

Amazon Apparel Recommendations

[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 [54]:

```
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_reivew', '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', 'title', 'formatted_price']]
```

```
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

	asin	brand	color	medium_image_url	product_type_name	title	formatted_price
0	B016I2TS4W	FNC7C	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	Minions Como Superheroes Ironman Long Sleeve R...	None
1	B01N49AI08	FIG Clothing	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	FIG Clothing Womens Izo Tunic	None
2	B01JDPCOHO	FIG Clothing	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	FIG Clothing Womens Won Top	None
3	B01N19U5H5	Focal18	None	https://images-na.ssl-images-amazon.com/images...	SHIRT	Focal18 Sailor Collar Bubble Sleeve Blouse Shi...	None
4	B004GSI2OS	FeatherLite	Onyx Black/Stone	https://images-na.ssl-images-amazon.com/images...	SHIRT	Featherlite Ladies' Long Sleeve Stain Resistan...	\$26.26

[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 print(data['color'].describe())
2
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 which have null values price ==
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, row 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)



These shirts exactly same except in color



:B00G278GZ6



:B00G278W6O



: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	formatted_price
4	B004GSI2OS	FeatherLite	Onyx Black/Stone	https://images-na.ssl-images-amazon.com/images...	SHIRT	Featherlite Ladies' Long Sleeve Stain Resistan...	\$26.26
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...	\$9.99
11	B001LOUGE4	Fitness Etc.	Black	https://images-na.ssl-images-amazon.com/images...	SHIRT	Ladies Cotton Tank 2x1 Ribbed Tank Top	\$11.99
15	B003BSRPB0	FeatherLite	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	FeatherLite Ladies' Moisture Free Mesh Sport S...	\$20.54
21	B014ICEDNA	FNC7C	Purple	https://images-na.ssl-images-amazon.com/images...	SHIRT	Supernatural Chibis Sam Dean And Castiel Short...	\$7.50

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[0])
```

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	formatted_price
61973	B06Y1KZ2WB	Éclair	Black/Pink	https://images-na.ssl-images-amazon.com/images...	SHIRT	Éclair Women's Printed Thin Strap Blouse Black...	\$24.99
133820	B010RV33VE	xiaoming	Pink	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Womens Sleeveless Loose Long T-shirts...	\$18.19
81461	B01DDSDLNS	xiaoming	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Women's White Long Sleeve Single Brea...	\$21.58
75995	B00X5LYO9Y	xiaoming	Red Anchors	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Stripes Tank Patch/Bear Sleeve Anchor...	\$15.91
151570	B00WPJG35K	xiaoming	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	xiaoming Sleeve Sheer Loose Tassel Kimono Woma...	\$14.32

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', 'Queen', 'of', 'Diamonds', 'W
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', 'The', 'Queen', 'of', 'Diamonds'
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 strings
24        count = 0
25
26        # itertools.zip_longest(a,b): will map the corresponding words in both strings, it will appened Non
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', 'd'), ('d', None)]
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 are considering it as those two
35        # if the number of words in which both strings differ are < 2 , we are considering it as those two
36        if (length - count) > 2: # number of words in which both sensences differ
37            # if both strings are differ by more than 2 words we include the 1st string index
38            stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[i]])
39
40        # if the comprision between is between num_data_points, num_data_points-1 strings and they dif
41        if j == num_data_points-1: stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[j]])
```

```

42
43      # start searching for similar apperals corresponds 2nd string
44      i = j
45      break
46  else:
47      j += 1
48 if previous_i == i:
49     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_apperial_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')

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 apparel's title
14     a = data['title'].loc[i].split()
15     # store the list of words of ith string in a, ex: a = ['tokidoki', 'The', 'Queen', 'of', 'Diamonds', 'W
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', 'The', 'Queen', 'of', 'Diamonds'
20
21         length = max(len(a), len(b))
22
23         # count is used to store the number of words that are matched in both strings
24         count = 0
25
26         # itertools.zip_longest(a,b): will map the corresponding words in both strings, it will append Non
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', 'd'), ('d', None) ]
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 are considering it as those two
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_apperal_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', 'ai n', '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', 'he r', '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', 'w ere', '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', 'w ouldn', 'we', 'do', 'having', 'yours', 'other', 'wasn', 'it', 'with', 'once', 'here', 'don', 'o', 'whom', 'th is', '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	formatted_price
4	B004GSI2OS	FeatherLite	Onyx Black/Stone	https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies long sleeve stain resistant...	\$26.26
6	B012YX2ZPI	HX-Kingdom Fashion T-shirts	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	womens unique 100 cotton special olympics wor...	\$9.99
15	B003BSRPB0	FeatherLite	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies moisture free mesh sport sh...	\$20.54
27	B014ICEJ1Q	FNC7C	Purple	https://images-na.ssl-images-amazon.com/images...	SHIRT	supernatural chibis sam dean castiel neck tshi...	\$7.39
46	B01NACPBG2	Fifth Degree	Black	https://images-na.ssl-images-amazon.com/images...	SHIRT	fifth degree womens gold foil graphic tees jun...	\$6.95

In [0]: 1 data.to_pickle('pickels/16k_apperial_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 [5]: 1 data = pd.read_pickle('pickles/16k_apperal_data_preprocessed')
2 data.head()
```

Out[5]:

	asin	brand	color	medium_image_url	product_type_name	title	formatted_price
4	B004GSI2OS	FeatherLite	Onyx Black/Stone	https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies long sleeve stain resistant...	\$26.26
6	B012YX2ZPI	HX-Kingdom Fashion T-shirts	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	womens unique 100 cotton special olympics wor...	\$9.99
15	B003BSRPB0	FeatherLite	White	https://images-na.ssl-images-amazon.com/images...	SHIRT	featherlite ladies moisture free mesh sport sh...	\$20.54
27	B014ICEJ1Q	FNC7C	Purple	https://images-na.ssl-images-amazon.com/images...	SHIRT	supernatural chibis sam dean castiel neck tshi...	\$7.39
46	B01NACPBG2	Fifth Degree	Black	https://images-na.ssl-images-amazon.com/images...	SHIRT	fifth degree womens gold foil graphic tees jun...	\$6.95

In [29]:

```
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
12 # plotting code to understand the algorithm's decision.
13 def plot_heatmap(keys, values, labels, url, text):
14     # keys: list of words of recommended title
15     # values: len(values) == len(keys), values(i) represents the occurrence of the word keys(i)
16     # labels: len(labels) == len(keys), the values of labels depends on the model we are using
17         # if model == 'bag of words': labels(i) = values(i)
18         # if model == 'tfidf weighted bag of words':labels(i) = tfidf(keys(i))
19         # if model == 'idf weighted bag of words':labels(i) = idf(keys(i))
20     # url : apparel's url
21
22     # we will devide the whole figure into two parts
23     gs = gridspec.GridSpec(2, 2, width_ratios=[4,1], height_ratios=[4,1])
24     fig = plt.figure(figsize=(25,3))
25
26     # 1st, plotting heat map that represents the count of commonly occurred words in title2
27     ax = plt.subplot(gs[0])
28     # it displays a cell in white color if the word is intersection(list of words of title1 and list of
29     ax = sns.heatmap(np.array([values]), annot=np.array([labels]))
30     ax.set_xticklabels(keys) # set that axis labels as the words of title
31     ax.set_title(text) # apparel title
32
33     # 2nd, plotting image of the the apparel
34     ax = plt.subplot(gs[1])
35     # we don't want any grid lines for image and no labels on x-axis and y-axis
36     ax.grid(False)
37     ax.set_xticks([])
38     ax.set_yticks([])
39
40     # we call dispaly_img based with paramete url
41     display_img(url, ax, fig)
```

```
42
43     # displays combine figure ( heat map and image together)
44     plt.show()
45
46 def plot_heatmap_image(doc_id, vec1, vec2, url, text, model):
47
48     # doc_id : index of the title1
49     # vec1 : input apparels's vector, it is of a dict type {word:count}
50     # vec2 : recommended apparels's vector, it is of a dict type {word:count}
51     # url : apparels image url
52     # text: title of recomonded apparel (used to keep title of image)
53     # model, it can be any of the models,
54         # 1. bag_of_words
55         # 2. tfidf
56         # 3. idf
57
58     # we find the common words in both titles, because these only words contribute to the distance between
59     intersection = set(vec1.keys()) & set(vec2.keys())
60
61     # we set the values of non intersecting words to zero, this is just to show the difference in heatmap
62     for i in vec2:
63         if i not in intersection:
64             vec2[i]=0
65
66     # for labeling heatmap, keys contains list of all words in title2
67     keys = list(vec2.keys())
68     # if ith word in intersection(list of words of title1 and list of words of title2): values(i)=count of
69     values = [vec2[x] for x in vec2.keys()]
70
71     # labels: len(labels) == len(keys), the values of labels depends on the model we are using
72         # if model == 'bag of words': labels(i) = values(i)
73         # if model == 'tfidf weighted bag of words':labels(i) = tfidf(keys(i))
74         # if model == 'idf weighted bag of words':labels(i) = idf(keys(i))
75
76     if model == 'bag_of_words':
77         labels = values
78     elif model == 'tfidf':
79         labels = []
80         for x in vec2.keys():
81             # tfidf_title_vectorizer.vocabulary_ it contains all the words in the corpus
82             # tfidf_title_features[doc_id, index_of_word_in_corpus] will give the tfidf value of word in gi
83             if x in tfidf_title_vectorizer.vocabulary_:
```

```

84         labels.append(tfidf_title_features[doc_id, tfidf_title_vectorizer.vocabulary_[x]])
85     else:
86         labels.append(0)
87 elif model == 'idf':
88     labels = []
89     for x in vec2.keys():
90         # idf_title_vectorizer.vocabulary_ it contains all the words in the corpus
91         # idf_title_features[doc_id, index_of_word_in_corpus] will give the idf value of word in given
92         if x in idf_title_vectorizer.vocabulary_:
93             labels.append(idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[x]])
94         else:
95             labels.append(0)
96
97 plot_heatmap(keys, values, labels, url, text)
98
99
100 # this function gets a list of words along with the frequency of each
101 # word given "text"
102 def text_to_vector(text):
103     word = re.compile(r'\w+')
104     words = word.findall(text)
105     # words stores list of all words in given string, you can try 'words = text.split()' this will also give
106     return Counter(words) # Counter counts the occurrence of each word in list, it returns dict type object
107
108
109
110 def get_result(doc_id, content_a, content_b, url, model):
111     text1 = content_a
112     text2 = content_b
113
114     # vector1 = dict{word11:#count, word12:#count, etc.}
115     vector1 = text_to_vector(text1)
116
117     # vector1 = dict{word21:#count, word22:#count, etc.}
118     vector2 = text_to_vector(text2)
119
120     plot_heatmap_image(doc_id, vector1, vector2, url, text2, model)

```

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

In [9]:

```
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 occured in that doc
12
13
```

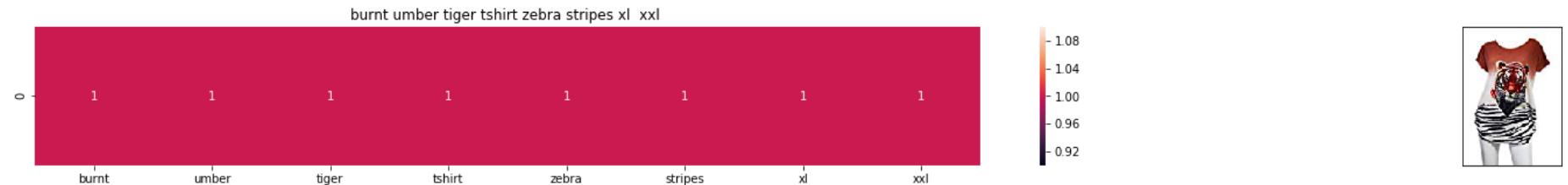
Out[9]: (16042, 12609)

In [10]:

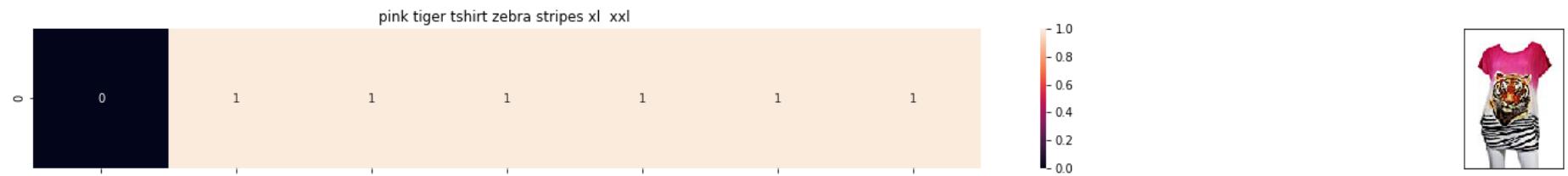
```

