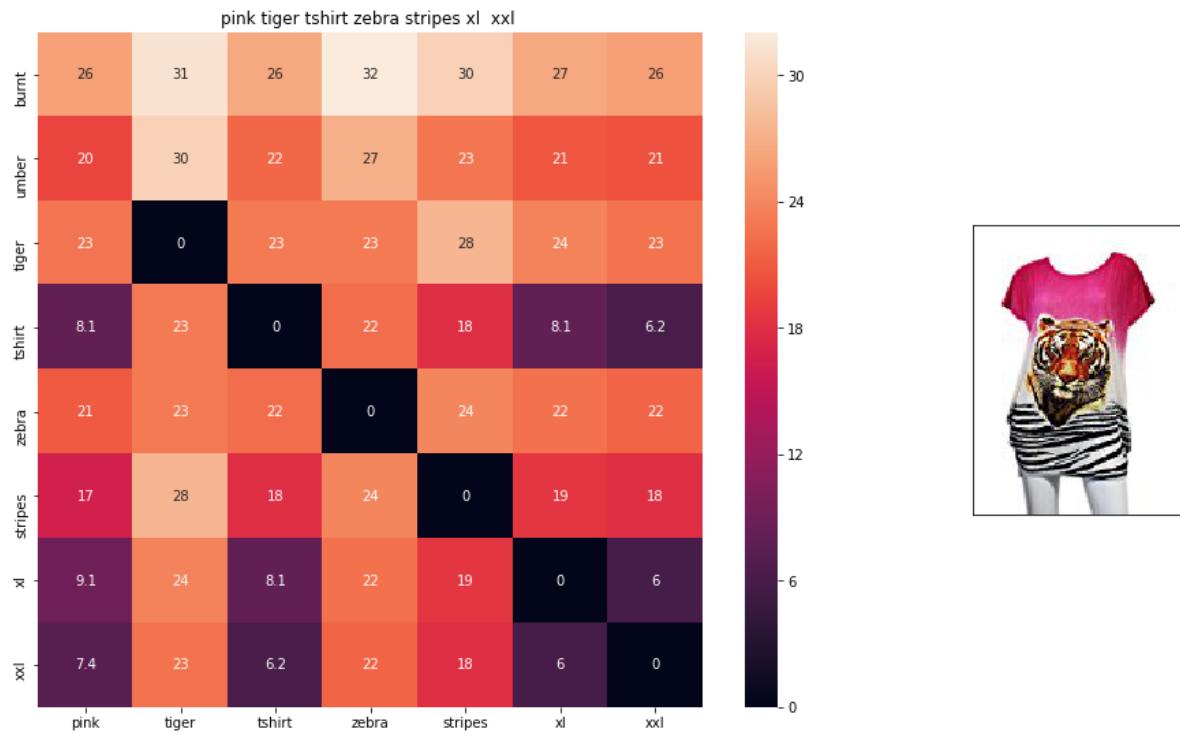
```



ASIN : B00JXQB5FQ

Brand : Si Row

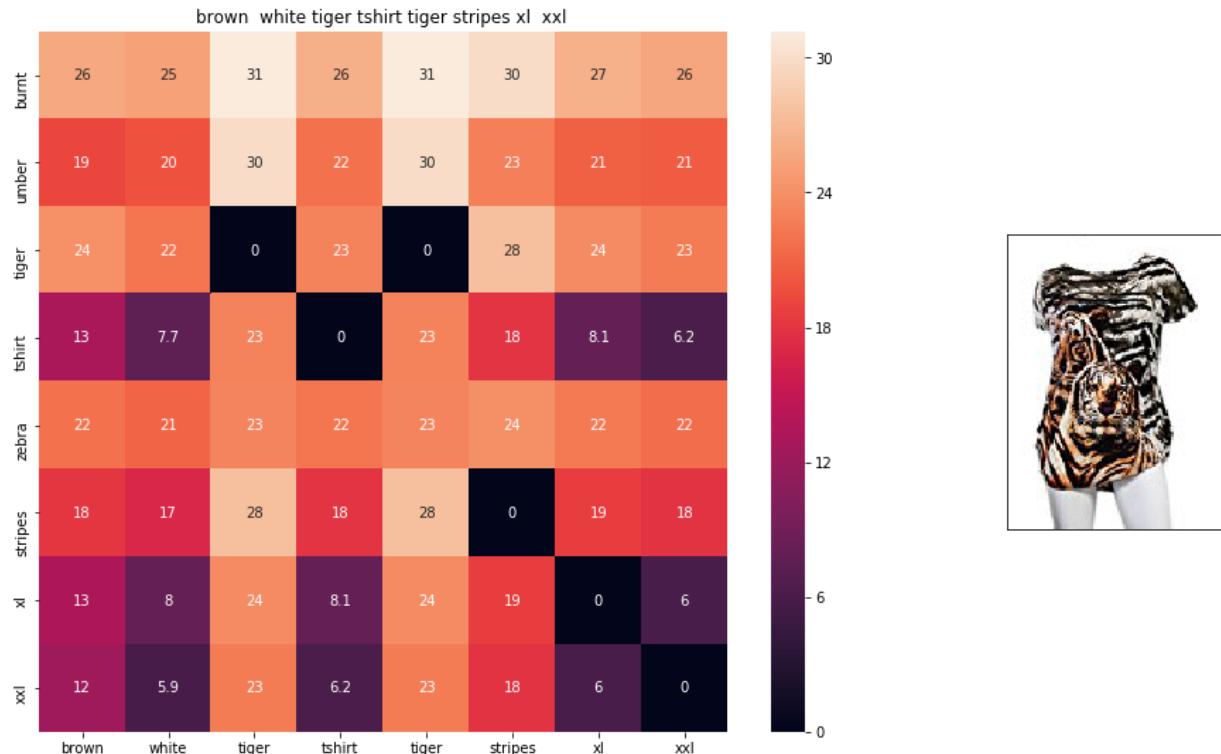
euclidean distance from input : 0.00390625



ASIN : B00JXQASS6

Brand : Si Row

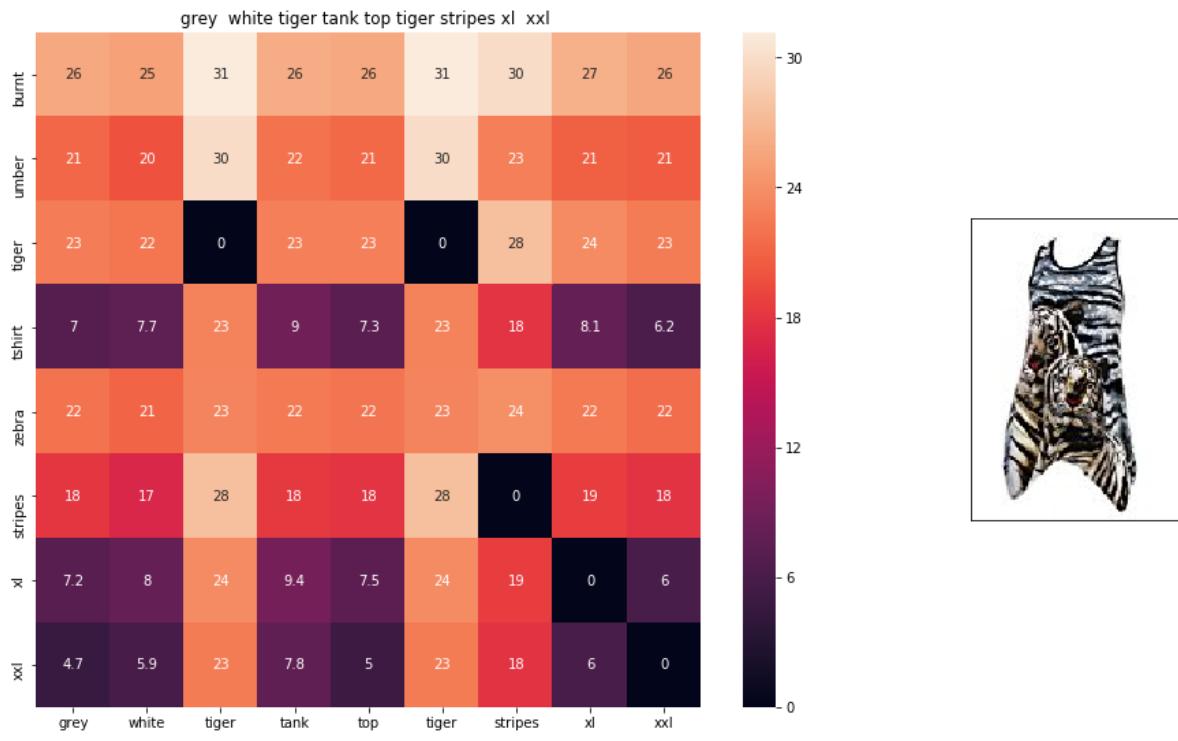
euclidean distance from input : 4.06389



ASIN : B00JXQCWT0

Brand : Si Row

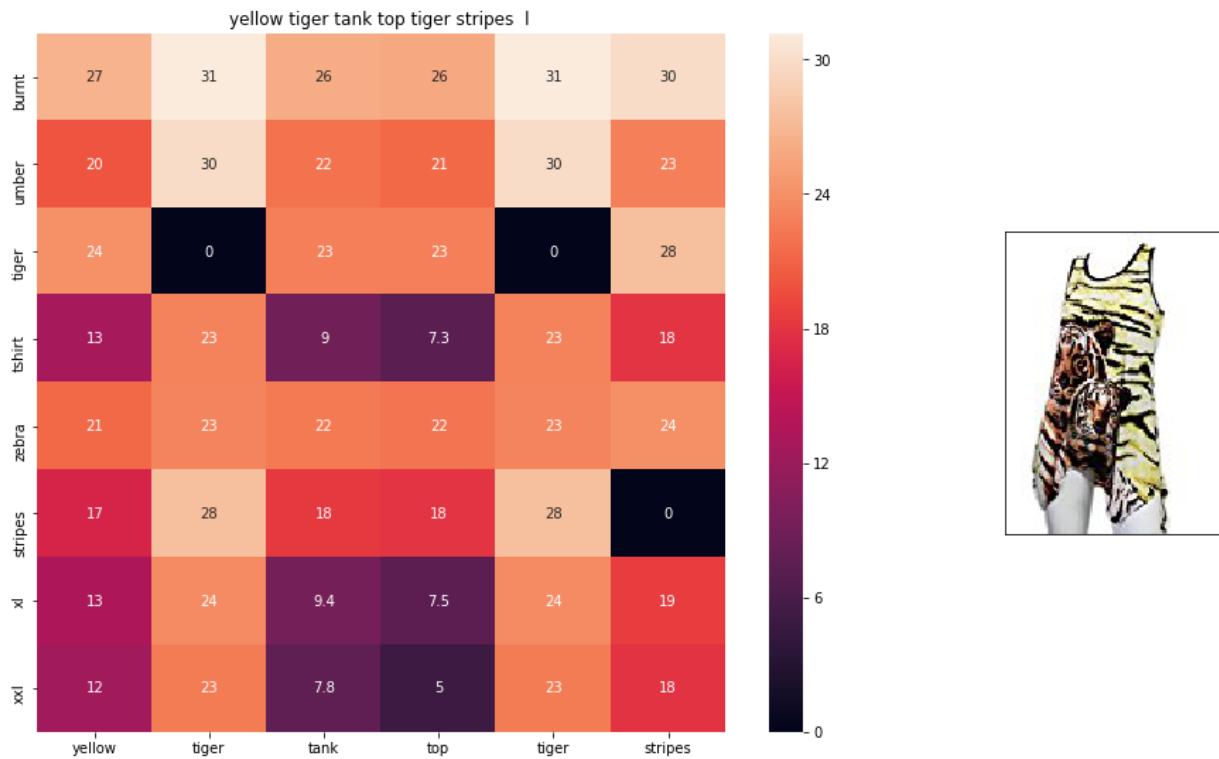
euclidean distance from input : 4.77094



ASIN : B00JXQAFZ2

Brand : Si Row

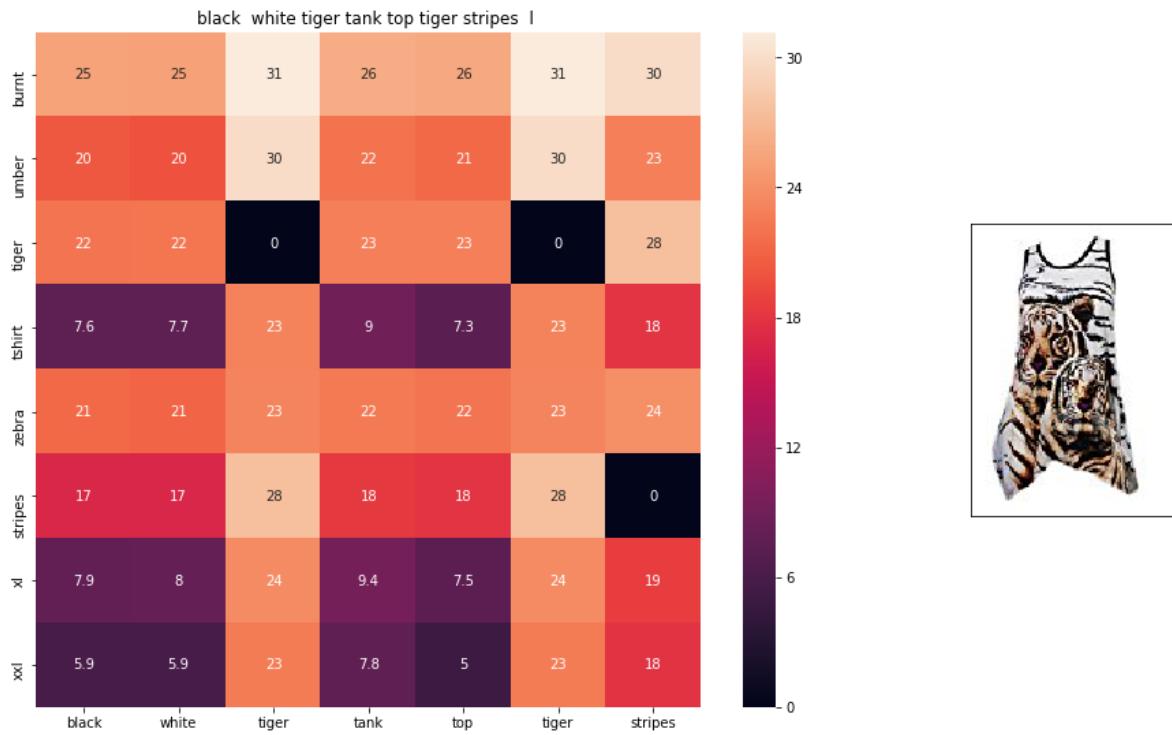
euclidean distance from input : 5.36016



ASIN : B00JXQAUWA

Brand : Si Row

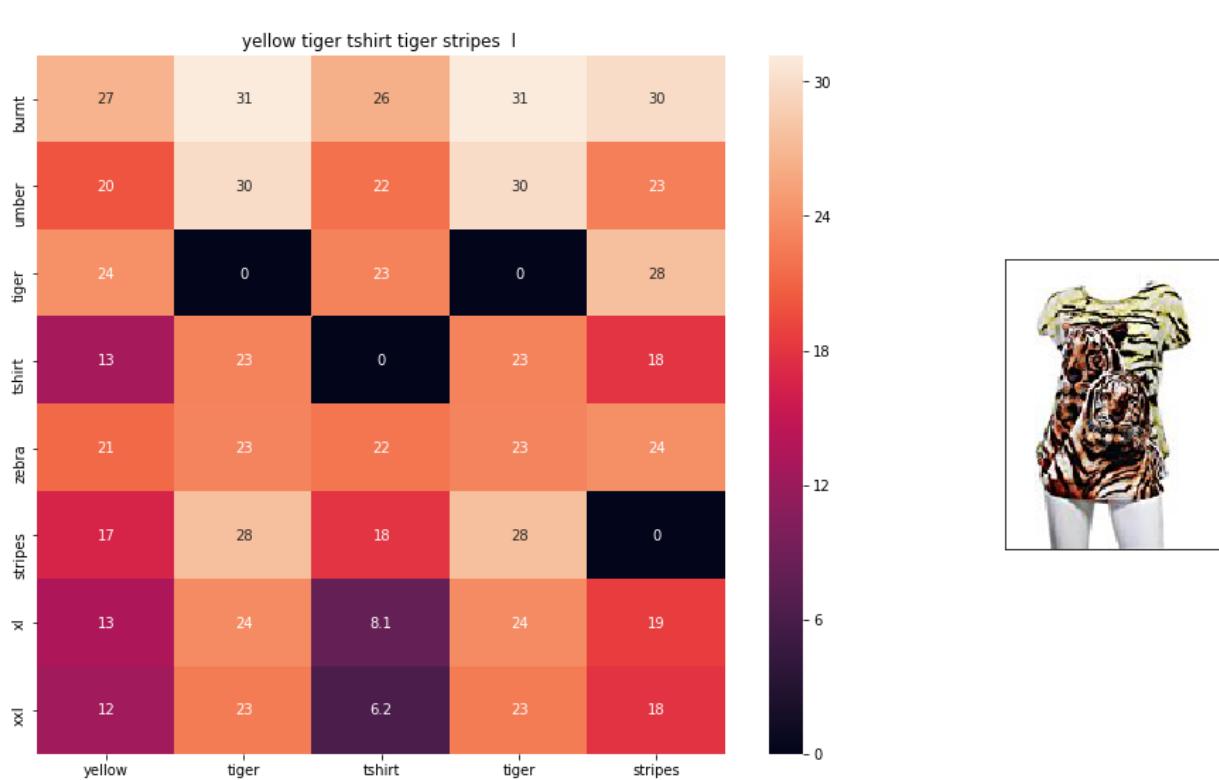
