

# AmzonFashionDiscovery

## Important Setsps to Process Amzon Fashion Discovery

1. Importing Python libraries
2. Data Acquisition
3. Data Cleaning
4. Data PreProcessing
5. Linear Algebra
6. Text Based Product Recomendataion
7. Image Based
8. A/B Testing

## 1. Importing Python Packages

In [1]:

```
#import all the necessary packages.

from PIL import Image
import requests
from io import BytesIO
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import warnings
from bs4 import BeautifulSoup
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
import nltk
import math
import time
import re
import os
import seaborn as sns
from collections import Counter
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.metrics import pairwise_distances
from matplotlib import gridspec
from scipy.sparse import hstack
import plotly
import plotly.figure_factory as ff
from plotly.graph_objs import Scatter, Layout

plotly.offline.init_notebook_mode(connected=True)
warnings.filterwarnings("ignore")
```

## 2. Data Acqusation

### Reading Amzon Data which is in JSON Format

- we have give a json file which consists of all information about the products
- loading the data using pandas' read\_json file.

In [3]:

```
data = pd.read_json('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/tops_fashion.json')
```

In [7]:

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

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

**Each product/item has 19 features in the raw dataset.**

In [8]:

```
data.columns
```

Out[8]:

```
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 [9]:

```
data=data[['asin','brand','color','product_type_name','medium_image_url','title','formatted_price']]
```

In [10]:

```
data.head()
```

Out[10]:

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

In [11]:

```
print ('Number of data points : ', data.shape[0], \
```

```
      '\n Number of features:', data.shape[1])
```

```
data.head()
```

```
Number of data points : 183138
```

```
Number of data points : 183138
Number of features: 7
```

Out[11]:

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

### 3. Data Understanding and Cleaning

- Very Important Stage in ML
- Often Overlooked Stage
- In this Stage we need to understand about data ie about each feature in data
- Remove unnecessary data

#### [5.1] Missing data for various features.

Basic stats for the feature: product\_type\_name

In [20]:

```
# We have total 72 unique type of product_type_names
print(data['product_type_name'].describe())

# 91.62% (167794/183138) of the products are shirts,
```

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

#### Names of different product types

In [21]:

```
# names of different product types
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']
```

## find the 10 most frequent product\_type\_names.

In [24]:

```
product_type_count = Counter(list(data['product_type_name']))  
product_type_count.most_common(10)
```

Out [24]:

```
[('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)]
```

In [33]:

```
product_type_count.elements()
```

Out [33]:

```
<itertools.chain at 0x7f58e71586a0>
```

## Basic stats for the feature: brand

In [34]:

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

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

In [35]:

```
brand_count = Counter(list(data['brand']))  
brand_count.most_common(10)
```

Out [35]:

```
[('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

The End.

```
In [36]:  
print(data['color'].describe())  
  
# we have 7380 unique colors  
# 7.2% of products are black in color  
# 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 [37]:

```
color_count = Counter(list(data['color']))  
color_count.most_common(10)
```

Out[37]:

```
[(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 [40]:

```
print(data['formatted_price'].describe())  
  
# 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 [41]:

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

Out[41]:

```
[(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 [42]:

```
print(data['title'].describe())
```

```

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

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

```

## Storing 180k data in 180k\_apparel\_data file in folder /home/bhargav/AAIC/DataSets/Amzon\_Fashion\_Discovery/pickels

In [44]:

```
data.to_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/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.

## Considering only products which have price information

In [45]:

```

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

```

Number of data points After eliminating price=NULL : 28395

## Consider products which have price information and color information

In [46]:

```

# consider products which have color information
# data['color'].isnull() => gives the information about the dataframe row's which have null values
# price == None|Null
data = data.loc[~data['color'].isnull()]
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.

## Storing 28k data in 28kk\_apparel\_data file in folder /home/bhargav/AAIC/DataSets/Amzon\_Fashion\_Discovery/pickels

In [50]:

```
data.to_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/28k_apparel_data')
```

## Wht we are considering 28k Data from 180k ??

- we are taking only 28k data from 180k data because
- if we conider all data it will take so much of time to process the data
- so we process the data on 28k data

## if we want to get all images from the medium-image-url ??

In [51]:

```
# You can download all these 28k images using this code below.  
# You do NOT need to run this code and hence it is commented.  
  
'''  
from PIL import Image  
import requests  
from io import BytesIO  
  
for index, row in images.iterrows():  
    url = row['large_image_url']  
    response = requests.get(url)  
    img = Image.open(BytesIO(response.content))  
    img.save('images/28k_images/'+row['asin']+'.jpeg')  
  
'''
```

Out[51]:

```
"\nfrom PIL import Image\nimport requests\nfrom io import BytesIO\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('images/28k_images/'+row['asin']+'.jpeg')\n\n"
```

### [5.2] Remove near duplicate items

#### [5.2.1] Understand about duplicates.

In [52]:

```
# read data from pickle file from previous stage  
data =  
pd.read_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/28k_apparel_data')  
  
# find number of products that have duplicate titles.  
print(sum(data.duplicated('title')))  
# we have 2325 products which have same title but different color
```

2325

- Example1: products like Shirts have Same Title for different sizes of Same shirts
- Example1: products like Shirts have Same Title for different Colours of Same shirts

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

In [ ]:

```
<table>  
<tr>  
<td> :B00AQ4GMCK</td>  
<td> :B00AQ4GMTS</td>  
</tr>
```

```
```
</tr>
<td> :B00AQ4GMLQ</td>
<td> :B00AQ4GN3I</td>
</tr>
</table>
```

**These shirts exactly same except in color**

In [ ]:

```
<table>
<tr>
<td> :B00G278GZ6</td>
<td> :B00G278W60</td>
</tr>
<tr>
<td> :B00G278Z2A</td>
<td> :B00G2786X8</td>
</tr>
</table>
```

**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 [57]:

```
# read data from pickle file from previous stage
data =
pd.read_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/28k_apparel_data')
```

In [58]:

```
data.head()
```

Out [58]:

	asin	brand	color	product_type_name	medium_image_url	title	formatted_price
4	B004GSI2OS	FeatherLite	Onyx Black/ Stone	SHIRT	https://images-na.ssl-images-amazon.com/images...	Featherlite Ladies' Long Sleeve Stain Resistan...	\$26.26
6	B012YX2ZPI	HX-Kingdom Fashion T-shirts	White	SHIRT	https://images-na.ssl-images-amazon.com/images...	Women's Unique 100% Cotton T - Special Olympic...	\$9.99
11	B001LOUGE4	Fitness Etc.	Black	SHIRT	https://images-na.ssl-images-amazon.com/images...	Ladies Cotton Tank 2x1 Ribbed Tank Top	\$11.99
15	B003BSRPB0	FeatherLite	White	SHIRT	https://images-na.ssl-images-amazon.com/images...	FeatherLite Ladies' Moisture Free Mesh Sport S...	\$20.54
21	B014ICEDNA	FNC7C	Purple	SHIRT	https://images-na.ssl-images-amazon.com/images...	Supernatural Chibis Sam Dean And Castiel Short...	\$7.50

## Remove All products with very few words in title

In [59]:

```
# Remove All products with very few words in title
data_sorted = data[data['title'].apply(lambda x: len(x.split())>4)]
print("After removal of products with short description:", data_sorted.shape[0])
```

After removal of products with short description: 27949

In [65]:

```
# Sort the whole data based on title (alphabetical order of title)
data_sorted.sort_values('title', inplace=True, ascending=False)
data_sorted.head()
```

Out[65]:

	asin	brand	color	product_type_name	medium_image_url	title	formatted_price
61973	B06Y1KZ2WB	Éclair	Black/Pink	SHIRT	https://images-na.ssl-images-amazon.com/images...	Éclair Women's Printed Thin Strap Blouse Black...	\$24.99
133820	B010RV33VE	xiaoming	Pink	SHIRT	https://images-na.ssl-images-amazon.com/images...	xiaoming Womens Sleeveless Loose Long T-shirts...	\$18.19
81461	B01DDSDLNS	xiaoming	White	SHIRT	https://images-na.ssl-images-amazon.com/images...	xiaoming Women's White Long Sleeve Single Brea...	\$21.58
75995	B00X5LYO9Y	xiaoming	Red Anchors	SHIRT	https://images-na.ssl-images-amazon.com/images...	xiaoming Stripes Tank Patch/Bear Sleeve Anchor...	\$15.91
151570	B00WPJG35K	xiaoming	White	SHIRT	https://images-na.ssl-images-amazon.com/images...	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 [66]:

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

In [73]:

```
import itertools
stage1_deduplicate_asins = []
i = 0
j = 0
num_data_points = data_sorted.shape[0]
while i < num_data_points and j < num_data_points:

    previous_i = i

    # store the list of words of ith string in a, ex: a = ['tokidoki', 'The', 'Queen', 'of', 'Diamonds', 'Women's', 'Shirt', 'X-Large']
```

```

a = data['title'].loc[indices[i]].split()

# search for the similar products sequentially
j = i+1
while j < num_data_points:

    # store the list of words of jth string in b, ex: b = ['tokidoki', 'The', 'Queen', 'of', 'Diamonds', 'Women's', 'Shirt', 'Small']
    b = data['title'].loc[indices[j]].split()

    # store the maximum length of two strings
    length = max(len(a), len(b))

    # count is used to store the number of words that are matched in both strings
    count = 0

    # itertools.zip_longest(a,b): will map the corresponding words in both strings, it will
    # appened None in case of unequal strings
    # example: a =['a', 'b', 'c', 'd']
    # b = ['a', 'b', 'd']
    # itertools.zip_longest(a,b): will give [('a', 'a'), ('b', 'b'), ('c', 'd'), ('d', None)]
    for k in itertools.zip_longest(a,b):
        if (k[0] == k[1]):
            count += 1

    # if the number of words in which both strings differ are > 2 , we are considering it as those two apperals are different
    # if the number of words in which both strings differ are < 2 , we are considering it as those two apperals are same, hence we are ignoring them
    if (length - count) > 2: # number of words in which both sensences differ
        # if both strings are differ by more than 2 words we include the 1st string index
        stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[i]])

    # if the comaprision between is between num_data_points, num_data_points-1 strings and they differ in more than 2 words we include both
    if j == num_data_points-1: stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[j]])

    # start searching for similar apperals corresponds 2nd string
    i = j
    break
else:
    j += 1
if previous_i == i:
    break

```

In [74]:

```
data = data.loc[data['asin'].isin(stage1_dedupe_asins)]
```

**We removed the dupliactes which differ only at the end.**

In [75]:

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

Number of data points : 17593

In [78]:

```
data.to_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/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 [80]:

```
data =  
pd.read_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/17k_apperial_data')
```

In [81]:

```
# This code snippet takes significant amount of time.  
# O(n^2) time.  
# Takes about an hour to run on a decent computer.  
  
indices = []  
for i, row in data.iterrows():  
    indices.append(i)  
  
stage2_dedupe_asins = []  
while len(indices)!=0:  
    i = indices.pop()  
    stage2_dedupe_asins.append(data['asin'].loc[i])  
    # consider the first apparel's title  
    a = data['title'].loc[i].split()  
    # store the list of words of ith string in a, ex: a = ['tokidoki', 'The', 'Queen', 'of', 'Diamonds', 'Women's', 'Shirt', 'X-Large']  
    for j in indices:  
  
        b = data['title'].loc[j].split()  
        # store the list of words of jth string in b, ex: b = ['tokidoki', 'The', 'Queen', 'of', 'Diamonds', 'Women's', 'Shirt', 'X-Large']  
  
        length = max(len(a), len(b))  
  
        # count is used to store the number of words that are matched in both strings  
        count = 0  
  
        # itertools.zip_longest(a,b): will map the corresponding words in both strings, it will  
        appened None in case of unequal strings  
        # example: a = ['a', 'b', 'c', 'd']  
        # b = ['a', 'b', 'd']  
        # itertools.zip_longest(a,b): will give [(a,a), (b,b), (c,d), (d, None)]  
        for k in itertools.zip_longest(a,b):  
            if (k[0]==k[1]):  
                count += 1  
  
            # if the number of words in which both strings differ are < 3 , we are considering it as those two apperals are same, hence we are ignoring them  
            if (length - count) < 3:  
                indices.remove(j)
```

```
-----  
KeyboardInterrupt                                     Traceback (most recent call last)  
<ipython-input-81-7cebe2e33e39> in <module>()  
  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', 'Women's', 'Shirt', 'X-Large']  
  20
```

```

~/anaconda3/lib/python3.6/site-packages/pandas/core/indexing.py in __getitem__(self, key)
 1476
 1477     maybe_callable = com._apply_if_callable(key, self.obj)
-> 1478     return self._getitem_axis(maybe_callable, axis=axis)
 1479
 1480     def _is_scalar_access(self, key):

```

```

~/anaconda3/lib/python3.6/site-packages/pandas/core/indexing.py in _getitem_axis(self, key, axis)
 1909
 1910     # fall thru to straight lookup
-> 1911     self._validate_key(key, axis)
 1912     return self._get_label(key, axis=axis)
 1913

```

```

~/anaconda3/lib/python3.6/site-packages/pandas/core/indexing.py in _validate_key(self, key, axis)
 1762     @Appender(_NDFrameIndexer._validate_key.__doc__)
 1763     def _validate_key(self, key, axis):
-> 1764         ax = self.obj_.get_axis(axis)
 1765
 1766         # valid for a label where all labels are in the index

```

KeyboardInterrupt:

In [ ]:

```

# from whole previous products we will consider only
# the products that are found in previous cell
data = data.loc[data['asin'].isin(stage2_dedupe_asins)]

```

In [ ]:

```

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

```

In [ ]:

```

data.to_pickle('pickels/16k_apperal_data')
# Storing these products in a pickle file
# candidates who wants to download these files instead
# of 180K they can download and use them from the Google Drive folder.

```

## 6. Text pre-processing

In [84]:

```

data=pd.read_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/16k_apperal_data')

# NLTK download stop words. [RUN ONLY ONCE]
# goto Terminal (Linux/Mac) or Command-Prompt (Window)
# In the temrinal, type these commands
# $python3
# $import nltk
# $nltk.download()

```

In [85]:

```

# we use the list of stop words that are downloaded from nltk lib.
stop_words = set(stopwords.words('english'))
print ('list of stop words:', stop_words)

def nlp_preprocessing(total_text, index, column):
    if type(total_text) is not int:
        string = ""
        for words in total_text.split():
            # remove the special chars in review like '#$@!%^&*()_+-~?>< etc.
            word = ("").join(e for e in words if e.isalnum())
            # Cover all letters to lower-case
            word = word.lower()
            # stop-word removal
            if not word in stop_words:
                string += word + " "

```

```

        string += word + ...
data[column][index] = string

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

In [86]:

```

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

13.01712000000009 seconds

In [87]:

```
data.head()
```

Out[87]:

	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 [89]:

```
data.to_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/16k_apperial_data_prep sed')
```

## Stemming

In [90]:

```

from nltk.stem.porter import *
stemmer = PorterStemmer()
print(stemmer.stem('arguing'))
print(stemmer.stem('fishing'))
```

```
# We tried using stemming on our titles and it didnot work very well.
```

```
argu  
fish
```

We tried using stemming on our titles and it didnot work very well.

## Text based product similarity

In [2]:

```
data =  
pd.read_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/16k_apperal_data_preproc  
sed')  
data.head()
```

Out [2]:

	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 [148]:

```
# Utility Functions which we will use through the rest of the workshop.  
  
#Display an image  
def display_img(url,ax,fig):  
    # we get the url of the apparel and download it  
    response = requests.get(url)  
    img = Image.open(BytesIO(response.content))  
    # we will display it in notebook  
    plt.imshow(img)  
  
#plotting code to understand the algorithm's decision.  
def plot_heatmap(keys, values, labels, url, text):  
    # keys: list of words of recommended title  
    # values: len(values) == len(keys), values(i) represents the occurence of the word  
    keys(i)  
    # labels: len(labels) == len(keys), the values of labels depends on the model we are using  
    # if model == 'bag of words': labels(i) = values(i)  
    # if model == 'tfidf weighted bag of words':labels(i) = tfidf(keys(i))  
    # if model == 'idf weighted bag of words':labels(i) = idf(keys(i))  
    # url : apparel's url  
  
    # we will devide the whole figure into two parts  
    gs = gridspec.GridSpec(2, 2, width_ratios=[4,1], height_ratios=[4,1])  
    fig = plt.figure(figsize=(25,3))  
  
    # 1st, plotting heat map that represents the count of commonly occurred words in title2  
    ax = plt.subplot(gs[0])  
    # it displays a cell in white color if the word is intersection(lis of words of title1 and  
    list of words of title2), in black if not  
    ax = sns.heatmap(np.array([values]), annot=np.array([labels]))  
    ax.set_xticklabels(keys) # set that axis labels as the words of title
```

```

        ax.set_title(text) # apparel title

        # 2nd, plotting image of the the apparel
        ax = plt.subplot(gs[1])
        # we don't want any grid lines for image and no labels on x-axis and y-axis
        ax.grid(False)
        ax.set_xticks([])
        ax.set_yticks([])

        # we call dispaly_img based with paramete url
        display_img(url, ax, fig)

        # displays combine figure ( heat map and image together)
        plt.show()

def plot_heatmap_image(doc_id, vec1, vec2, url, text, model):

    # doc_id : index of the title1
    # vec1 : input apparels's vector, it is of a dict type {word:count}
    # vec2 : recommended apparels's vector, it is of a dict type {word:count}
    # url : apparels image url
    # text: title of recomonded apparel (used to keep title of image)
    # model, it can be any of the models,
    # 1. bag_of_words
    # 2. tfidf
    # 3. idf

    # we find the common words in both titles, because these only words contribute to the distance
    between two title vec's
    intersection = set(vec1.keys()) & set(vec2.keys())

    # we set the values of non intersecting words to zero, this is just to show the difference in
    heatmap
    for i in vec2:
        if i not in intersection:
            vec2[i]=0

    # for labeling heatmap, keys contains list of all words in title2
    keys = list(vec2.keys())
    # if ith word in intersection(lis of words of title1 and list of words of title2):
    values(i)=count of that word in title2 else values(i)=0
    values = [vec2[x] for x in vec2.keys()]

    # labels: len(labels) == len(keys), the values of labels depends on the model we are using
    # if model == 'bag of words': labels(i) = values(i)
    # if model == 'tfidf weighted bag of words':labels(i) = tfidf(keys(i))
    # if model == 'idf weighted bag of words':labels(i) = idf(keys(i))

    if model == 'bag_of_words':
        labels = values
    elif model == 'tfidf':
        labels = []
        for x in vec2.keys():
            # tfidf_title_vectorizer.vocabulary_ it contains all the words in the corpus
            # tfidf_title_features[doc_id, index_of_word_in_corpus] will give the tfidf value of w
            rd in given document (doc_id)
            if x in tfidf_title_vectorizer.vocabulary_:
                labels.append(tfidf_title_features[doc_id, tfidf_title_vectorizer.vocabulary_[x]])
            else:
                labels.append(0)
    elif model == 'idf':
        labels = []
        for x in vec2.keys():
            # idf_title_vectorizer.vocabulary_ it contains all the words in the corpus
            # idf_title_features[doc_id, index_of_word_in_corpus] will give the idf value of word
            in given document (doc_id)
            if x in idf_title_vectorizer.vocabulary_:
                labels.append(idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[x]])
            else:
                labels.append(0)

    plot_heatmap(keys, values, labels, url, text)

# this function gets a list of wrods along with the frequency of each
# word given "text"
def text_to_vector(text):

```

```

word_ = re.compile(r'\w+')
words = word_.findall(text)
# words stores list of all words in given string, you can try 'words = text.split()' this will
also gives same result
return Counter(words) # Counter counts the occurrence of each word in list, it returns dict
type object {word1:count}

def get_result(doc_id, content_a, content_b, url, model):
    text1 = content_a
    text2 = content_b

    # vector1 = dict{word1:#count, word2:#count, etc.}
    vector1 = text_to_vector(text1)

    # vector1 = dict{word21:#count, word22:#count, etc.}
    vector2 = text_to_vector(text2)

    plot_heatmap_image(doc_id, vector1, vector2, url, text2, model)

# get_result(indices[i],data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data['m
edium_image_url'].loc[df_indices[i]], 'bag_of_words')

```

In [24]:

```

def text_to_vector(text):
    word = re.compile(r'\w+')
    print(word)
    words = word.findall(text)
    print(words)
    # words stores list of all words in given string, you can try 'words = text.split()' this will
also gives same result
    return Counter(words) #

text_to_vector("apple is a bad bad boy bad is of another tyoe of boy")

re.compile('\w+')
['apple', 'is', 'a', 'bad', 'bad', 'boy', 'bad', 'is', 'of', 'another', 'tyoe', 'of', 'boy']

```

Out[24]:

```

Counter({'apple': 1,
         'is': 2,
         'a': 1,
         'bad': 3,
         'boy': 2,
         'of': 2,
         'another': 1,
         'tyoe': 1})

```

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

In [4]:

```

from sklearn.feature_extraction.text import CountVectorizer
title_vectorizer = CountVectorizer()
title_features = title_vectorizer.fit_transform(data['title'])
title_features.get_shape() # get number of rows and columns in feature matrix.
# title_features.shape = #data_points * #words_in_corpus
# CountVectorizer().fit_transform(corpus) returns
# the a sparse matrix of dimensions #data_points * #words_in_corpus

# What is a sparse vector?

# title_features[doc_id, index_of_word_in_corpus] = number of times the word occurred in that doc

```

Out[4]:

```
(16042, 12609)
```

In [8]:

```
title_features
```

Out[8]:

```
<16042x12609 sparse matrix of type '<class 'numpy.int64'>'  
with 147545 stored elements in Compressed Sparse Row format>
```

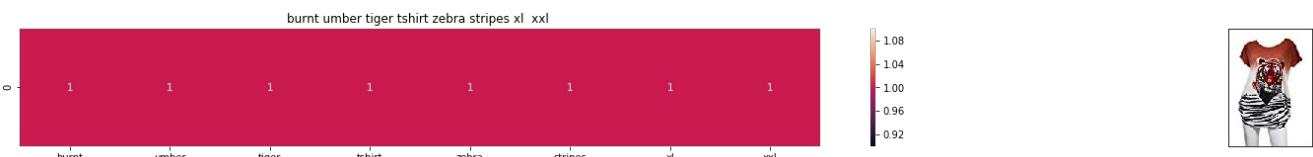
In [19]:

```
def bag_of_words_model(doc_id, num_results):  
    pairwise_dist = pairwise_distances(title_features,title_features[doc_id])  
    print(pairwise_dist.shape)  
    #Top 20 shortest pairwiseDistances Index  
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]  
    #print(indices)  
    #Top 20 shortest pairwiseDistances values  
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]  
    #print(pdists)  
    df_indices = list(data.index[indices])  
    #print(df_indices)  
    print(data.index[[0,1,2,3,4,5]])  
  
bag_of_words_model(12566, 20)
```

```
(16042, 1)  
Int64Index([4, 6, 15, 27, 46, 60], dtype='int64')
```

In [151]:

```
def bag_of_words_model(doc_id, num_results):  
    # doc_id: apparel's id in given corpus  
  
    # pairwise_dist will store the distance from given input apparel to all remaining apparels  
    # the metric we used here is cosine, the coside distance is mesured as  $K(X, Y) = \langle X, Y \rangle / (\|X\| * \|Y\|)$   
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity  
    pairwise_dist = pairwise_distances(title_features,title_features[doc_id])  
  
    # np.argsort will return indices of the smallest distances  
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]  
    #pdists will store the smallest distances  
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]  
  
    #data frame indices of the 9 smallest distace's  
    df_indices = list(data.index[indices])  
  
    for i in range(0,len(indices)):  
        # we will pass 1. doc_id, 2. title1, 3. title2, url, model  
        get_result(indices[i],data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data['medium_image_url'].loc[df_indices[i]], 'bag_of_words')  
        print('ASIN :',data['asin'].loc[df_indices[i]])  
        print ('Brand:', data['brand'].loc[df_indices[i]])  
        print ('Title:', data['title'].loc[df_indices[i]])  
        print ('Euclidean similarity with the query image :', pdists[i])  
        print('='*60)  
  
    #call the bag-of-words model for a product to get similar products.  
    bag_of_words_model(12566, 20) # change the index if you want to.  
    # In the output heat map each value represents the count value  
    # of the label word, the color represents the intersection  
    # with inputs title.  
  
    #try 12566  
    #try 931
```