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 remaining apparels
5     # the metric we used here is cosine, the coside distance is mesured as K(X, Y) = <X, Y> / (||X|| * ||Y||)
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'].loc[df_indices[i]], data['med
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
=====
```



[8.5] TF-IDF based product similarity

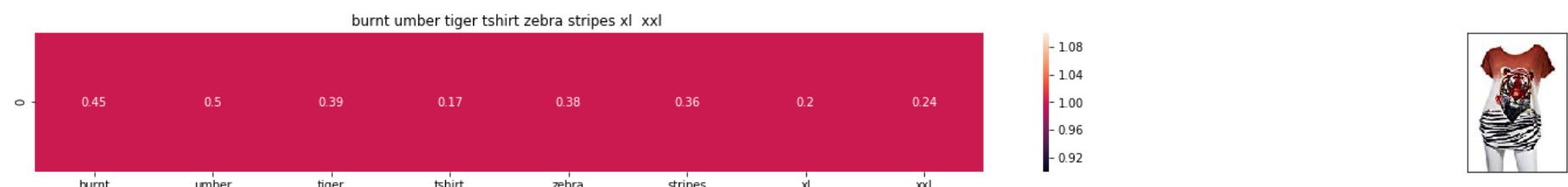
```
In [11]: 1 tfidf_title_vectorizer = TfidfVectorizer(min_df = 0)
2 tfidf_title_features = tfidf_title_vectorizer.fit_transform(data['title'])
3 # tfidf_title_features.shape = #data_points * #words_in_corpus
4 # CountVectorizer().fit_transform(courpus) returns the a sparase matrix of dimensions #data_points * #words
5 # tfidf_title_features[doc_id, index_of_word_in_corpus] = tfidf values of the word in given doc
```

In [12]:

```

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 remaining apparels
5     # the metric we used here is cosine, the coside distance is mesured as K(X, Y) = <X, Y> / (||X|| * ||Y||)
6     # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(tfidf_title_features, tfidf_title_features[doc_id])
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'].loc[df_indices[i]], data['me
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 word, the color represents the

```

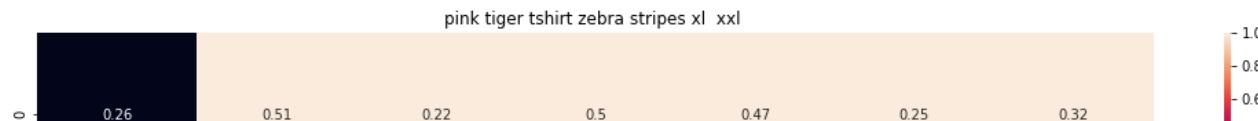


ASIN : B00JXQB5FQ

BRAND : Si Row

Eucliden distance from the given image : 0.0





[8.5] IDF based product similarity

```
In [13]: 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 dimensions #data_points * #words
6 # idf_title_features[doc_id, index_of_word_in_corpus] = number of times the word occured in that doc
```



```
In [0]: 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 [0]: 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 the idf values of the word
9     # idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0] will return all documents in
10    for j in idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0]:
11
12        # we replace the count values of word i in document j with idf_value of word i
13        # idf_title_features[doc_id, index_of_word_in_courpus] = idf value of word
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 remaining apparels
5     # the metric we used here is cosine, the coside distance is mesured as K(X, Y) = <X, Y> / (||X|| * ||Y||)
6     # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(idf_title_features,idf_title_features[doc_id])
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'].loc[df_indices[i]], data['med
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 word, the color represents the i

```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from the given image : 0.0





[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 [14]:

```
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', binary=True)
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 [15]:

```
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 of word i
8         # if m_name == 'weighted', we will multiply each w2v[word] with the idf(word)
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.vocabulary_[i]] * model[i])
14            elif m_name == 'avg':
15                vec.append(model[i])
16        else:
17            # if the word in our courpus is not there in the google word2vec corpus, we are just ignoring i
18            vec.append(np.zeros(shape=(300,)))
19    # we will return a numpy array of shape (#number of words in title * 300 ) 300 = len(w2v_model[word])
20    # each row represents the word2vec representation of each word (weighted/avg) in given sentance
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 length 300 corresponds to each
25     # vec2 = np.array(#number_of_words_title2 * 300), each row is a vector of length 300 corresponds to each
26
27     final_dist = []
28     # for each vector in vec1 we caluclate the distance(euclidean) to all vectors in vec2
29     for i in vec1:
30         dist = []
31         for j in vec2:
32             # np.linalg.norm(i-j) will result the euclidean distance between vectors i, j
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 title2)
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     # sentance1 : title1, input apparel
```

```
42 # sentance2 : 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(weighted/avg) of length 300 cor
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(weighted/avg) of length 300 cor
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 title2)
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 displays image of apparel
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 [16]:

```

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 of given sentance
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 of word i
11    # if m_name == 'weighted', we will multiply each w2v[word] with the idf(word)
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_])
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 [17]:

```
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 corresponds to a doc
9 w2v_title = np.array(w2v_title)
10
```

In [18]:

```
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,-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_indices[i]], data['medium_image_'
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 words i, j
```

burnt umber tiger tshirt zebra stripes xl xxl



[9.4] IDF weighted Word2Vec for product similarity

In [19]:

```
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 corresponds to a doc
8 w2v_title_weight = np.array(w2v_title_weight)
```

In [20]:

```
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 remaining apparels
5     # the metric we used here is cosine, the coside distance is mesured as  $K(X, Y) = \langle X, Y \rangle / (\|X\| * \|Y\|)$ 
6     # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
7     pairwise_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1, -1))
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_indices[i]], data['medium_image_'])
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 words i, j
```



[9.6] Weighted similarity using brand and color.

```
In [21]: 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)).tocsr()
20
```