euclidean distance from input : 5.68952



ASIN : B00JXQA094

Brand : Si Row

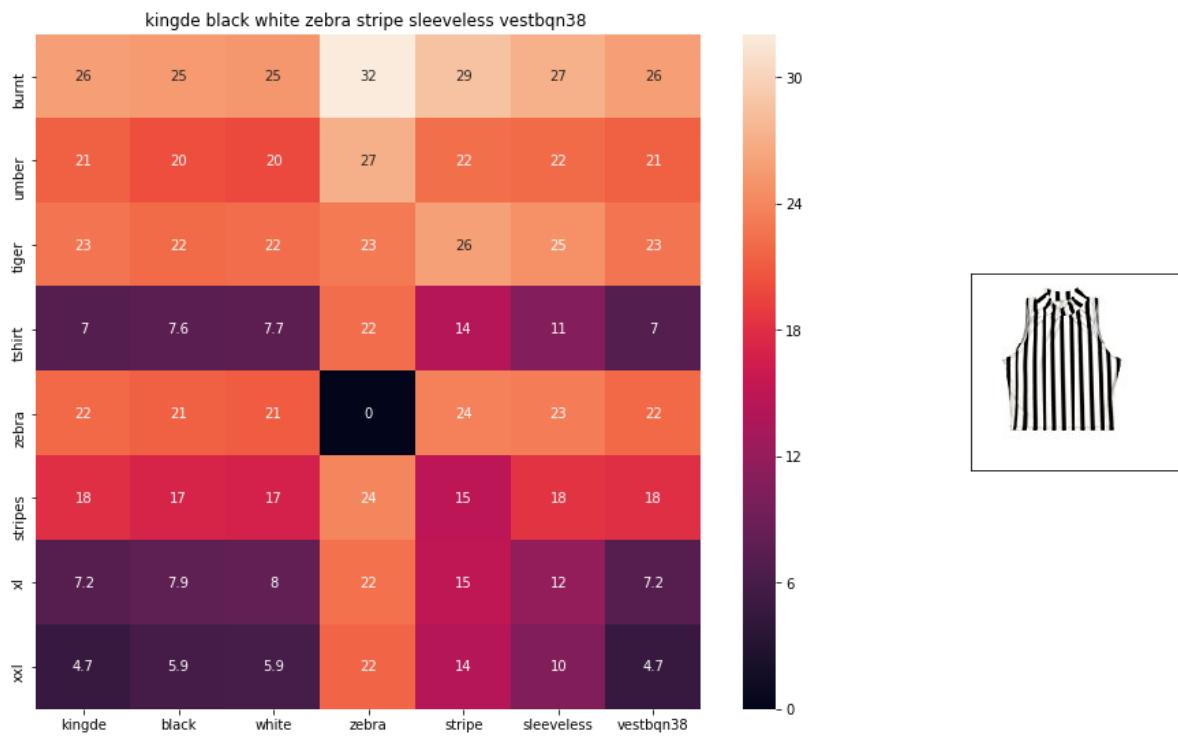
euclidean distance from input : 5.69302



ASIN : B00JXQCUIC

Brand : Si Row

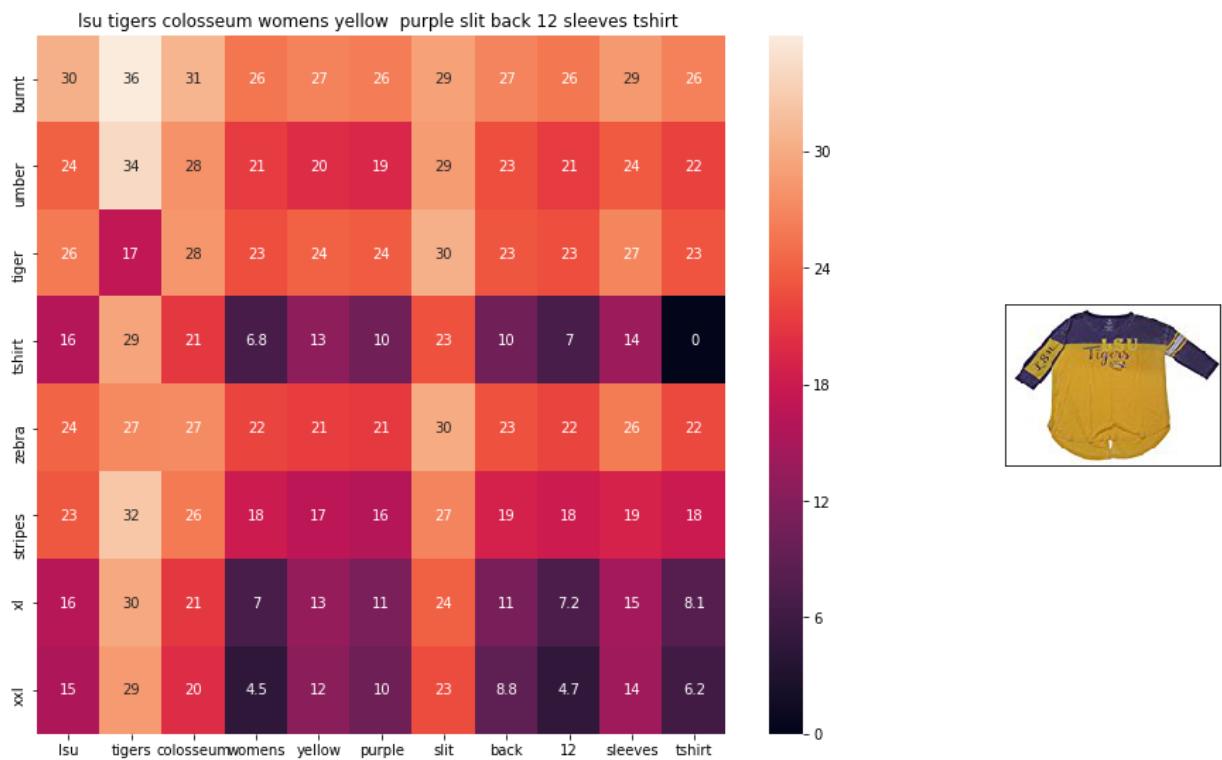
euclidean distance from input : 5.89344



ASIN : B015H41F6G

Brand : KINGDE

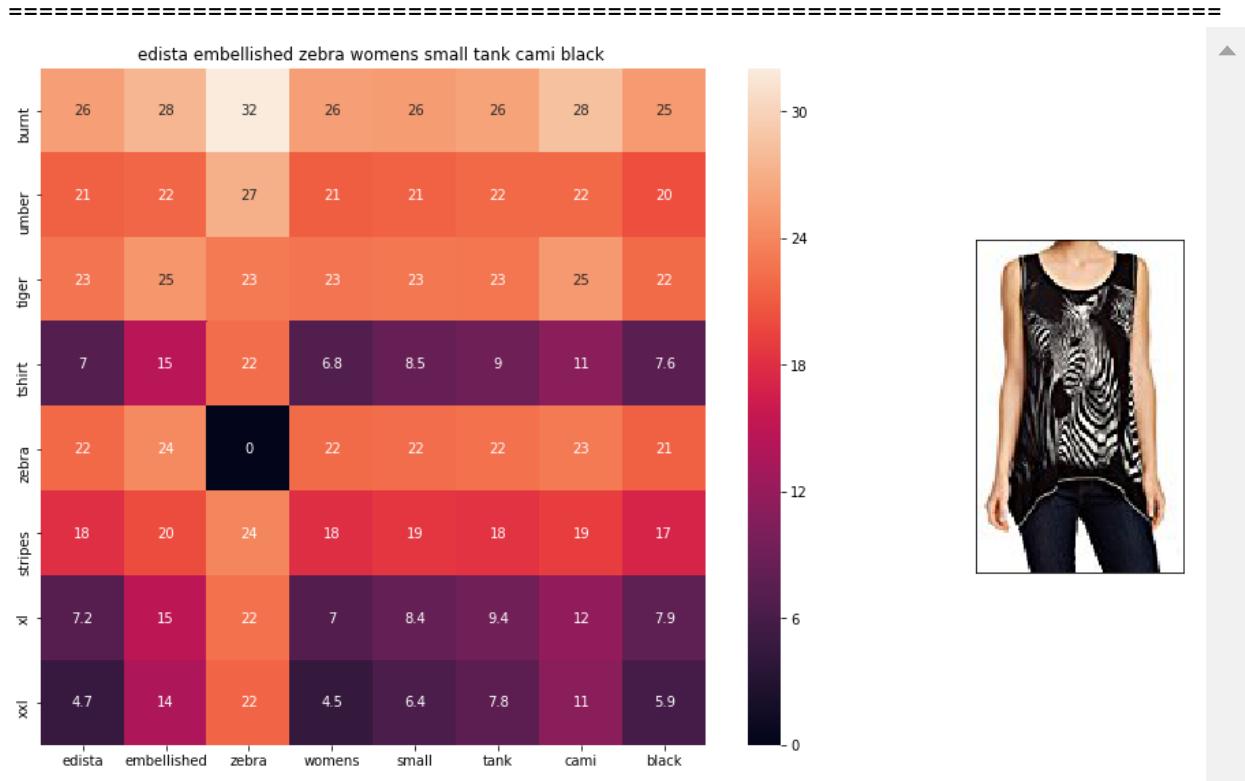
euclidean distance from input : 6.13299



ASIN : B073R5Q8HD

Brand : Colosseum

euclidean distance from input : 6.25671



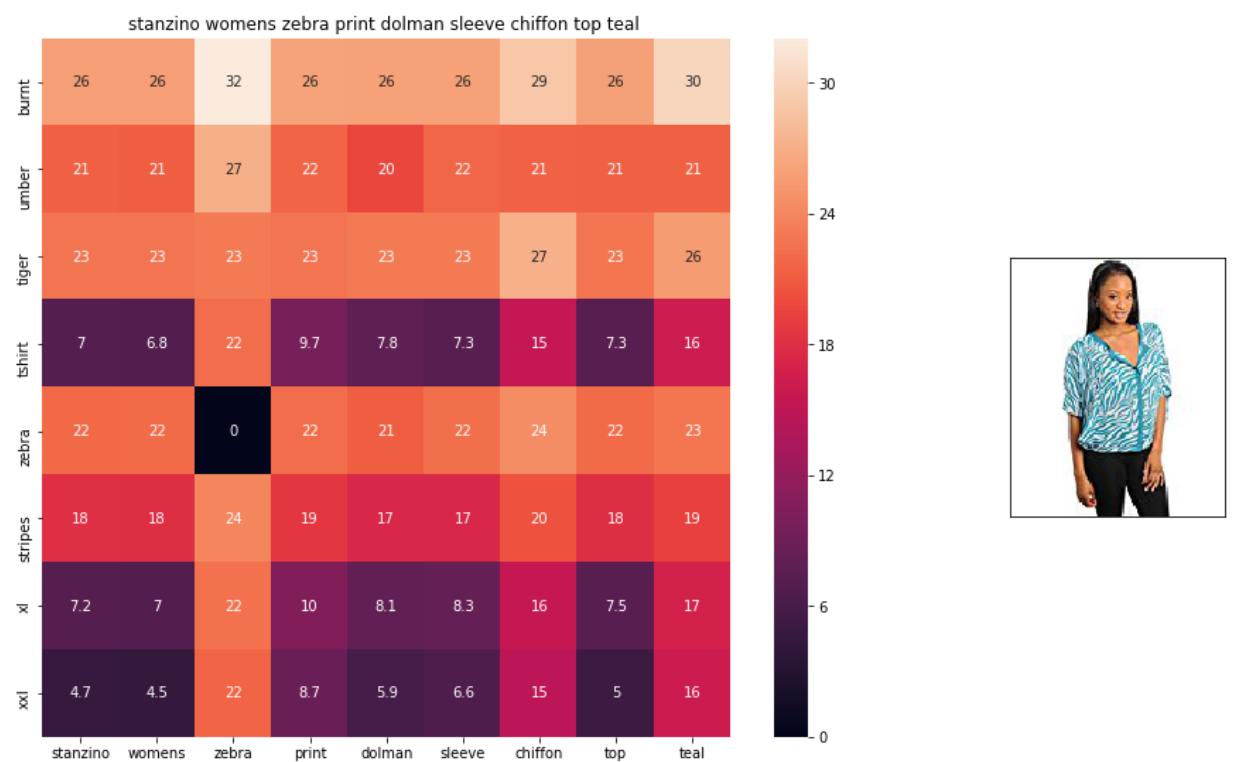
ASIN : B074P8MD22

Brand : Edista

euclidean distance from input : 6.3922

=====