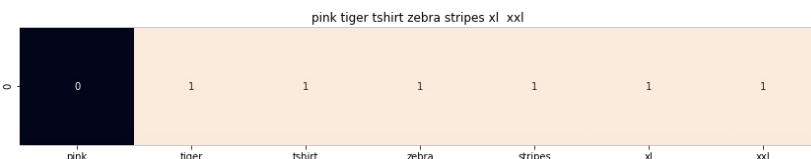
ASIN : B00JXQB5FQ

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 : 1.7320508075688772

=====



ASIN : B00JXQCWTO

Brand: Si Row

Title: brown white tiger tshirt tiger stripes xl xxl

Euclidean similarity with the query image : 2.449489742783178

=====



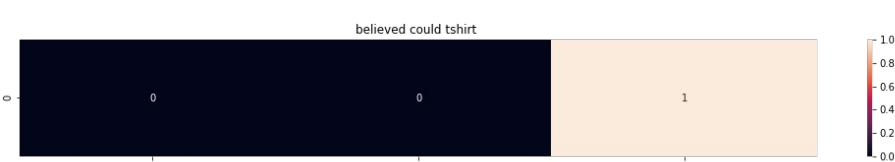
ASIN : B00JXQCUIC

Brand: Si Row

Title: yellow tiger tshirt tiger stripes l

Euclidean similarity with the query image : 2.6457513110645907

=====



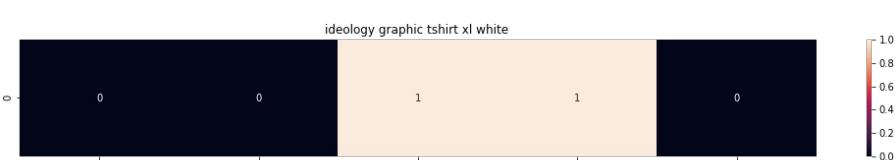
ASIN : B07568NZX4

Brand: Rustic Grace

Title: believed could tshirt

Euclidean similarity with the query image : 3.0

=====

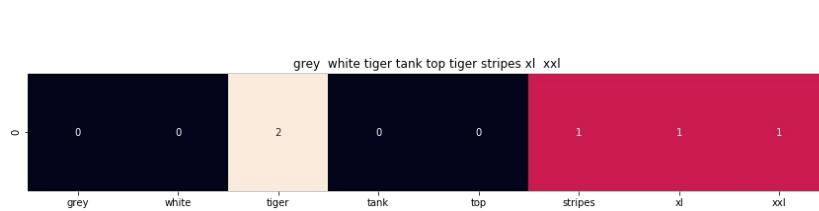


ASIN : B01NB0NKRO

Brand: Ideology

Title: ideology graphic tshirt xl white

Euclidean similarity with the query image : 3.0



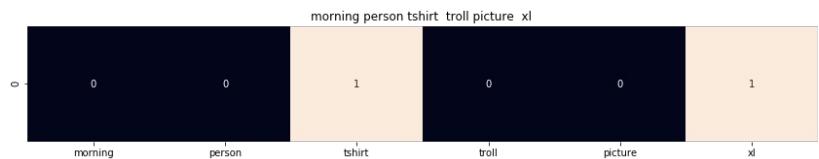
ASIN : B00JXQAFZ2

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 : 3.1622776601683795

=====



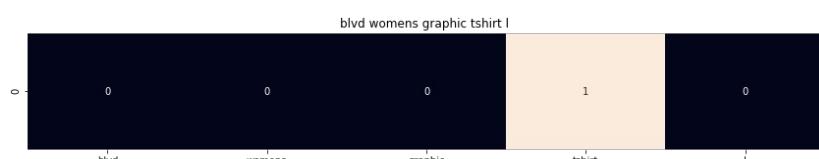
ASIN : B01KVZUB6G

Brand: Merona

Title: merona green gold stripes

Euclidean similarity with the query image : 3.1622776601683795

=====



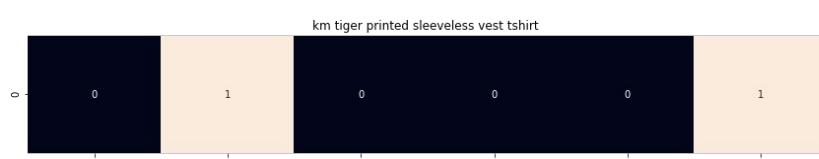
ASIN : B0733R2CJK

Brand: BLVD

Title: blvd womens graphic tshirt l

Euclidean similarity with the query image : 3.1622776601683795

=====



ASIN : B012VQLT6Y

Brand: KM T-shirt

Title: km tiger printed sleeveless vest tshirt

Euclidean similarity with the query image : 3.1622776601683795

=====



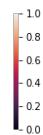
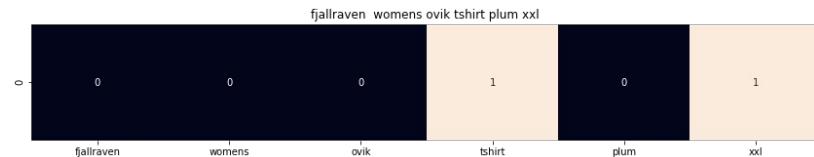


ASIN : B00JXQC8L6

Brand: Si Row

Title: blue peacock print tshirt l

Euclidean similarity with the query image : 3.1622776601683795



ASIN : B06XC3CZF6

Brand: Fjallraven

Title: fjallraven womens ovik tshirt plum xxl

Euclidean similarity with the query image : 3.1622776601683795

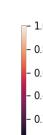
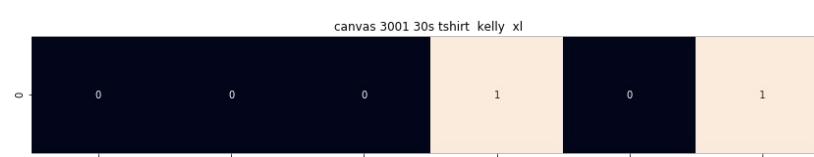


ASIN : B005IT8OBA

Brand: Hetalia

Title: hetalia us girl tshirt

Euclidean similarity with the query image : 3.1622776601683795

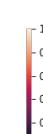


ASIN : B0088PN0LA

Brand: Red House

Title: canvas 3001 30s tshirt kelly xl

Euclidean similarity with the query image : 3.1622776601683795

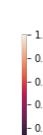


ASIN : B06X99V6WC

Brand: Brunello Cucinelli

Title: brunello cucinelli tshirt women white xl

Euclidean similarity with the query image : 3.1622776601683795



ASIN : B06Y1JPW1Q

Brand: Xhilaration

Title: xhilaration womens lace tshirt salmon xxl

```
Euclidean similarity with the query image : 3.1622776601683795
=====

```



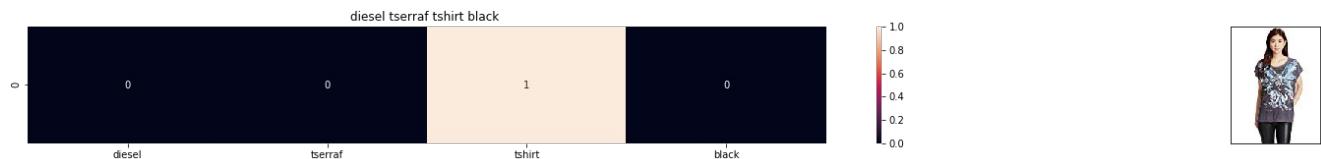
ASIN : B06X6GX6WG

Brand: Animal

Title: animal oceania tshirt yellow

```
Euclidean similarity with the query image : 3.1622776601683795
=====

```



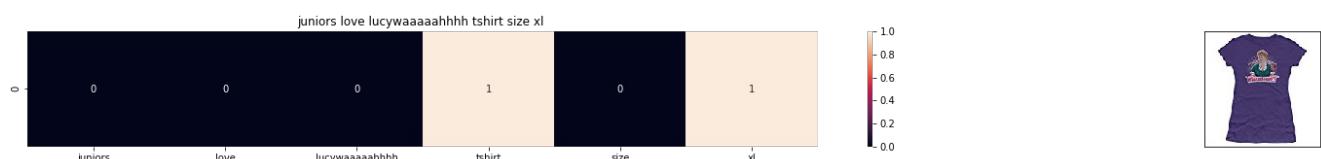
ASIN : B017X8PW9U

Brand: Diesel

Title: diesel tserraf tshirt black

```
Euclidean similarity with the query image : 3.1622776601683795
=====

```



ASIN : B00IAA4JIQ

Brand: I Love Lucy

Title: juniors love lucywaaaaahhhh tshirt size xl

```
Euclidean similarity with the query image : 3.1622776601683795
=====

```

## [8.5] TF-IDF based product similarity

In [96]:

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

In [97]:

```
def tfidf_model(doc_id, num_results):
    # doc_id: apparel's id in given corpus

    # pairwise_dist will store the distance from given input apparel to all remaining apparels
    # the metric we used here is cosine, the coside distance is mesured as K(X, Y) = <X, Y> / (||X|| * ||Y||)
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
    pairwise_dist = pairwise_distances(tfidf_title_features,tfidf_title_features[doc_id])

    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

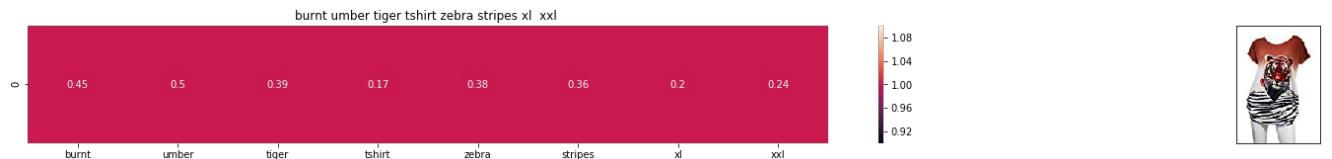
    #data frame indices of the 9 smallest distace's
```

```

df_indices = list(data.index[indices])

for i in range(0,len(indices)):
    # we will pass 1. doc_id, 2. title1, 3. title2, url, model
    get_result(indices[i], data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], 
data['medium_image_url'].loc[df_indices[i]], 'tfidf')
    print('ASIN :',data['asin'].loc[df_indices[i]])
    print('BRAND :',data['brand'].loc[df_indices[i]])
    print ('Eucliden distance from the given image :', pdists[i])
    print('='*125)
tfidf_model(12566, 20)
# in the output heat map each value represents the tfidf values of the label word, the color represents the intersection with inputs title

```



ASIN : B00JXQB5FQ

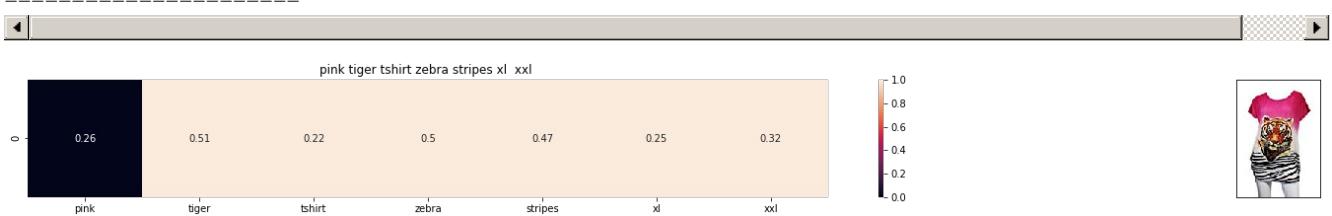
BRAND : Si Row

Eucliden distance from the given image : 0.0

---



---



ASIN : B00JXQASS6

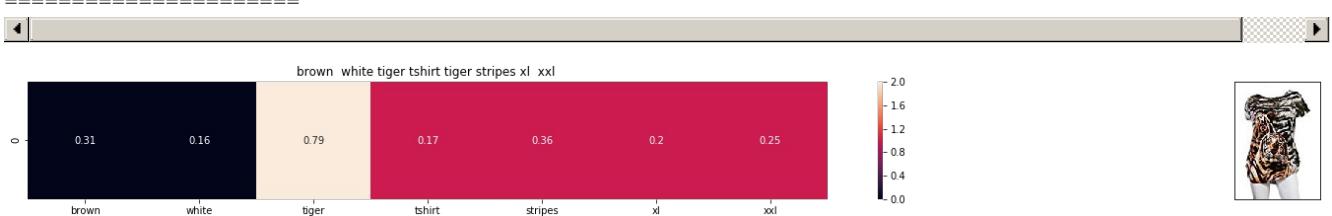
BRAND : Si Row

Eucliden distance from the given image : 0.7536331912451361

---



---



ASIN : B00JXQCWT0

BRAND : Si Row

Eucliden distance from the given image : 0.9357643943769645

---



---



ASIN : B00JXQAFZ2

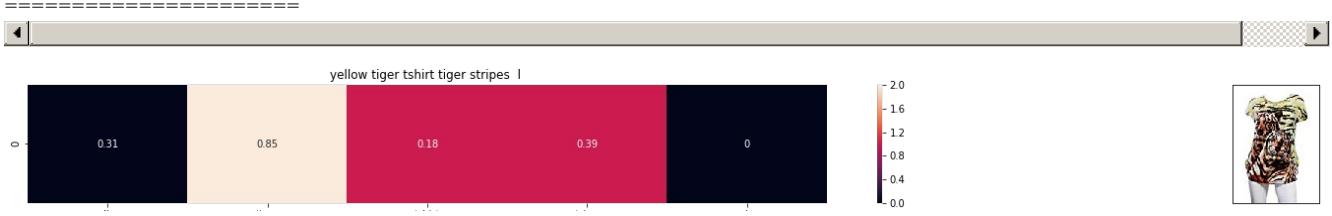
BRAND : Si Row

Eucliden distance from the given image : 0.9586153524200749

---



---



yellow

tiger

tshirt

stripes

l

ASIN : B00JXQCUIC

BRAND : Si Row

Eucliden distance from the given image : 1.000074961446881



ASIN : B00JXQA094

BRAND : Si Row

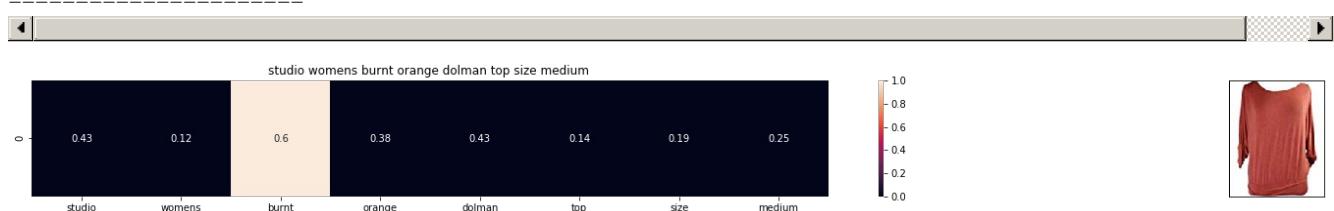
Eucliden distance from the given image : 1.023215552457452



ASIN : B00JXQAUWA

BRAND : Si Row

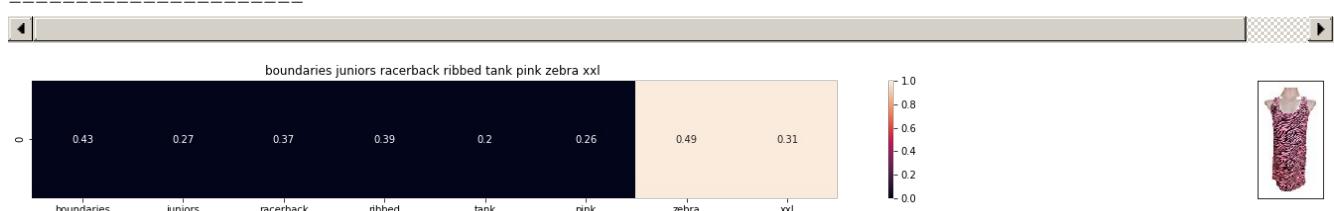
Eucliden distance from the given image : 1.031991846303421



ASIN : B06XSCVFT5

BRAND : Studio M

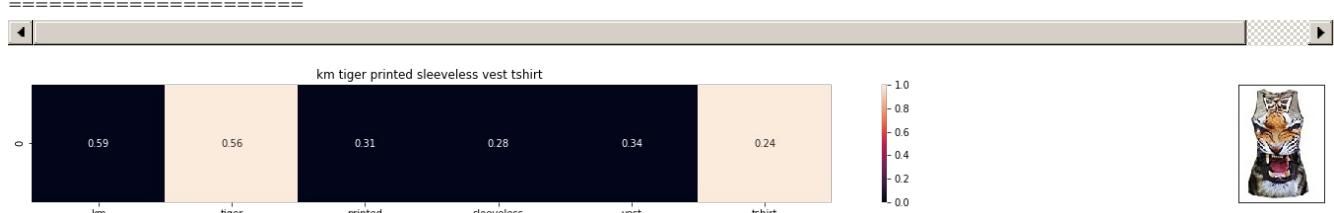
Eucliden distance from the given image : 1.2106843670424716



ASIN : B06Y2GTYPM

BRAND : No Boundaries

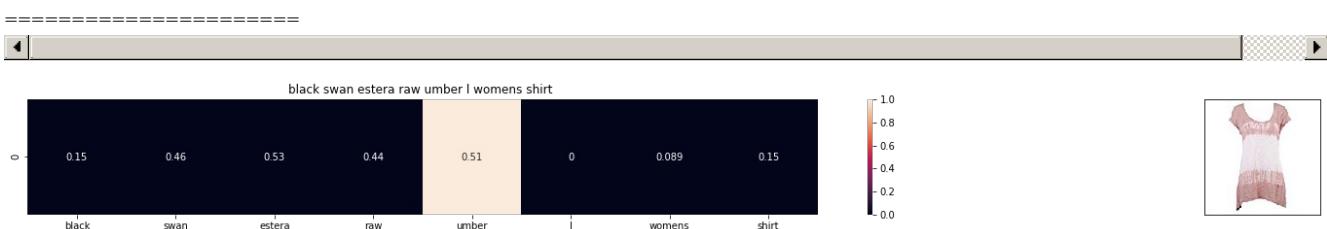
Eucliden distance from the given image : 1.212168381072083



ASIN : B012VQLT6Y

BRAND : KM T-shirt

Eucliden distance from the given image : 1.219790640280982

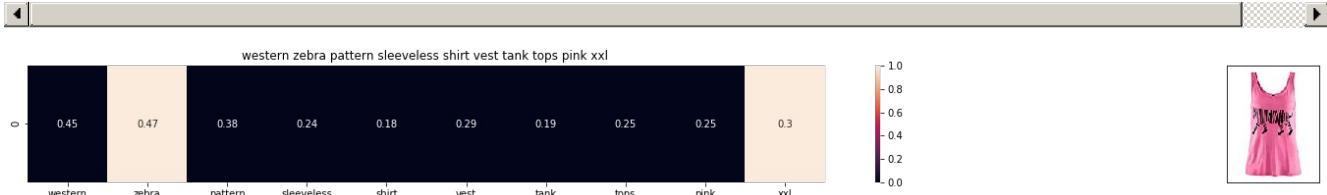


ASIN : B06Y1VN8WQ

BRAND : Black Swan

Eucliden distance from the given image : 1.2206849659998316

=====

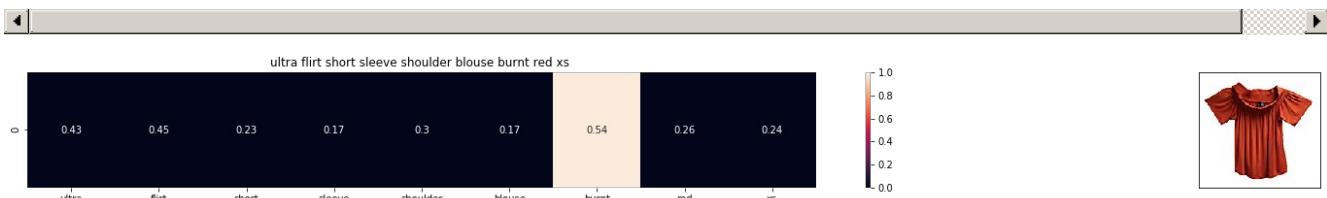


ASIN : B00Z6HEXWI

BRAND : Black Temptation

Eucliden distance from the given image : 1.221281392120943

=====

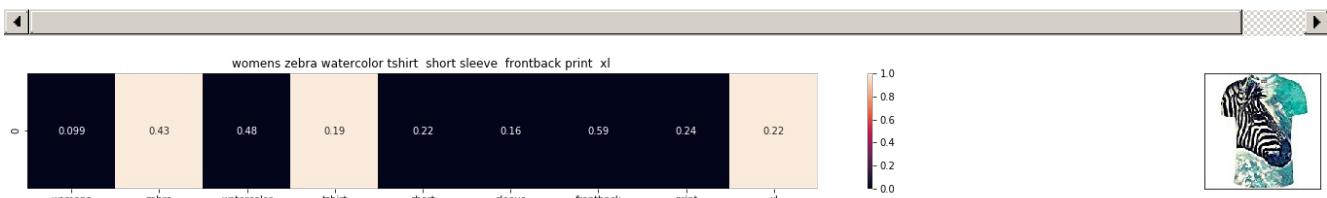


ASIN : B074TR12BH

BRAND : Ultra Flirt

Eucliden distance from the given image : 1.2313364094597743

=====

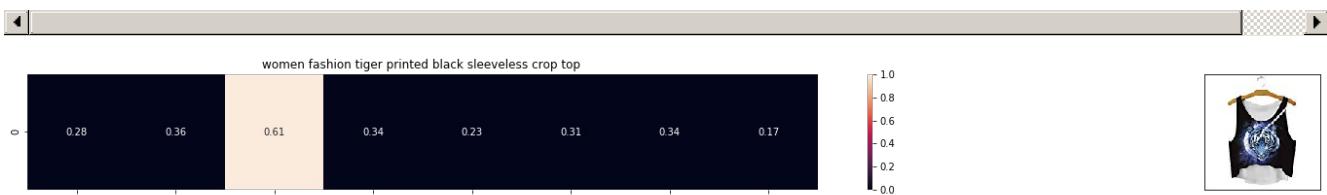


ASIN : B072R2JXKW

BRAND : WHAT ON EARTH

Eucliden distance from the given image : 1.2318451972624518

=====



ASIN : B074T8ZYGX

BRAND : MKP Crop Top

Eucliden distance from the given image : 1.2340607457359425

=====

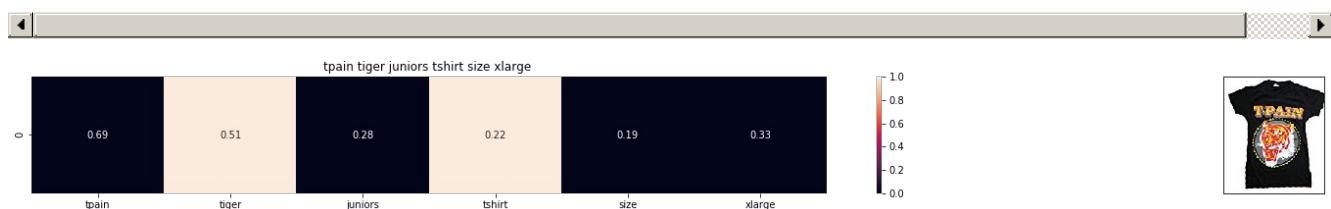




ASIN : B071ZDF6T2

BRAND : Mossimo

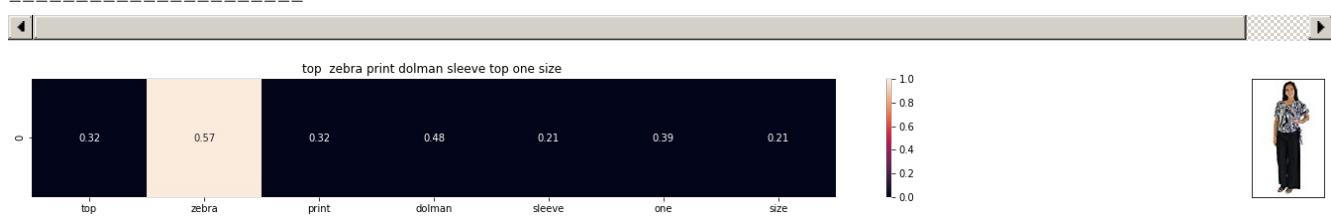
Euclidean distance from the given image : 1.2352785577664824



ASIN : B01K0H02OG

BRAND : Tultex

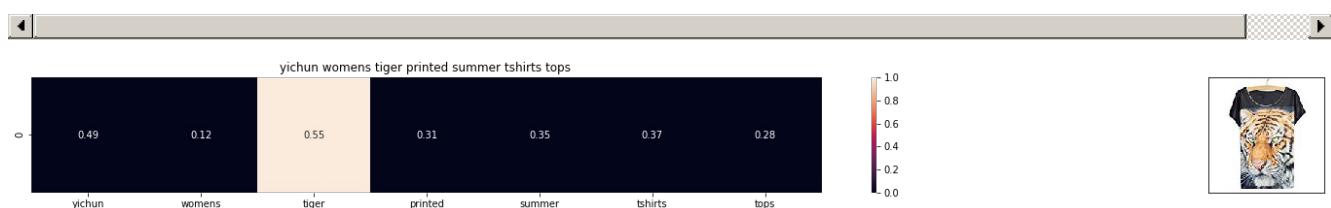
Euclidean distance from the given image : 1.236457298812782



ASIN : B00H8A6ZLI

BRAND : Vivian's Fashions

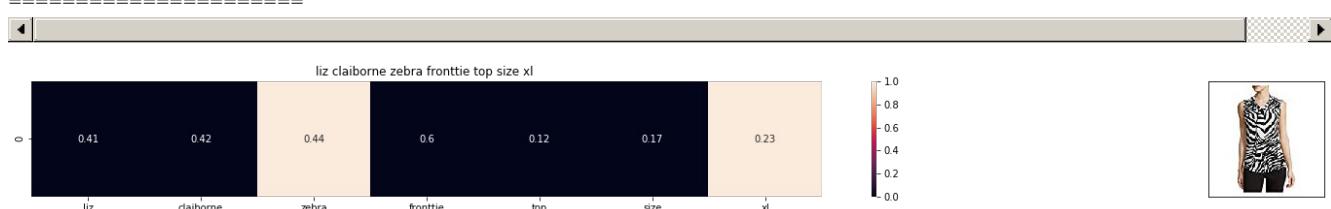
Euclidean distance from the given image : 1.24996155052848



ASIN : B010NN9RXO

BRAND : YICHUN

Euclidean distance from the given image : 1.25354614208561



ASIN : B06XBY5QXL

BRAND : Liz Claiborne

Euclidean distance from the given image : 1.2538832938357722

## [8.5] IDF based product similarity

In [98]:

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

# idf_title_features.shape = #data_points * #words_in_corpus
```

```
# CountVectorizer().fit_transform(courpus) returns the a sparase matrix of dimensions #data_points  
* #words_in_corpus  
# idf_title_features[doc_id, index_of_word_in_corpus] = number of times the word occured in that doc
```

In [99]:

```
def nContaining(word):
    # return the number of documents which had the given word
    return sum(1 for blob in data['title'] if word in blob.split())

def idf(word):
    #  $idf = \log(\#number\ of\ docs / \#number\ of\ docs\ which\ had\ the\ given\ word)$ 
    return math.log(data.shape[0] / (nContaining(word)))
```

In [100]:

```

# we need to convert the values into float
idf_title_features = idf_title_features.astype(np.float)

for i in idf_title_vectorizer.vocabulary_.keys():
    # for every word in whole corpus we will find its idf value
    idf_val = idf(i)

    # to calculate idf_title_features we need to replace the count values with the idf values of the word
    # idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0] will return all documents in which the word i present
    for j in idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0]:
        # we replace the count values of word i in document j with idf_value of word i
        # idf_title_features[doc_id, index_of_word_in_corpus] = idf value of word
        idf_title_features[j,idf_title_vectorizer.vocabulary_[i]] = idf_val

```

In [101]:

```

def idf_model(doc_id, num_results):
    # doc_id: apparel's id in given corpus

    # pairwise_dist will store the distance from given input apparel to all remaining apparels
    # the metric we used here is cosine, the coside distance is measured as  $K(X, Y) = \langle X, Y \rangle / (\|X\| * \|Y\|)$ 
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
    pairwise_dist = pairwise_distances(idf_title_features,idf_title_features[doc_id])