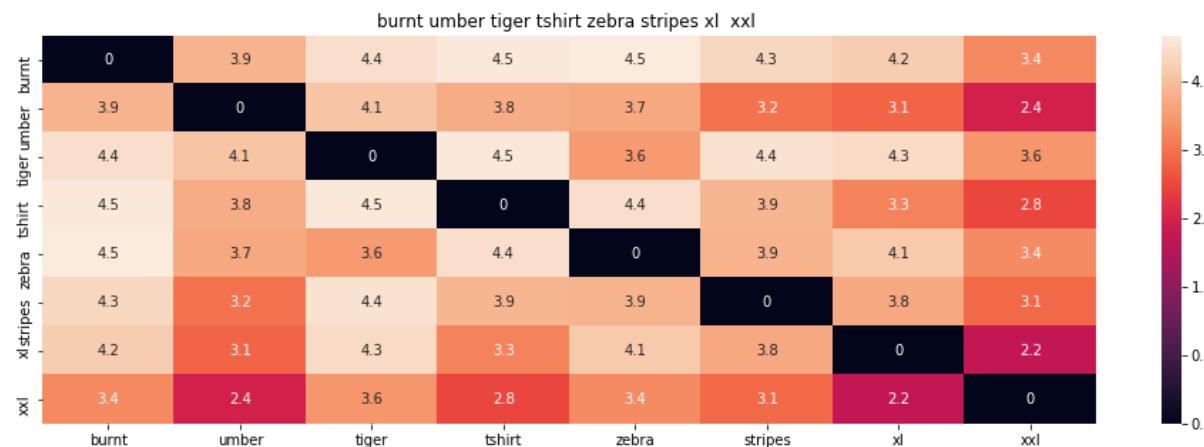
In [22]:

```
1 def heat_map_w2v_brand(sentance1, sentance2, url, doc_id1, doc_id2, df_id1, df_id2, model):
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(weighted/avg) of length 300 cor
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 cor
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 title2)
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], types[doc_id1]], # input apparel
23                  [data['asin'].loc[df_id2], brands[doc_id2], colors[doc_id2], types[doc_id2]]] # recommended a
24
25    colorscale = [[0, '#1d004d'], [.5, '#f2e5ff'], [1, '#f2e5d1']] # to color the headings of each column
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 lins 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 [23]:

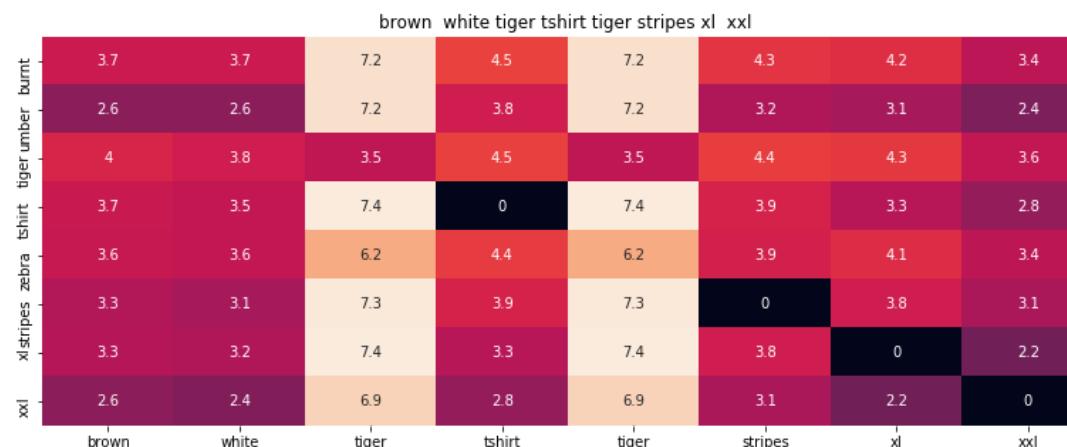
```
1 def idf_w2v_brand(doc_id, w1, w2, num_results):
2     # doc_id: apparel's id in given corpus
3     # w1: weight for w2v features
4     # w2: weight for brand and color features
5
6     # pairwise_dist will store the distance from given input apparel to all remaining apparels
7     # the metric we used here is cosine, the coside distance is mesured as  $K(X, Y) = \langle X, Y \rangle / (\|X\| * \|Y\|)$ 
8     # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
9     idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1, -1))
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[df_indices[i]], data['medium_'
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 words i, j
```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from input : 0.0



ASIN : B00JXQCWT0

Brand : Si Row

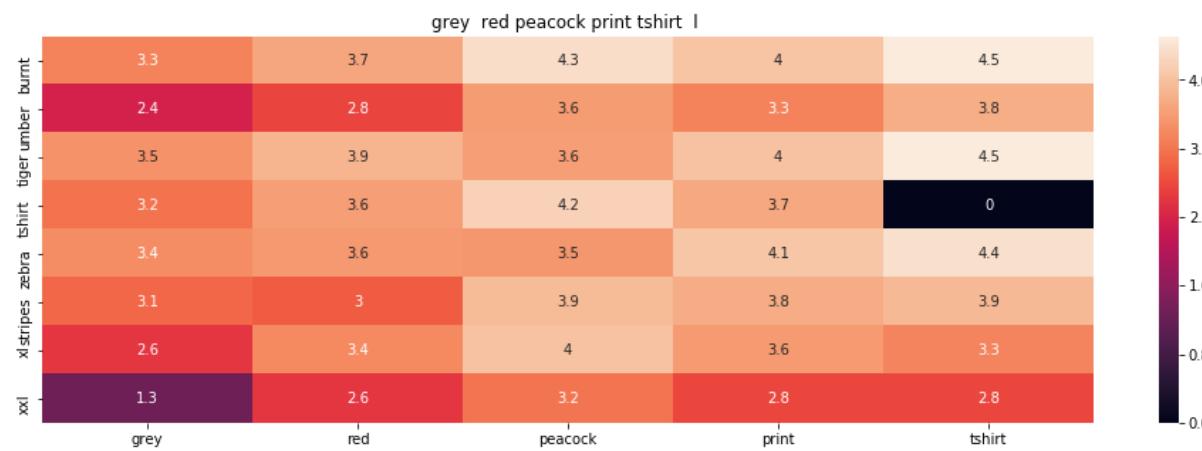
euclidean distance from input : 0.6528799057006835



ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from input : 1.0017030956167843



ASIN : B00JXQCFRS

Brand : Si Row

euclidean distance from input : 1.2372833253760007

=====

=====

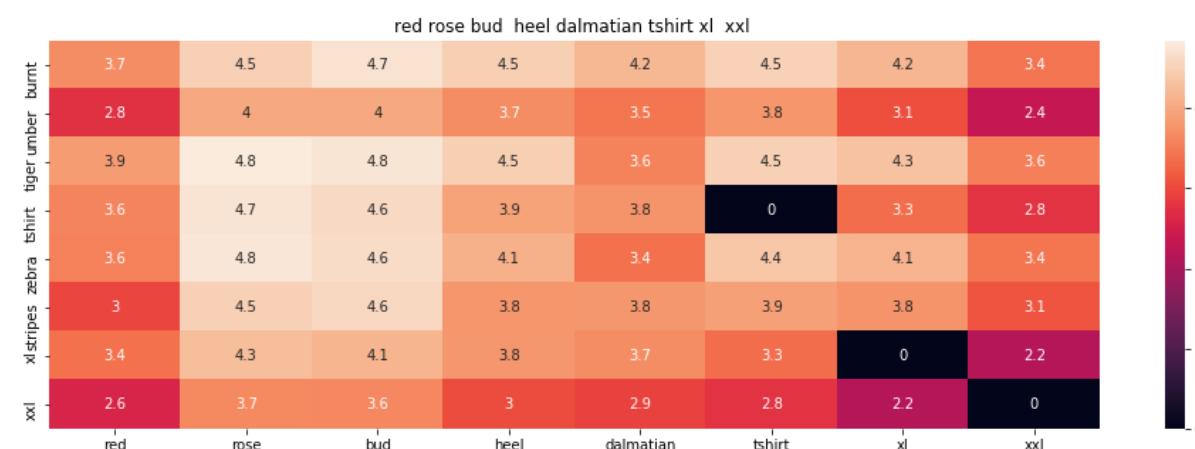


ASIN : B00JXQBBMI

Brand : Si Row