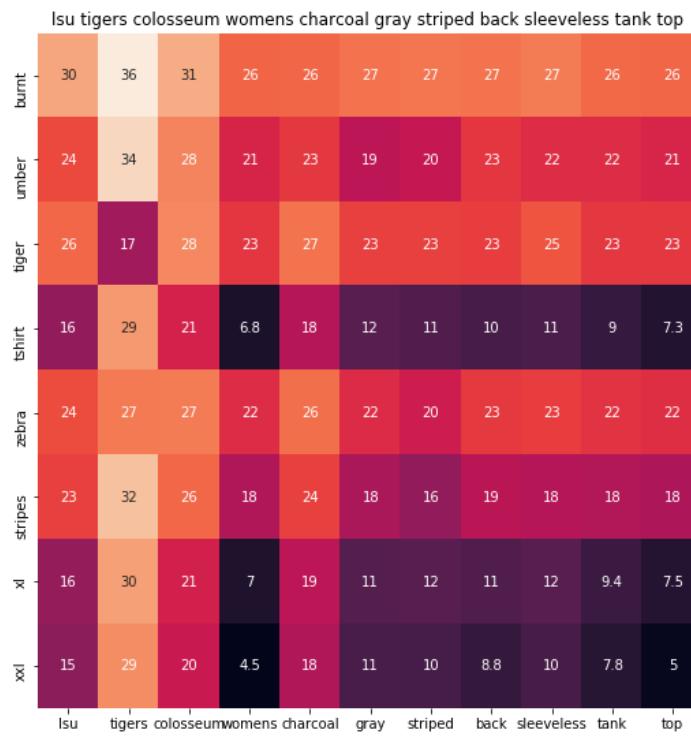
=====



ASIN : B00C0I3U3E

Brand : Stanzino

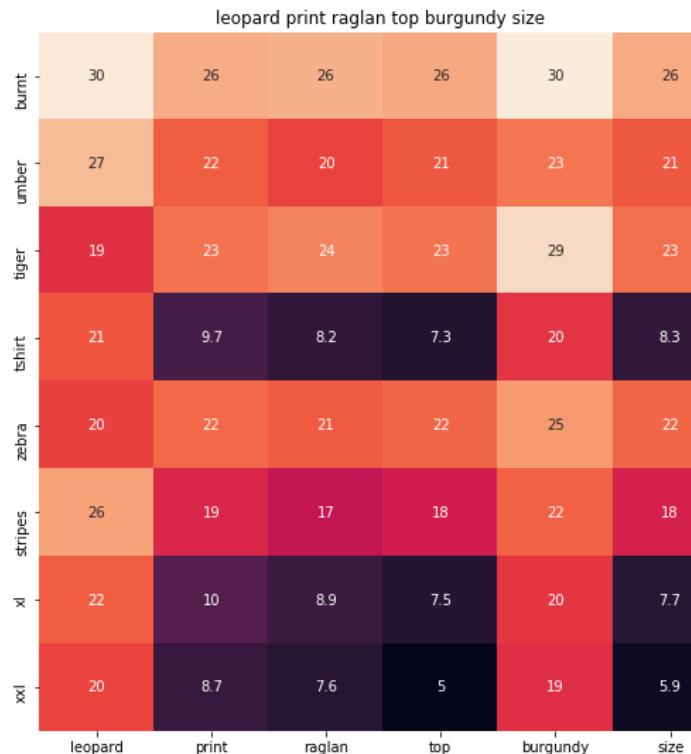
euclidean distance from input : 6.4149



ASIN : B073R4ZM7Y

Brand : Colosseum

euclidean distance from input : 6.45096



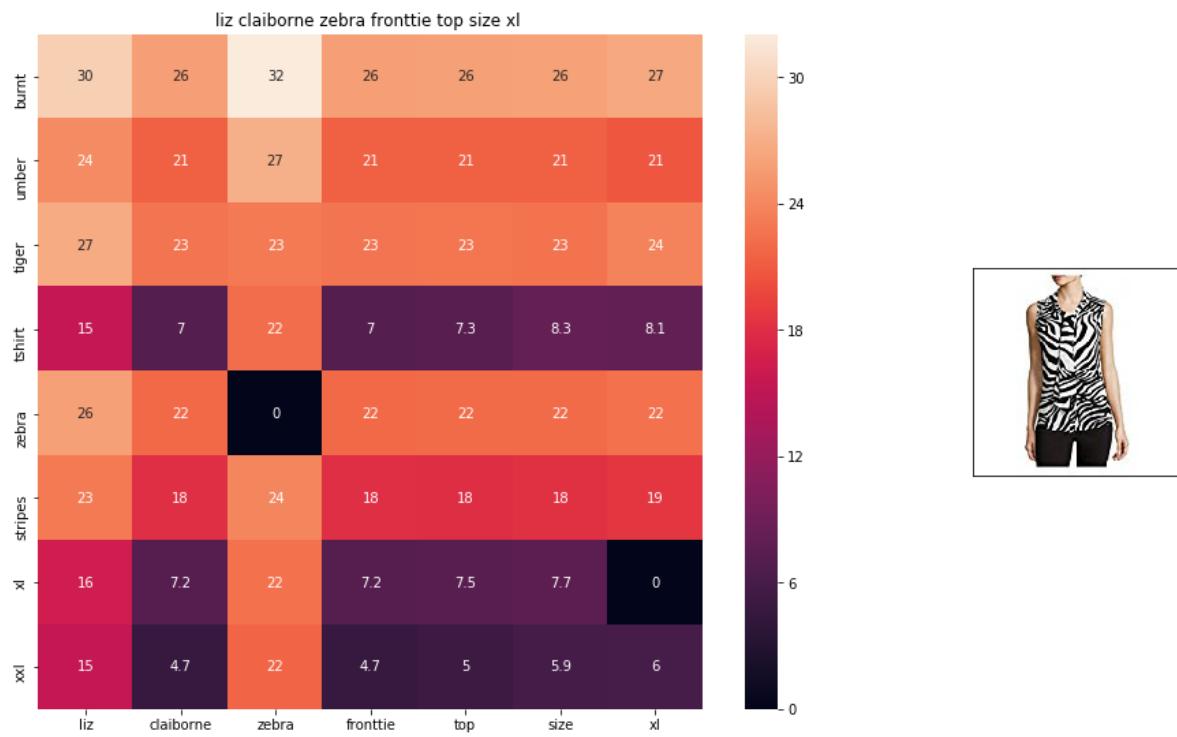
ASIN : B01C60RLDQ

Brand : 1 Mad Fit

euclidean distance from input : 6.46341

=====

=====



ASIN : B06XBY5QXL

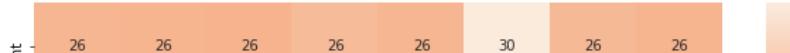
Brand : Liz Claiborne

euclidean distance from input : 6.53922

=====

=====

merona womens printed blouse brown leopard print xxl



ASIN : B071YF3WDD