    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    #data frame indices of the 9 smallest distace's
    df_indices = list(data.index[indices])

    for i in range(0,len(indices)):
        get_result(indices[i],data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data['medium_image_url'].loc[df_indices[i]], 'idf')
        print('ASIN :',data['asin'].loc[df_indices[i]])
        print('Brand :',data['brand'].loc[df_indices[i]])
        print ('euclidean distance from the given image :', pdists[i])
        print('='*125)

idf_model(12566,20)
# in the output heat map each value represents the idf values of the label word, the color represents the intersection with inputs title

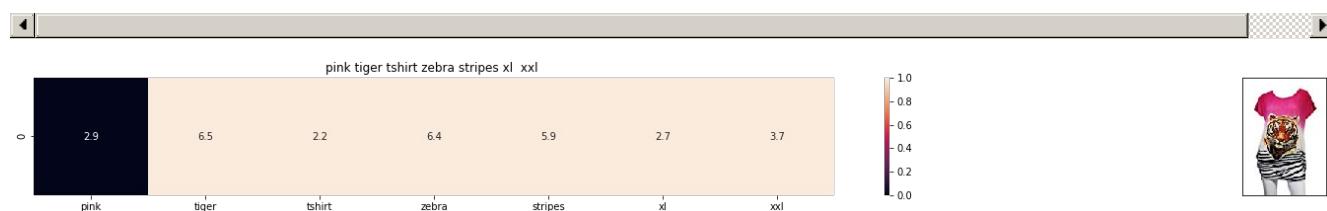
```



ASIN : B00JXQB5FQ

Brand : Si Row

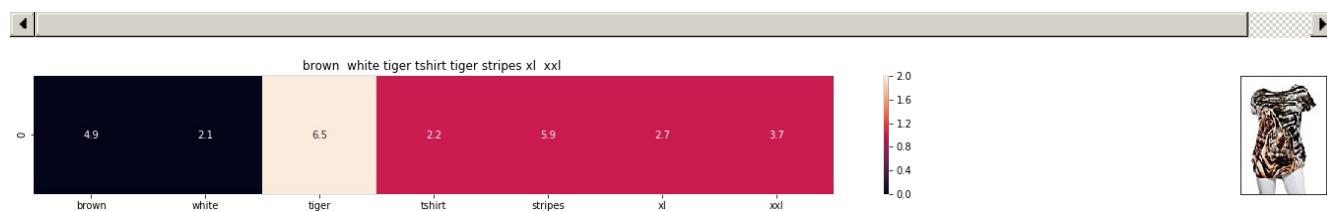
euclidean distance from the given image : 0.0



ASIN : B00JXQASS6

Brand : Si Row

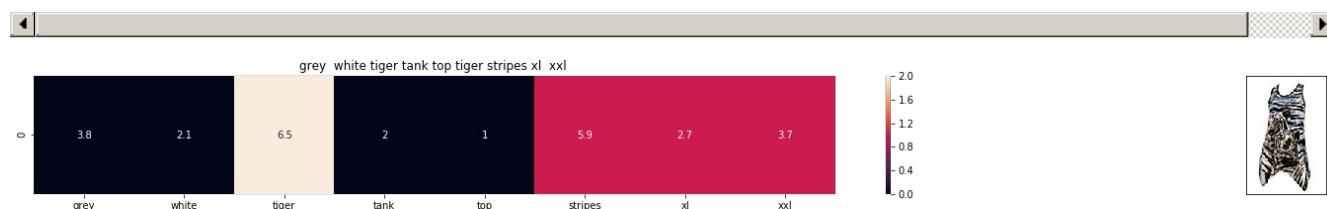
euclidean distance from the given image : 12.20507131122177



ASIN : B00JXQCWT0

Brand : Si Row

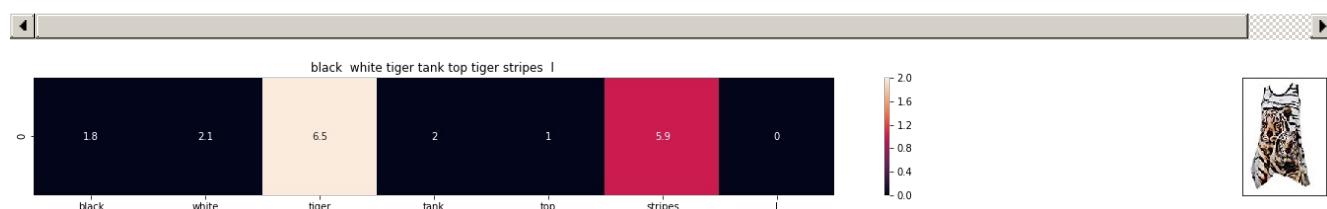
euclidean distance from the given image : 14.468362685603465



ASIN : B00JXQAFZ2

Brand : Si Row

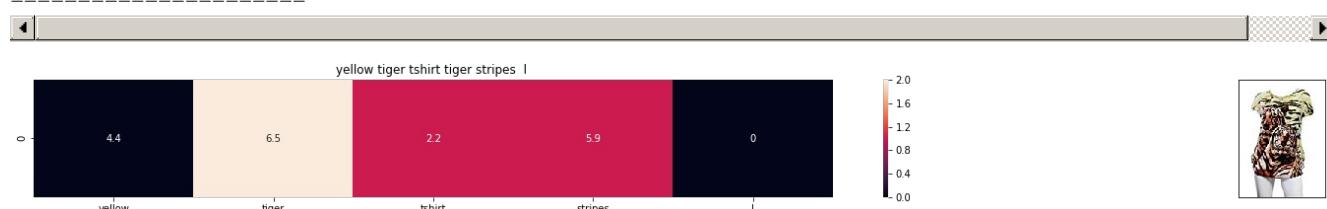
euclidean distance from the given image : 14.486832924778964



ASIN : B00JXQAO94

Brand : Si Row

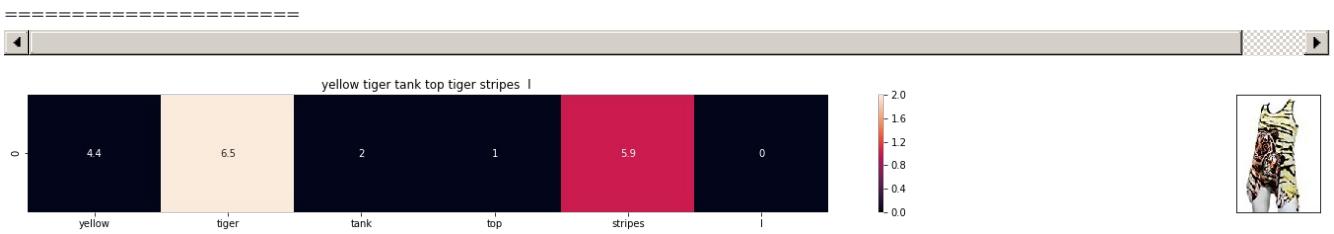
euclidean distance from the given image : 14.833392966672909



ASIN : B00JXOCUIC

Brand : Si Row

euclidean distance from the given image : 14.898744516719225

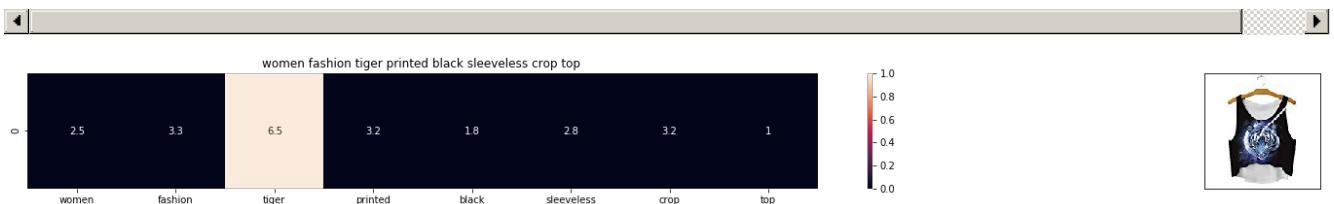


ASIN : B00JXQAUWA

Brand : Si Row

euclidean distance from the given image : 15.224458287343769

=====

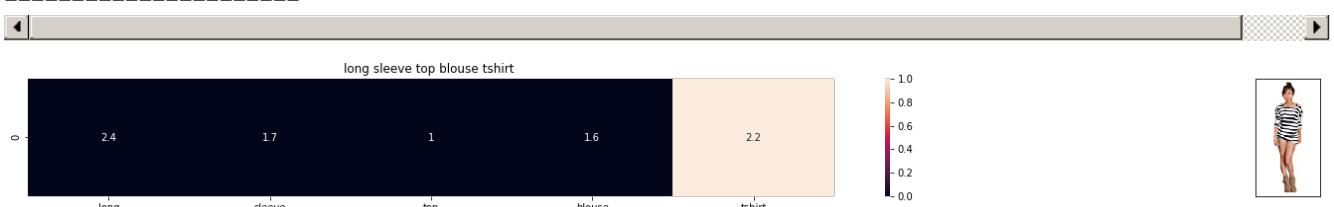


ASIN : B074T8ZYGX

Brand : MKP Crop Top

euclidean distance from the given image : 17.080812955631995

=====



ASIN : B00KF2N5PU

Brand : Vietsbay

euclidean distance from the given image : 17.090168125645416

=====

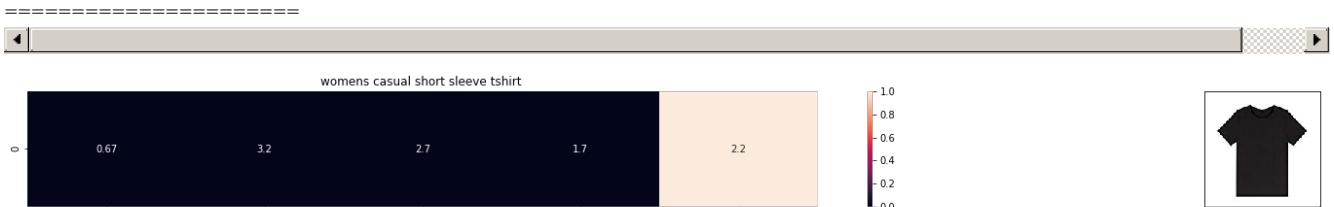


ASIN : B00JPOZ9GM

Brand : Sofra

euclidean distance from the given image : 17.153215337562703

=====



ASIN : B074T9KG9Q

Brand : Rain

euclidean distance from the given image : 17.33671523874989

=====





ASIN : B00H8A6ZLI

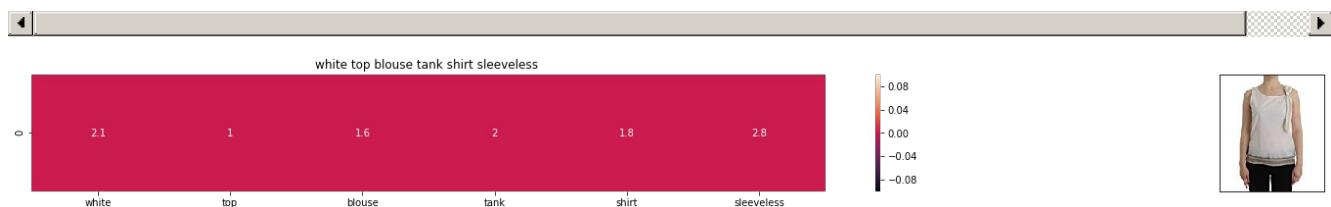
Brand : Vivian's Fashions

euclidean distance from the given image : 17.410075941001253

---



---



ASIN : B074G5G5RK

Brand : ERMANNO SCERVINO

euclidean distance from the given image : 17.539921335459557

---



---



ASIN : B06XSCVFT5

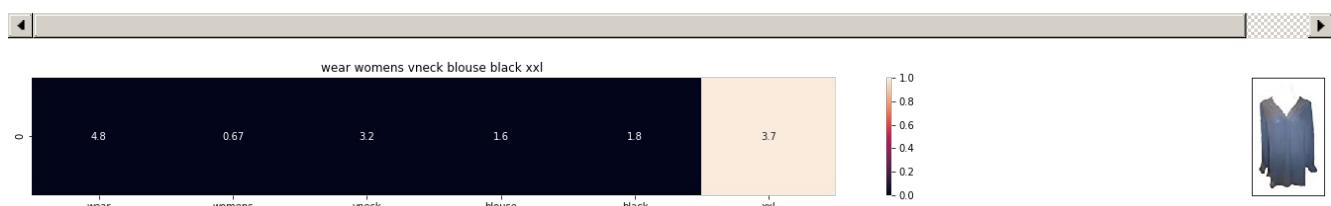
Brand : Studio M

euclidean distance from the given image : 17.61275854366134

---



---



ASIN : B06Y6FH453

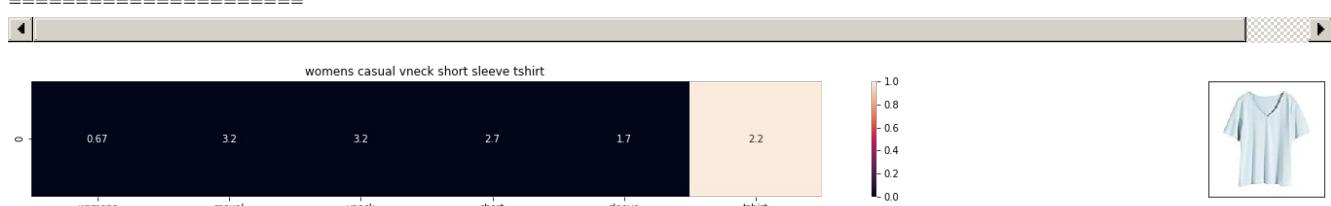
Brand : Who What Wear

euclidean distance from the given image : 17.623745282500135

---



---



ASIN : B074V45DCX

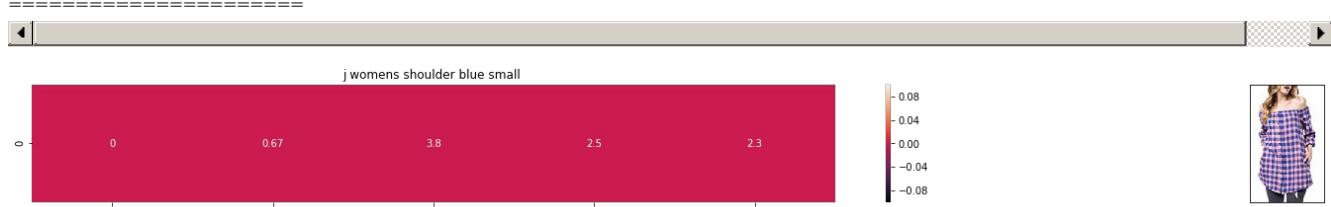
Brand : Rain

euclidean distance from the given image : 17.634342496835046

---



---



ASIN : B07583CQFT

Brand : Very J

euclidean distance from the given image : 17.63753712743611

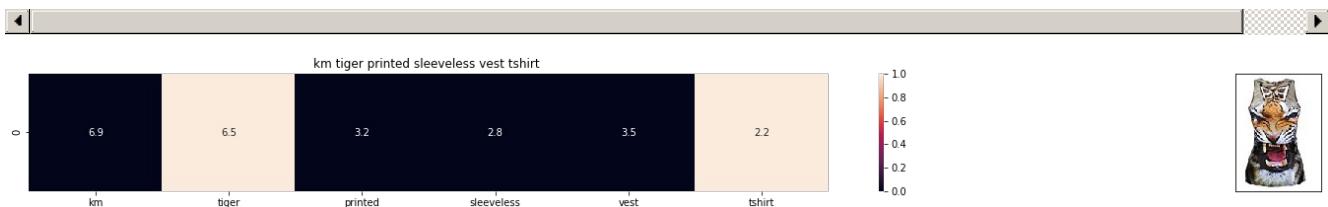
euclidean distance from the given image : 17.0010512745011



ASIN : B073GJGVBN

Brand : Ivan Levi

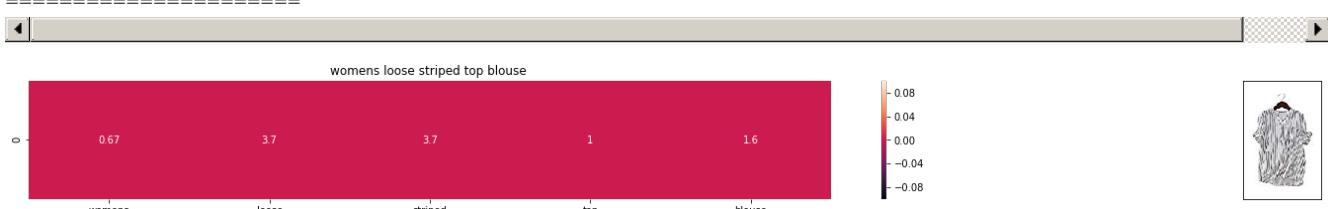
euclidean distance from the given image : 17.7230738913371



ASIN : B012VQLT6Y

Brand : KM T-shirt

euclidean distance from the given image : 17.762588561202364



ASIN : B00ZZMYBRG

Brand : HP-LEISURE

euclidean distance from the given image : 17.779536864674238

## [9] Text Semantics based product similarity

In [102]:

```
# credits: https://www.kaggle.com/c/word2vec-nlp-tutorial#part-2-word-vectors
# Custom Word2Vec using your own text data.
# Do NOT RUN this code.
# It is meant as a reference to build your own Word2Vec when you have
# lots of data.

'''
# Set values for various parameters
num_features = 300      # Word vector dimensionality
min_word_count = 1       # Minimum word count
num_workers = 4           # Number of threads to run in parallel
context = 10              # Context window size

downsampling = 1e-3      # Downsample setting for frequent words

# Initialize and train the model (this will take some time)
from gensim.models import word2vec
print ("Training model...")
model = word2vec.Word2Vec(sen_corpus, workers=num_workers,
                           size=num_features, min_count = min_word_count,
                           window = context)

'''
```

In [102]:

```
'\n# Set values for various parameters\nnum_features = 300      # Word vector dimensionality\n\nmin_word_count = 1      # Minimum word count\n\nnum_workers = 4          # Number of threads to run in parallel\ncontext = 10            # Context window size\n\ndownsampling = 1e-3     # Downsample setting for frequent words\n\n# Initialize and train the model (this will take some time)\nfrom gensim.models import Word2Vec\nprint ("Training model...")\nmodel = Word2Vec(sen_corpus, workers=num_workers,\nsize=num_features, min_count = min_word_count,\n                  window = context)\n'\n
```

In [104]:

```
from gensim.models import Word2Vec\nfrom gensim.models import KeyedVectors\nimport pickle\n\n# in this project we are using a pretrained model by google\n# its 3.3G file, once you load this into your memory\n# it occupies ~9Gb, so please do this step only if you have >12G of ram\n# we will provide a pickle file which contains a dict ,\n# and it contains all our corpus words as keys and model[word] as values\n# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"\n# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTtSS21pQmM/edit\n# it's 1.9GB in size.\n\n'''\nmodel = KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)\n'''\n\n#if you do NOT have RAM >= 12GB, use the code below.\nwith open('/home/bhargav/AAIC/DataSets/Amazon_Fashion_Discovery/word2vec_model', 'rb') as handle:\n    model = pickle.load(handle)
```

In [105]:

```
# Utility functions\n\n\ndef get_word_vec(sentence, doc_id, m_name):\n    # sentence : title of the apparel\n    # doc_id: document id in our corpus\n    # m_name: model information it will take two values\n        # if m_name == 'avg', we will append the model[i], w2v representation of word i\n        # if m_name == 'weighted', we will multiply each w2v[word] with the idf(word)\n    vec = []\n    for i in sentence.split():\n        if i in vocab:\n            if m_name == 'weighted' and i in idf_title_vectorizer.vocabulary_:\n                vec.append(idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[i]] * model[i])\n            elif m_name == 'avg':\n                vec.append(model[i])\n            else:\n                # if the word in our corpus is not there in the google word2vec corpus, we are just\n                # ignoring it\n                vec.append(np.zeros(shape=(300,)))\n        # we will return a numpy array of shape (#number of words in title * 300 ) 300 =\n        # len(w2v_model[word])\n        # each row represents the word2vec representation of each word (weighted/avg) in given\n        # sentance\n    return np.array(vec)\n\n\ndef get_distance(vec1, vec2):\n    # vec1 = np.array(#number_of_words_title1 * 300), each row is a vector of length 300\n    # corresponds to each word in give title\n    # vec2 = np.array(#number_of_words_title2 * 300), each row is a vector of length 300\n    # corresponds to each word in give title\n\n    final_dist = []\n    # for each vector in vec1 we calculate the distance(euclidean) to all vectors in vec2\n    for i in vec1:\n        dist = []\n        for j in vec2:\n            # np.linalg.norm(i-j) will result the euclidean distance between vectors i, j\n            dist.append(np.linalg.norm(i-j))\n        final_dist.append(np.array(dist))
```

```

# final_dist = np.array(#number of words in title1 * #number of words in title2)
# final_dist[i,j] = euclidean distance between vectors i, j
return np.array(final_dist)

def heat_map_w2v(sentence1, sentence2, url, doc_id1, doc_id2, model):
    # sentance1 : title1, input apparel
    # sentance2 : title2, recommended apparel
    # url: apparel image url
    # doc_id1: document id of input apparel
    # doc_id2: document id of recommended apparel
    # model: it can have two values, 1. avg 2. weighted

    #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
    s1_vec = get_word_vec(sentence1, doc_id1, model)
    #s2_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
    s2_vec = get_word_vec(sentence2, doc_id2, model)

    # s1_s2_dist = np.array(#number of words in title1 * #number of words in title2)
    # s1_s2_dist[i,j] = euclidean distance between words i, j
    s1_s2_dist = get_distance(s1_vec, s2_vec)

    # devide whole figure into 2 parts 1st part displays heatmap 2nd part displays image of apparel
1
    gs = gridspec.GridSpec(2, 2, width_ratios=[4,1],height_ratios=[2,1])
    fig = plt.figure(figsize=(15,15))

    ax = plt.subplot(gs[0])
    # plotting the heap map based on the pairwise distances
    ax = sns.heatmap(np.round(s1_s2_dist,4), annot=True)
    # set the x axis labels as recommended apparels title
    ax.set_xticklabels(sentence2.split())
    # set the y axis labels as input apparels title
    ax.set_yticklabels(sentence1.split())
    # set title as recommended apparels title
    ax.set_title(sentence2)

    ax = plt.subplot(gs[1])
    # we remove all grids and axis labels for image
    ax.grid(False)
    ax.set_xticks([])
    ax.set_yticks([])
    display_img(url, ax, fig)

    plt.show()

```

In [106]:

```

# vocab = stores all the words that are there in google w2v model
# vocab = model.wv.vocab.keys() # if you are using Google word2Vec

vocab = model.keys()
# this function will add the vectors of each word and returns the avg vector of given sentance
def build_avg_vec(sentence, num_features, doc_id, m_name):
    # sentace: its title of the apparel
    # num_features: the lenght of word2vec vector, its values = 300
    # m_name: model information it will take two values
        # if m_name == 'avg', we will append the model[i], w2v representation of word i
        # if m_name == 'weighted', we will multiply each w2v[word] with the idf(word)

    featureVec = np.zeros((num_features,), dtype="float32")
    # we will intialize a vector of size 300 with all zeros
    # we add each word2vec(wordi) to this fatureVec
    nwords = 0

    for word in sentence.split():
        nwords += 1
        if word in vocab:
            if m_name == 'weighted' and word in idf_title_vectorizer.vocabulary_:
                featureVec = np.add(featureVec, idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[word]] * model[word])
            elif m_name == 'avg':

```

```
        featureVec = np.add(featureVec, model[word])
if(nwords>0):
    featureVec = np.divide(featureVec, nwords)
# returns the avg vector of given sentence, its of shape (1, 300)
return featureVec
```

### [9.2] Average Word2Vec product similarity.

In [107]:

```
doc_id = 0
w2v_title = []
# for every title we build a avg vector representation
for i in data['title']:
    w2v_title.append(build_avg_vec(i, 300, doc_id,'avg'))
    doc_id += 1

# w2v_title = np.array(# number of doc in courpus * 300), each row corresponds to a doc
w2v_title = np.array(w2v_title)
```

In [108]:

```

def avg_w2v_model(doc_id, num_results):
    # doc_id: apparel's id in given corpus

    # dist(x, y) = sqrt(dot(x, x) - 2 * dot(x, y) + dot(y, y))
    pairwise_dist = pairwise_distances(w2v_title, w2v_title[doc_id].reshape(1,-1))

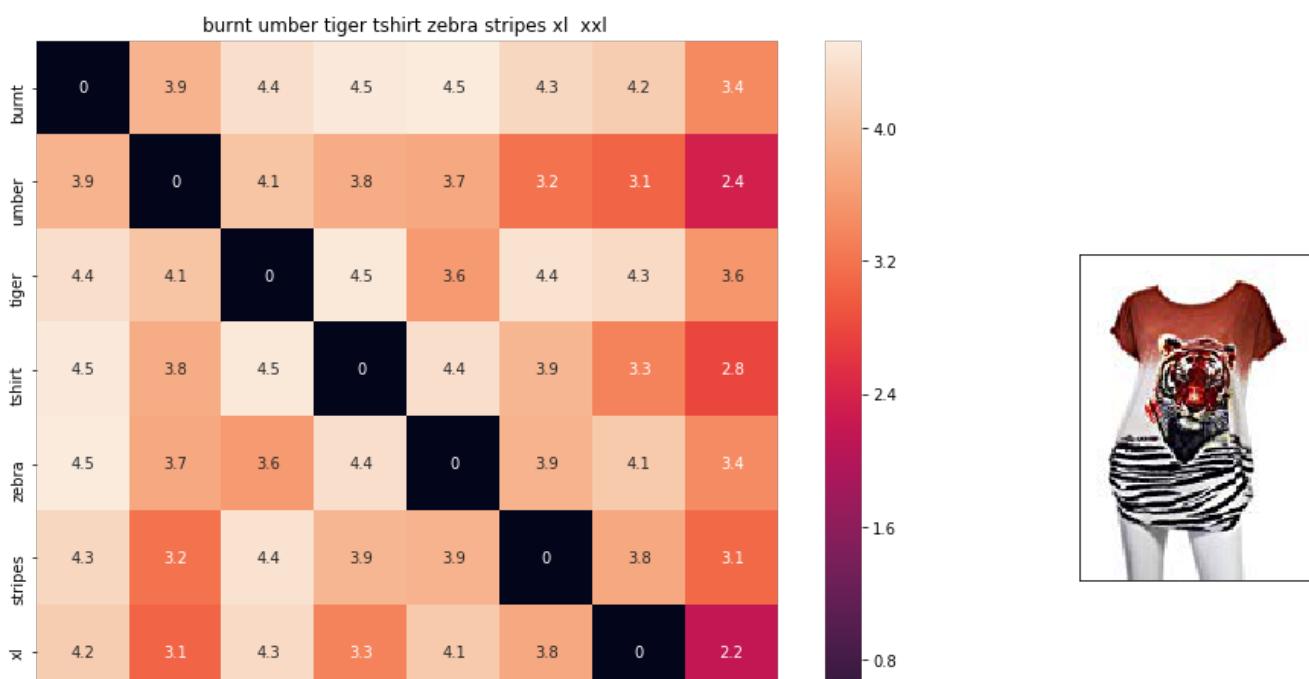
    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