euclidean distance from input : 1.2686225892920162

三

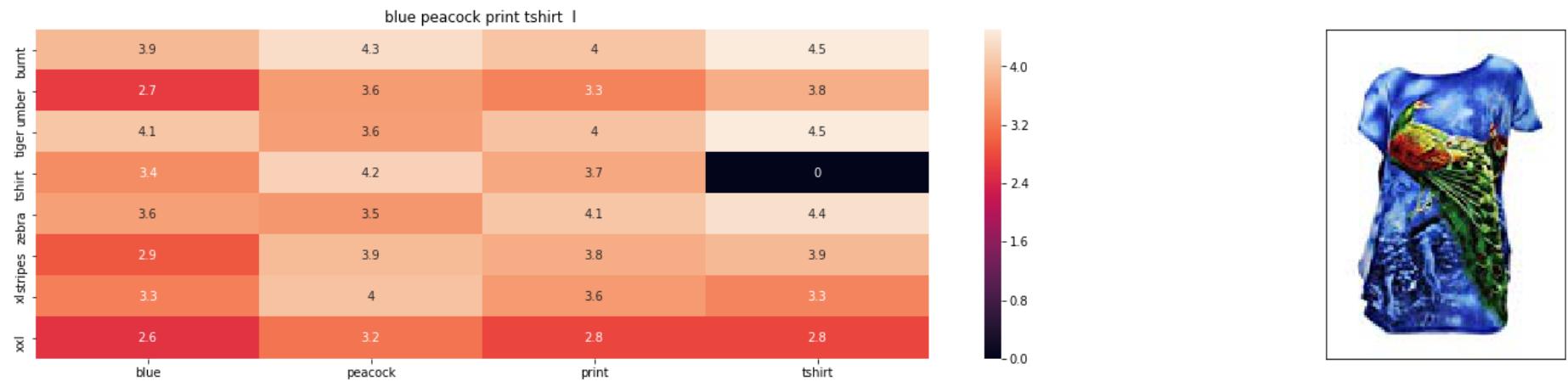


ASIN : B00JXOABB0

Brand : Si Row

euclidean distance from input : 1.2733484746832517

Sacredan distar



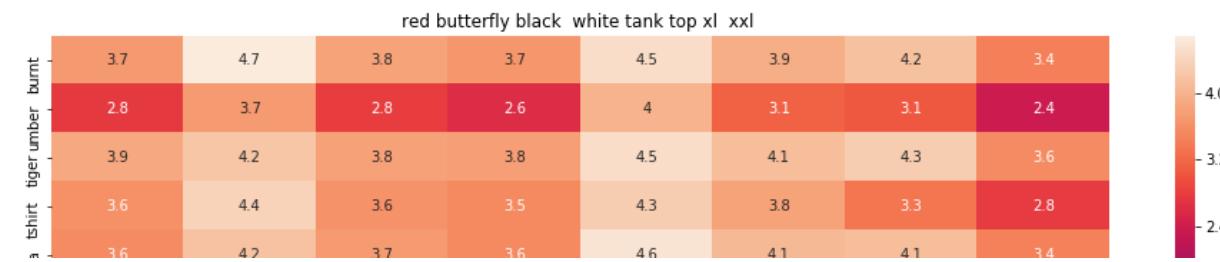
ASIN : B00JXQC8L6

Brand : Si Row

euclidean distance from input : 1.2818687440771726

=====

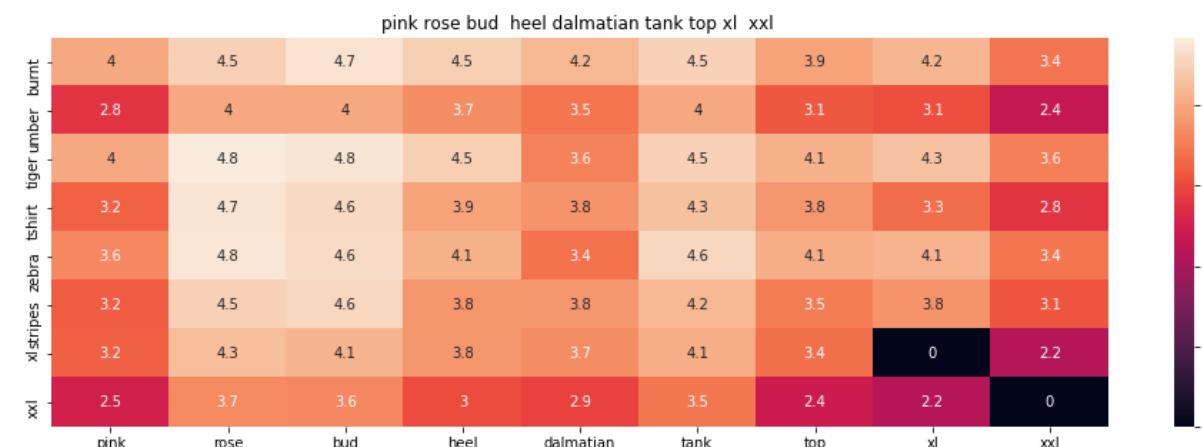
=====



ASIN : B00JV63CW2

Brand : Si Row

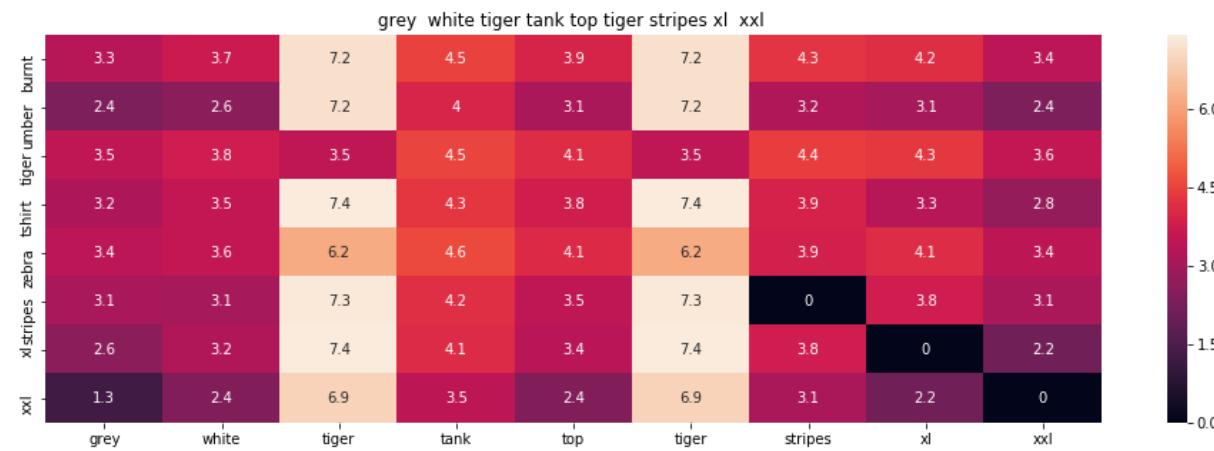
euclidean distance from input : 1.2947488786596921



ASIN : B00JXQAX2C

Brand : Si Row

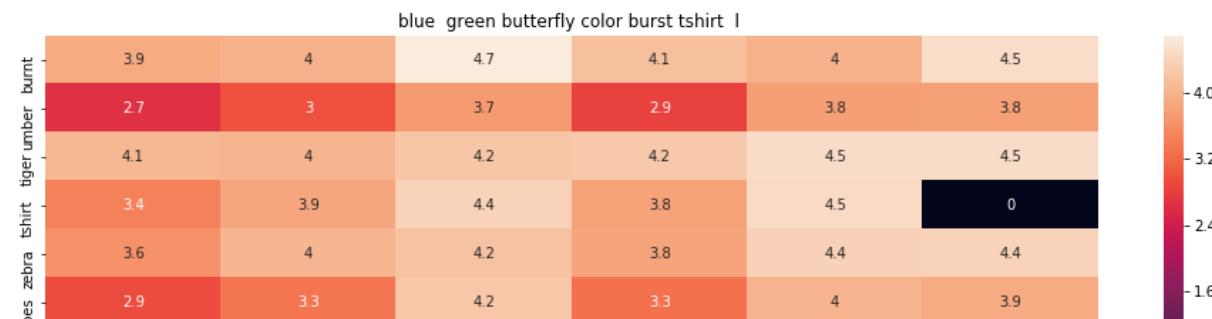
euclidean distance from input : 1.3236358167547848



ASIN : B00JXQAFZ2

Brand : Si Row

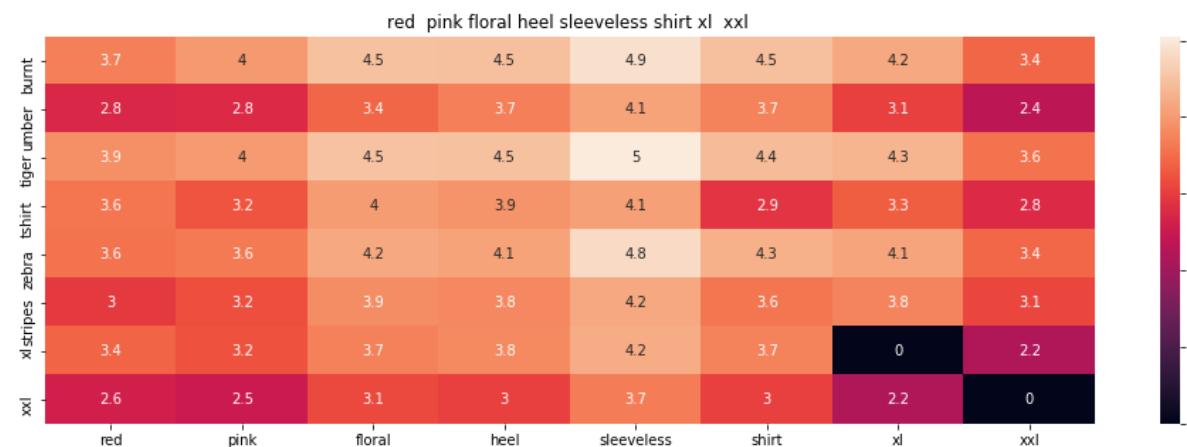
euclidean distance from input : 1.3293567659277585



ASIN : B00JXQC0C8

Brand : Si Row

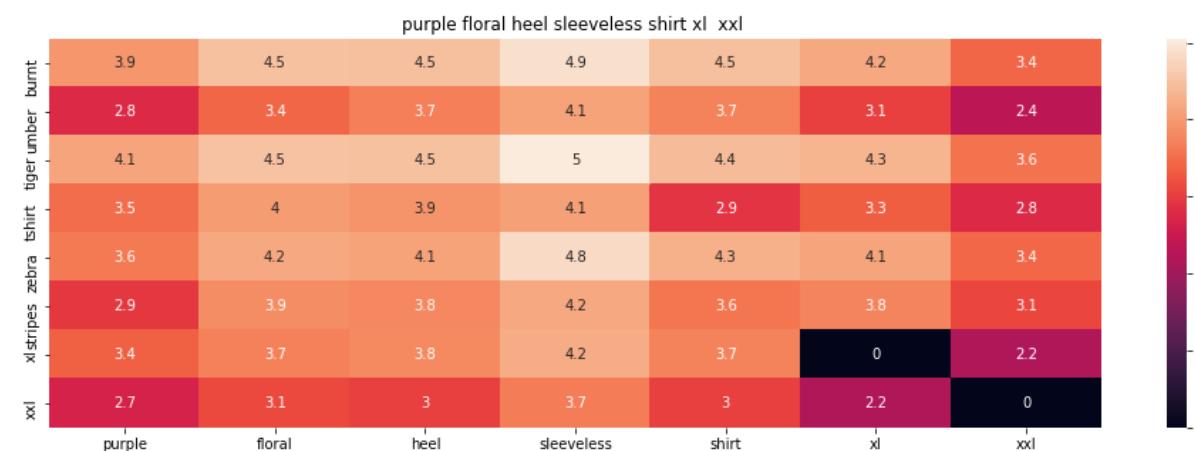
euclidean distance from input : 1.34993848818728



ASIN : B00JV63QQE

Brand : Si Row

euclidean distance from input : 1.3649898530859617

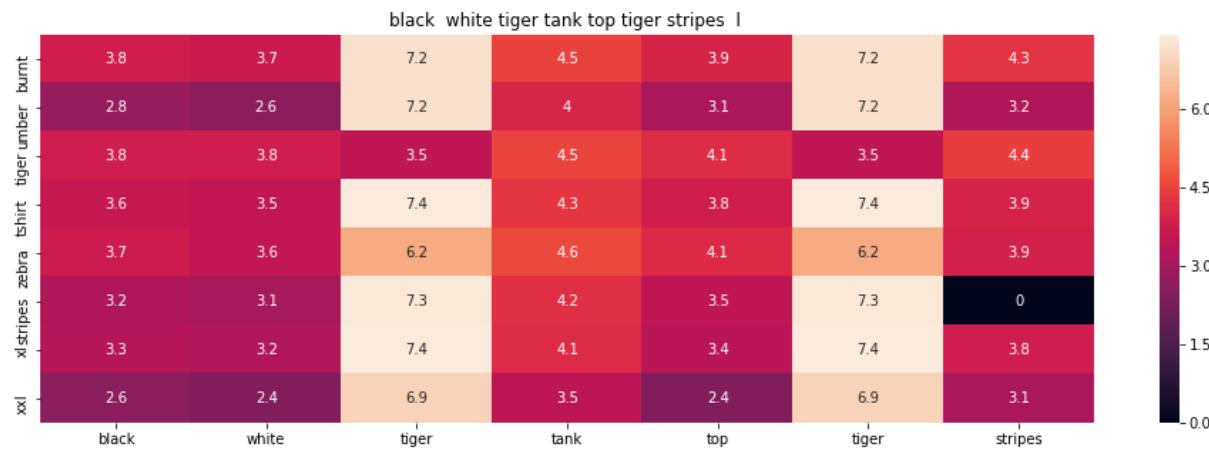


ASIN : B00JV63VC8

Brand : Si Row

euclidean distance from input : 1.4083902837652829

=====



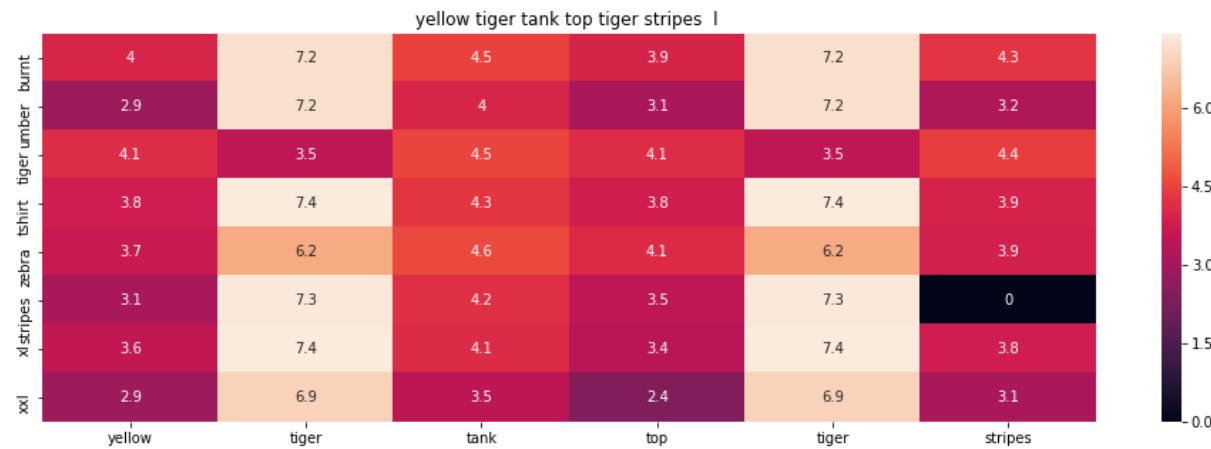
ASIN : B00JXQAO94

Brand : Si Row

euclidean distance from input : 1.5019501687903074

=====

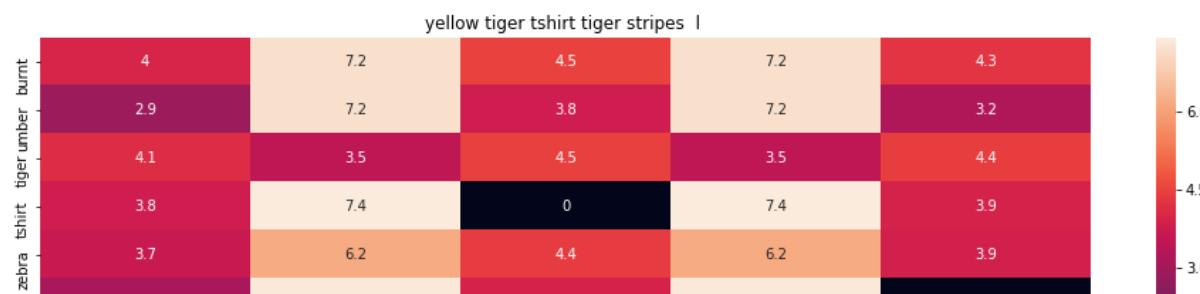
=====



ASIN : B00JXQAUWA

Brand : Si Row

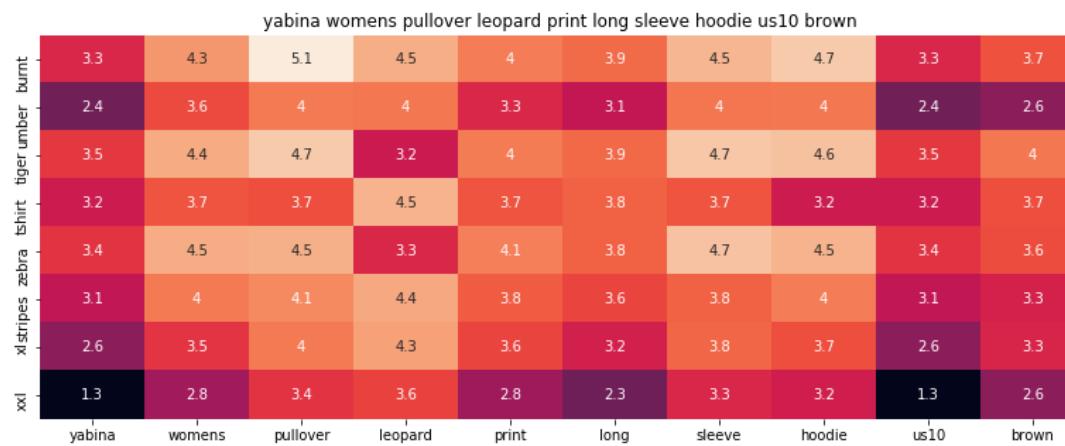
euclidean distance from input : 1.5693789483923581



ASIN : B00JXQCUIC

Brand : Si Row

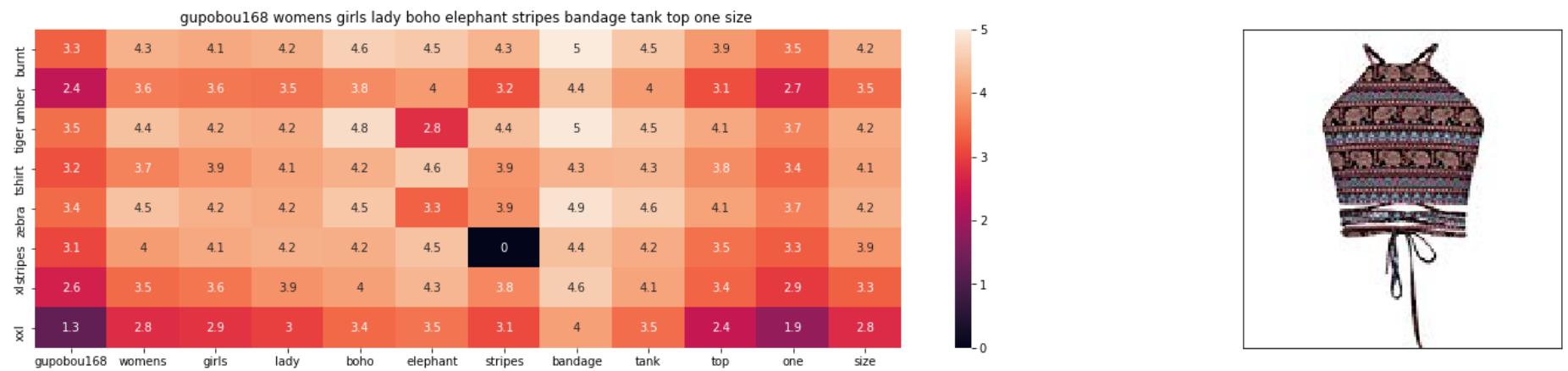
euclidean distance from input : 1.633365059079614



ASIN : B01KJUM6JI

Brand : YABINA

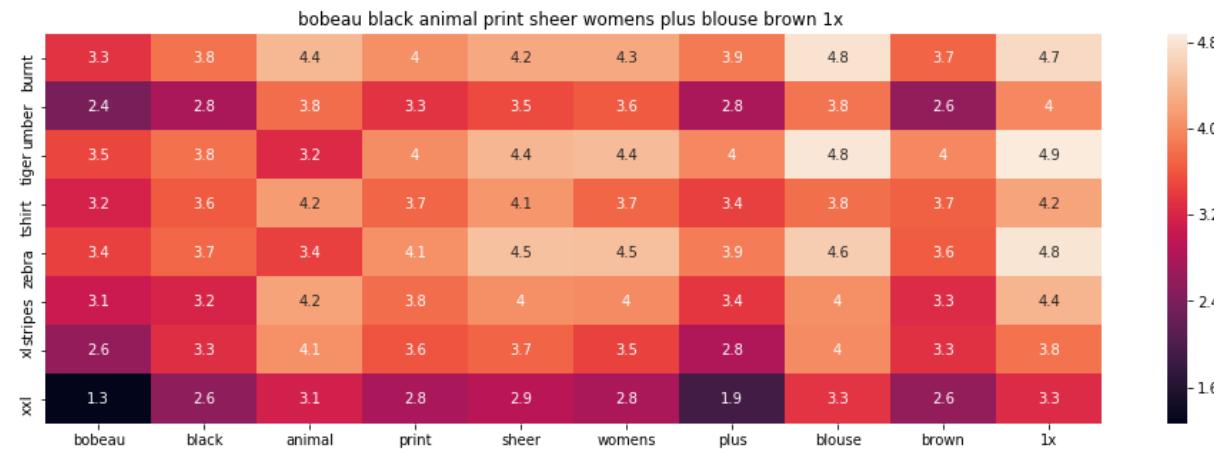
euclidean distance from input : 1.7323767736314868



ASIN : B01ER18406

Brand : GuPoBoU168

euclidean distance from input : 1.7352904394029711

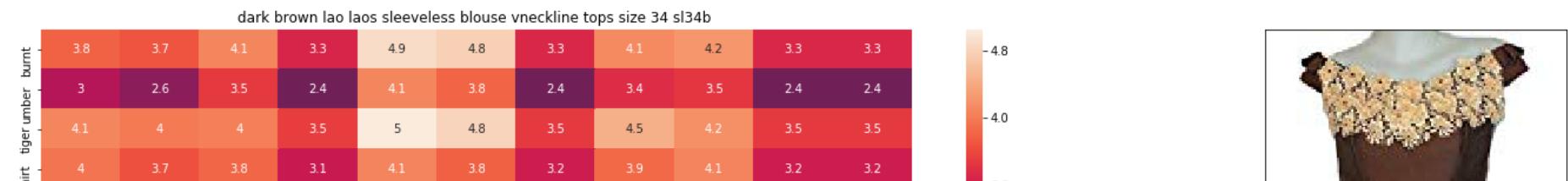


ASIN : B074MH886R

Brand : Bobeau

euclidean distance from input : 1.7428852155565355

=====

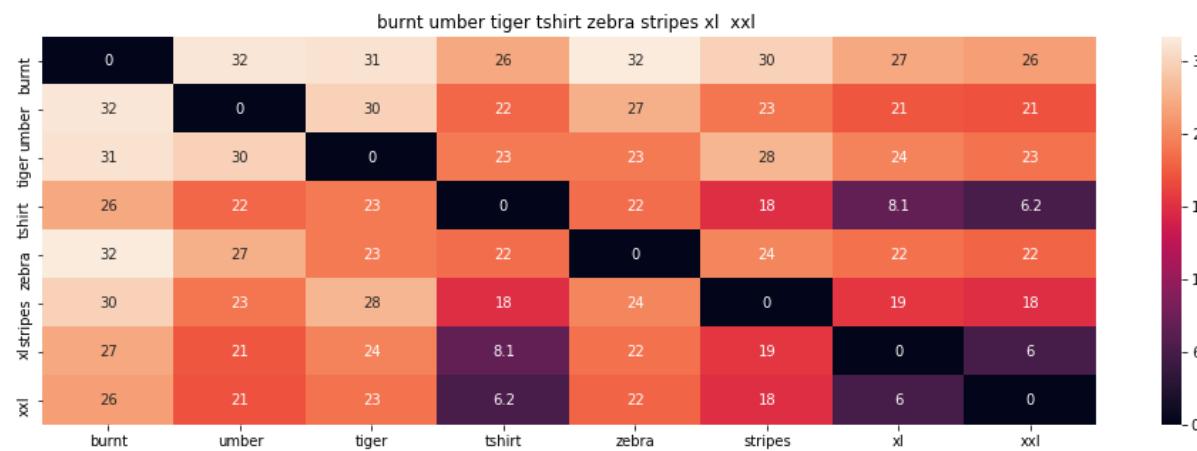


In [0]:

```

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

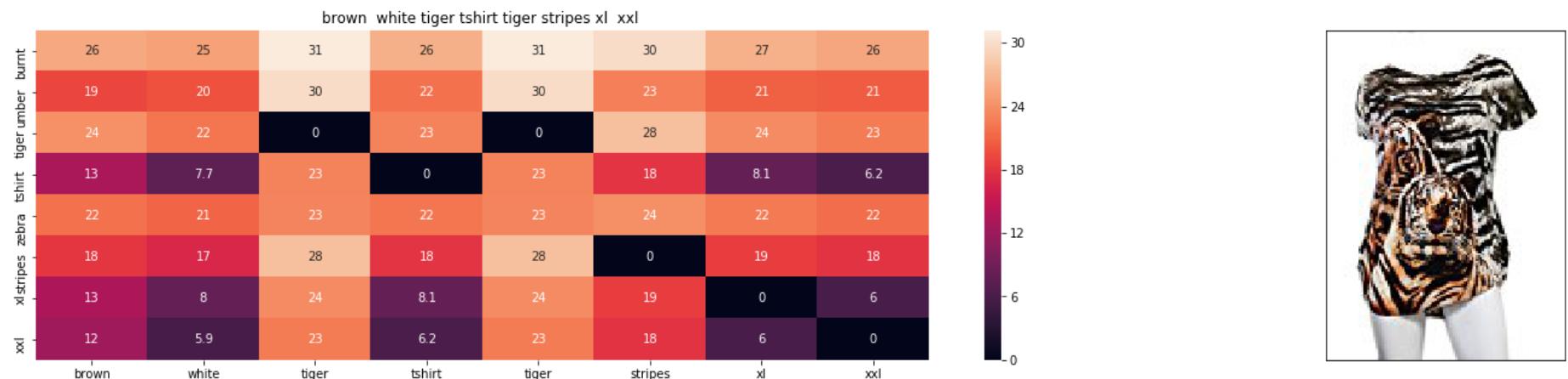
```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from input : 0.000355113636364



ASIN : B00JXQCWTO

Brand : Si Row

euclidean distance from input : 0.433722027865

=====

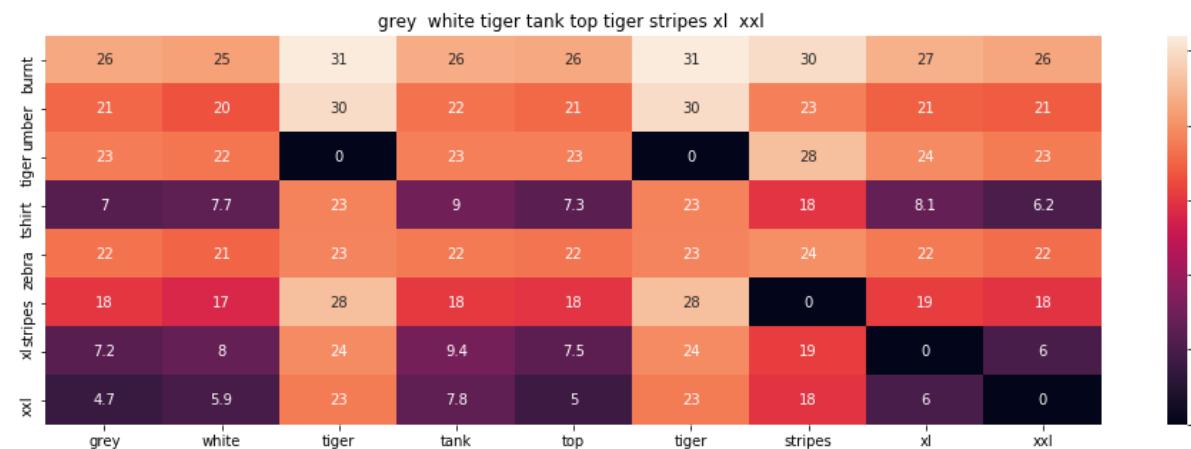
=====



ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from input : 1.65509310669



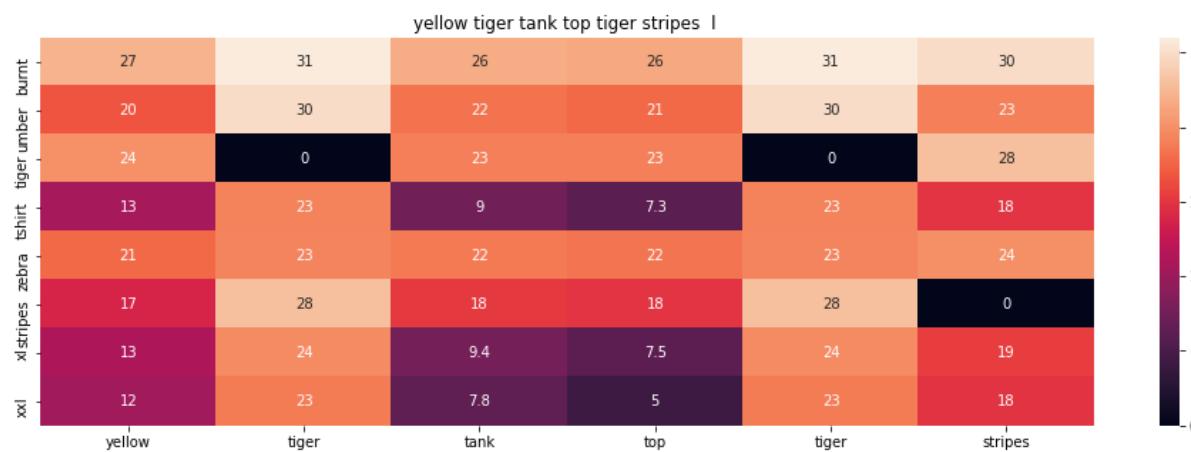
ASIN : B00JXQAFZ2

Brand : Si Row

euclidean distance from input : 1.77293604103

=====

=====



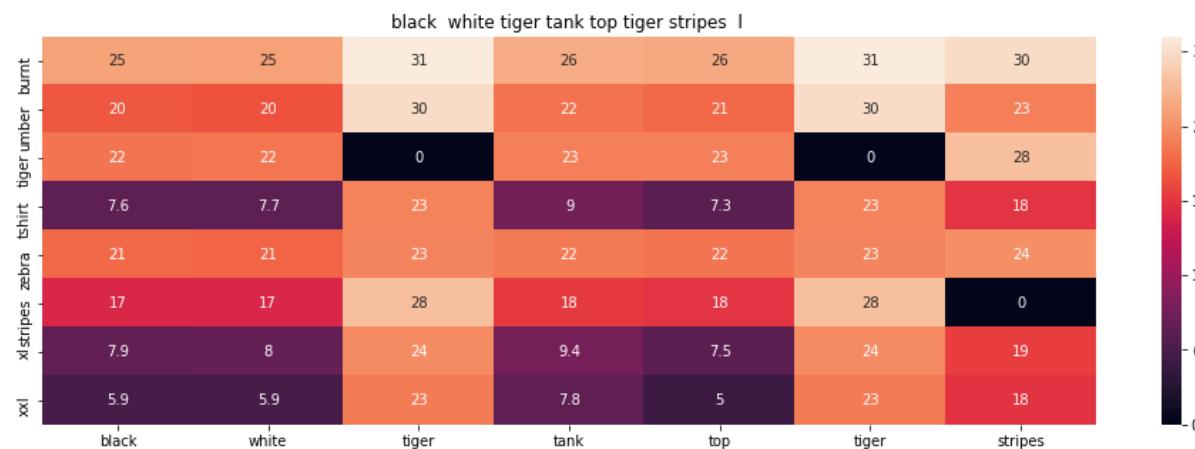
ASIN : B00JXQAUWA

Brand : Si Row

euclidean distance from input : 1.80287808538

=====

=====



ASIN : B00JXQAO94

Brand : Si Row

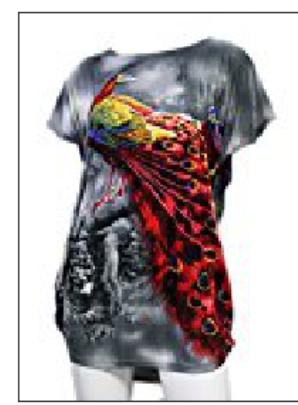
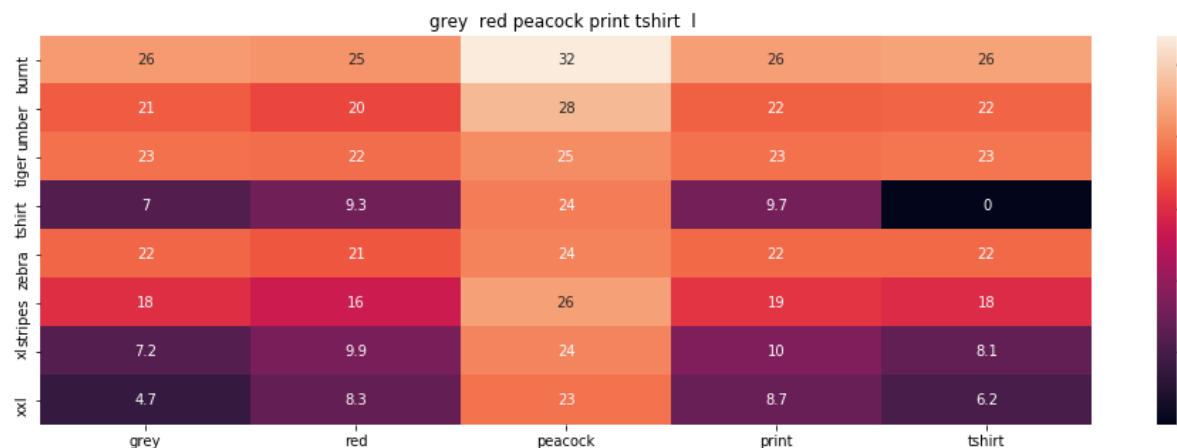
euclidean distance from input : 1.80319609241



ASIN : B00JXQCUIC

Brand : Si Row

euclidean distance from input : 1.82141619628



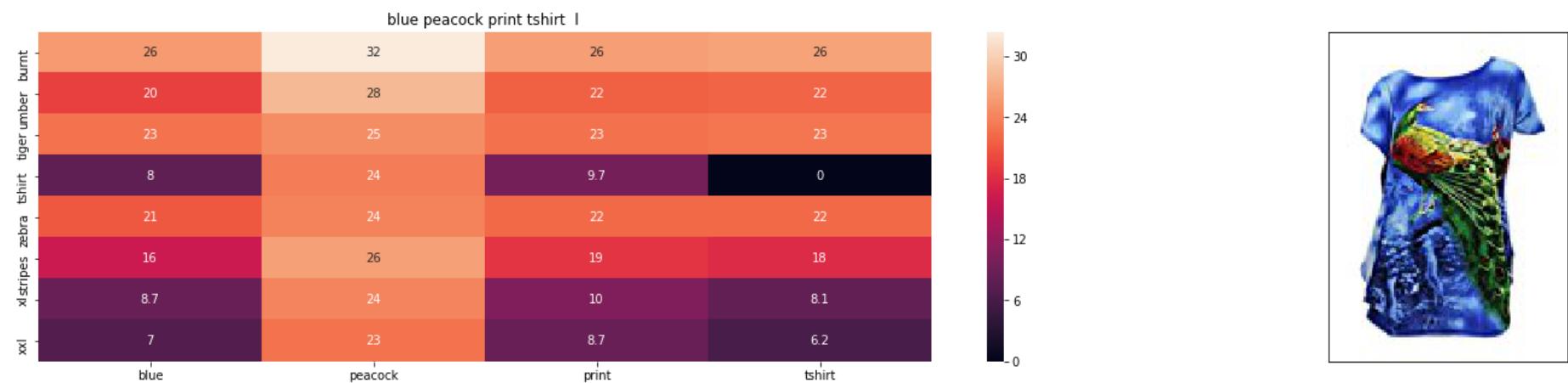
ASIN : B00JXQCFRS

Brand : Si Row

euclidean distance from input : 1.90777685025

=====

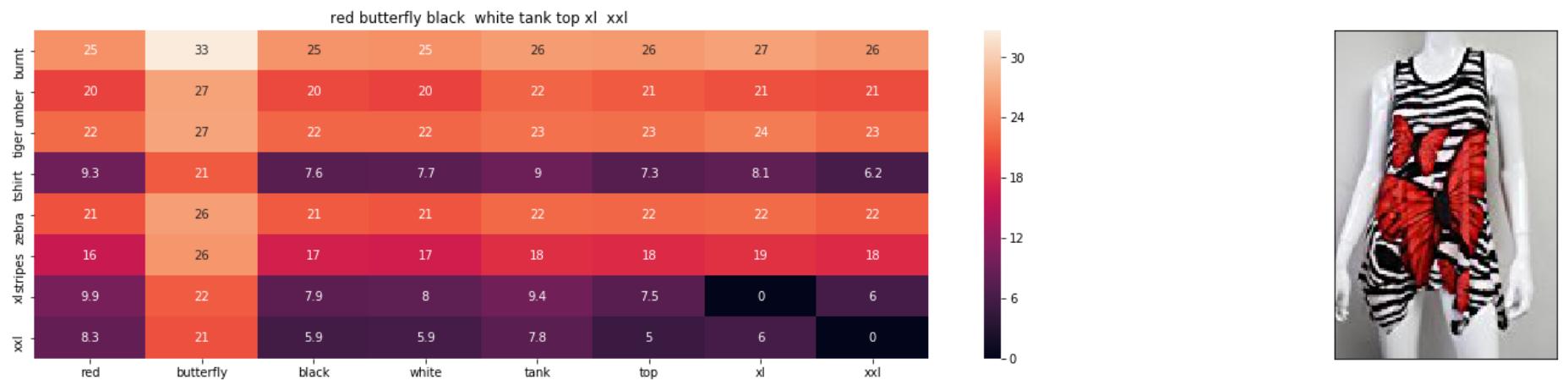
=====



ASIN : B00JXQC8L6

Brand : Si Row

euclidean distance from input : 1.92142937433



ASIN : B00JV63CW2

Brand : Si Row

euclidean distance from input : 1.93646323497

=====

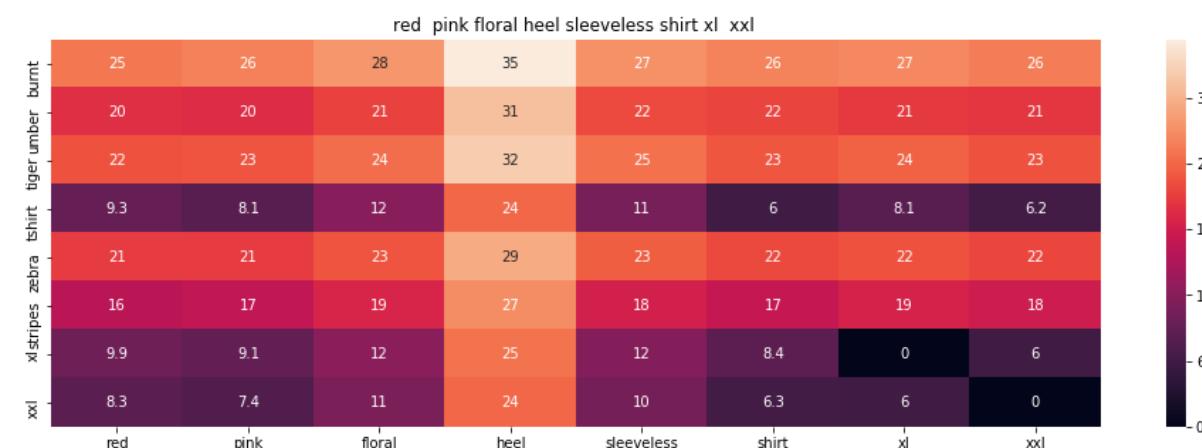
=====



ASIN : B00JXQBBMI

Brand : Si Row

euclidean distance from input : 1.95670381059



ASIN : B00JV63QQE