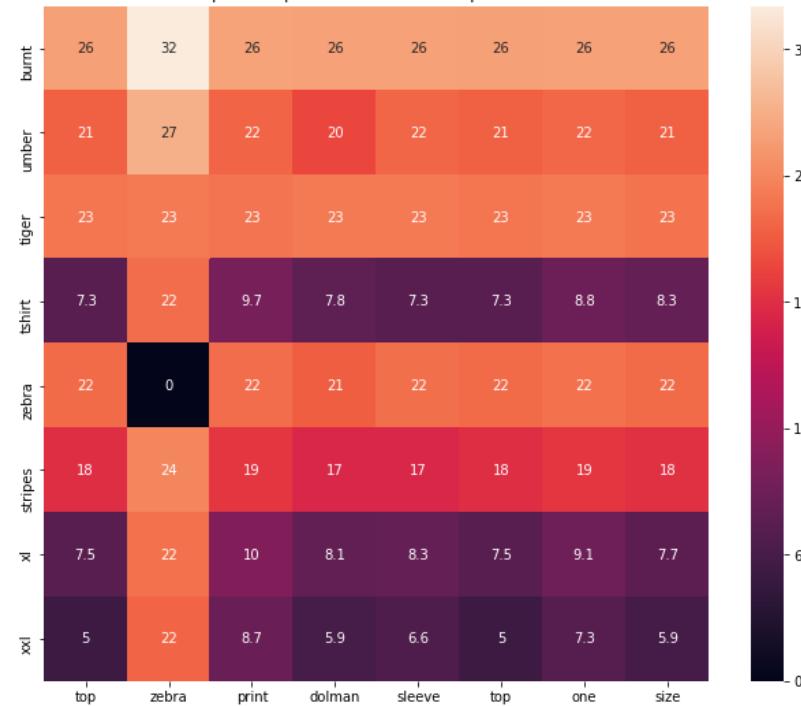
Brand : Merona

euclidean distance from input : 6.5755

=====

=====

top zebra print dolman sleeve top one size



ASIN : B00H8A6ZLI

Brand : Vivian's Fashions

euclidean distance from input : 6.63821

=====

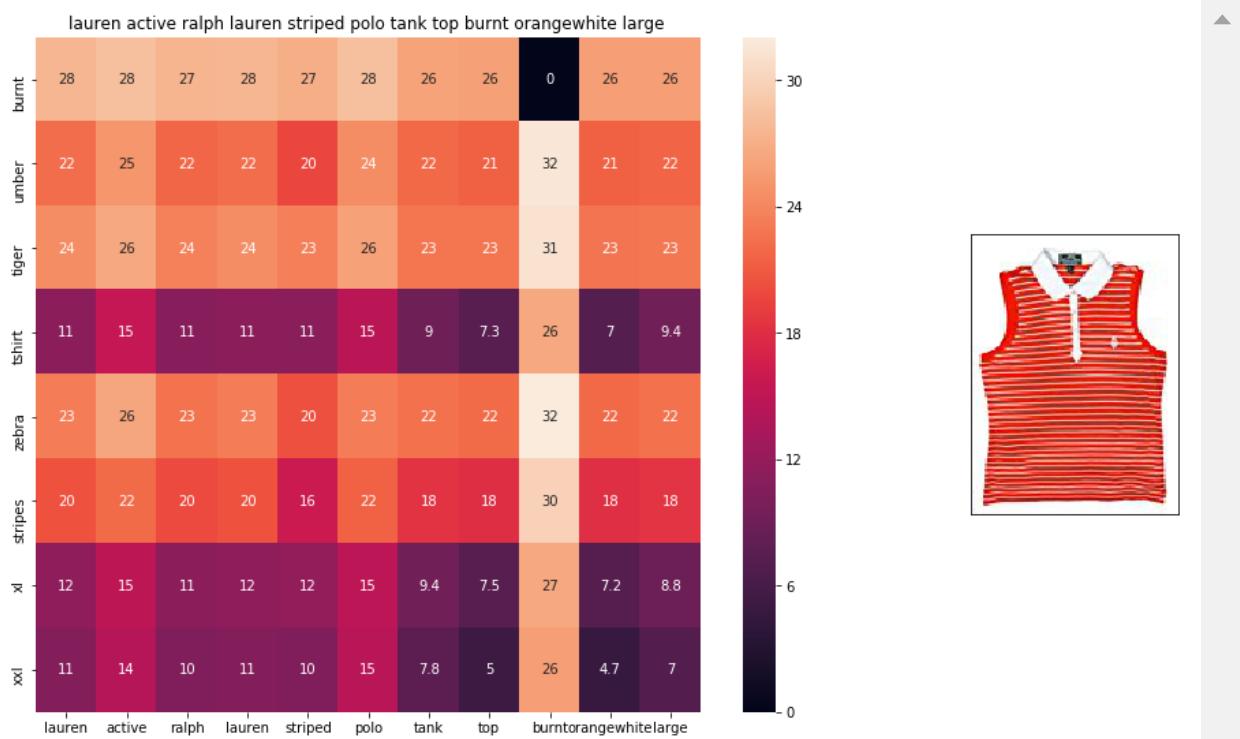
=====



ASIN : B00Z6HEXWI

Brand : Black Temptation

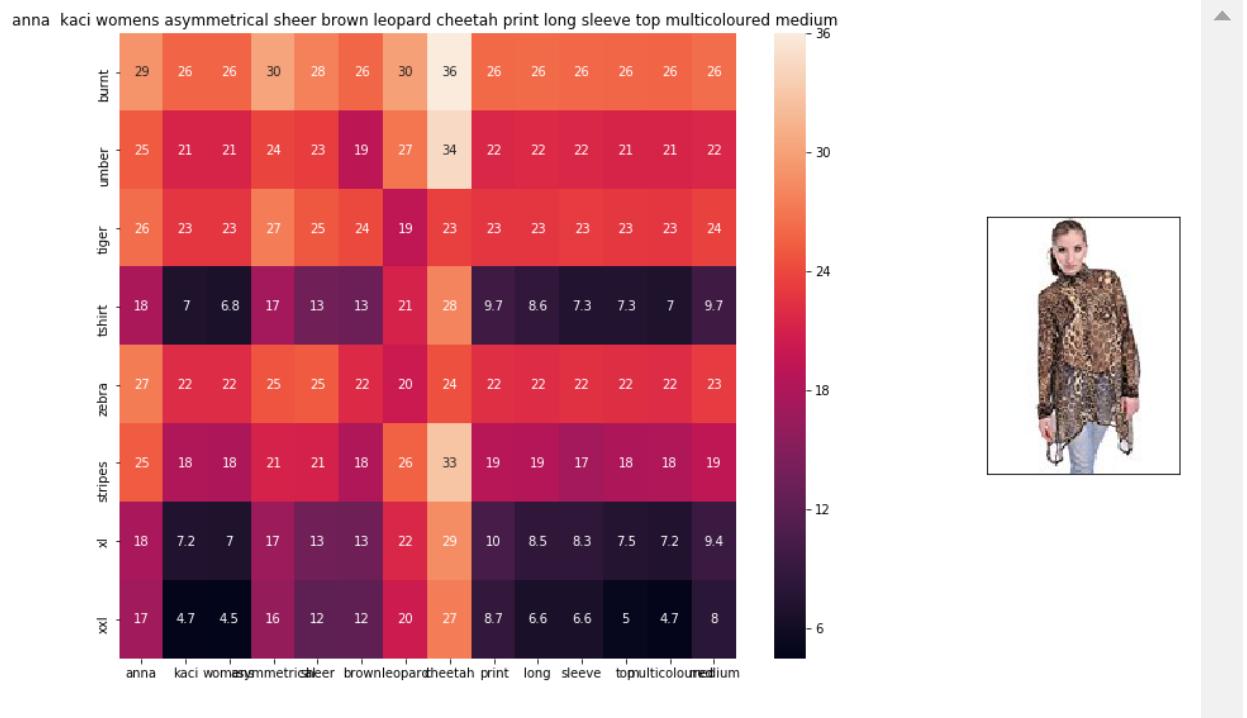
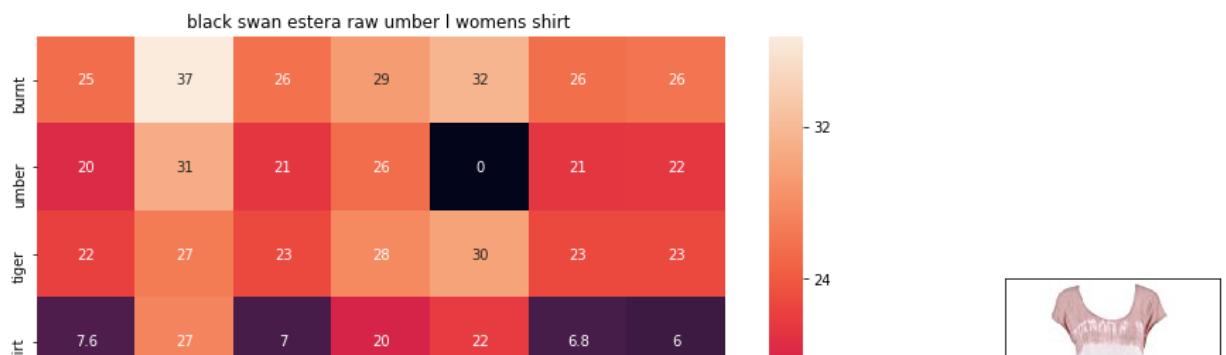
euclidean distance from input : 6.66074



ASIN : B00ILGH50Y

Brand : Ralph Lauren Active

euclidean distance from input : 6.68391



[9.6] Weighted similarity using brand and color.

In [0]:

```
1 # some of the brand values are empty.
2 # Need to replace Null with string "NULL"
3 data['brand'].fillna(value="Not given", inplace=True )
4
5 # replace spaces with hyphen
6 brands = [x.replace(" ", "-") for x in data['brand'].values]
7 types = [x.replace(" ", "-") for x in data['product_type_name'].values]
8 colors = [x.replace(" ", "-") for x in data['color'].values]
9
10 brand_vectorizer = CountVectorizer()
11 brand_features = brand_vectorizer.fit_transform(brands)
12
13 type_vectorizer = CountVectorizer()
14 type_features = type_vectorizer.fit_transform(types)
15
16 color_vectorizer = CountVectorizer()
17 color_features = color_vectorizer.fit_transform(colors)
18
19 extra_features = hstack((brand_features, type_features, color_features)).tocs
```

In [27]:

```

1 def heat_map_w2v_brand(sentance1, sentance2, url, doc_id1, doc_id2, df_id1, d
2
3     # sentance1 : title1, input apparel
4     # sentance2 : title2, recommended apparel
5     # url: apparel image url
6     # doc_id1: document id of input apparel
7     # doc_id2: document id of recommended apparel
8     # df_id1: index of document1 in the data frame
9     # df_id2: index of document2 in the data frame
10    # model: it can have two values, 1. avg 2. weighted
11
12    #s1_vec = np.array(#number_of_words_title1 * 300), each row is a vector(w
13    s1_vec = get_word_vec(sentance1, doc_id1, model)
14    #s2_vec = np.array(#number_of_words_title2 * 300), each row is a vector(w
15    s2_vec = get_word_vec(sentance2, doc_id2, model)
16
17    # s1_s2_dist = np.array(#number of words in title1 * #number of words in
18    # s1_s2_dist[i,j] = euclidean distance between words i, j
19    s1_s2_dist = get_distance(s1_vec, s2_vec)
20
21    data_matrix = [['Asin','Brand', 'Color', 'Product type'],
22                  [data['asin'].loc[df_id1],brands[doc_id1], colors[doc_id1], ty
23                   [data['asin'].loc[df_id2],brands[doc_id2], colors[doc_id2], ty
24
25    colorscale = [[0, '#1d004d'], [.5, '#f2e5ff'], [1, '#f2e5d1']] # to color t
26
27    # we create a table with the data_matrix
28    table = ff.create_table(data_matrix, index=True, colorscale=colorscale)
29    # plot it with plotly
30    plotly.offline.iplot(table, filename='simple_table')
31
32    # devide whole figure space into 25 * 1:10 grids
33    gs = gridspec.GridSpec(25, 15)
34    fig = plt.figure(figsize=(25,5))
35
36    # in first 25*10 grids we plot heatmap
37    ax1 = plt.subplot(gs[:, :-5])
38    # plotting the heap map based on the pairwise distances
39    ax1 = sns.heatmap(np.round(s1_s2_dist,6), annot=True)
40    # set the x axis labels as recommended apparels title
41    ax1.set_xticklabels(sentance2.split())
42    # set the y axis labels as input apparels title
43    ax1.set_yticklabels(sentance1.split())
44    # set title as recommended apparels title
45    ax1.set_title(sentance2)
46
47    # in Last 25 * 10:15 grids we display image
48    ax2 = plt.subplot(gs[:, 10:16])
49    # we dont display grid lines and axis labels to images
50    ax2.grid(False)
51    ax2.set_xticks([])
52    ax2.set_yticks([])
53
54    # pass the url it display it
55    display_img(url, ax2, fig)
56

```

57 | plt.show()

In [28]:

```

1 def idf_w2v_brand(doc_id, w1, w2, num_results):
2     # doc_id: apparel's id in given corpus
3     # w1: weight for w2v features
4     # w2: weight for brand and color features
5
6     # pairwise_dist will store the distance from given input apparel to all r
7     # the metric we used here is cosine, the coside distance is mesured as K(
8     # http://scikit-Learn.org/stable/modules/metrics.html#cosine-similarity
9     idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_
10    ex_feat_dist = pairwise_distances(extra_features, extra_features[doc_id])
11    pairwise_dist = (w1 * idf_w2v_dist + w2 * ex_feat_dist)/float(w1 + w2)
12
13    # np.argsort will return indices of 9 smallest distances
14    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
15    #pdists will store the 9 smallest distances
16    pdists = np.sort(pairwise_dist.flatten())[0:num_results]
17
18    #data frame indices of the 9 smallest distace's
19    df_indices = list(data.index[indices])
20
21
22    for i in range(0, len(indices)):
23        heat_map_w2v_brand(data['title'].loc[df_indices[0]], data['title'].loc[
24            print('ASIN :', data['asin'].loc[df_indices[i]])
25            print('Brand :', data['brand'].loc[df_indices[i]])
26            print('euclidean distance from input :', pdists[i])
27            print('='*125)
28
29 idf_w2v_brand(12566, 5, 5, 20)
30 # in the give heat map, each cell contains the euclidean distance between wor

```

NameError

Traceback (most recent call last)

```

<ipython-input-28-822d8df49624> in <module>()
      27         print('*'*125)
      28
----> 29 idf_w2v_brand(12566, 5, 5, 20)
      30 # in the give heat map, each cell contains the euclidean distance bet
ween words i, j

<ipython-input-28-822d8df49624> in idf_w2v_brand(doc_id, w1, w2, num_results)
      21
      22     for i in range(0, len(indices)):
----> 23         heat_map_w2v_brand(data['title'].loc[df_indices[0]], data['tit
le'].loc[df_indices[i]], data['medium_image_url'].loc[df_indices[i]], indices
[0], indices[i], df_indices[0], df_indices[i], 'weighted')
      24         print('ASIN :', data['asin'].loc[df_indices[i]])
      25         print('Brand :', data['brand'].loc[df_indices[i]])

<ipython-input-27-ac9ce1dbaed2> in heat_map_w2v_brand(sentance1, sentance2, u
rl, doc_id1, doc_id2, df_id1, df_id2, model)
      11
      12     #s1_vec = np.array(#number_of_words_title1 * 300), each row is a
vector(weighted/avg) of length 300 corresponds to each word in give title
----> 13     s1_vec = get_word_vec(sentance1, doc_id1, model)

```

```
14      #s2_vec = np.array(#number_of_words_title2 * 300), each row is a  
vector(weighted/avg) of length 300 corresponds to each word in give title  
15      s2_vec = get_word_vec(sentance2, doc_id2, model)
```

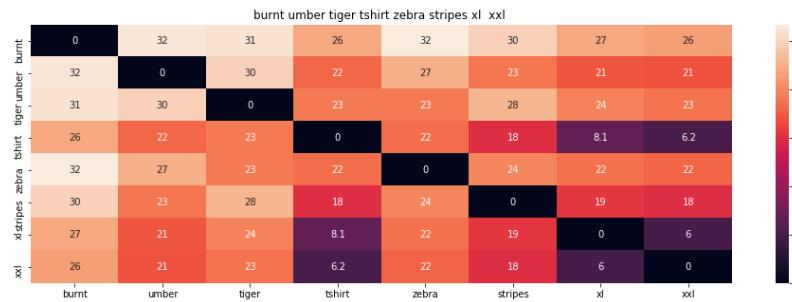
NameError: name 'get_word_vec' is not defined

In [0]:

```

1 # brand and color weight =50
2 # title vector weight = 5
3
4 idf_w2v_brand(12566, 5, 50, 20)