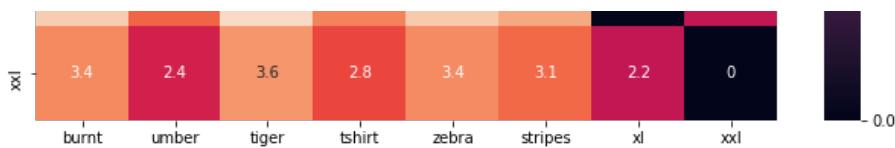
    #data frame indices of the 9 smallest distance's
    df_indices = list(data.index[indices])

    for i in range(0, len(indices)):
        heat_map_w2v(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data['medium_image_url'].loc[df_indices[i]], indices[0], indices[i], 'avg')
        print('ASIN :', data['asin'].loc[df_indices[i]])
        print('BRAND :', data['brand'].loc[df_indices[i]])
        print ('euclidean distance from given input image :', pdists[i])
        print('='*125)

avg_w2v_model(12566, 20)
# in the give heat map, each cell contains the euclidean distance between words i, j

```





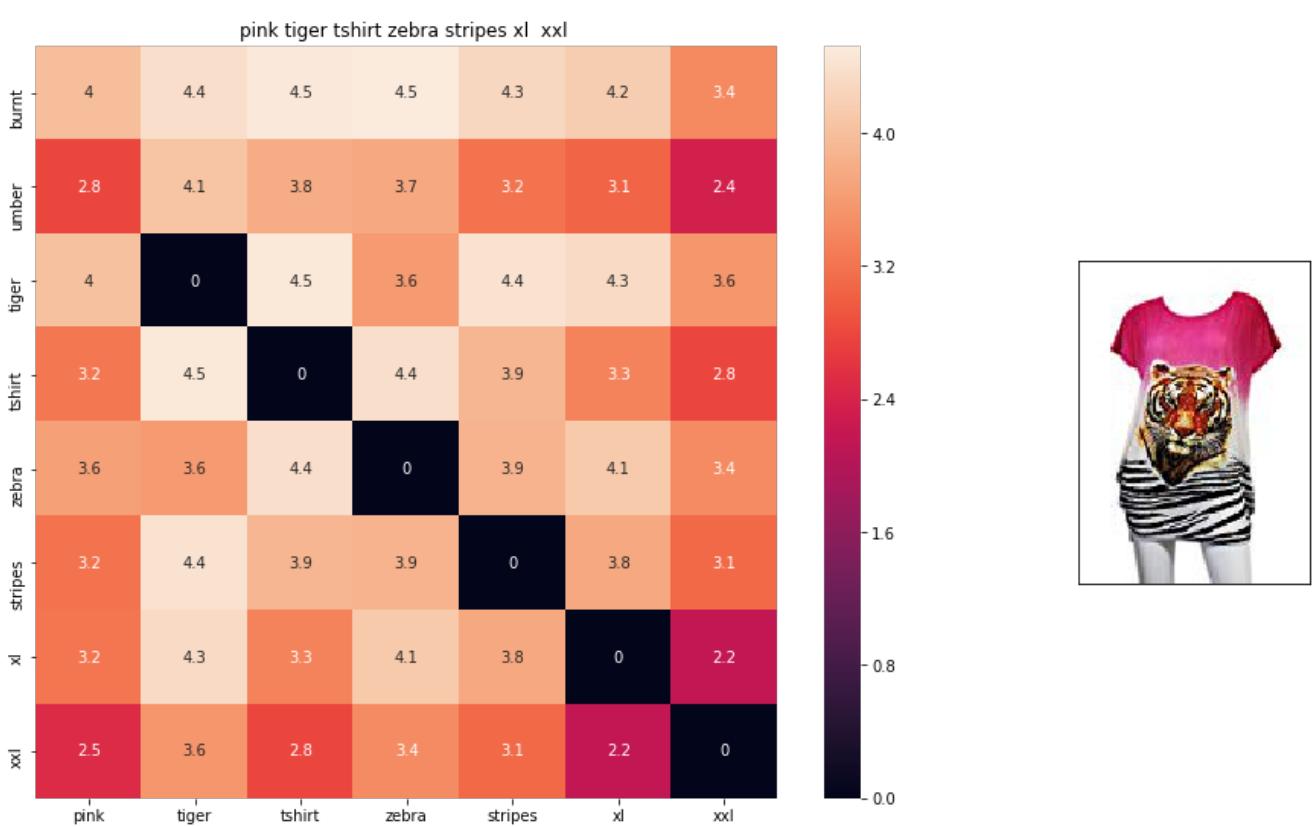
ASIN : B00JXQB5FQ

BRAND : Si Row

euclidean distance from given input image : 0.00069053395

-----

-----



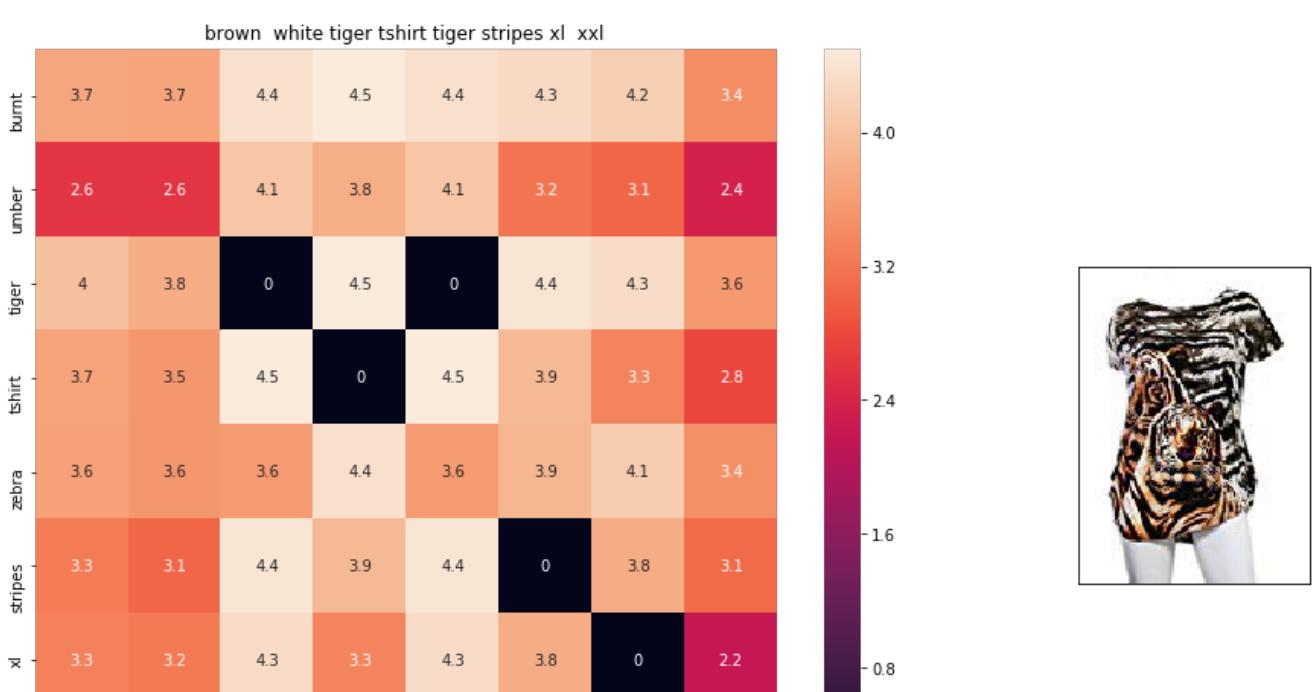
ASIN : B00JXQASS6

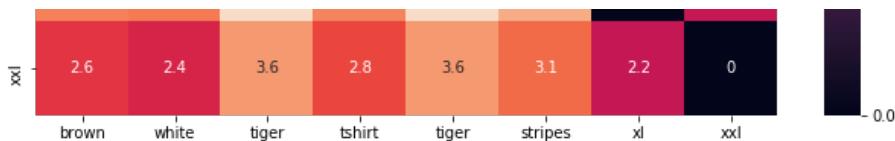
BRAND : Si Row

euclidean distance from given input image : 0.5891932

=====

=====





ASIN : B00JXQCWT0

BRAND : Si Row

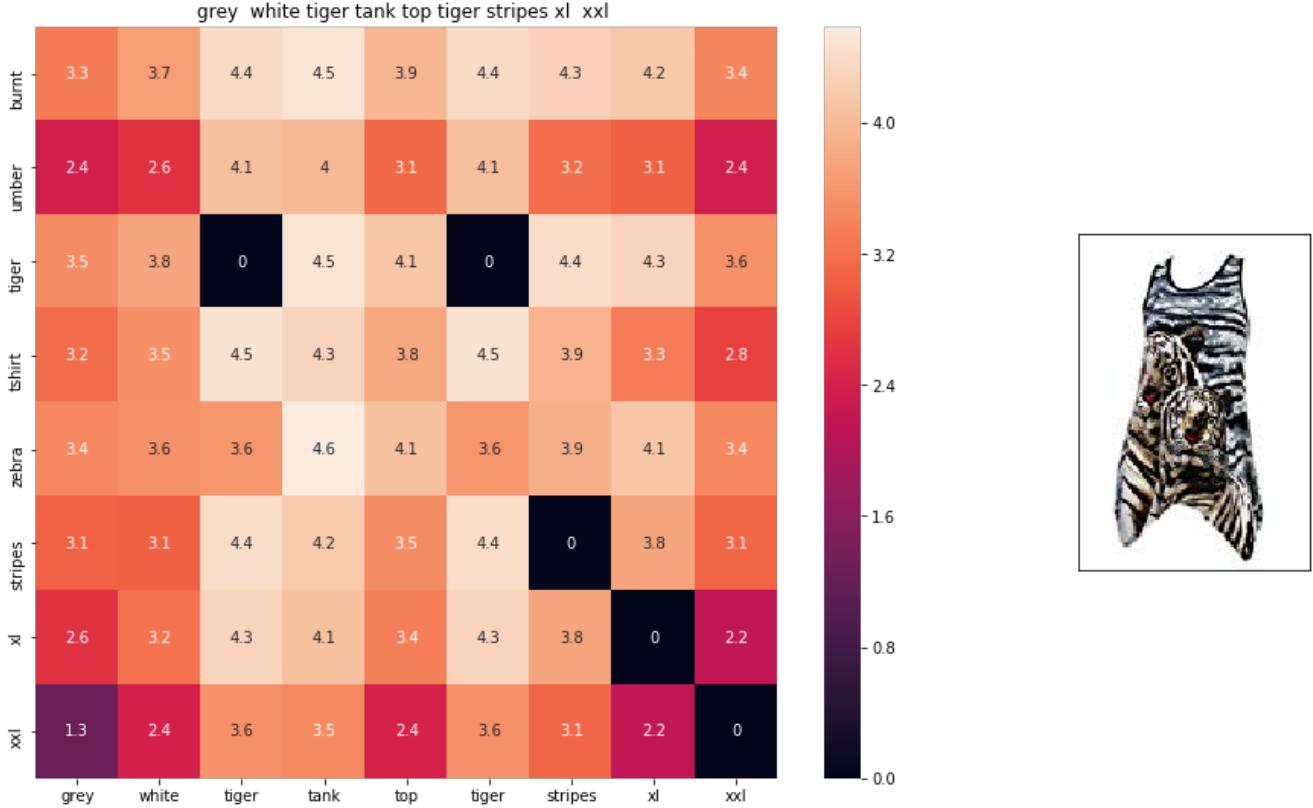
euclidean distance from given input image : 0.7003439

=====

=====

=====

=====



ASIN : B00JXQAFZ2

BRAND : Si Row

euclidean distance from given input image : 0.8928398

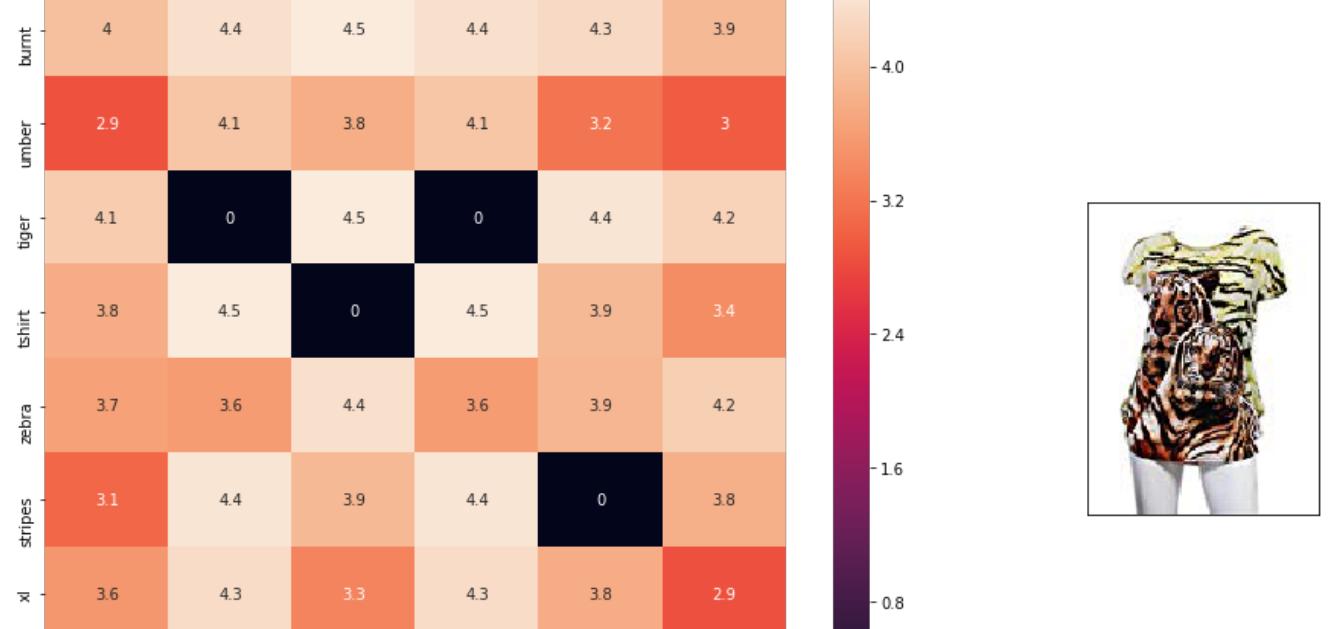
=====

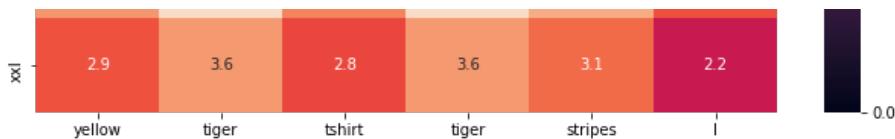
=====

=====

=====

=====





ASIN : B00JXQCUIC

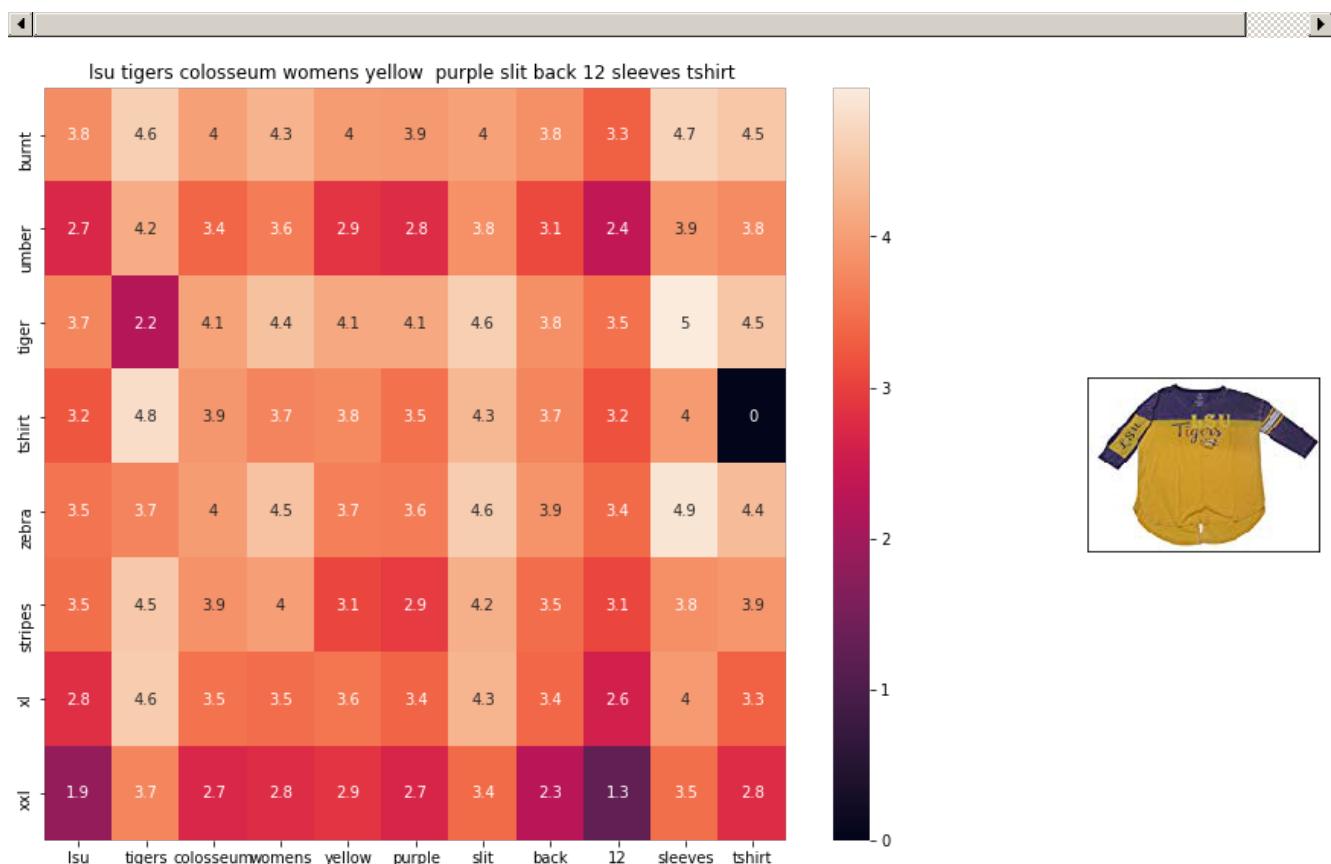
BRAND : Si Row

euclidean distance from given input image : 0.95601267

---



---



ASIN : B073R5Q8HD

BRAND : Colosseum

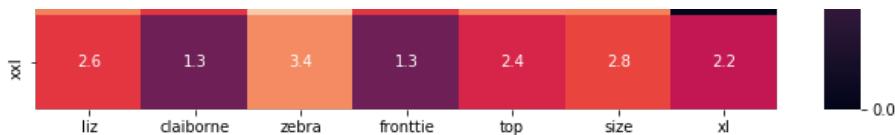
euclidean distance from given input image : 1.0229691

---



---





ASIN : B06XBY5QXL

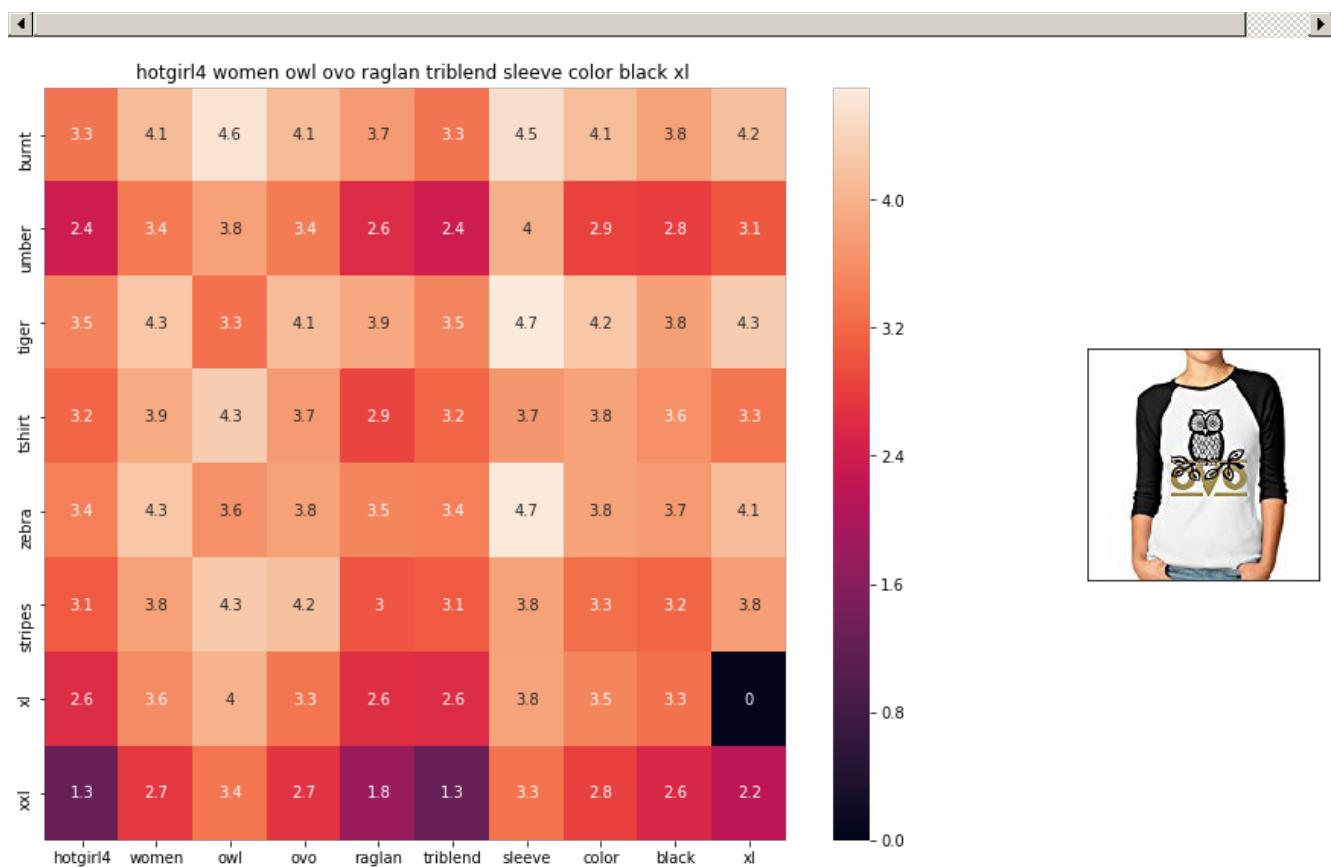
BRAND : Liz Claiborne

euclidean distance from given input image : 1.0669324

---



---



ASIN : B01L8L73M2

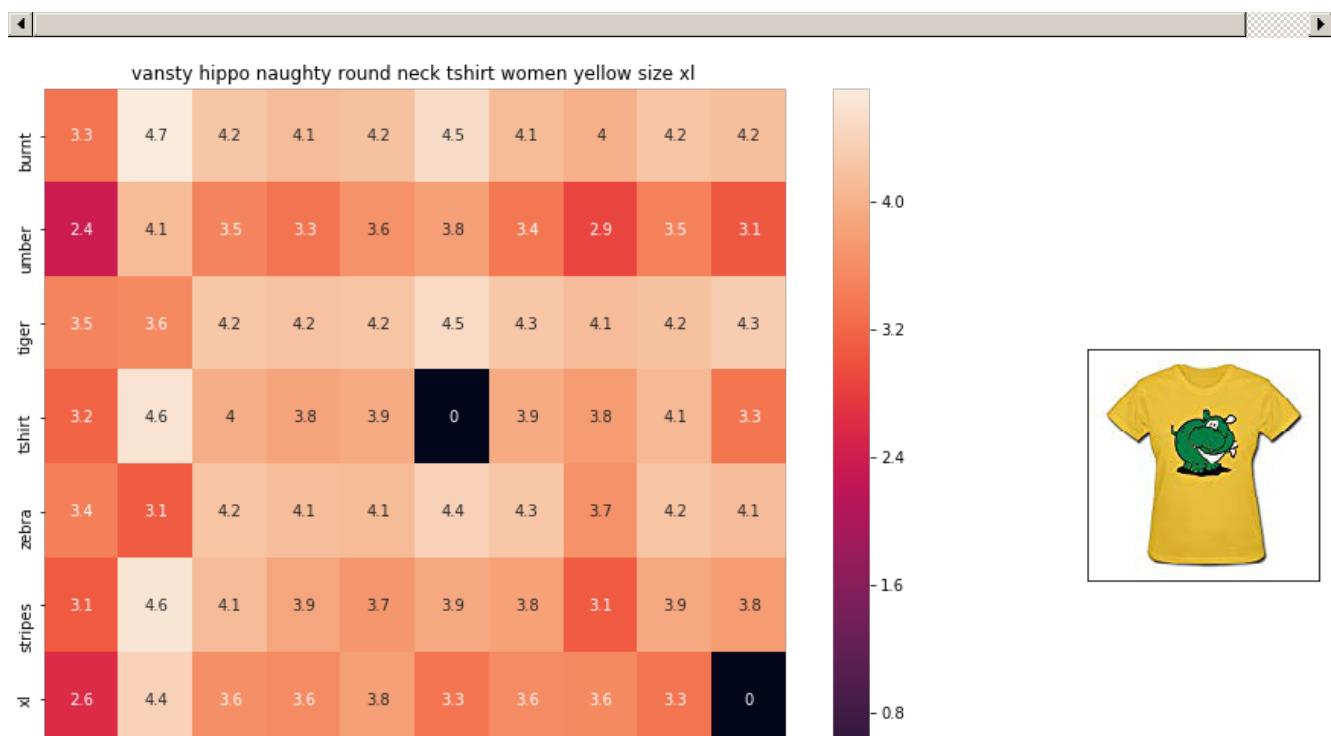
BRAND : Hotgirl4 Raglan Design

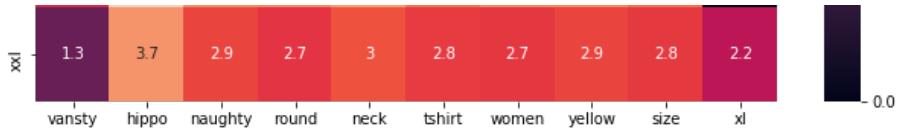
euclidean distance from given input image : 1.0731405

---



---

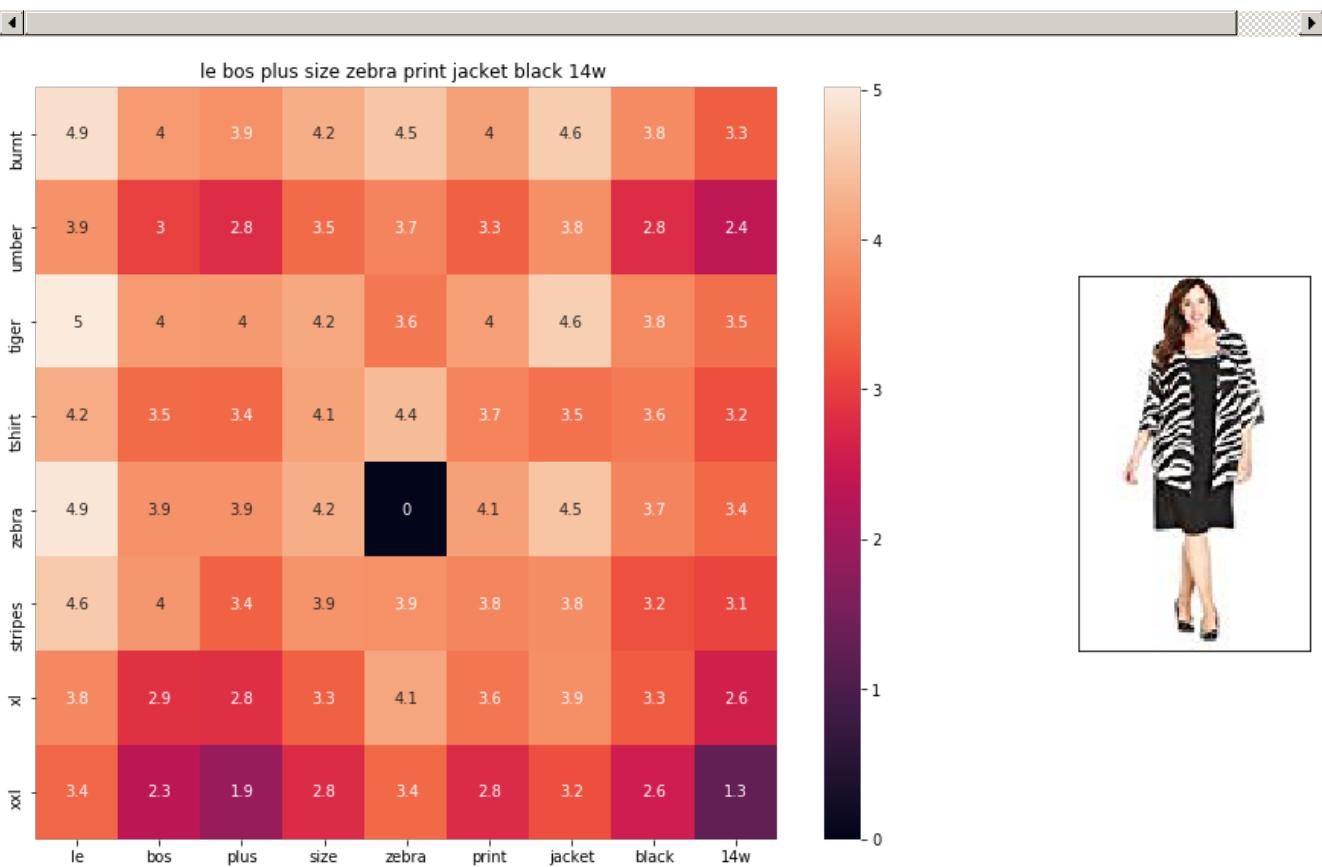