Brand : Si Row

euclidean distance from input : 1.9806066343

=====

=====



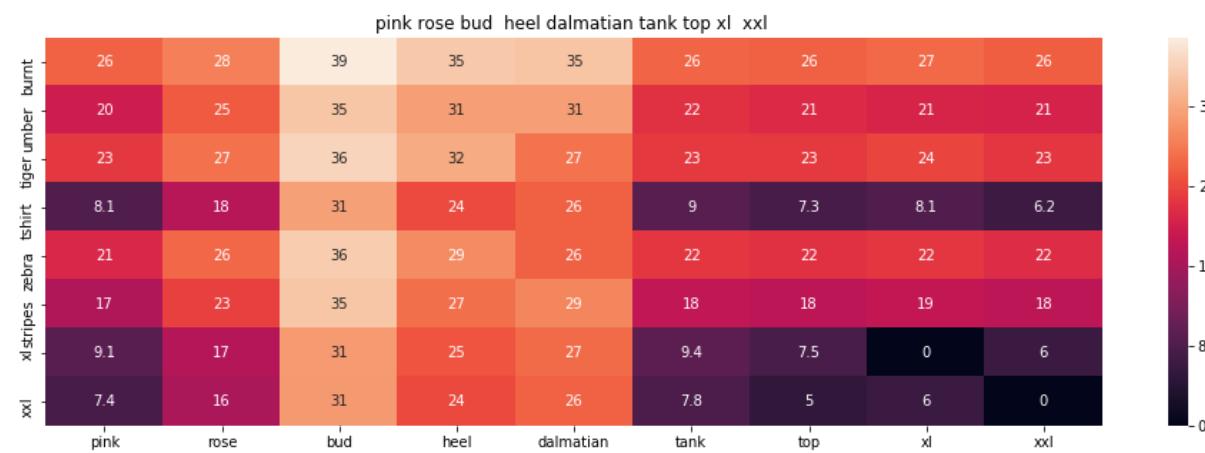
ASIN : B00JV63VC8

Brand : Si Row

euclidean distance from input : 2.01218559992

=====

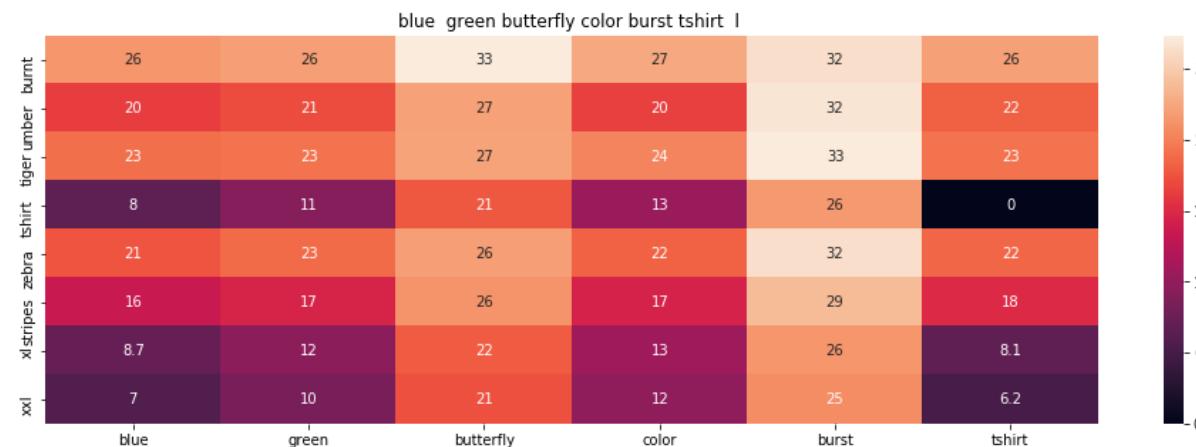
=====



ASIN : B00JXQAX2C

Brand : Si Row

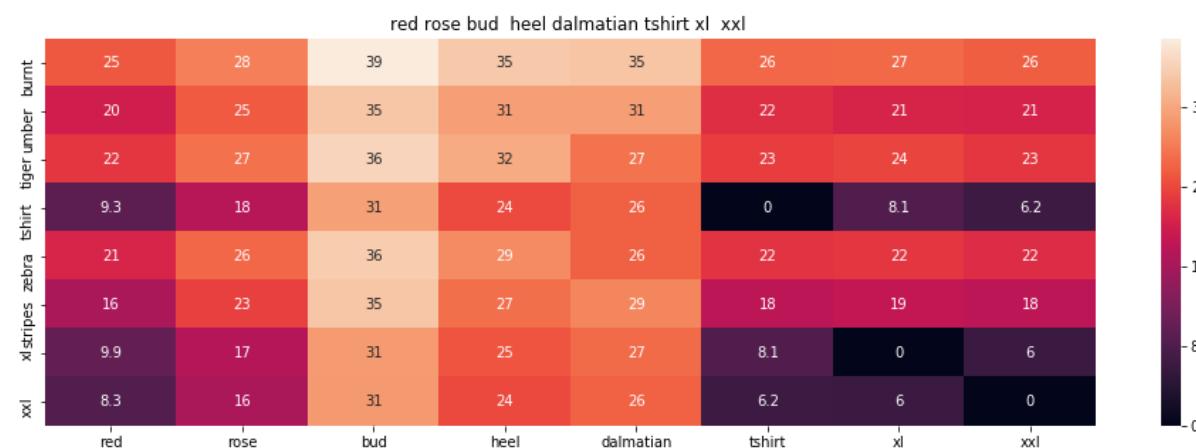
euclidean distance from input : 2.01335178755



ASIN : B00JXQC0C8

Brand : Si Row

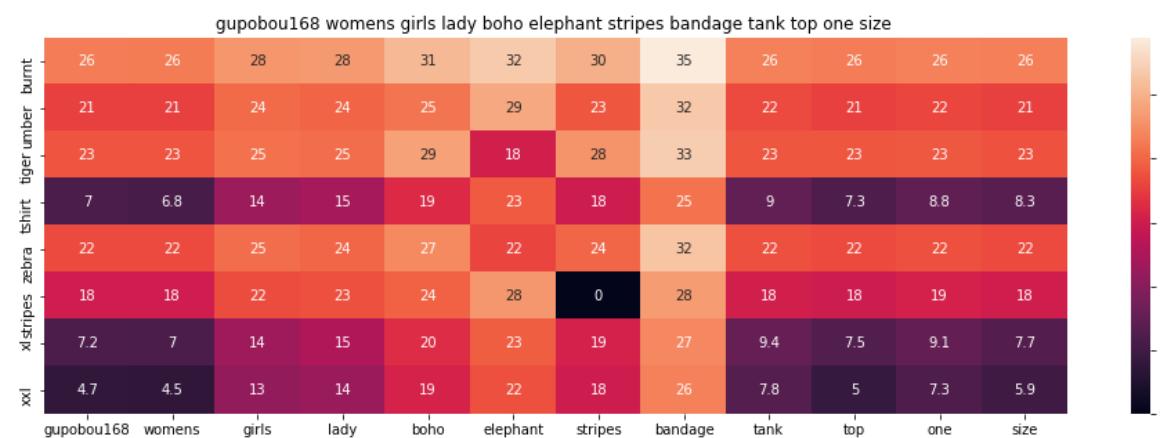
euclidean distance from input : 2.01388334827



ASIN : B00JXQABBO

Brand : Si Row

euclidean distance from input : 2.0367257555

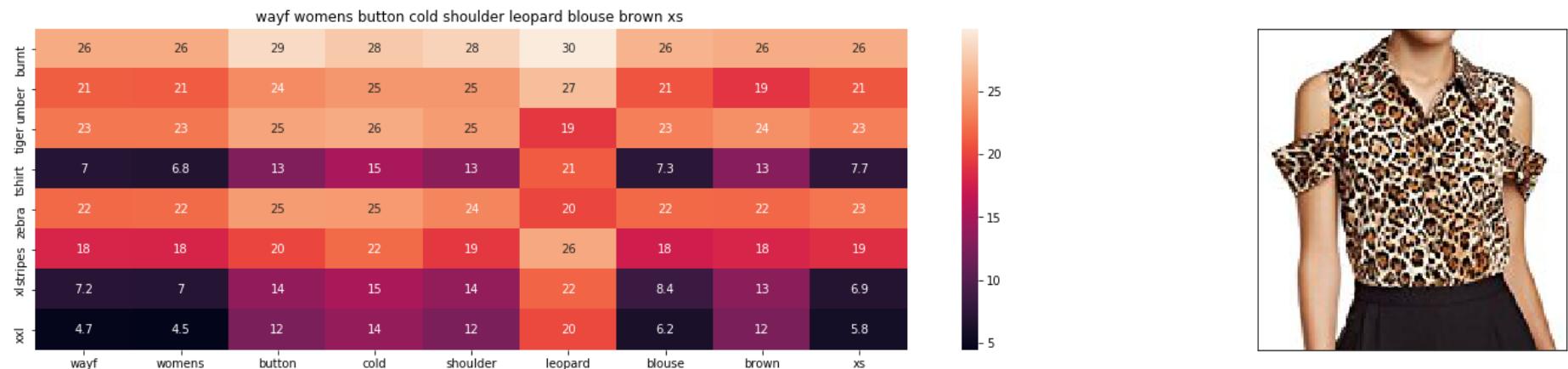
=====
=====


ASIN : B01ER18406

Brand : GuPoBoU168

euclidean distance from input : 2.65620416778

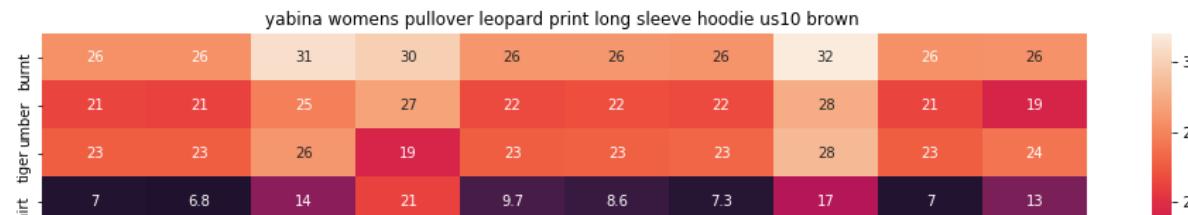
=====
=====



ASIN : B01LZ7BQ4H

Brand : WAYF

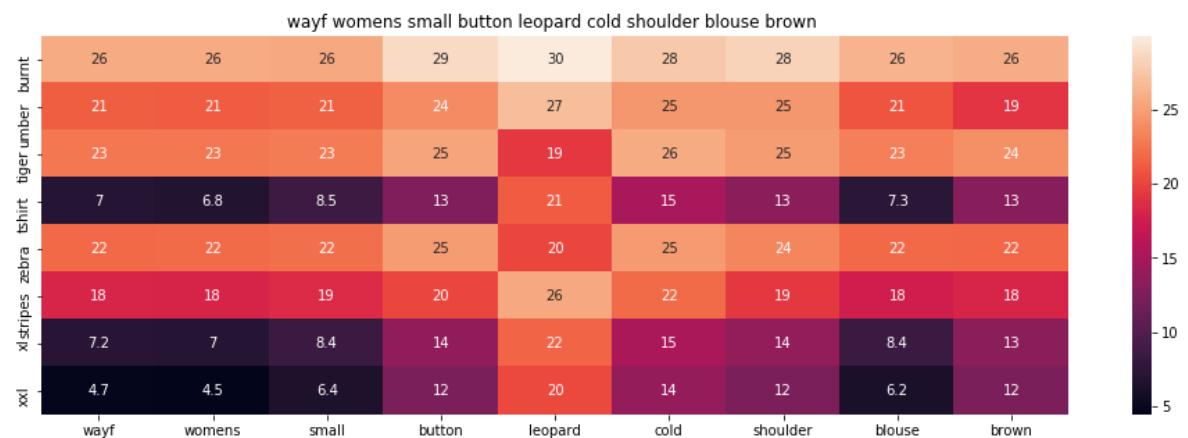
euclidean distance from input : 2.6849067823



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-classification-models-using-very-little-data
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 output
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_samples // batch_size)  
49     bottleneck_features_train = bottleneck_features_train.reshape((16042,25088))  
50  
51     np.save(open('16k_data_cnn_features.npy', 'wb'), bottleneck_features_train)  
52     np.save(open('16k_data_cnn_feature_asins.npy', 'wb'), np.array(asins))  
53  
54  
55 save_bottlebeck_features()  
56  
57 ...
```