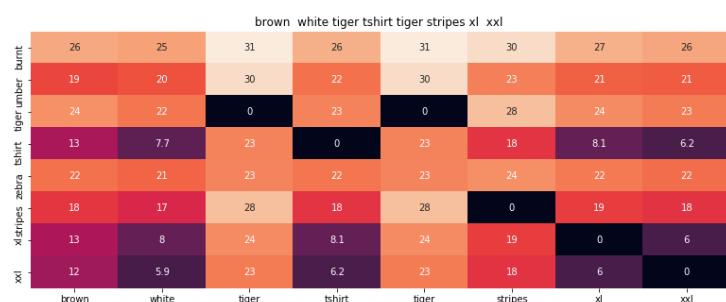
```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from input : 0.000355113636364



ASIN : B00JXQCWT0

Brand : Si Row

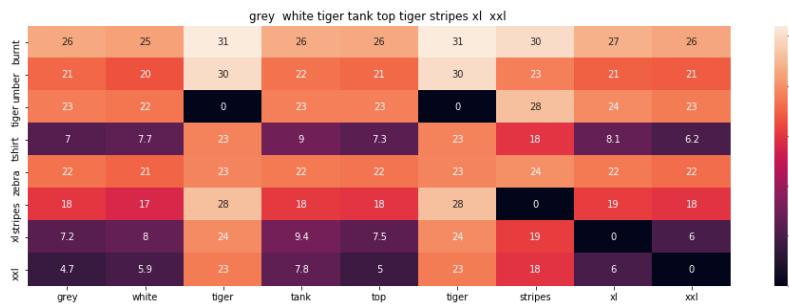
euclidean distance from input : 0.433722027865



ASIN : B00JXQASS6

Brand : Si Row

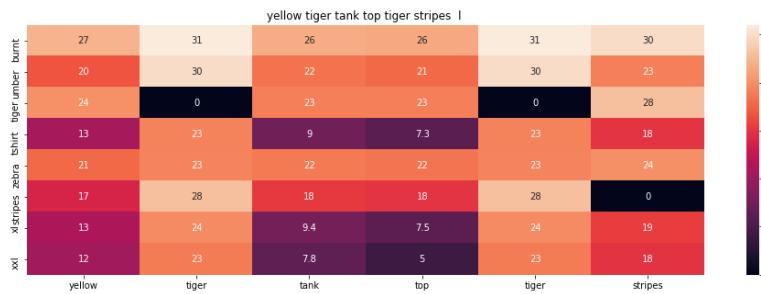
euclidean distance from input : 1.65509310669



ASIN : B00JXQAFZ2

Brand : Si Row

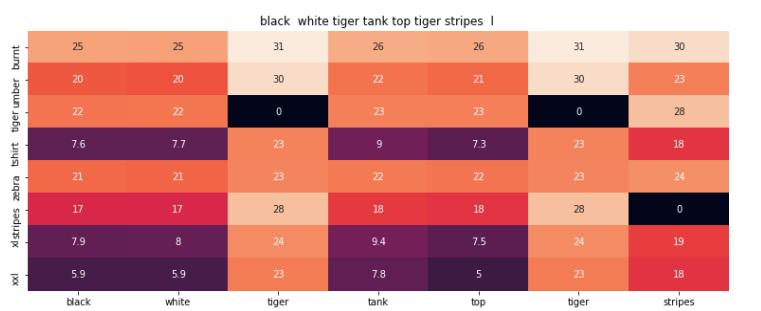
euclidean distance from input : 1.77293604103



ASIN : B00JXQAUWA

Brand : Si Row

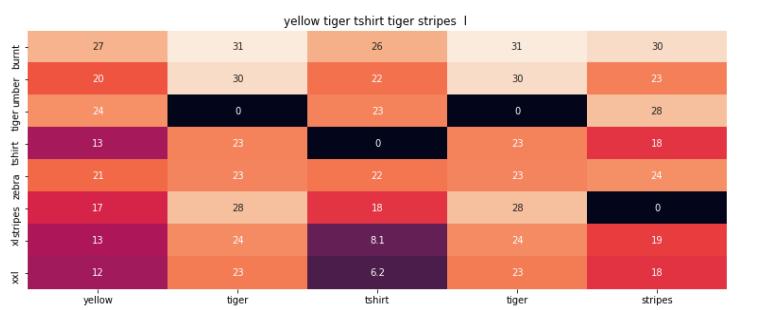
euclidean distance from input : 1.80287808538



ASIN : B00JXQA094

Brand : Si Row

euclidean distance from input : 1.80319609241



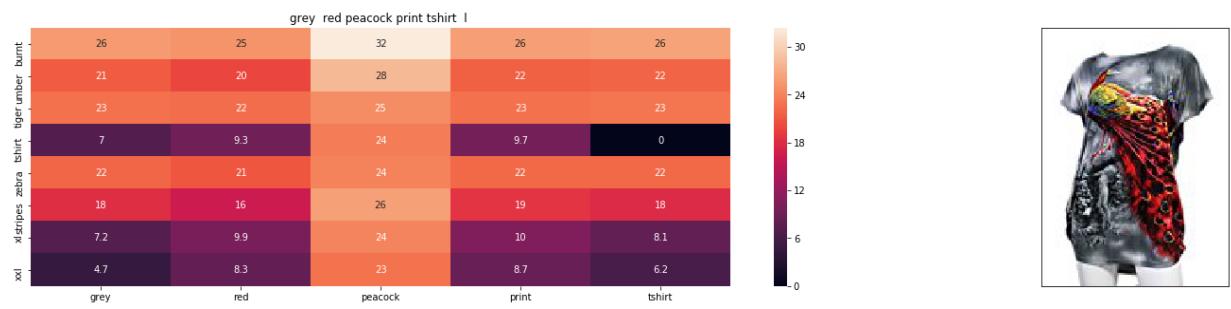
ASIN : B00JXQCUIC

Brand : Si Row

euclidean distance from input : 1.82141619628

=====

=====



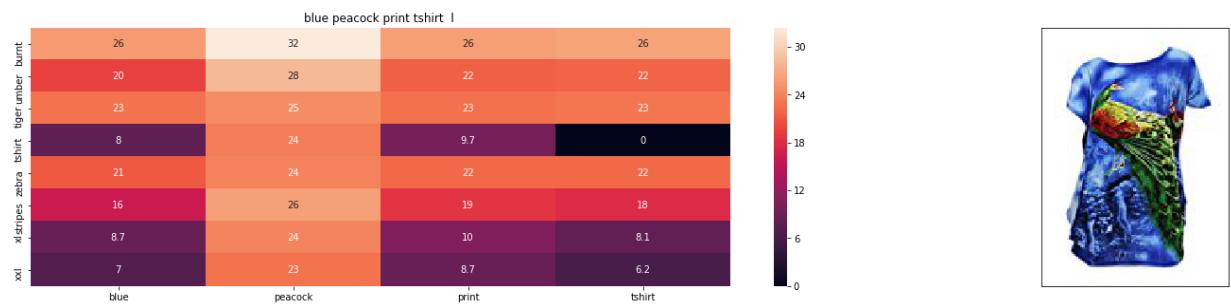
ASIN : B00JXQCFRS

Brand : Si Row

euclidean distance from input : 1.90777685025

=====

=====



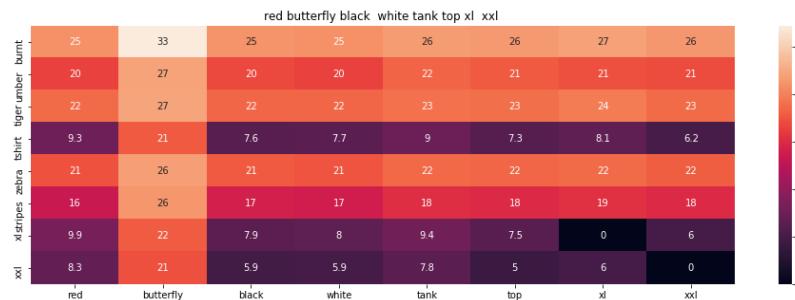
ASIN : B00JXQC8L6

Brand : Si Row

euclidean distance from input : 1.92142937433

=====

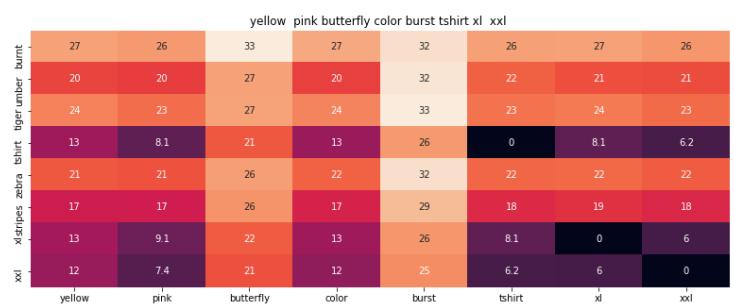
=====



ASIN : B00JV63CW2

Brand : Si Row

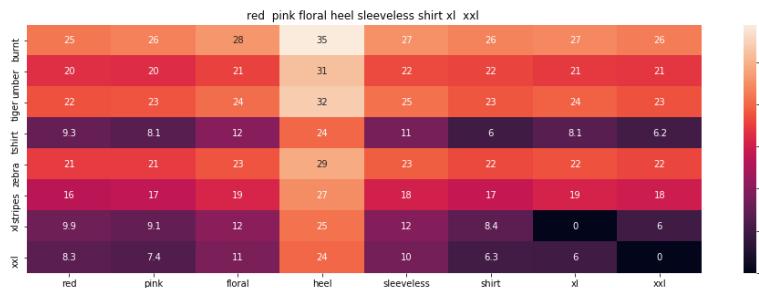
euclidean distance from input : 1.93646323497



ASIN : B00JXQBBMI

Brand : Si Row

euclidean distance from input : 1.95670381059



ASIN : B00JV63QQE

Brand : Si Row

euclidean distance from input : 1.9806066343



ASIN : B00JV63VC8

Brand : Si Row

euclidean distance from input : 2.01218559992



ASIN : B00JXQAX2C

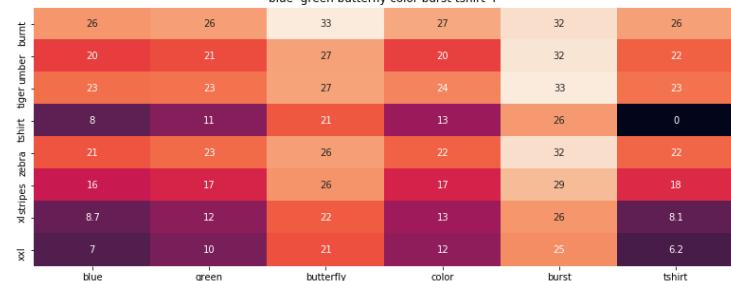
Brand : Si Row

euclidean distance from input : 2.01335178755

=====

=====

blue green butterfly color burst tshirt l



ASIN : B00JXQC0C8

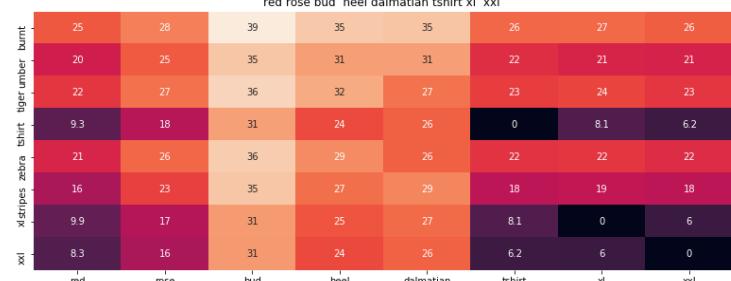
Brand : Si Row

euclidean distance from input : 2.01388334827

=====

=====

red rose bud heel dalmatian tshirt xl xxl

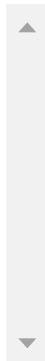
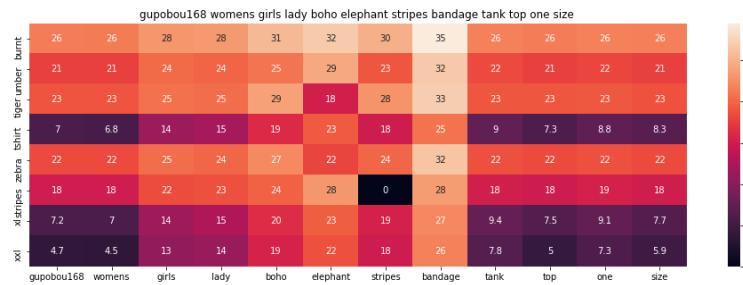


ASIN : B00JXQABB0

Brand : Si Row

euclidean distance from input : 2.0367257555

=====



ASIN : B01ER18406

Brand : GuPoBoU168

euclidean distance from input : 2.65620416778

=====

=====



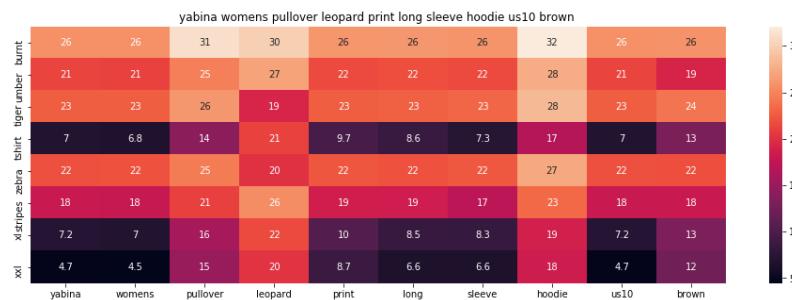
ASIN : B01LZ7BQ4H

Brand : WAYF

euclidean distance from input : 2.6849067823

=====

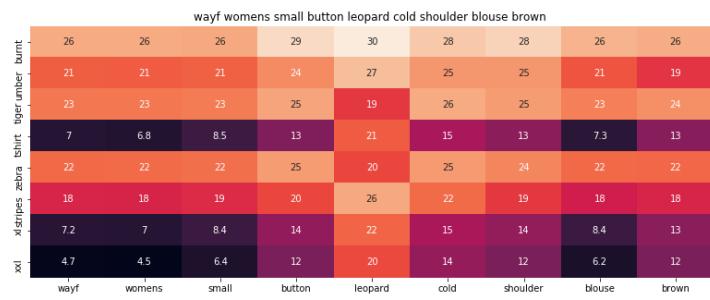
=====



ASIN : B01KJUM6JI

Brand : YABINA

euclidean distance from input : 2.68583819266



ASIN : B01M06V4X1

Brand : WAYF

euclidean distance from input : 2.69476194865

[10.2] Keras and Tensorflow to extract features

In [0]:

```
1 import numpy as np
2 from keras.preprocessing.image import ImageDataGenerator
3 from keras.models import Sequential
4 from keras.layers import Dropout, Flatten, Dense
5 from keras import applications
6 from sklearn.metrics import pairwise_distances
7 import matplotlib.pyplot as plt
8 import requests
9 from PIL import Image
10 import pandas as pd
11 import pickle
```

Using TensorFlow backend.

In [0]:

```
1 # https://gist.github.com/fchollet/f35fbc80e066a49d65f1688a7e99f069
2 # Code reference: https://blog.keras.io/building-powerful-image-classificatio
3
4
5
6 # This code takes 40 minutes to run on a modern GPU (graphics card)
7 # like Nvidia 1050.
8 # GPU (Nvidia 1050): 0.175 seconds per image
9
10 # This codse takes 160 minutes to run on a high end i7 CPU
11 # CPU (i7): 0.615 seconds per image.
12
13 #Do NOT run this code unless you want to wait a few hours for it to generate
14
15 # each image is converted into 25088 Length dense-vector
16
17
18 ...
19 # dimensions of our images.
20 img_width, img_height = 224, 224
21
22 top_model_weights_path = 'bottleneck_fc_model.h5'
23 train_data_dir = 'images2/'
24 nb_train_samples = 16042
25 epochs = 50
26 batch_size = 1
27
28
29 def save_bottlebeck_features():
30
31     #Function to compute VGG-16 CNN for image feature extraction.
32
33     asins = []
34     datagen = ImageDataGenerator(rescale=1. / 255)
35
36     # build the VGG16 network
37     model = applications.VGG16(include_top=False, weights='imagenet')
38     generator = datagen.flow_from_directory(
39         train_data_dir,
40         target_size=(img_width, img_height),
41         batch_size=batch_size,
42         class_mode=None,
43         shuffle=False)
44
45     for i in generator.filenames:
46         asins.append(i[2:-5])
47
48     bottleneck_features_train = model.predict_generator(generator, nb_train_s
49     bottleneck_features_train = bottleneck_features_train.reshape((16042,2508
50
51     np.save(open('16k_data_cnn_features.npy', 'wb'), bottleneck_features_trai
52     np.save(open('16k_data_cnn_feature_asins.npy', 'wb'), np.array(asins))
53
54
55 save_bottlebeck_features()
56
```

[10.3] Visual features based product similarity.