ASIN : B01EJS5H06

BRAND : Vansty

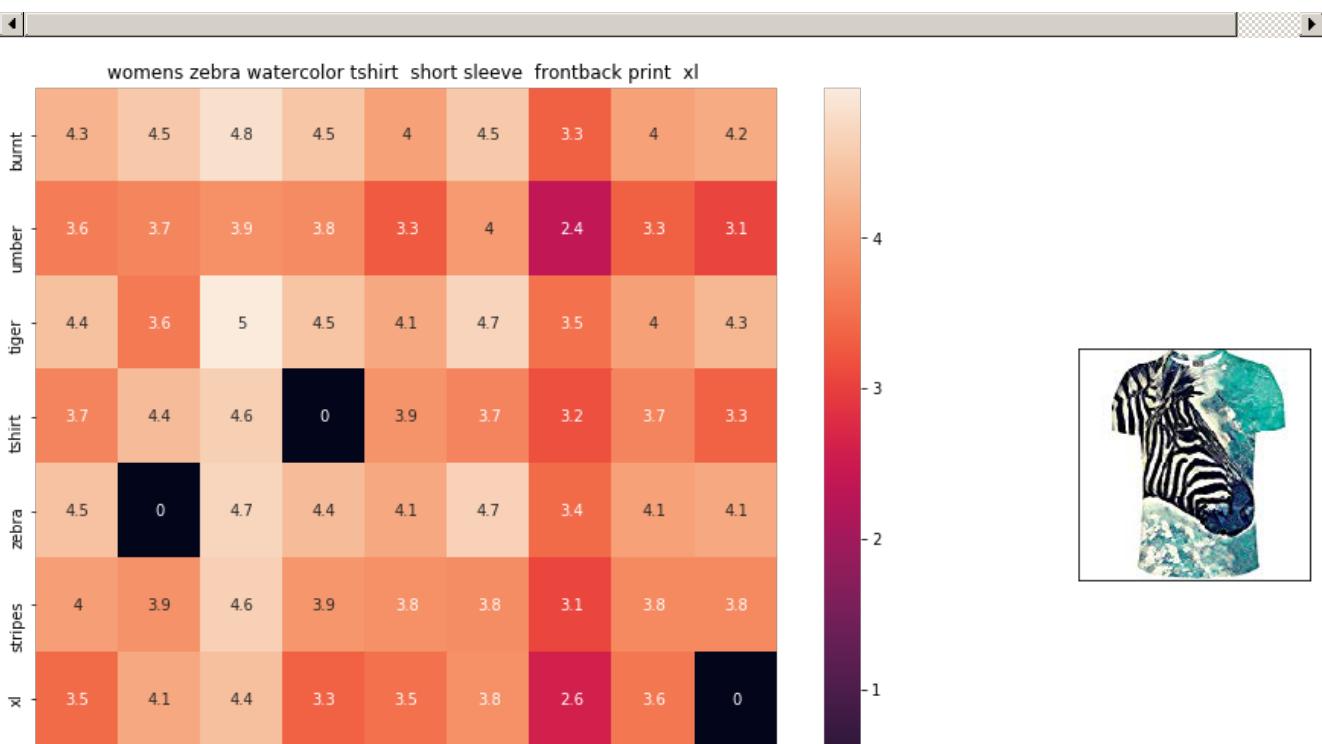
euclidean distance from given input image : 1.075719

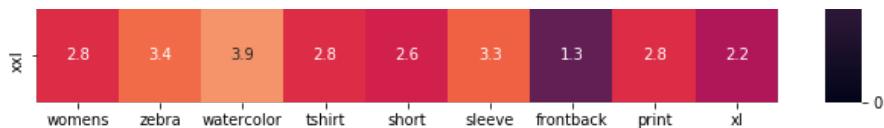


ASIN : B01B01XRK8

BRAND : Le Bos

euclidean distance from given input image : 1.0839964





ASIN : B072R2JXKW

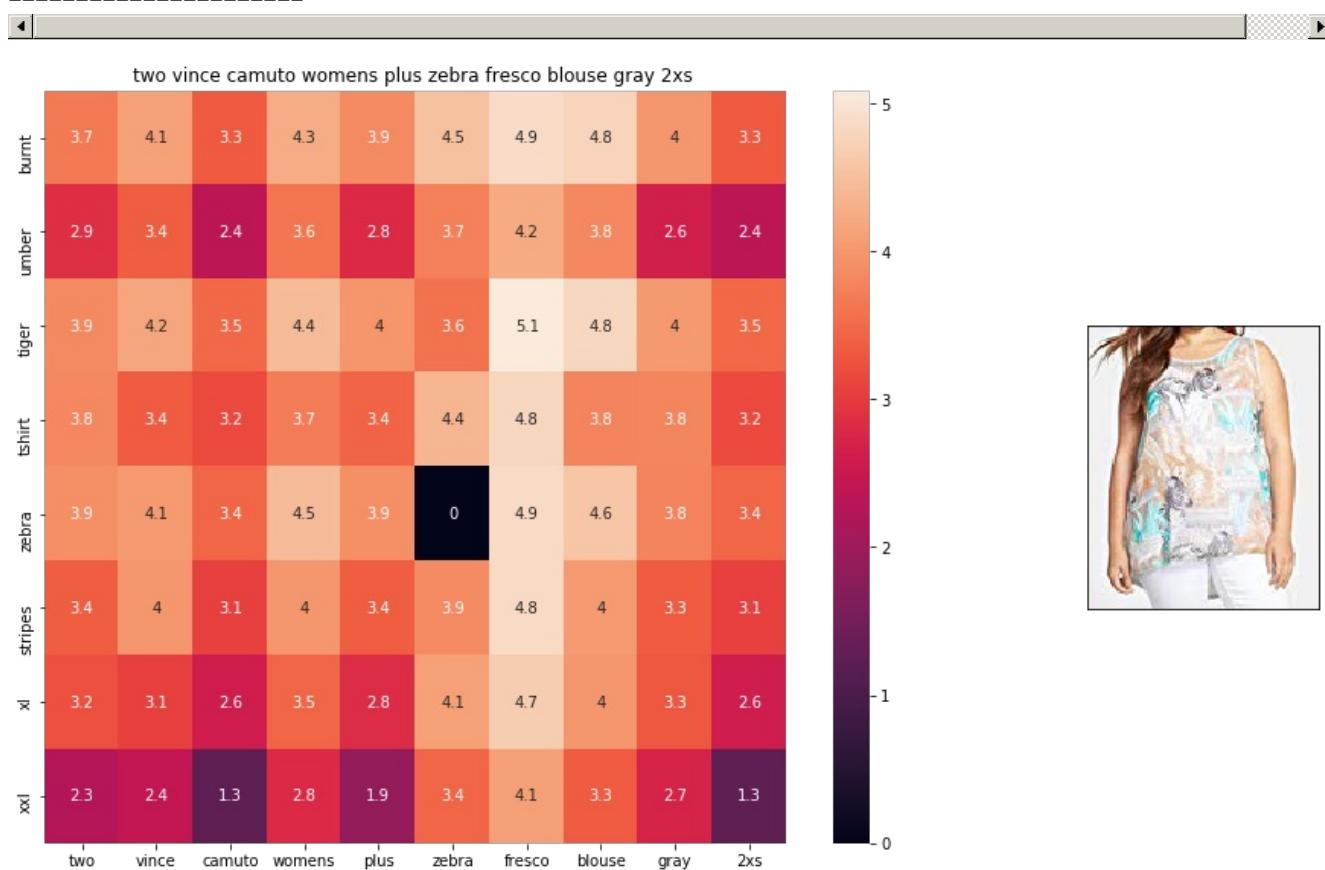
BRAND : WHAT ON EARTH

euclidean distance from given input image : 1.0842218

---



---



ASIN : B074MJRGW6

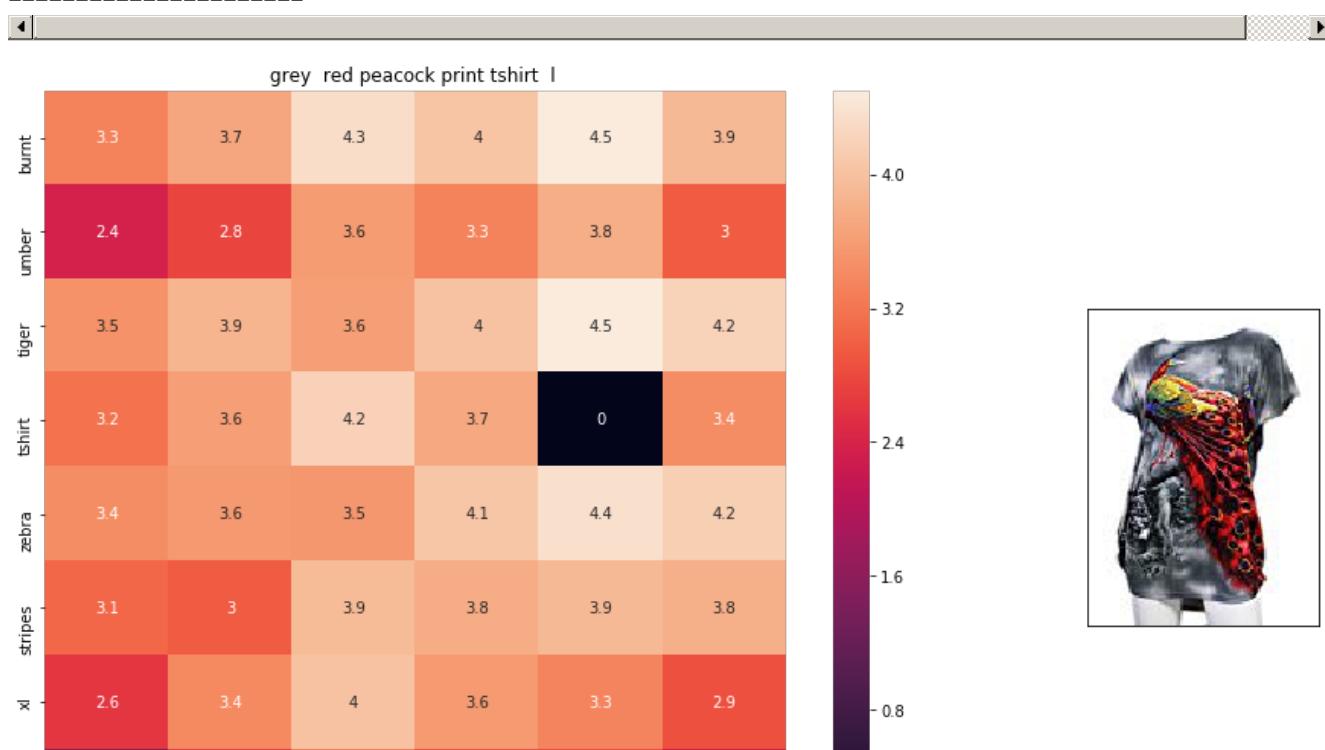
BRAND : Two by Vince Camuto

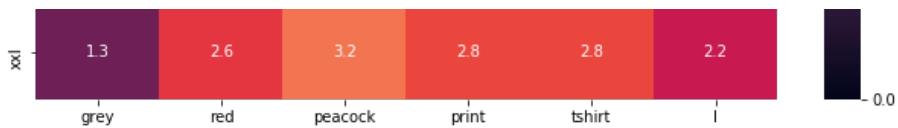
euclidean distance from given input image : 1.0895038

---



---

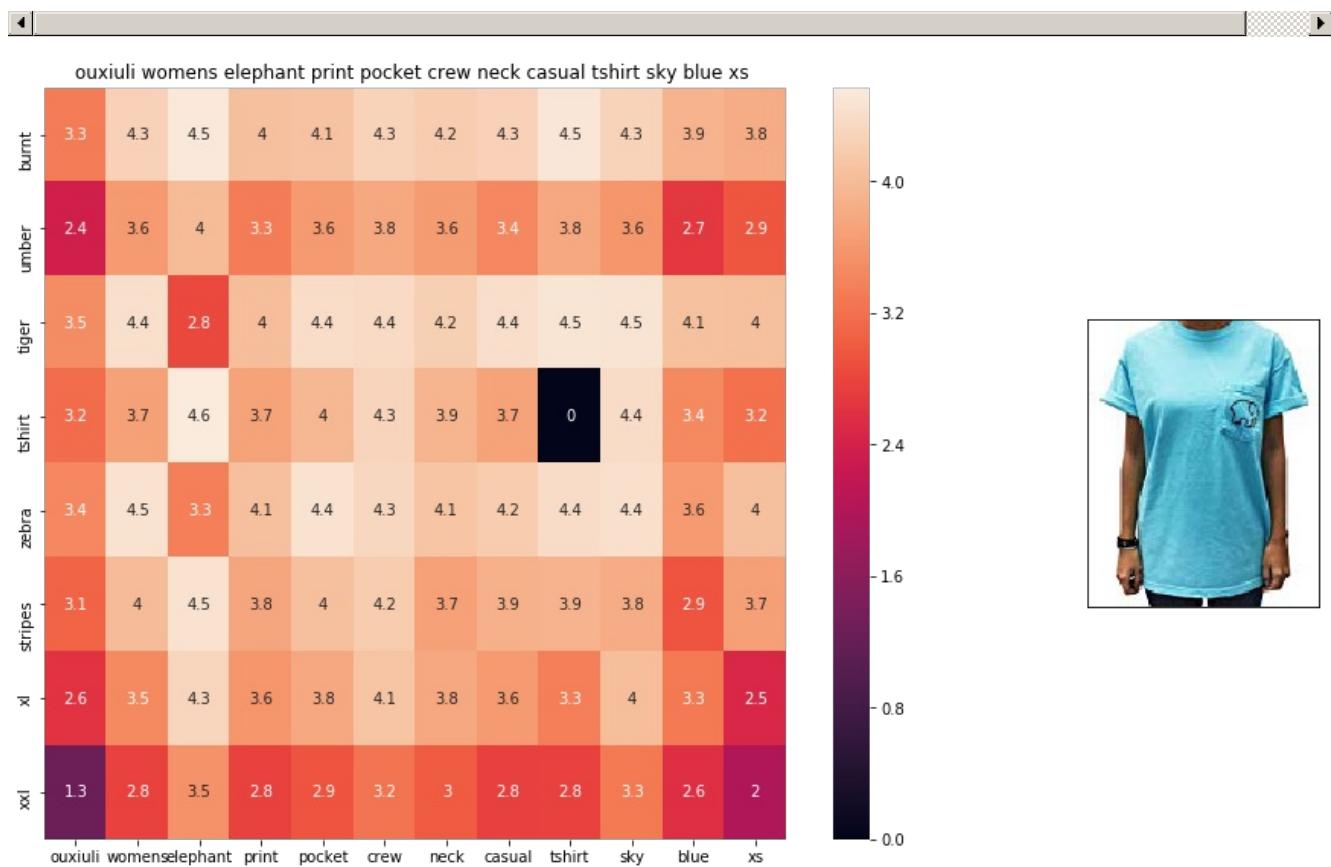




ASIN : B00JXQCFRS

BRAND : Si Row

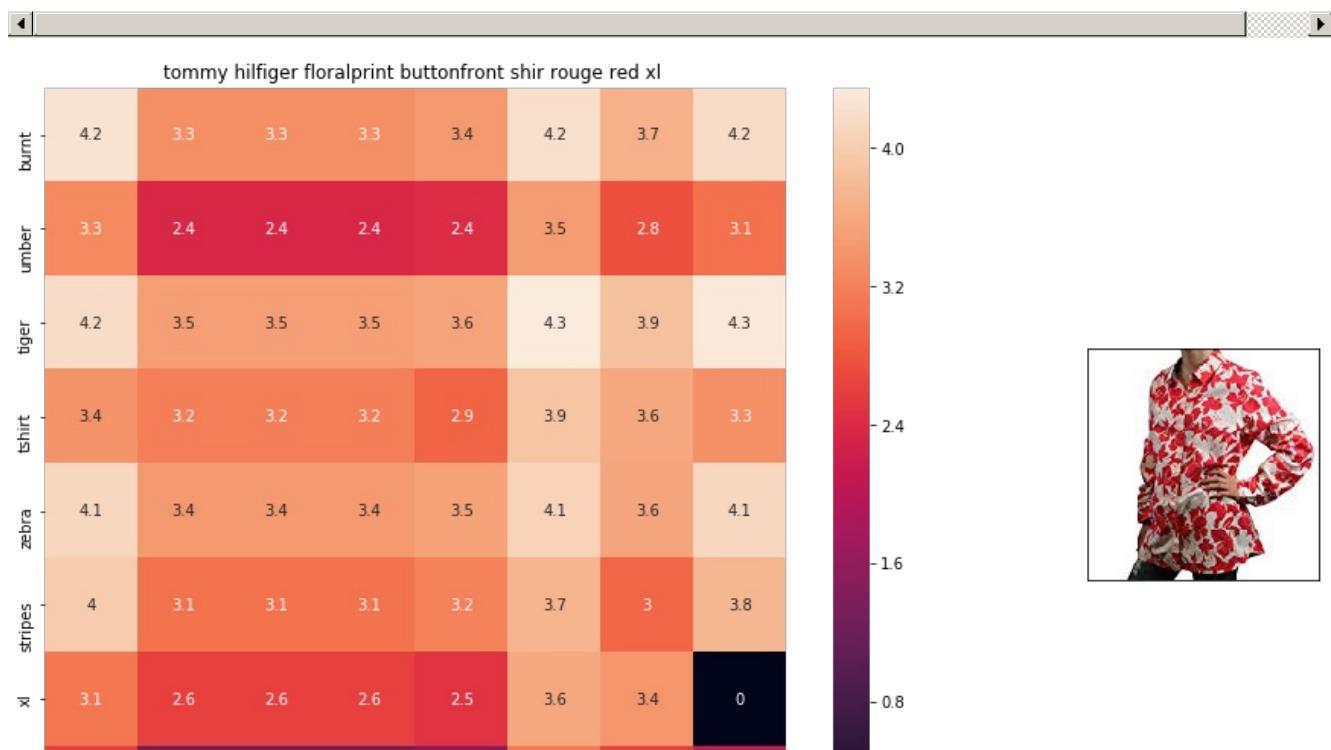
euclidean distance from given input image : 1.0900588

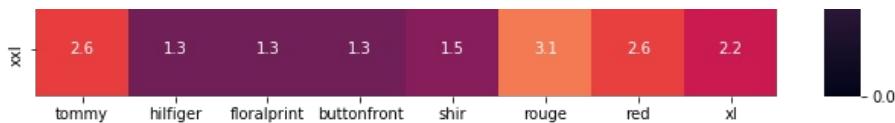


ASIN : B01I53HU6K

BRAND : ouxiuli

euclidean distance from given input image : 1.0920112





ASIN : B0711NGTQM

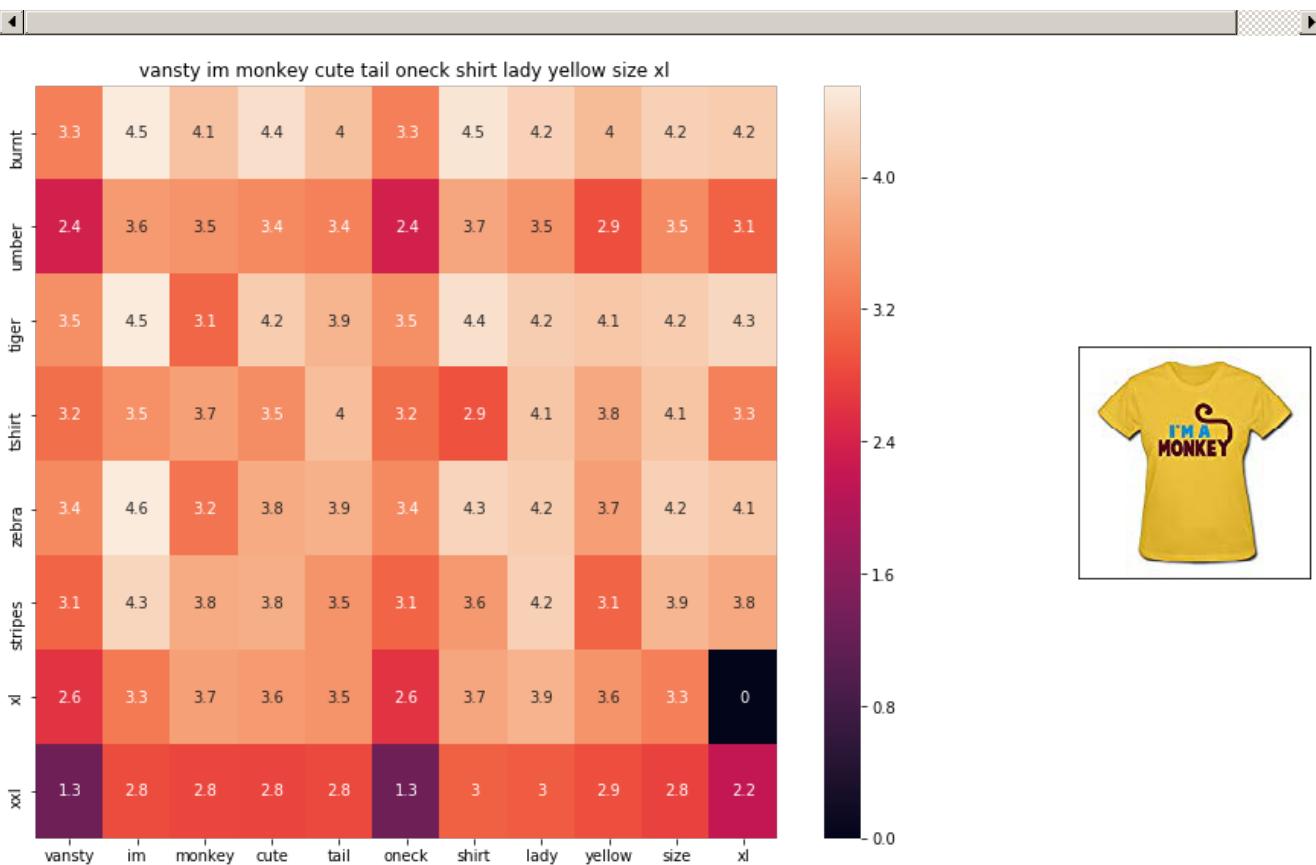
BRAND : THILFIGER RTW

euclidean distance from given input image : 1.0923418

---



---



ASIN : B01EFSLO8Y

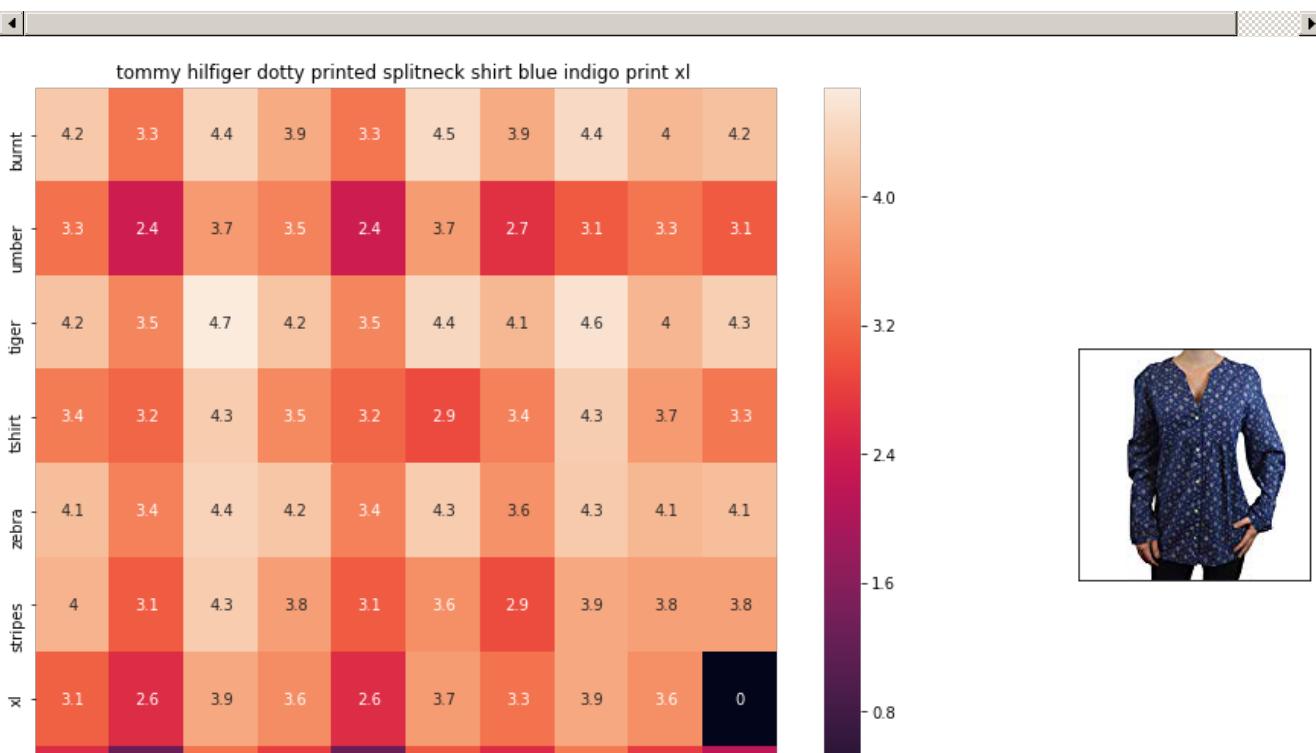
BRAND : Vansty

euclidean distance from given input image : 1.0934006

---



---





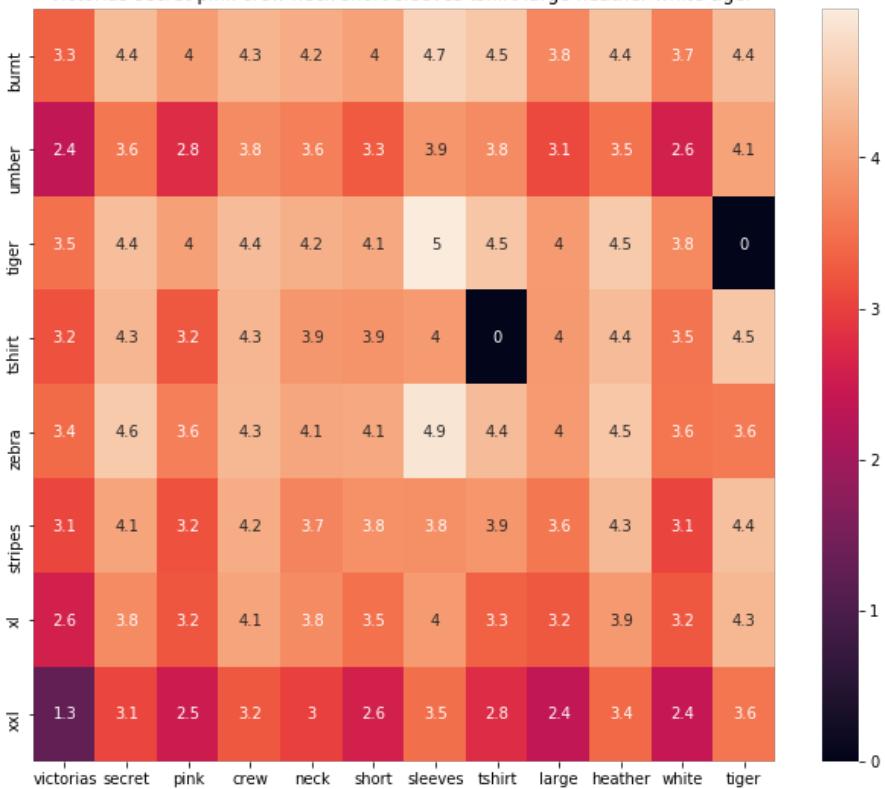
ASIN : B0716TVWQ4

BRAND : THILFIGER RTW

euclidean distance from given input image : 1.0942026

=====  
=====

victorias secret pink crew neck short sleeves tshirt large heather white tiger



ASIN : B0716MVPGV

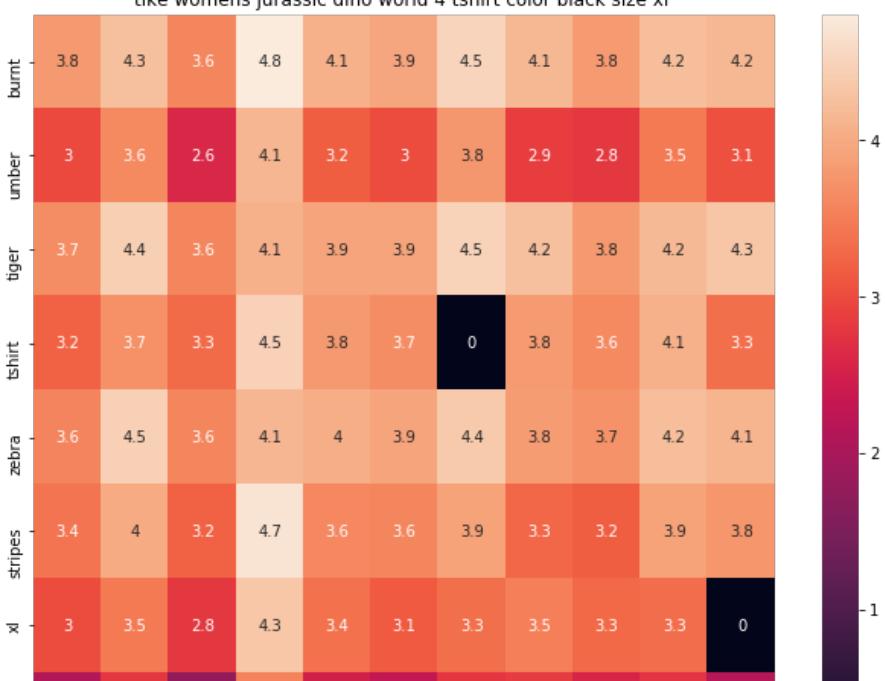
BRAND : V.Secret

euclidean distance from given input image : 1.0948305

=====  
=====

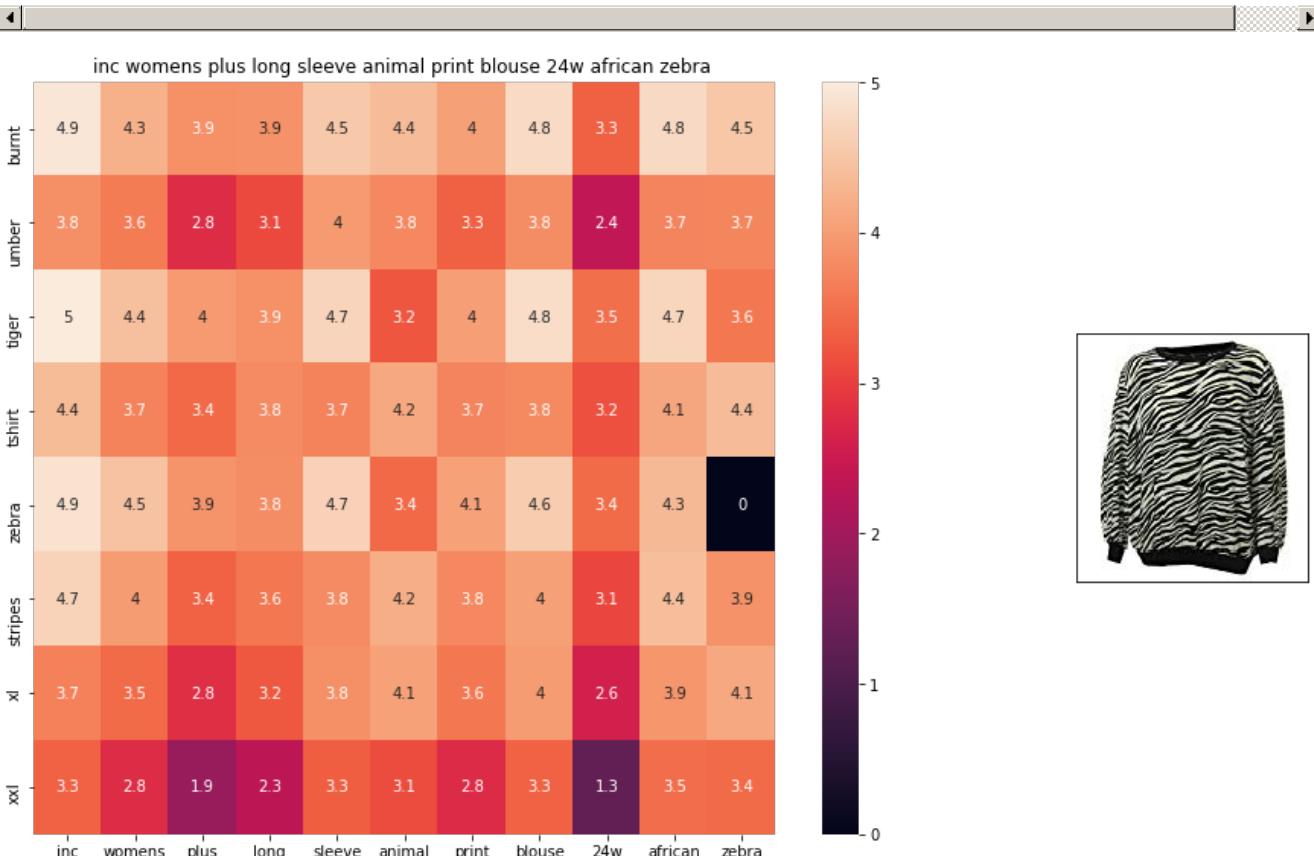
◀

tike womens jurassic dino world 4 tshirt color black size xl





ASIN : B016OPN40I  
 BRAND : TIKE Fashions  
 euclidean distance from given input image : 1.0951275



ASIN : B018WDJCUA  
 BRAND : INC - International Concepts Woman  
 euclidean distance from given input image : 1.0966892

## [9.4] IDF weighted Word2Vec for product similarity

In [109]:

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

In [110]:

```
def weighted_w2v_model(doc_id, num_results):
    # doc_id: apparel's id in given corpus

    # pairwise_dist will store the distance from given input apparel to all remaining apparels
    # the metric we used here is cosine, the coside distance is mesured as K(X, Y) = <X, Y> / (||X|| * ||Y||)
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
    pairwise_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1,-1))
```

```

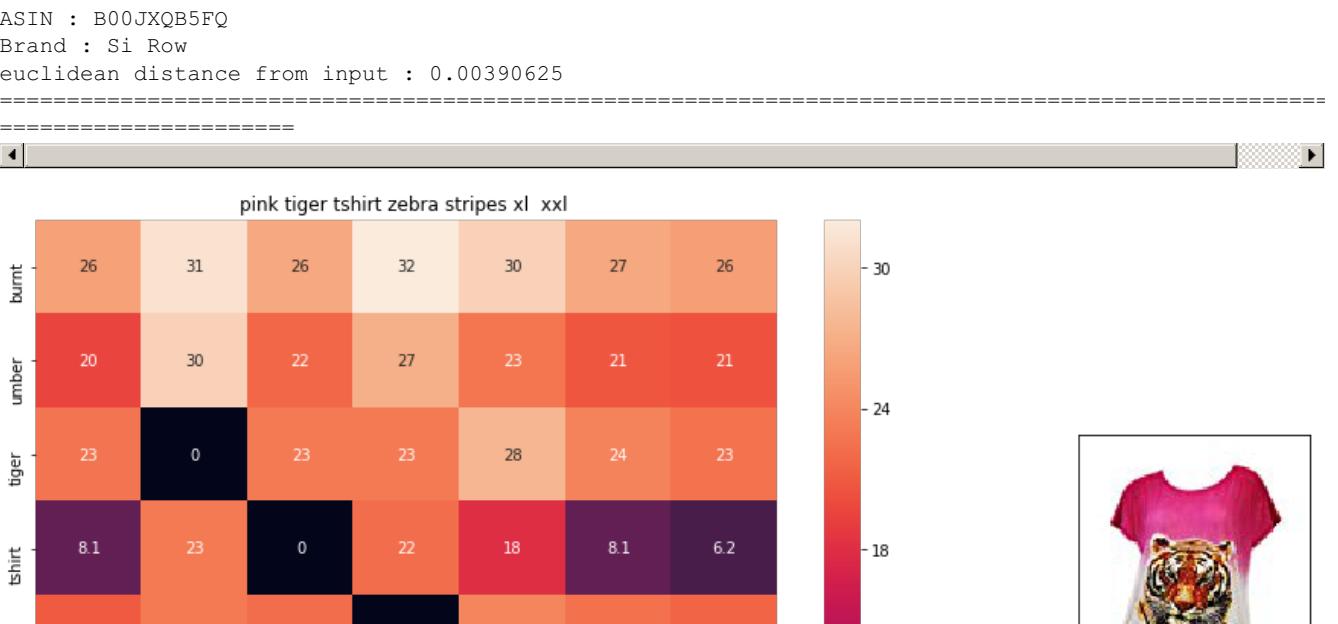
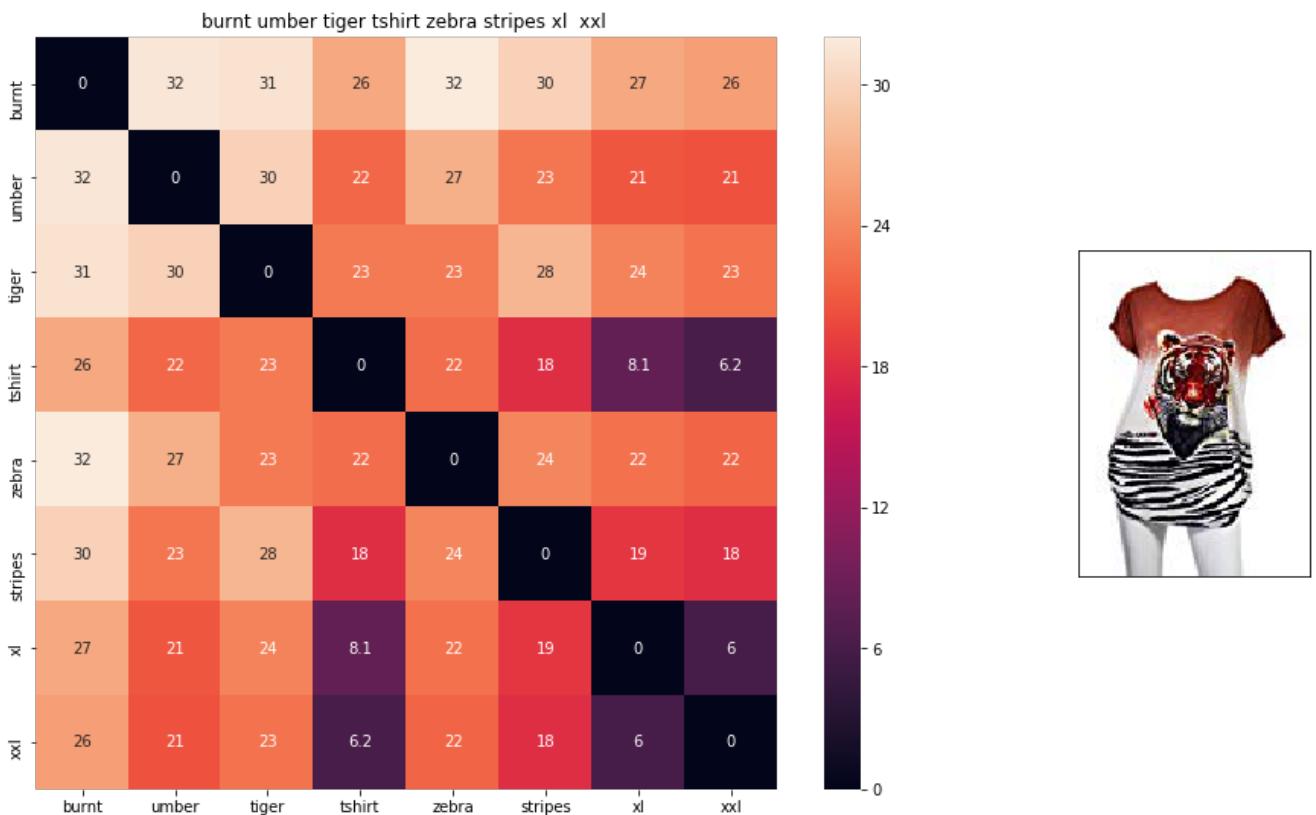
# np.argsort will return indices of 9 smallest distances
indices = np.argsort(pairwise_dist.flatten())[0:num_results]
#pdists will store the 9 smallest distances
pdists = np.sort(pairwise_dist.flatten())[0:num_results]

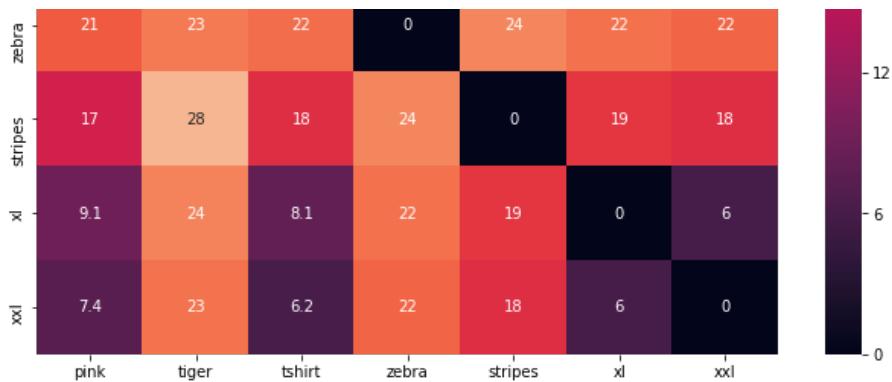
#data frame indices of the 9 smallest distance's
df_indices = list(data.index[indices])

for i in range(0, len(indices)):
    heat_map_w2v(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data['medium_image_url'].loc[df_indices[i]], indices[0], indices[i], 'weighted')
    print('ASIN :', data['asin'].loc[df_indices[i]])
    print('Brand :', data['brand'].loc[df_indices[i]])
    print('euclidean distance from input :', pdists[i])
    print('='*125)