[10.3] Visual features based product similarity.

In [56]:

```
1 #load the features and corresponding ASINS info.
2 bottleneck_features_train = np.load('pickles/16k_data_cnn_features.npy')
3 asins = np.load('pickles/16k_data_cnn_feature_asins.npy')
4 asins = list(asins)
5
6 # load the original 16K dataset
7 data = pd.read_pickle('pickles/16k_apperial_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_features_train[doc_id].reshape(1, -1))
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, 10)
31
```



Product Title: burnt umber tiger tshirt zebra stripes xl xxl

Euclidean Distance from input image: 4.2649613e-06

Amazon Url: www.amazon.com/dp/B00JXQB5FQ



Product Title: pink tiger tshirt zebra stripes xl xxl

Euclidean Distance from input image: 30.05017

Amazon Url: www.amazon.com/dp/B00JXQASS6



Product Title: yellow tiger tshirt tiger stripes l

Euclidean Distance from input image: 41.261116

Amazon Url: www.amazon.com/dp/B00JXQCUIC



Product Title: brown white tiger tshirt tiger stripes xl xxl

Euclidean Distance from input image: 44.000156

Amazon Url: www.amzon.com/dp/B00JXQCWTO



Product Title: kawaii pastel tops tees pink flower design

Euclidean Distance from input image: 47.38248

Amazon Url: www.amzon.com/dp/B071FCWD97



Product Title: womens thin style tops tees pastel watermelon print

Euclidean Distance from input image: 47.71842

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



Product Title: kawaii pastel tops tees baby blue flower design

Euclidean Distance from input image: 47.90206

Amazon Url: www.amazon.com/dp/B071SRCY9W



Product Title: edv cheetah run purple multi xl

Euclidean Distance from input image: 48.046482

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



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

Euclidean Distance from input image: 48.101837

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



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

Euclidean Distance from input image: 48.118866

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

In []:

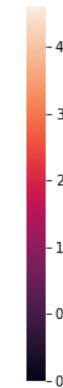
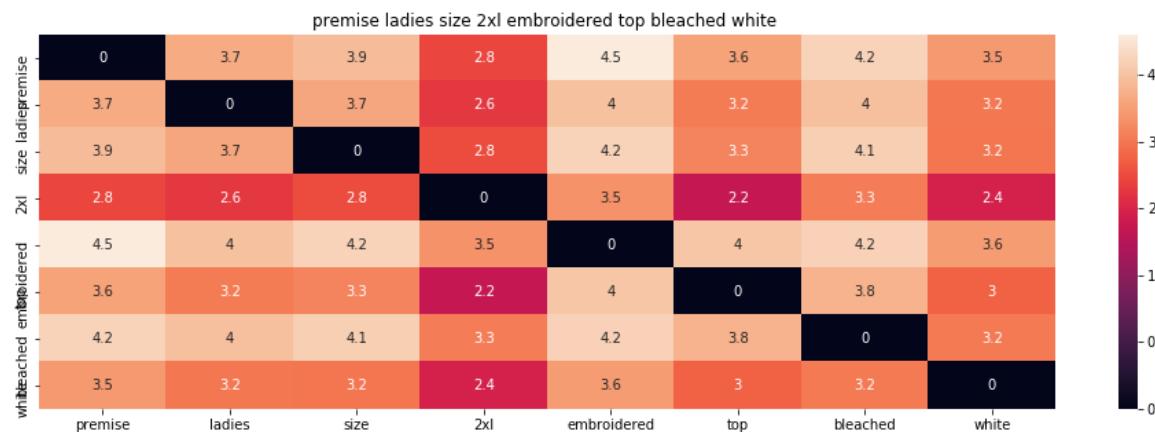
1

Assignment :

Combine idfw2v, bow, tfidf, image feature and give weights to get recommendation.

In [59]:

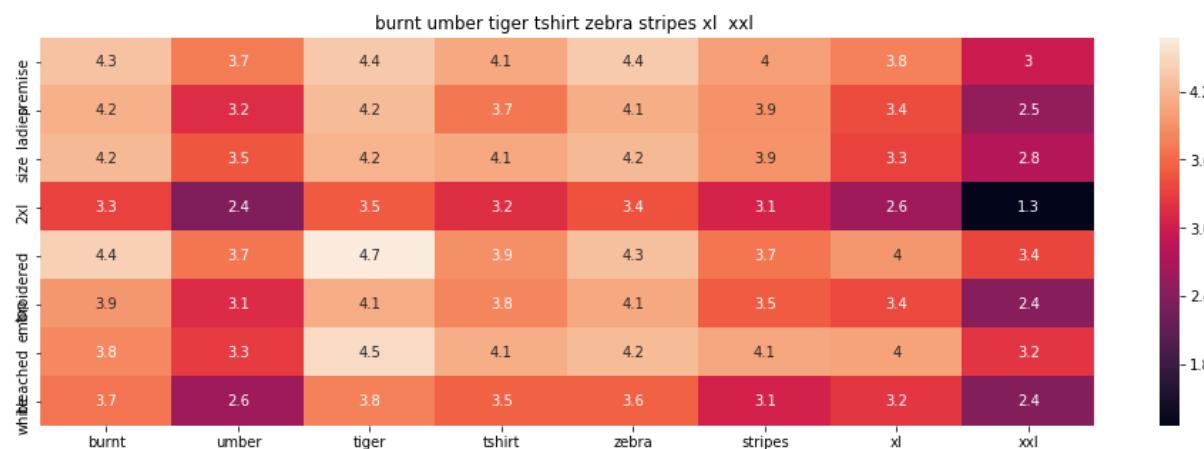
```
1 from PIL import Image
2 def combined(doc_id, w1, w2, w3, num_results):
3     # doc_id: apparel's id in given corpus
4     # w1: weight for w2v features
5     # w2: weight for brand and color features
6
7     # pairwise_dist will store the distance from given input apparel to all remaining apparels
8     # the metric we used here is cosine, the coside distance is mesured as  $K(X, Y) = \langle X, Y \rangle / (\|X\| * \|Y\|)$ 
9     # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
10    doc_id_asin = asins.index(df_asins[doc_id])
11    image_dist = pairwise_distances(bottleneck_features_train, bottleneck_features_train[doc_id_asin].reshape(1, -1))
12
13
14    idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1, -1))
15    ex_feat_dist = pairwise_distances(extra_features, extra_features[doc_id])
16
17    pairwise_dist = (w1 * idf_w2v_dist + w2 * ex_feat_dist + w3 * image_dist)/float(w1 + w2 + w3)
18
19    # np.argsort will return indices of 9 smallest distances
20    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
21    #pdists will store the 9 smallest distances
22    pdists = np.sort(pairwise_dist.flatten())[0:num_results]
23
24    #data frame indices of the 9 smallest distace's
25    df_indices = list(data.index[indices])
26
27
28    for i in range(0, len(indices)):
29        heat_map_w2v_brand(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data['medium_')
30        print('ASIN :', data['asin'].loc[df_indices[i]])
31        print('Brand :', data['brand'].loc[df_indices[i]])
32        print('euclidean distance from input :', pdists[i])
33        print('='*125)
34
35 combined(12566, 10, 12, 1, 10)
36 # in the give heat map, each cell contains the euclidean distance between words i, j
```