```
In [2]: 1 b_features = np.load('16k_data_cnn_features.npy')
```

```
In [5]: 1 df_cnn = pd.DataFrame(b_features)
```

```
Out[5]: (25088,
```

In [0]:

```

1 #Load the features and corresponding ASINS info.
2 bottleneck_features_train = np.load('16k_data_cnn_features.npy')
3 asins = np.load('16k_data_cnn_feature_asins.npy')
4 asins = list(asins)
5
6 # Load the original 16K dataset
7 data = pd.read_pickle('pickels/16k_apperal_data_preprocessed')
8 df_asins = list(data['asin'])
9
10
11 from IPython.display import display, Image, SVG, Math, YouTubeVideo
12
13
14 #get similar products using CNN features (VGG-16)
15 def get_similar_products_cnn(doc_id, num_results):
16     doc_id = asins.index(df_asins[doc_id])
17     pairwise_dist = pairwise_distances(bottleneck_features_train, bottleneck_
18
19     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
20     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
21
22     for i in range(len(indices)):
23         rows = data[['medium_image_url','title']].loc[data['asin']==asins[indices[i]]]
24         for idx, row in rows.iterrows():
25             display(Image(url=row['medium_image_url'], embed=True))
26             print('Product Title: ', row['title'])
27             print('Euclidean Distance from input image:', pdists[i])
28             print('Amazon Url: www.amazon.com/dp/' + asins[indices[i]])
29
30 get_similar_products_cnn(12566, 20)
31

```



Product Title: burnt umber tiger tshirt zebra stripes xl xxl
 Euclidean Distance from input image: 0.0
 Amazon Url: www.amazon.com/dp/B00JXQB5FQ



Product Title: pink tiger tshirt zebra stripes xl xxl
Euclidean Distance from input image: 30.0501
Amazon Url: www.amazon.com/dp/B00JXQASS6



Product Title: yellow tiger tshirt tiger stripes l
Euclidean Distance from input image: 41.2611
Amazon Url: www.amazon.com/dp/B00JXQCUIC



Product Title: brown white tiger tshirt tiger stripes xl xxl
Euclidean Distance from input image: 44.0002
Amazon Url: www.amazon.com/dp/B00JXQCWTO



Product Title: kawaii pastel tops tees pink flower design
Euclidean Distance from input image: 47.3825
Amazon Url: www.amazon.com/dp/B071FCWD97



Product Title: womens thin style tops tees pastel watermelon print
Euclidean Distance from input image: 47.7184

Amazon Url: www.amazon.com/dp/B01JUNHBRM



Product Title: kawaii pastel tops tees baby blue flower design

Euclidean Distance from input image: 47.9021

Amazon Url: www.amazon.com/dp/B071SBCY9W



Product Title: edv cheetah run purple multi xl

Euclidean Distance from input image: 48.0465

Amazon Url: www.amazon.com/dp/B01CUPYBM0



Product Title: danskin womens vneck loose performance tee xsmall pink ombre

Euclidean Distance from input image: 48.1019

Amazon Url: www.amazon.com/dp/B01F7PHXY8



Product Title: summer alpaca 3d pastel casual loose tops tee design

Euclidean Distance from input image: 48.1189

Amazon Url: www.amazon.com/dp/B01I80A93G



Product Title: miss chievous juniors striped peplum tank top medium shadowpeach

Euclidean Distance from input image: 48.1313

Amazon Url: www.amazon.com/dp/B0177DM70S



Product Title: red pink floral heel sleeveless shirt xl xxl

Euclidean Distance from input image: 48.1695

Amazon Url: www.amazon.com/dp/B00JV63QQE



Product Title: moana logo adults hot v neck shirt black xxl

Euclidean Distance from input image: 48.2568

Amazon Url: www.amazon.com/dp/B01LX6H43D



Product Title: abaday multicolor cartoon cat print short sleeve longline shirt large

Euclidean Distance from input image: 48.2657

Amazon Url: www.amazon.com/dp/B01CR57YY0



Product Title: kawaii cotton pastel tops tees peach pink cactus design
Euclidean Distance from input image: 48.3626
Amazon Url: www.amazon.com/dp/B071WYLBZS



Product Title: chicago chicago 18 shirt women pink
Euclidean Distance from input image: 48.3836
Amazon Url: www.amazon.com/dp/B01GXAZTRY



Product Title: yichun womens tiger printed summer tshirts tops
Euclidean Distance from input image: 48.4493
Amazon Url: www.amazon.com/dp/B010NN9RX0



Product Title: nancy lopez whimsy short sleeve whiteblacklemon drop xs
Euclidean Distance from input image: 48.4788
Amazon Url: www.amazon.com/dp/B01MPX6IDX



Product Title: womens tops tees pastel peach ice cream cone print
Euclidean Distance from input image: 48.558
Amazon Url: www.amazon.com/dp/B0734GRKZL



Product Title: uswomens mary j blige without tshirts shirt
Euclidean Distance from input image: 48.6144
Amazon Url: www.amazon.com/dp/B01M0XXFKK

Assignments

so we will concatenate the idf w2v features, one hot encoded brand and color features and the image based feature that we got using vgc16 model and then use them to find the similarity with weighted features

In [1]:

```
1 #import all the necessary packages.  
2  
3 from PIL import Image  
4 import requests  
5 from io import BytesIO  
6 import matplotlib.pyplot as plt  
7 import numpy as np  
8 import pandas as pd  
9 import warnings  
10 from bs4 import BeautifulSoup  
11 from nltk.corpus import stopwords  
12 from nltk.tokenize import word_tokenize  
13 import nltk  
14 import math  
15 import time  
16 import re  
17 import os  
18 import seaborn as sns  
19 from collections import Counter  
20 from sklearn.feature_extraction.text import CountVectorizer  
21 from sklearn.feature_extraction.text import TfidfVectorizer  
22 from sklearn.metrics.pairwise import cosine_similarity  
23 from sklearn.metrics import pairwise_distances  
24 from matplotlib import gridspec  
25 from scipy.sparse import hstack  
26 import plotly  
27 import plotly.figure_factory as ff  
28 from plotly.graph_objs import Scatter, Layout  
29  
30 plotly.offline.init_notebook_mode(connected=True)  
31 warnings.filterwarnings("ignore")  
32  
33 from IPython.core.display import display, HTML  
34 display(HTML("<style>.container { width:100% !important; }</style>"))
```

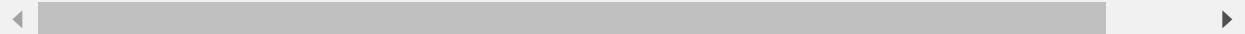
Loading the data

In [2]:

```
1 data = pd.read_pickle('pickels/16k_apperal_data_preprocessed')
2 data.head()
```

Out[2]:

	asin	brand	color	medium_image_url	product_type_name	title	format
4	B004GSI2OS	FeatherLite	Onyx Black/ Stone	https://images-na.ssl- images- amazon.com/images...		SHIRT	featherlite ladies long sleeve stain resistant...
6	B012YX2ZPI	HX- Kingdom Fashion T- shirts	White	https://images-na.ssl- images- amazon.com/images...		SHIRT	womens unique 100 cotton special olympics wor...
15	B003BSRPB0	FeatherLite	White	https://images-na.ssl- images- amazon.com/images...		SHIRT	featherlite ladies moisture free mesh sport sh...
27	B014ICEJ1Q	FNC7C	Purple	https://images-na.ssl- images- amazon.com/images...		SHIRT	supernatural chibis sam dean castiel neck tshi...
46	B01NACPBG2	Fifth Degree	Black	https://images-na.ssl- images- amazon.com/images...		SHIRT	fifth degree womens gold foil graphic tees jun...



IDF Weighted W2V features

In [3]:

```
1 idf_title_vectorizer = CountVectorizer()
2 idf_title_features = idf_title_vectorizer.fit_transform(data['title'])
```

In [4]:

```
1 # idf_title_features.shape = #data_points * #words_in_corpus
2 # CountVectorizer().fit_transform(courpus) returns the a sparase matrix of di
3 # idf_title_features[doc_id, index_of_word_in_corpus] = number of times the w
4 def nContaining(word):
5     # return the number of documents which has the given word
6     return sum(1 for blob in data['title'] if word in blob.split())
7
8 def idf(word):
9     # idf = log(#number of docs / #number of docs which had the given word)
10    return math.log(data.shape[0] / (nContaining(word)))
```

In [5]:

```

1 # we need to convert the values into float
2 idf_title_features = idf_title_features.astype(np.float)
3
4 for i in idf_title_vectorizer.vocabulary_.keys():
5     # for every word in whole corpus we will find its idf value
6     idf_val = idf(i)
7
8     # to calculate idf_title_features we need to replace the count values with
9     # idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0]
10    for j in idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonze
11
12        # we replace the count values of word i in document j with idf_value
13        # idf_title_features[doc_id, index_of_word_in_corpus] = idf value of
14        idf_title_features[j,idf_title_vectorizer.vocabulary_[i]] = idf_val #

```

Building the idf weighted w2v

In [6]:

```

1 import pickle
2 #if you do NOT have RAM >= 12GB, use the code below.
3 with open('word2vec_model', 'rb') as handle:
4     model = pickle.load(handle)

```

In [7]:

```

1 # vocab = stores all the words that are there in google w2v model
2 # vocab = model.wv.vocab.keys() # if you are using Google word2Vec
3
4 vocab = model.keys()
5 # this function will add the vectors of each word and returns the avg vector
6 def build_avg_vec(sentence, num_features, doc_id, m_name):
7     # sentence: its title of the apparel
8     # num_features: the length of word2vec vector, its values = 300
9     # m_name: model information it will take two values
10    # if m_name == 'avg', we will append the model[i], w2v representation
11    # if m_name == 'weighted', we will multiply each w2v[word] with the i
12
13    featureVec = np.zeros((num_features,), dtype="float32")
14    # we will initialize a vector of size 300 with all zeros
15    # we add each word2vec(wordi) to this featureVec
16    nwords = 0
17
18    for word in sentence.split():
19        nwords += 1
20        if word in vocab:
21            if m_name == 'weighted' and word in idf_title_vectorizer.vocabul
22                featureVec = np.add(featureVec, idf_title_features[doc_id, id
23            elif m_name == 'avg':
24                featureVec = np.add(featureVec, model[word])
25        if(nwords>0):
26            featureVec = np.divide(featureVec, nwords)
27        # returns the avg vector of given sentence, its of shape (1, 300)
28
29    return featureVec

```

In [8]: 1 import datetime as dt

In [9]: 1 start = dt.datetime.now()
 2 doc_id = 0
 3 w2v_title_weight = []
 4 # for every title we build a weighted vector representation
 5 for i in data['title']:
 6 w2v_title_weight.append(build_avg_vec(i, 300, doc_id, 'weighted'))
 7 doc_id += 1
 8 # w2v_title = np.array(# number of doc in courpus * 300), each row correspond
 9 w2v_title_weight = np.array(w2v_title_weight)
 10 print('time taken for execution is:',dt.datetime.now() - start)

time taken for execution is: 0:00:04.014132

In [10]: 1 w2v_title_weight.shape

Out[10]: (16042, 300)

ONe hot encoding for brand and color

In [13]: 1 # some of the brand values are empty.
 2 # Need to replace Null with string "NULL"
 3 data['brand'].fillna(value="Not given", inplace=True)
 4
 5 # replace spaces with hyphen
 6 brands = [x.replace(" ", "-") for x in data['brand'].values]
 7 #types = [x.replace(" ", "-") for x in data['product_type_name'].values]
 8 colors = [x.replace(" ", "-") for x in data['color'].values]
 9
 10 brand_vectorizer = CountVectorizer()
 11 brand_features = brand_vectorizer.fit_transform(brands)
 12
 13 color_vectorizer = CountVectorizer()
 14 color_features = color_vectorizer.fit_transform(colors)
 15
 16 extra_features = hstack((brand_features,color_features)).tocsr()

In [14]:

```

1 ## final features being concatenated
2 from IPython.display import display, SVG, Math, YouTubeVideo
3
4 #Display an image
5 def display_img(url,ax,fig):
6     # we get the url of the apparel and download it
7     response = requests.get(url)
8     img = Image.open(BytesIO(response.content))
9     # we will display it in notebook
10    plt.imshow(img)
11
12 # Utility functions
13
14 def get_word_vec(sentence, doc_id, m_name):
15     # sentence : title of the apparel
16     # doc_id: document id in our corpus
17     # m_name: model information it will take two values
18     # if m_name == 'avg', we will append the model[i], w2v representation
19     # if m_name == 'weighted', we will multiply each w2v[word] with the i
20     vec = []
21     for i in sentence.split():
22         if i in vocab:
23             if m_name == 'weighted' and i in idf_title_vectorizer.vocabulary_
24                 vec.append(idf_title_features[doc_id, idf_title_vectorizer.voc
25             elif m_name == 'avg':
26                 vec.append(model[i])
27         else:
28             # if the word in our courpus is not there in the google word2vec
29             vec.append(np.zeros(shape=(300,)))
30     # we will return a numpy array of shape (#number of words in title * 300
31     # each row represents the word2vec representation of each word (weighted)
32     return np.array(vec)
33
34 def get_distance(vec1, vec2):
35     # vec1 = np.array(#number_of_words_title1 * 300), each row is a vector of
36     # vec2 = np.array(#number_of_words_title2 * 300), each row is a vector of
37
38     final_dist = []
39     # for each vector in vec1 we caluclate the distance(euclidean) to all vec
40     for i in vec1:
41         dist = []
42         for j in vec2:
43             # np.linalg.norm(i-j) will result the euclidean distance between
44             dist.append(np.linalg.norm(i-j))
45         final_dist.append(np.array(dist))
46     # final_dist = np.array(#number of words in title1 * #number of words in
47     # final_dist[i,j] = euclidean distance between vectors i, j
48     return np.array(final_dist)
49
50 def create_table(df_id1, df_id2, doc_id1, doc_id2):
51     import plotly
52     import plotly.figure_factory as ff
53     from plotly.graph_objs import Scatter, Layout
54     plotly.offline.init_notebook_mode(connected=False)
55
56     data_matrix = [['Asin', 'Brand', 'Color'],

```