weighted_w2v_model(12566, 20)
#931
#12566
# in the give heat map, each cell contains the euclidean distance between words i, j

```





ASIN : B00JXQASS6

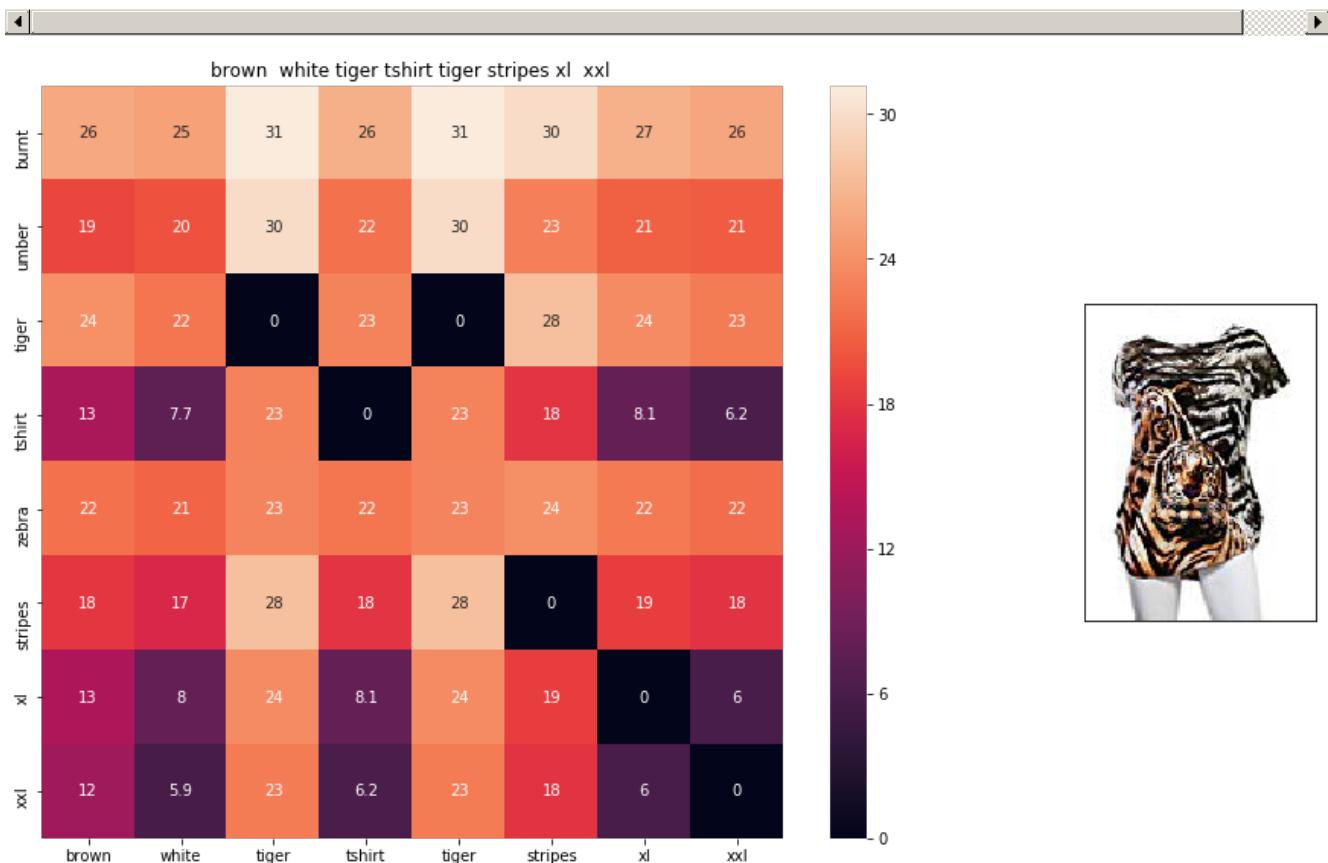
Brand : Si Row

euclidean distance from input : 4.0638885

---



---



ASIN : B00JXQCWTO

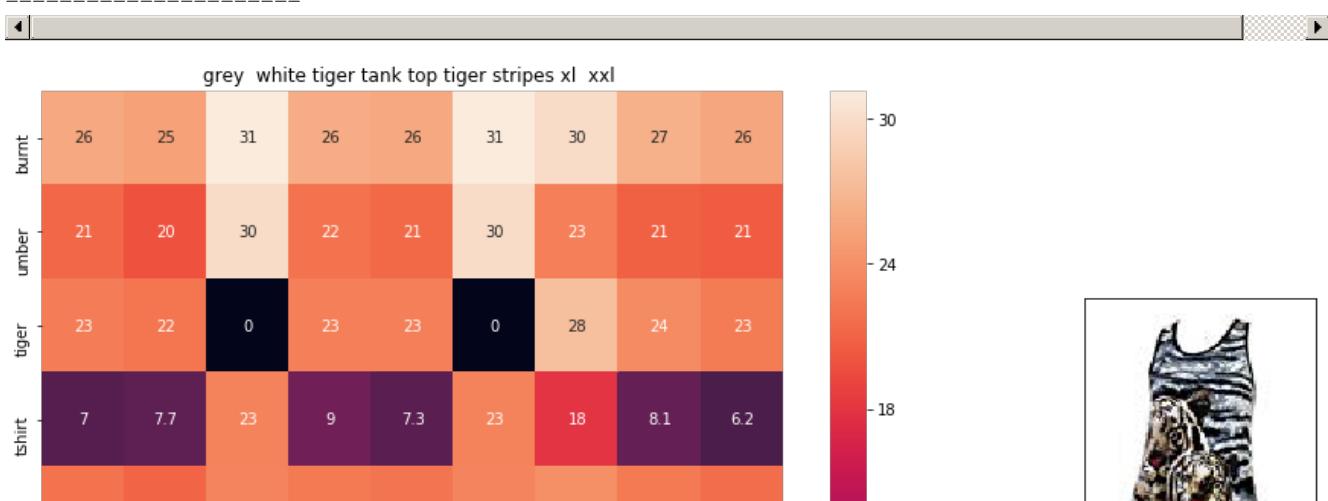
Brand : Si Row

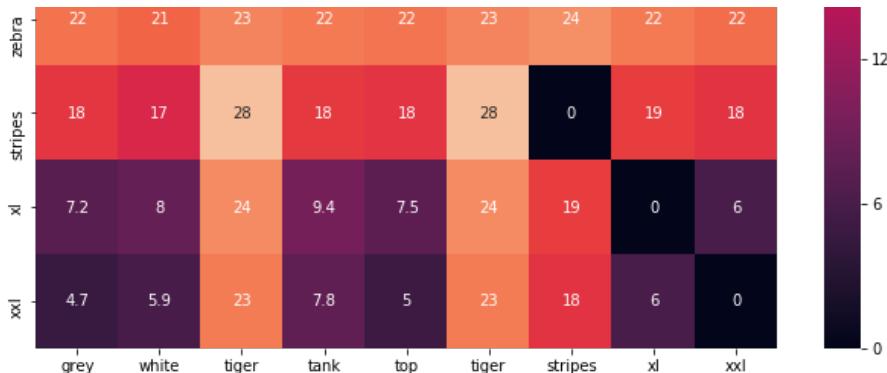
euclidean distance from input : 4.770942

---



---





ASIN : B00JXQAFZ2

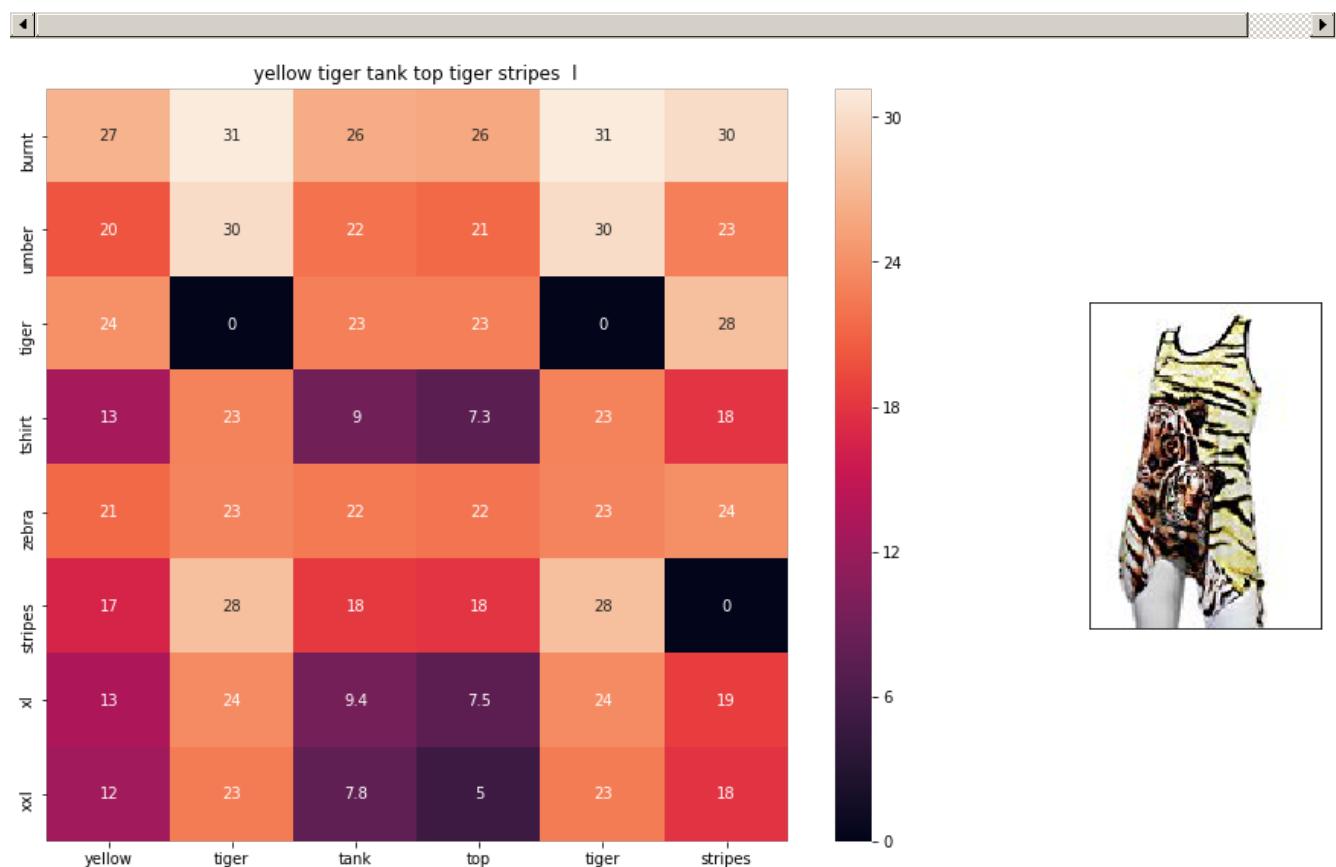
Brand : Si Row

euclidean distance from input : 5.360161

---



---



ASIN : B00JXQAUWA

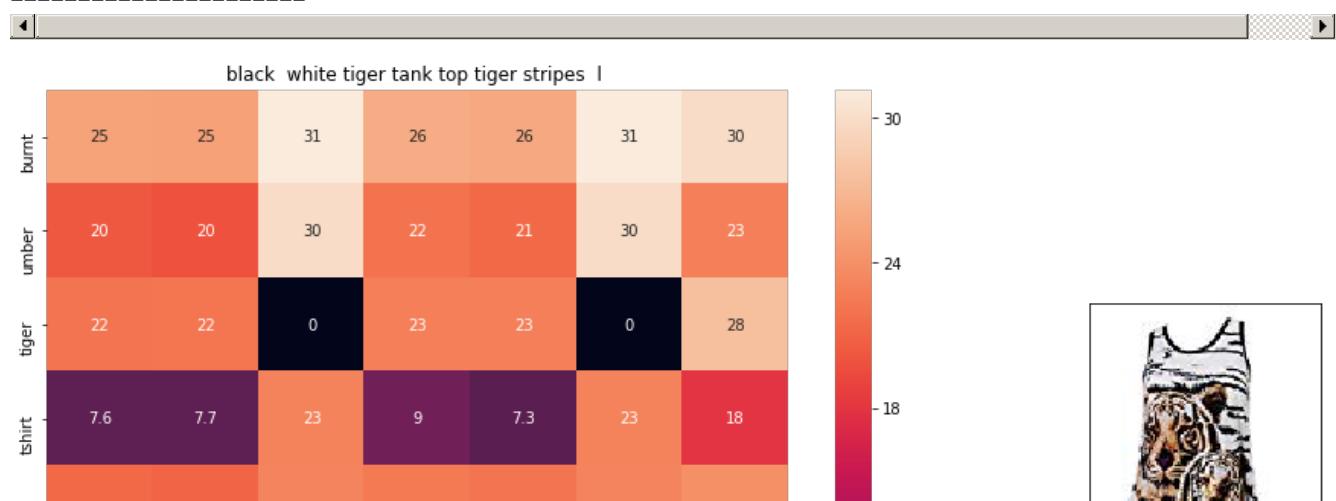
Brand : Si Row

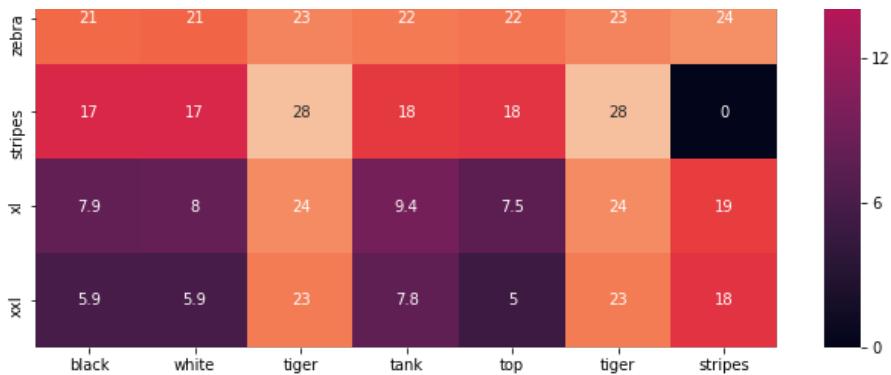
euclidean distance from input : 5.6895237

---



---





ASIN : B00JXQAO94

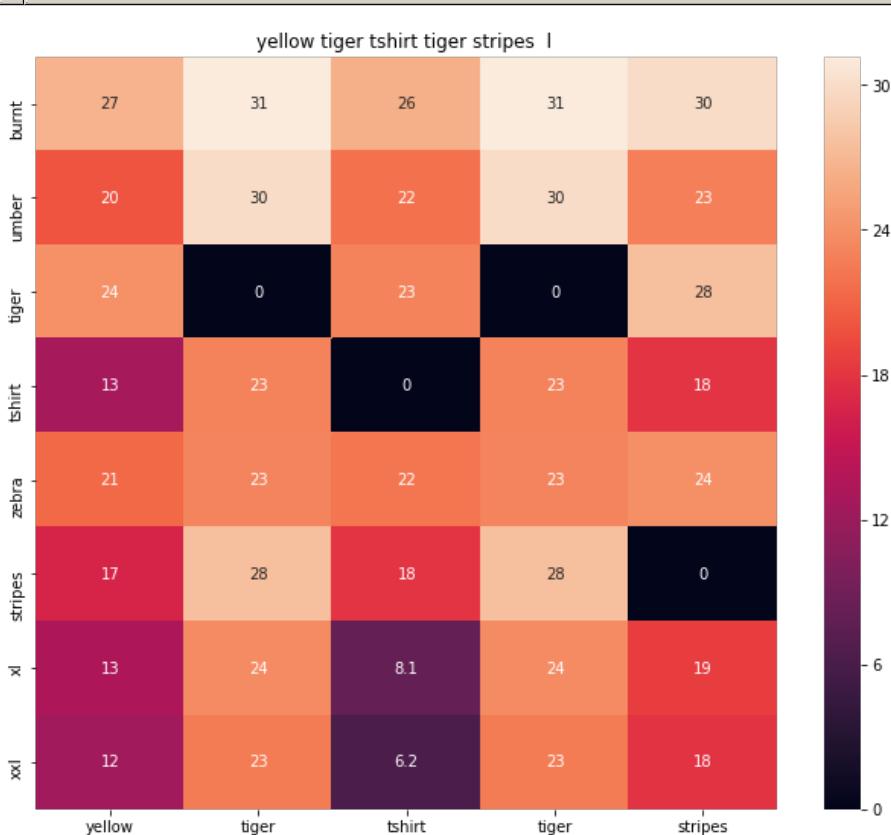
Brand : Si Row

euclidean distance from input : 5.6930213

---



---



ASIN : B00JXQCUIC

Brand : Si Row

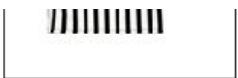
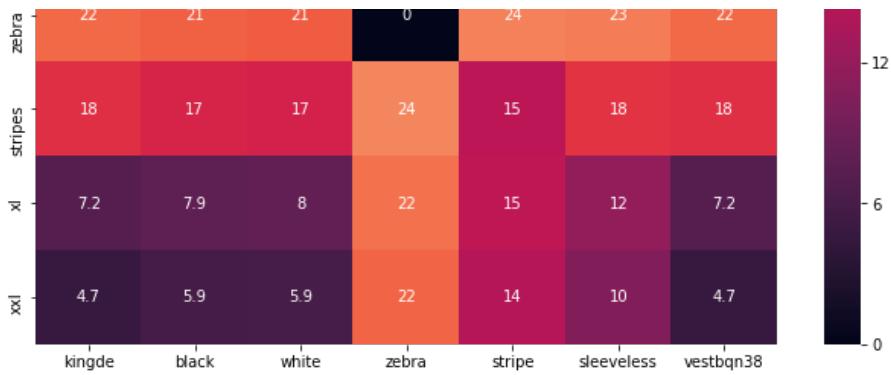
euclidean distance from input : 5.8934426

---



---





ASIN : B015H41F6G

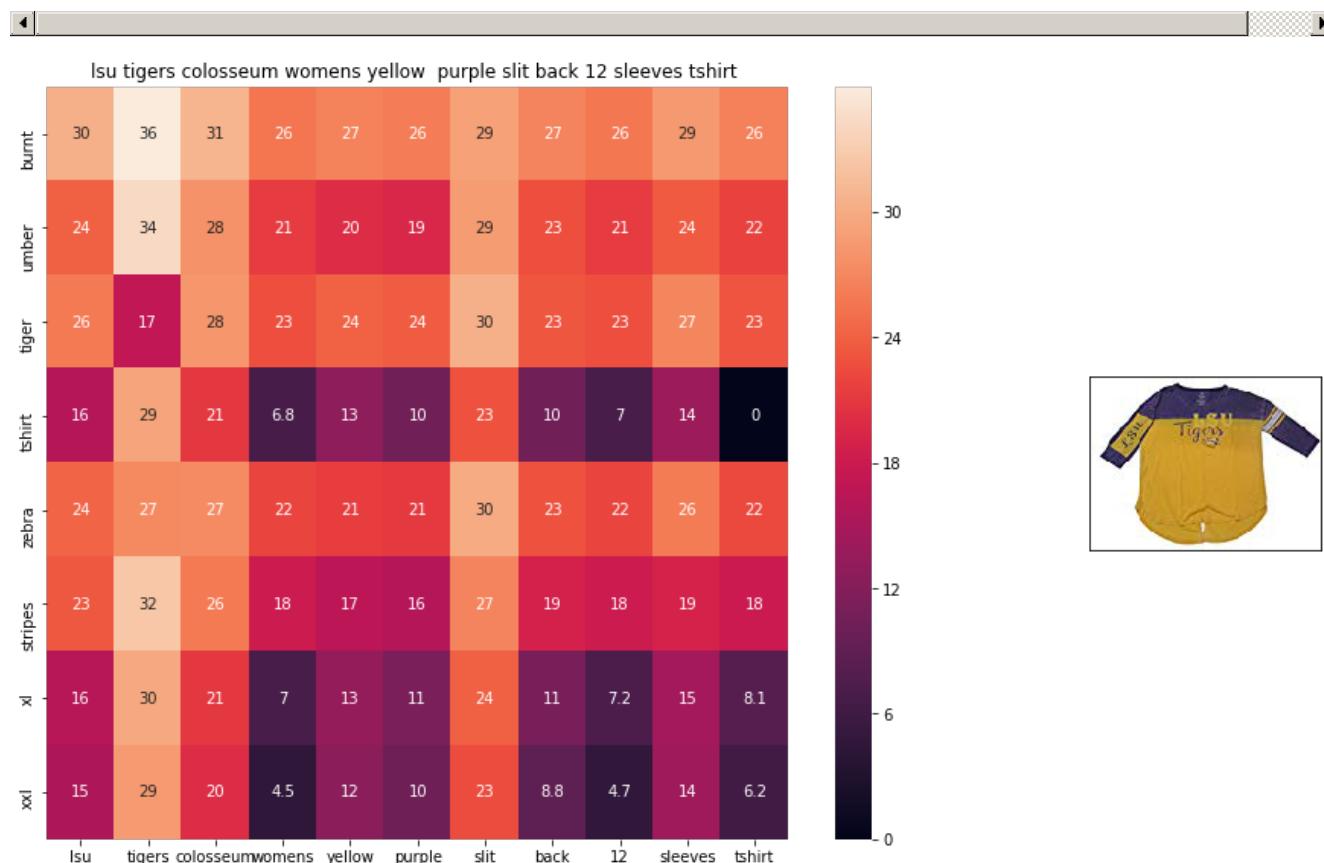
Brand : KINGDE

euclidean distance from input : 6.13299

---



---



ASIN : B073R5Q8HD

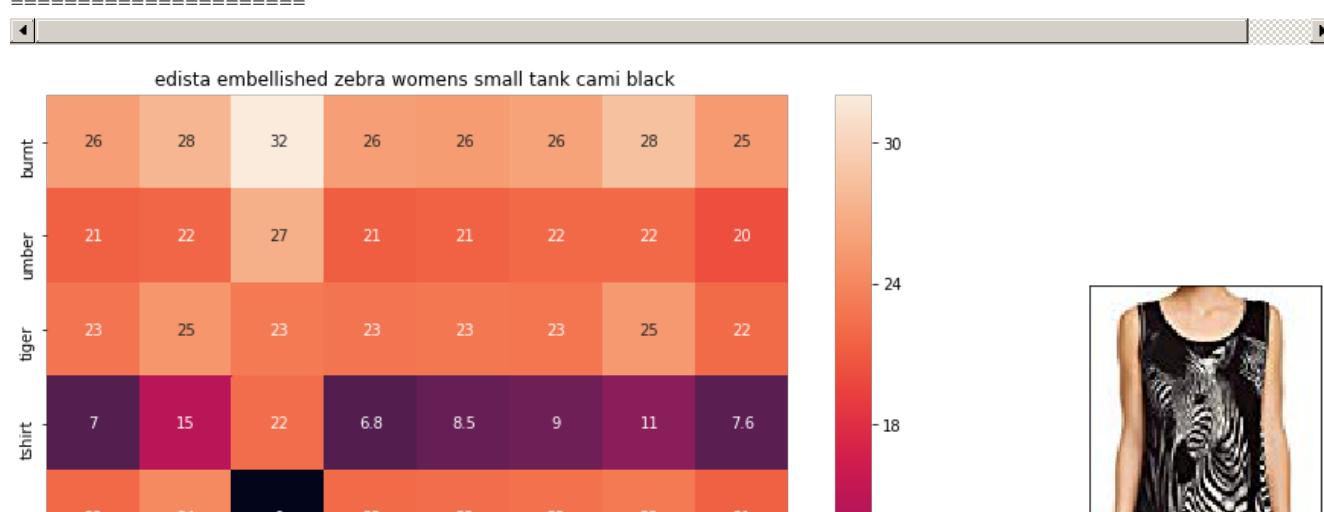
Brand : Colosseum

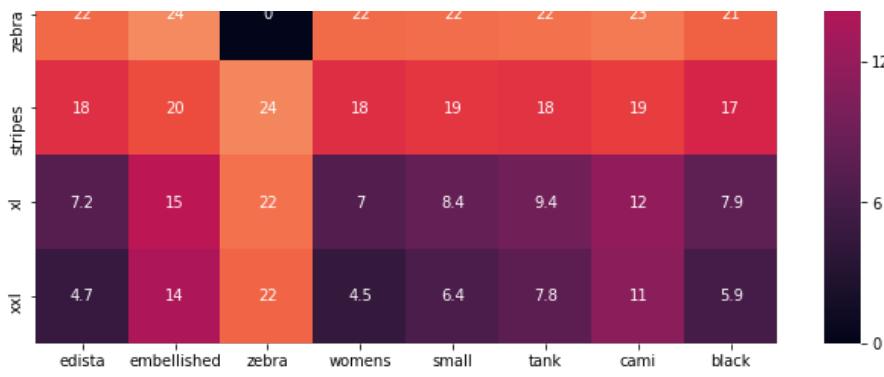
euclidean distance from input : 6.2567067

---



---





ASIN : B074P8MD22

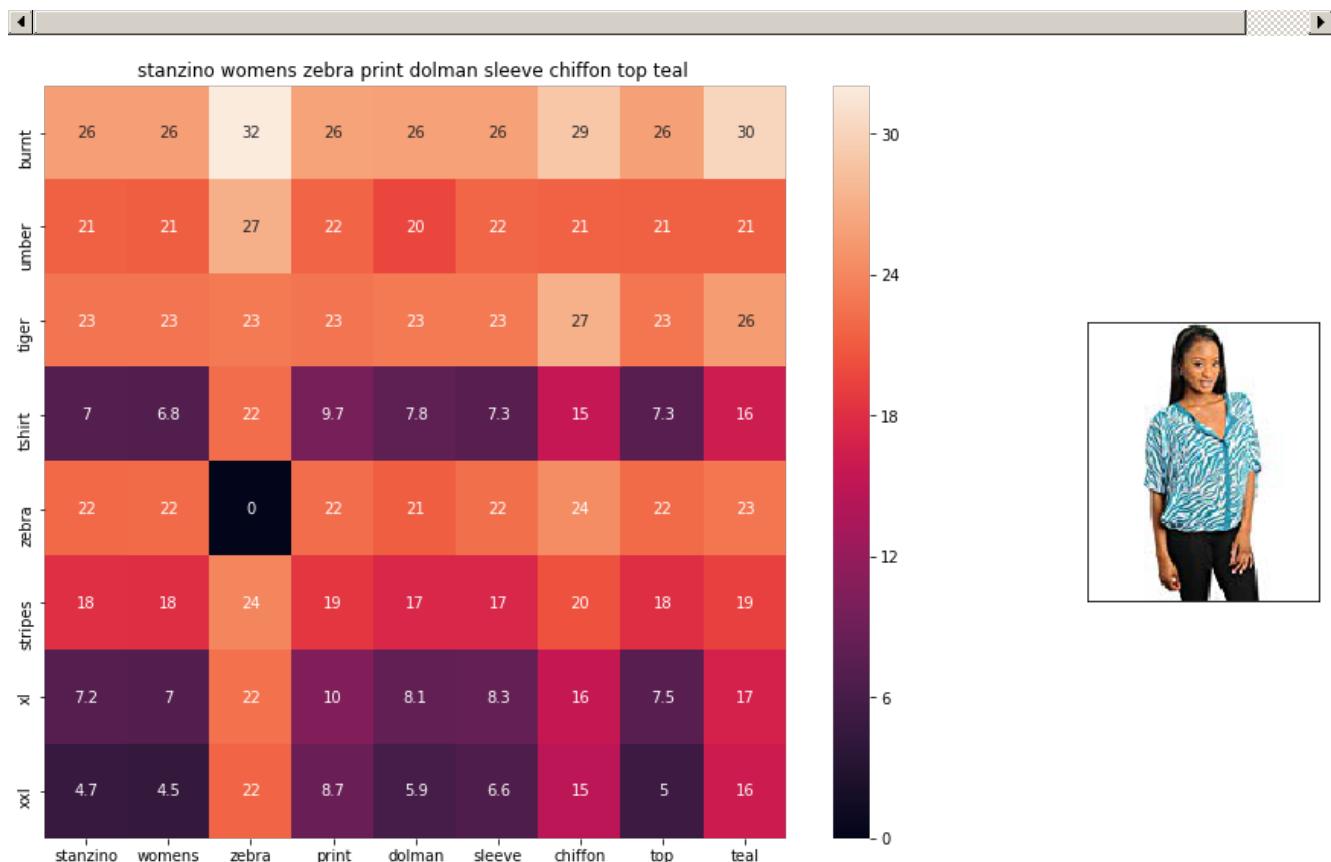
Brand : Edista

euclidean distance from input : 6.3922043

---



---



ASIN : B00C0I3U3E

Brand : Stanzino

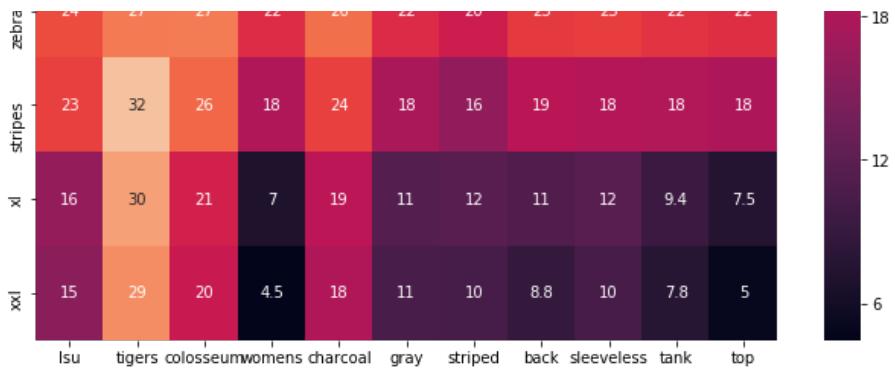
euclidean distance from input : 6.414901

---



---





ASIN : B073R4ZM7Y

Brand : Colosseum

euclidean distance from input : 6.45096

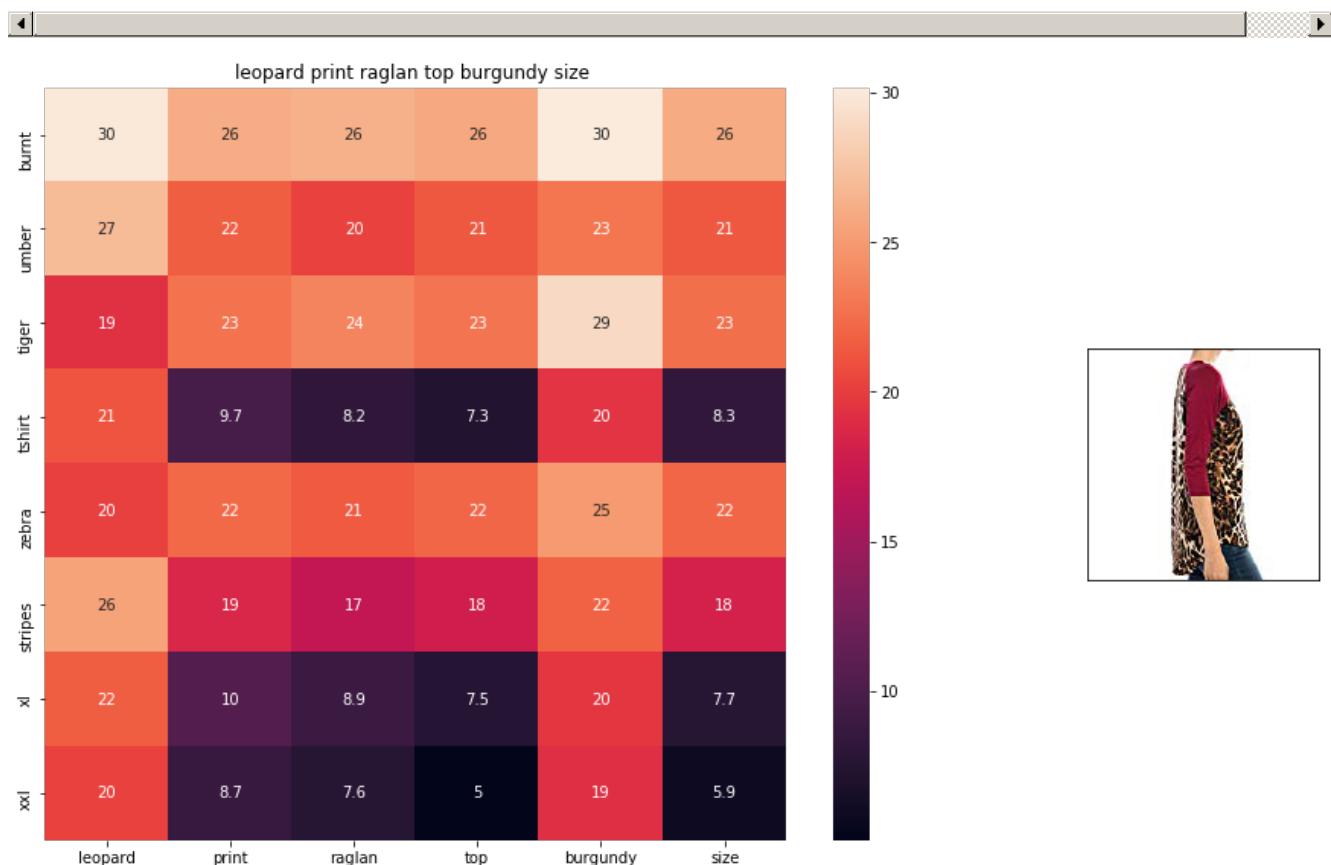
---



---



---



ASIN : B01C60RLDQ

Brand : 1 Mad Fit

euclidean distance from input : 6.4634094

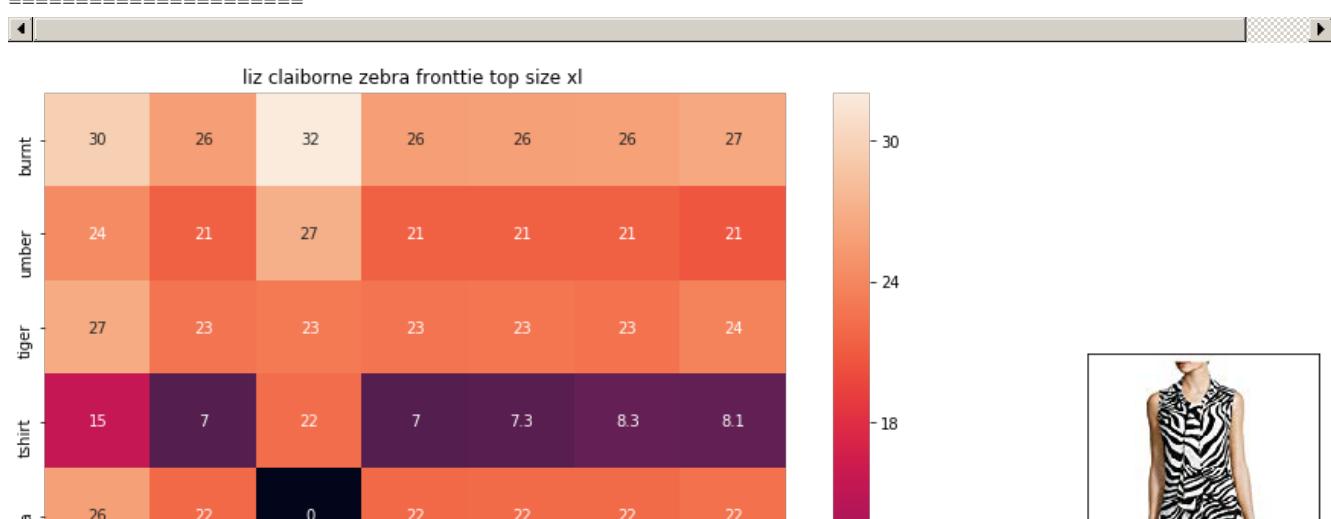
---

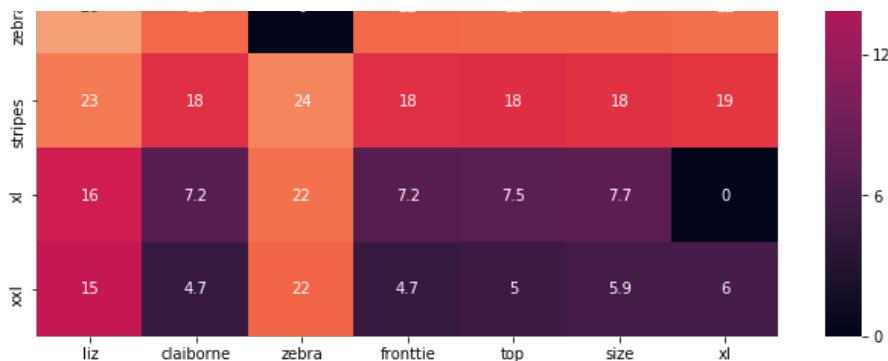


---



---





ASIN : B06XBY5QXL

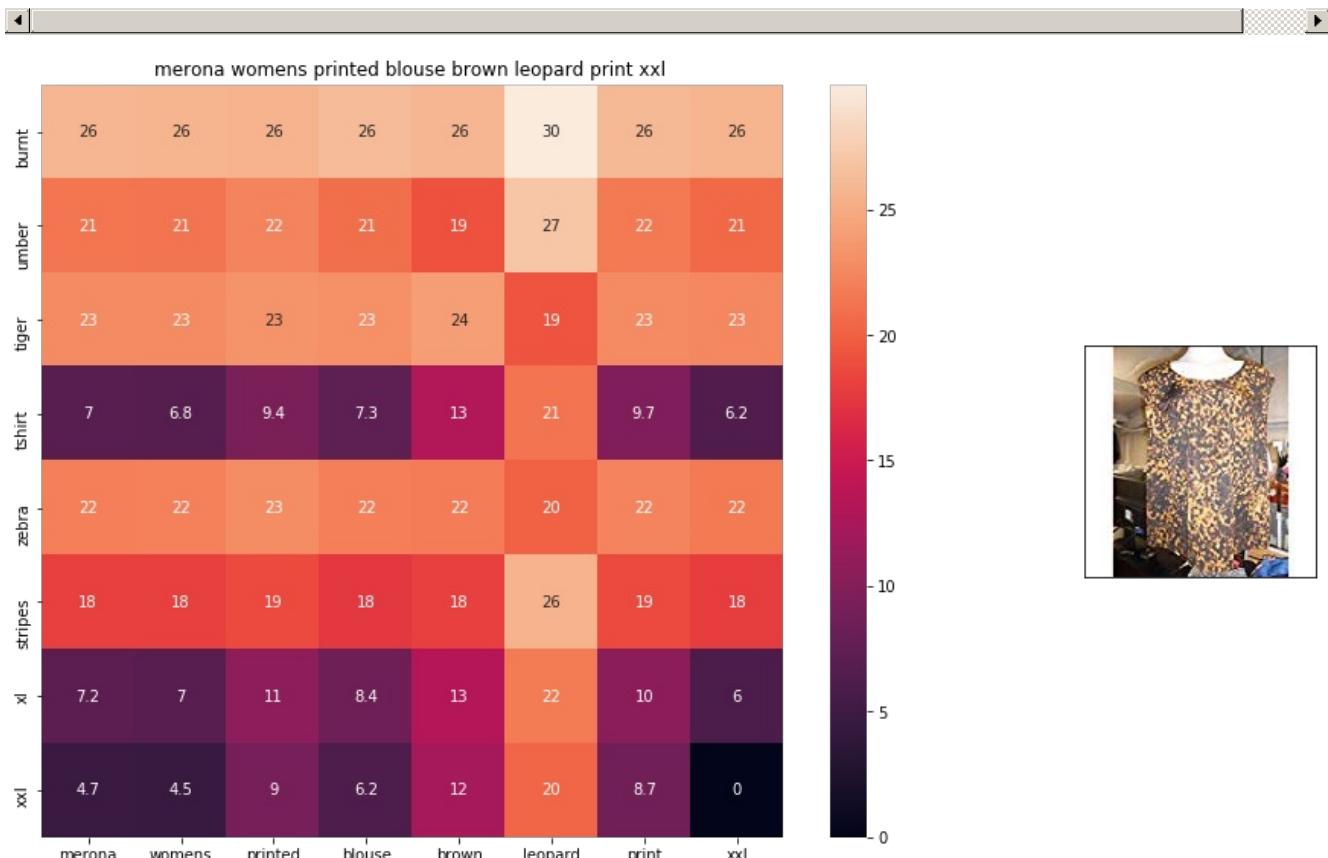
Brand : Liz Claiborne

euclidean distance from input : 6.5392237

---



---



ASIN : B071YF3WDD

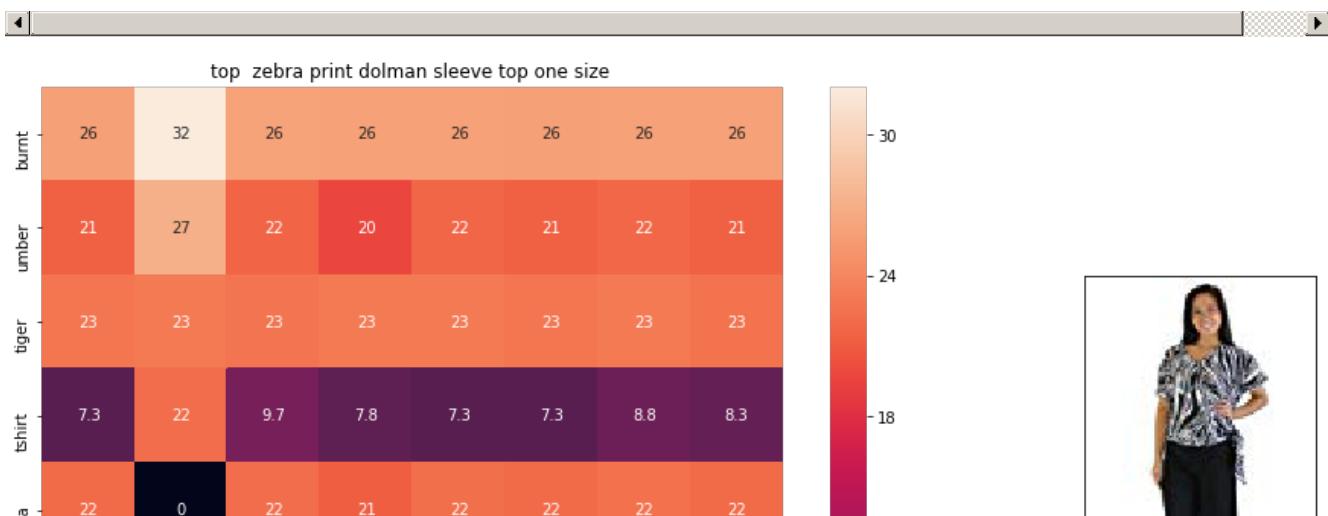
Brand : Merona

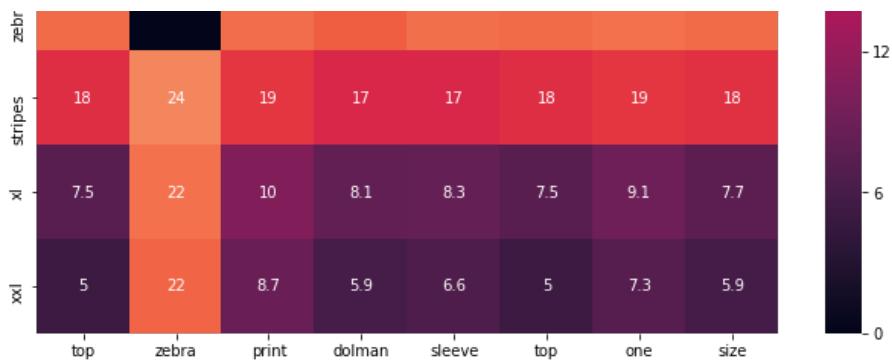
euclidean distance from input : 6.575504

---



---





ASIN : B00H8A6ZLI

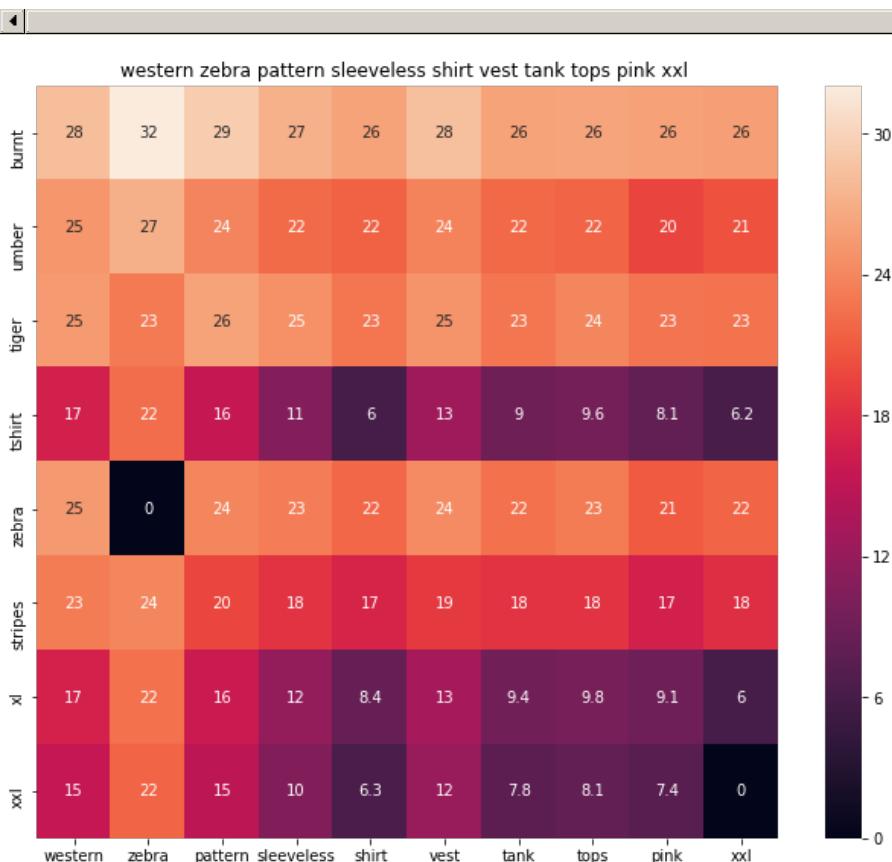
Brand : Vivian's Fashions

euclidean distance from input : 6.6382155

---



---



ASIN : B00Z6HEXWI

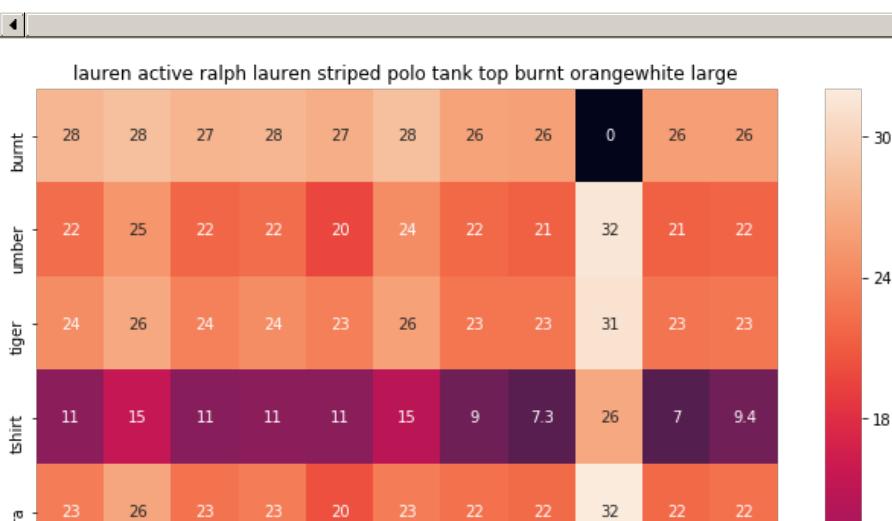
Brand : Black Temptation

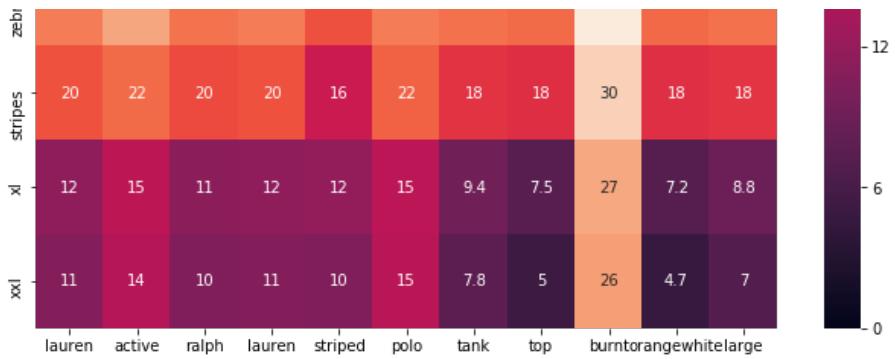
euclidean distance from input : 6.660737

---



---





ASIN : B00ILGH5OY

Brand : Ralph Lauren Active

euclidean distance from input : 6.6839056

---



---



ASIN : B06Y1VN8WQ

Brand : Black Swan

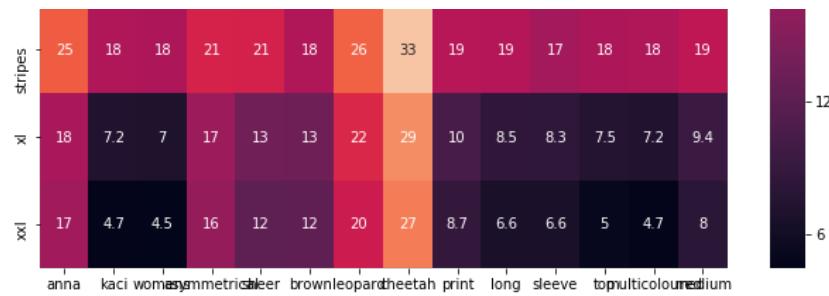
euclidean distance from input : 6.7057643

---



---