ASIN : B01M0IDUCV

Brand : Premise

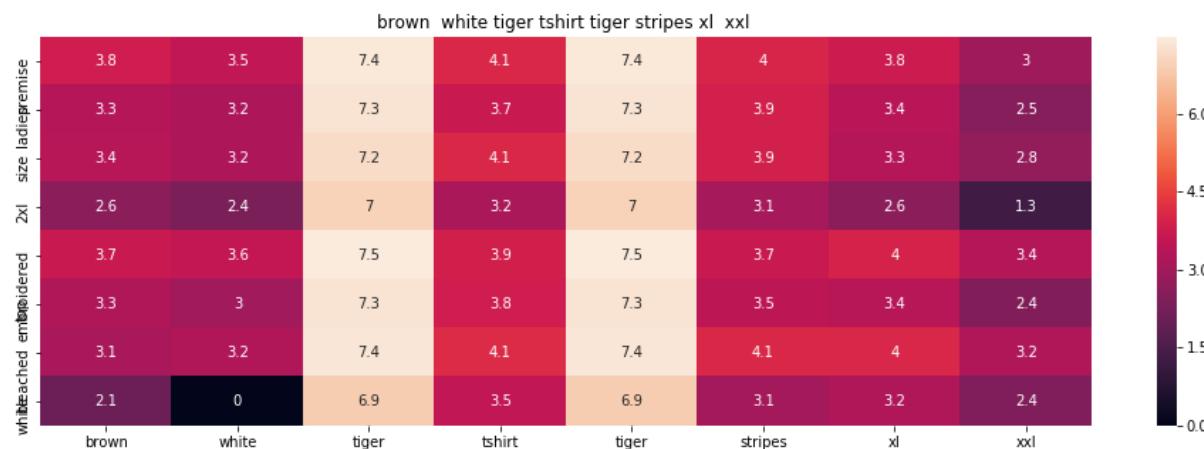
euclidean distance from input : 2.0446252441916406



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from input : 2.5679238360861074



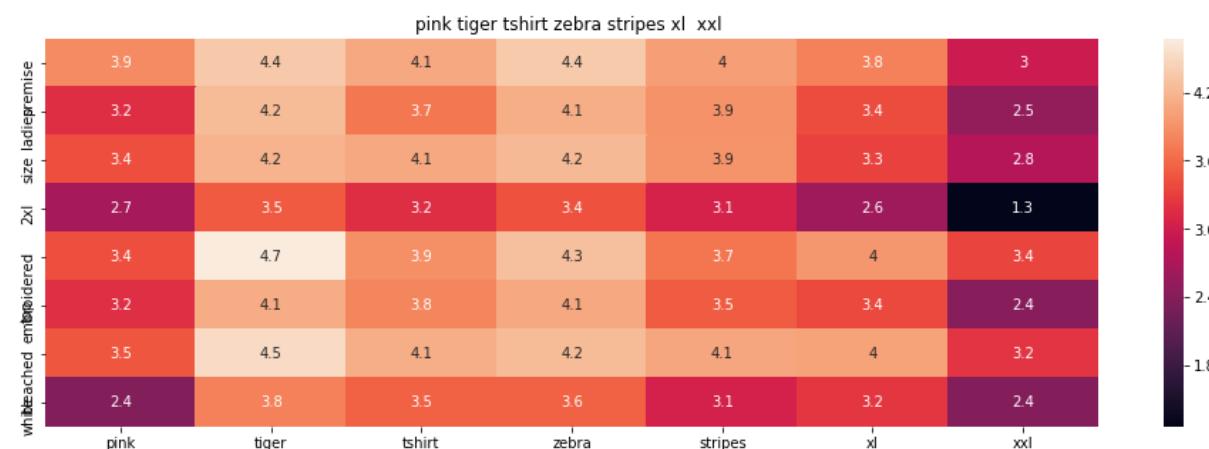
ASIN : B00JXQCWT0

Brand : Si Row

euclidean distance from input : 3.1388397216796875

=====

=====



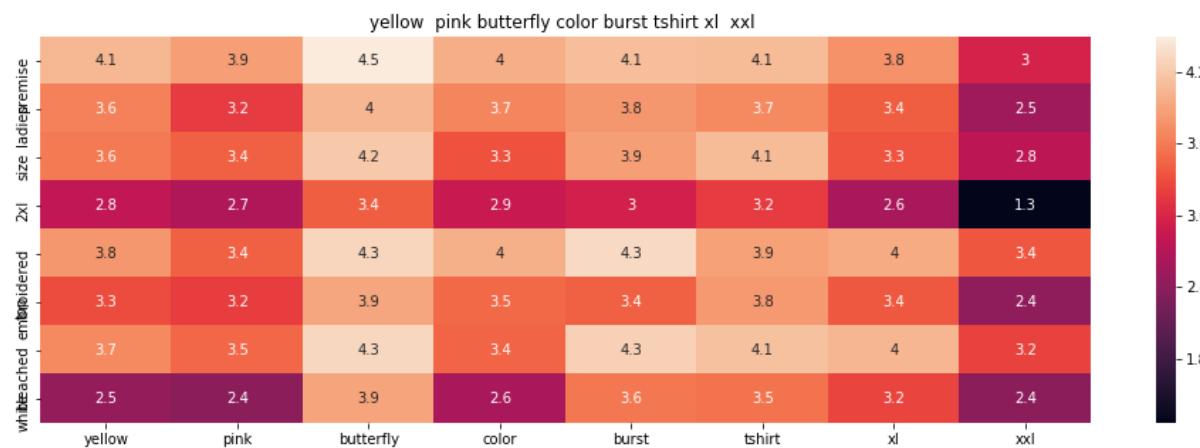
ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from input : 3.3450818810067866

=====

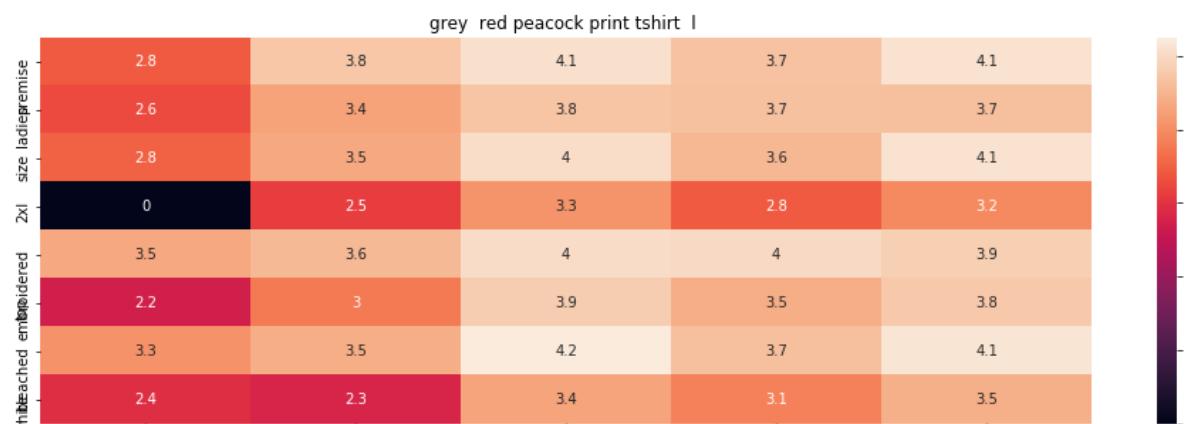
=====



ASIN : B00JXQBBMI

Brand : Si Row

euclidean distance from input : 3.361690554600059



ASIN : B00JXQCFRS

Brand : Si Row

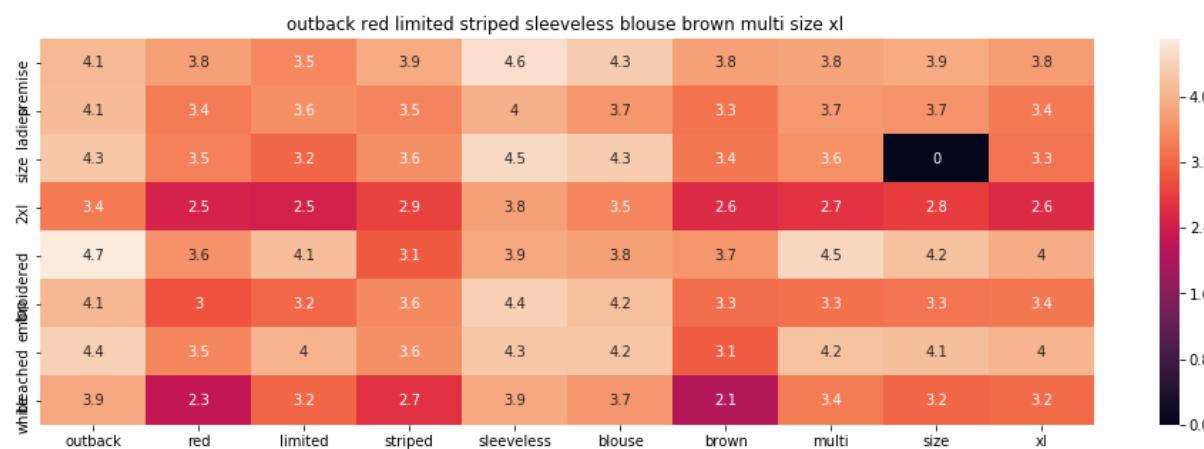
euclidean distance from input : 3.607642497375619



ASIN : B01N4NQ7LX

Brand : CeCe by Cynthia Steffe

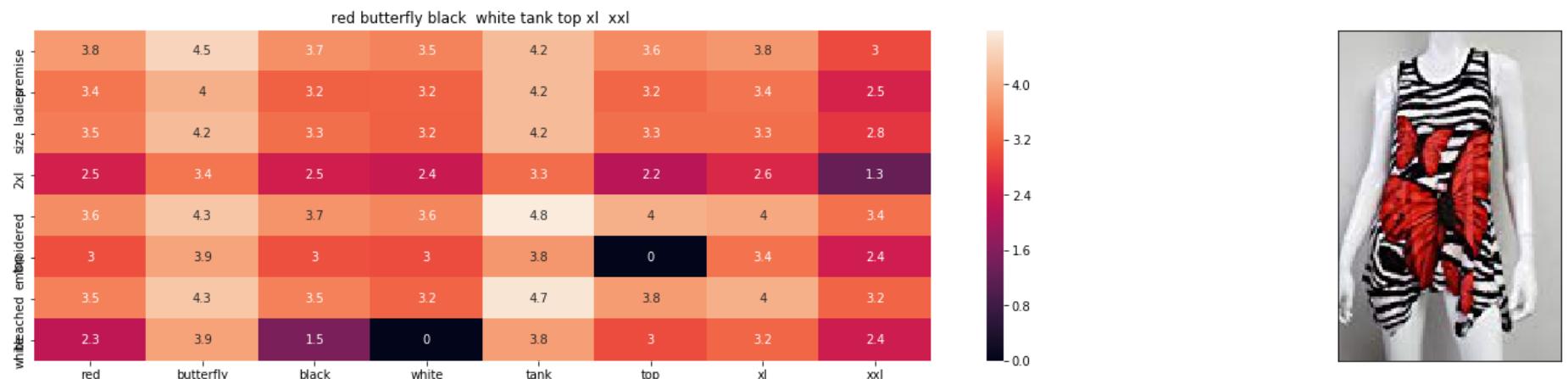
euclidean distance from input : 3.629845199338683



ASIN : B01IU645VU

Brand : Outback Red

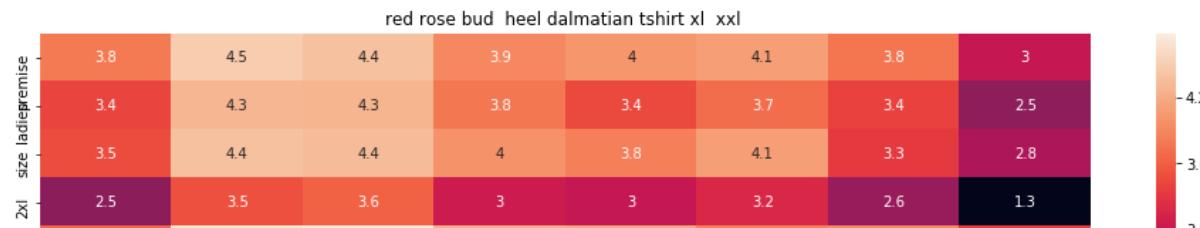
euclidean distance from input : 3.739560704480477



ASIN : B00JV63CW2

Brand : Si Row

euclidean distance from input : 3.7469922274940806



ASIN : B00JXQABBO

Brand : Si Row

euclidean distance from input : 3.74840121910042

=====

=====

In []:

1