```

57 [data['asin'].loc[df_id1],brands[doc_id1], colors[doc_id1]
58 [data['asin'].loc[df_id2],brands[doc_id2], colors[doc_id2]
59
60 colorscale = [[0, '#1d004d'], [.5, '#f2e5ff'], [1, '#f2e5d1']] # to color
61
62 # we create a table with the data_matrix
63 table = ff.create_table(data_matrix, index=True, colorscale=colorscale)
64 # plot it with plotly
65 plotly.offline.iplot(table, filename='simple_table')
66
67
68 def heat_map_idf_brand_col_img(sentance1, sentance2, url, doc_id1, doc_id2, model):
69
70     # sentance1 : title1, input apparel
71     # sentance2 : title2, recommended apparel
72     # url: apparel image url
73     # doc_id1: document id of input apparel
74     # doc_id2: document id of recommended apparel
75     # df_id1: index of document1 in the data frame
76     # df_id2: index of document2 in the data frame
77     # model: it can have two values, 1. avg 2. weighted
78
79     #s1_vec = np.array(#number_of_words_title1 * 300), each row is a vector(1*300)
80     s1_vec = get_word_vec(sentance1, doc_id1, model)
81     #s2_vec = np.array(#number_of_words_title2 * 300), each row is a vector(1*300)
82     s2_vec = get_word_vec(sentance2, doc_id2, model)
83
84     # s1_s2_dist = np.array(#number of words in title1 * #number of words in title2)
85     # s1_s2_dist[i,j] = euclidean distance between words i, j
86     s1_s2_dist = get_distance(s1_vec, s2_vec)
87
88     #Create a table
89     #create_table(df_id1, df_id2, doc_id1, doc_id2)
90     '''Extremely high ram consumption. Try to avoid running the above line'''
91
92
93     # devide whole figure space into 25 * 1:10 grids
94     gs = gridspec.GridSpec(25, 15)
95     fig = plt.figure(figsize=(25,5))
96
97     # in first 25*10 grids we plot heatmap
98     ax1 = plt.subplot(gs[:, :-5])
99     # plotting the heap map based on the pairwise distances
100    ax1 = sns.heatmap(np.round(s1_s2_dist,6), annot=True)
101    # set the x axis labels as recommended apparels title
102    ax1.set_xticklabels(sentance2.split())
103    # set the y axis labels as input apparels title
104    ax1.set_yticklabels(sentance1.split())
105    # set title as recommended apparels title
106    ax1.set_title(sentance2)
107
108    # in last 25 * 10:15 grids we display image
109    ax2 = plt.subplot(gs[:, 10:16])
110    # we dont display grid lines and axis labels to images
111    ax2.grid(False)
112    ax2.set_xticks([])
113    ax2.set_yticks([])
```

```

114
115      # pass the url it display it
116      display_img(url, ax2, fig)
117
118      plt.show()

```

CNN features and concatenating all the features

In [15]:

```

1  cnn_features = np.load('16k_data_cnn_features.npy')
2  cnn_features_asins = list(np.load('16k_data_cnn_feature_asins.npy'))

```

In [16]:

```

1  def idf_w2v_final(doc_id, w1, w2,w3, num_results):
2      # doc_id: apparel's id in given corpus
3      # w1: weight for w2v features
4      # w2: weight for brand and color features
5      #w3 :weight for cnn features
6
7
8
9      # pairwise_dist will store the distance from given input apparel to all r
10     # the metric we used here is cosine, the coside distance is mesured as K(
11     idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id])
12     ex_feat_dist = pairwise_distances(extra_features, extra_features[doc_id])
13     cnn_dist = pairwise_distances(cnn_features,cnn_features[doc_id].reshape(1,-1))
14     pairwise_dist = (w1 * idf_w2v_dist + w2 * ex_feat_dist) + w3 * cnn_dist
15
16     # np.argsort will return indices of 9 smallest distances
17     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
18     #pdists will store the 9 smallest distances
19     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
20
21     #data frame indices of the 9 smallest distace's
22     df_indices = list(data.index[indices])
23
24
25     for i in range(0, len(indices)):
26         heat_map_idf_brand_col_img(data['title'].loc[df_indices[i]],data['title'])
27         print('ASIN : ',data['asin'].loc[df_indices[i]])
28         print('Brand : ',data['brand'].loc[df_indices[i]])
29         print('Color : ',data['color'].loc[df_indices[i]])
30         print('Product Type : ',data['product_type_name'].loc[df_indices[i]])
31         print('Euclidean distance from input : ', pdists[i])
32
33         print('*'*125)

```

Giving equal weights to all

In [19]: 1 idf_w2v_final(12566, 1, 1, 1, 5)



ASIN : B00JXQB5FQ

Brand : Si Row

Color : Brown

Product Type : TOYS_AND_GAMES

Euclidean distance from input : 0.018637641333



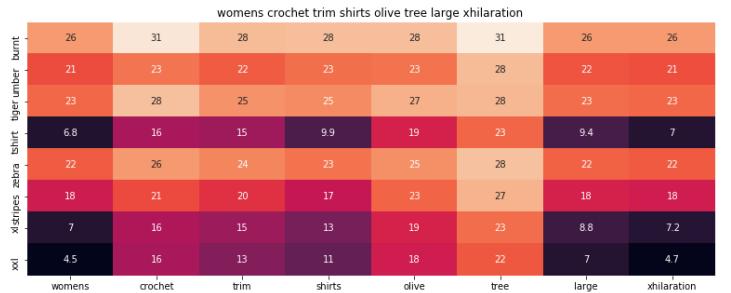
ASIN : B00JXQASS6

Brand : Si Row

Color : Pink

Product Type : TOYS_AND_GAMES

Euclidean distance from input : 21.6480869297



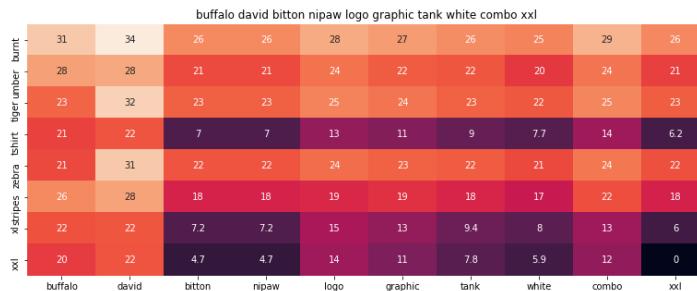
ASIN : B06XBHNM7J

Brand : Xhilaration

Color : Olive Tree

Product Type : SHIRT

Euclidean distance from input : 22.7769724431



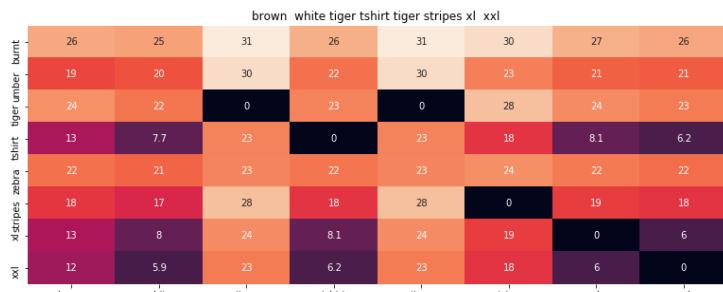
ASIN : B018H5AZXQ

Brand : Buffalo

Color : White Combo

Product Type : SHIRT

Euclidean distance from input : 23.0580307546



ASIN : B00JXQCWTO

Brand : Si Row

Color : Brown

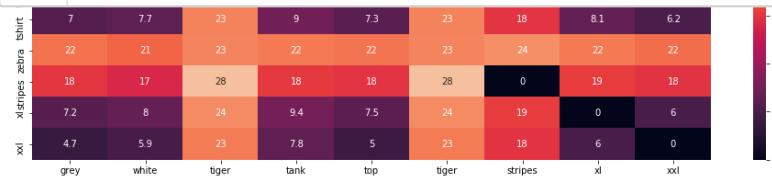
Product Type : TOYS_AND_GAMES

Euclidean distance from input : 23.2869067192

Givning more weight to brand+color features

In [20]:

```
1 ## Giving more weightage to brand an color features
2 idf_w2v_final(12566, 1, 20, 1, 10)
```



ASIN : B00JXQAFZ2

Brand : Si Row

Color : Grey

Product Type : TOYS_AND_GAMES

Euclidean distance from input : 36.0956416202



Giving more weightage to idf w22v features

In [21]:

```

1 ## Giving more weightage to idf-w2v features
2
3 idf_w2v_final(12566, 20,1,1, 10)

```



ASIN : B00JXQAFZ2

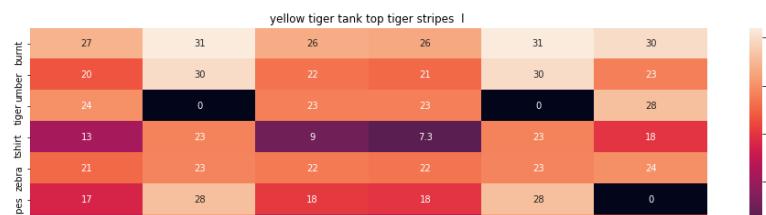
Brand : Si Row

Color : Grey

Product Type : TOYS_AND_GAMES

Euclidean distance from input : 111.06863966

```
=====
=====
```



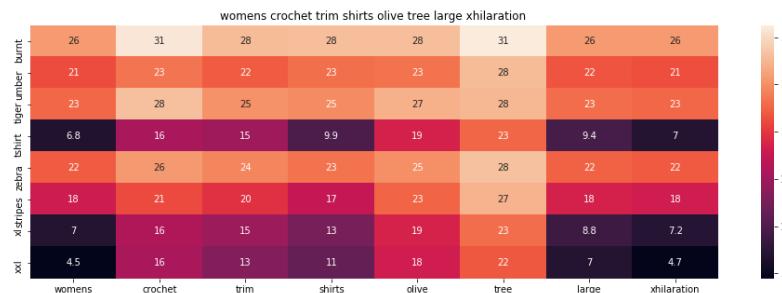
Giving more weight to cnn based features

In [22]:

```

1 ## Giving more weightage to cnn features
2 idf_w2v_final(12566, 1,1,20, 10)

```



ASIN : B06XBHNM7J

Brand : Xhilaration

Color : Olive Tree

Product Type : SHIRT

Euclidean distance from input : 44.1881395879

```
=====
=====
```



Observation and conclusion

There is a wide range of suggestions that pop up when we use equal weights for each of our four features. There is a healthy combination of Brand names, Colors and Textual similarity with the title. In general, what we observe mostly is that the suggestions are mostly related to animal names like zebra, tiger and a few other animals.

Giving more weight to IDF vectors has resulted in most of the recommendations belonging to the same brand and almost 90% of the top 20 recommendations had tiger and zebra prints specifically!

Giving more weights to brand names has resulted in suggestions which belong to the input brand, in our case we have Si Row. Almost 90% of the recommendations have the same brand name.

Here we have given more weights to colors and not to our surprise we see that most of the recommendations belonged to the same color, in our case the input color is Brown. Here again, 90% of the recommendations belonged to apparels having the same color as the query data point, i.e. Brown color.

Lastly, we have tried giving more weights to the CNN features and see how the system behaves! As we can see based solely on the CNN features, the model has picked up some really nice and diverse collections for recommendation purpose. Our query data point had a tiger image, it was white in color. The recommendations that were shown also had more or less the same features. We can see lots of diversity here including zebra stripes, tiger images, white images etc!