```
ASIN : B00KSNTY7Y
Brand : Anna-Kaci
euclidean distance from input : 6.706125
=====
```

## [9.6] Weighted similarity using brand and color.

In [135]:

```
# some of the brand values are empty.
# Need to replace Null with string "NULL"
data['brand'].fillna(value="Not given", inplace=True)

# replace spaces with hyphen
brands = [x.replace(" ", "-") for x in data['brand'].values]
types = [x.replace(" ", "-") for x in data['product_type_name'].values]
colors = [x.replace(" ", "-") for x in data['color'].values]

brand_vectorizer = CountVectorizer()
brand_features = brand_vectorizer.fit_transform(brands)

type_vectorizer = CountVectorizer()
type_features = type_vectorizer.fit_transform(types)

color_vectorizer = CountVectorizer()
color_features = color_vectorizer.fit_transform(colors)

extra_features = hstack((brand_features, type_features, color_features)).tocsr()
```

In [125]:

```
print(brand_features.shape,color_features.shape)
print(w2v_title_weight.shape)
print(w2v_title.shape)
```

```
(16042, 3835) (16042, 1845)
(16042, 300)
(16042, 300)
```

In [140]:

```
def heat_map_w2v_brand(sentance1, sentance2, url, doc_id1, doc_id2, df_id1, df_id2, model):

    # sentance1 : title1, input apparel
    # sentance2 : title2, recommended apparel
    # url: apparel image url
    # doc_id1: document id of input apparel
    # doc_id2: document id of recommended apparel
    # df_id1: index of document1 in the data frame
    # df_id2: index of document2 in the data frame
    # model: it can have two values, 1. avg 2. weighted

    #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
    s1_vec = get_word_vec(sentance1, doc_id1, model)
    #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
    s2_vec = get_word_vec(sentance2, doc_id2, model)
```

```

# s1_s2_dist = np.array(#number of words in title1 * #number of words in title2)
# s1_s2_dist[i,j] = euclidean distance between words i, j
s1_s2_dist = get_distance(s1_vec, s2_vec)

data_matrix = [['Asin', 'Brand', 'Color', 'Product type'],
               [data['asin'].loc[df_id1], brands[doc_id1], colors[doc_id1], types[doc_id1]], # input
apparel's features
               [data['asin'].loc[df_id2], brands[doc_id2], colors[doc_id2], types[doc_id2]]] # recommended apparel's features

colorscale = [[0, '#1d004d'], [.5, '#f2e5ff'], [1, '#f2e5d1']] # to color the headings of each column

# we create a table with the data_matrix
table = ff.create_table(data_matrix, index=True, colorscale=colorscale)
# plot it with plotly
plotly.offline.iplot(table, filename='simple_table')

# devide whole figure space into 25 * 1:10 grids
gs = gridspec.GridSpec(25, 15)
fig = plt.figure(figsize=(25,5))

# in first 25*10 grids we plot heatmap
ax1 = plt.subplot(gs[:, :-5])
# plotting the heap map based on the pairwise distances
ax1 = sns.heatmap(np.round(s1_s2_dist,6), annot=True)
# set the x axis labels as recommended apparels title
ax1.set_xticklabels(sentance2.split())
# set the y axis labels as input apparels title
ax1.set_yticklabels(sentance1.split())
# set title as recommended apparels title
ax1.set_title(sentance2)

# in last 25 * 10:15 grids we display image
ax2 = plt.subplot(gs[:, 10:16])
# we dont display grid lens and axis labels to images
ax2.grid(False)
ax2.set_xticks([])
ax2.set_yticks([])

# pass the url it display it
display_img(url, ax2, fig)

plt.show()

```

In [141]:

```

def idf_w2v_brand(doc_id, w1, w2, num_results):
    # doc_id: apparel's id in given corpus
    # w1: weight for w2v features
    # w2: weight for brand and color features

    # pairwise_dist will store the distance from given input apparel to all remaining apparels
    # the metric we used here is cosine, the coside distance is mesured as  $K(X, Y) = \langle X, Y \rangle / (\|X\| * \|Y\|)$ 
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
    idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1,-1))
    ex_feat_dist = pairwise_distances(extra_features, extra_features[doc_id])
    pairwise_dist = (w1 * idf_w2v_dist + w2 * ex_feat_dist)/float(w1 + w2)

    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    #data frame indices of the 9 smallest distace's
    df_indices = list(data.index[indices])

    for i in range(0, len(indices)):
        heat_map_w2v_brand(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]], data[
'medium_image_url'].loc[df_indices[i]], indices[0], indices[i], df_indices[0], df_indices[i], 'weigh
ted')
        print('ASIN :', data['asin'].loc[df_indices[i]])
        print('Brand :', data['brand'].loc[df_indices[i]])
        print('euclidean distance from input :', pdists[i])

```

```

print('euclidean distance from input . , pairs[1])
print('*125)

idf_w2v_brand(12566, 5, 5, 20)
# in the give heat map, each cell contains the euclidean distance between words i, j

```

```

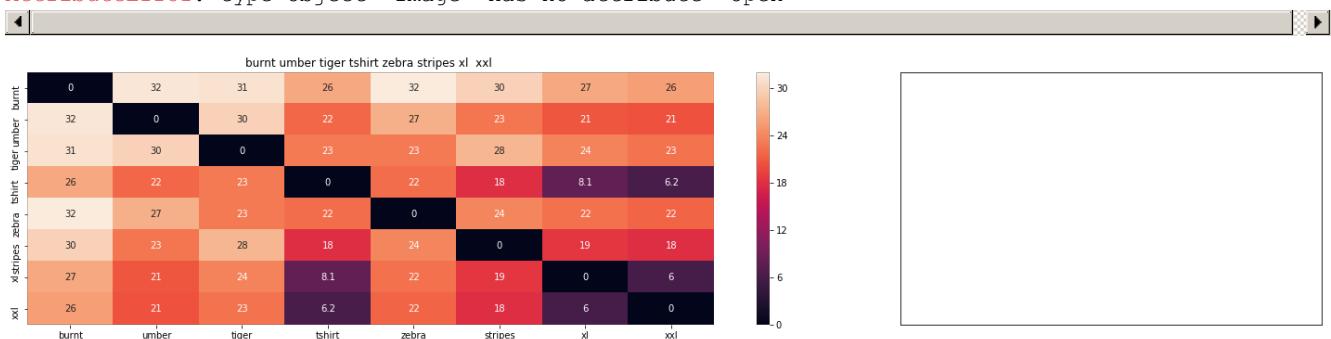
-----
AttributeError                                 Traceback (most recent call last)
<ipython-input-141-822d8df49624> in <module>()
    27         print('*125)
    28
---> 29 idf_w2v_brand(12566, 5, 5, 20)
    30 # in the give heat map, each cell contains the euclidean distance between words i, j

<ipython-input-141-822d8df49624> 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['title'].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-140-ac9ce1dbaed2> in heat_map_w2v_brand(sentance1, sentance2, url, doc_id1,
doc_id2, df_id1, df_id2, model)
    53
    54     # pass the url it display it
---> 55     display_img(url, ax2, fig)
    56
    57     plt.show()

<ipython-input-93-f15034a92765> in 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)
```

AttributeError: type object 'Image' has no attribute 'open'



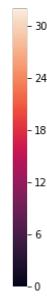
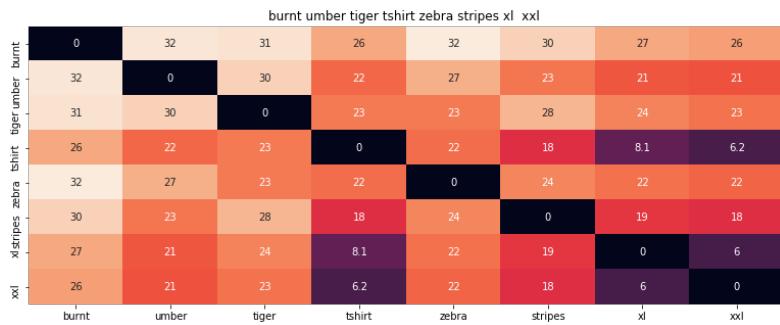
In [114]:

```

# brand and color weight =50
# title vector weight = 5

idf_w2v_brand(12566, 5, 50, 20)

```



ASIN : B00JXQB5FQ

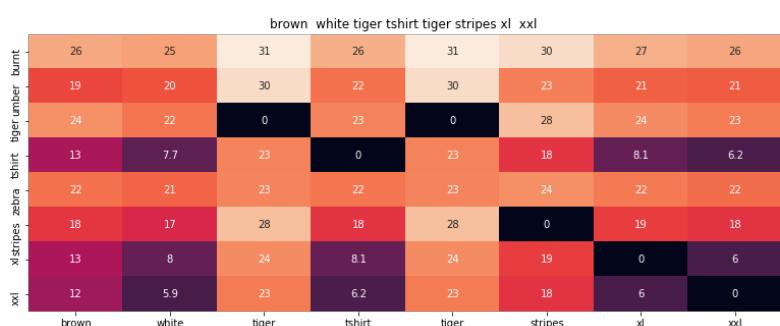
Brand : Si Row

euclidean distance from input : 0.0003551136363636364

---



---



ASIN : B00JXQCWT0

Brand : Si Row

euclidean distance from input : 0.43372202786532316

---



---



ASIN : B00JXQASS6

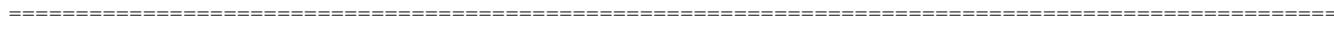
Brand : Si Row

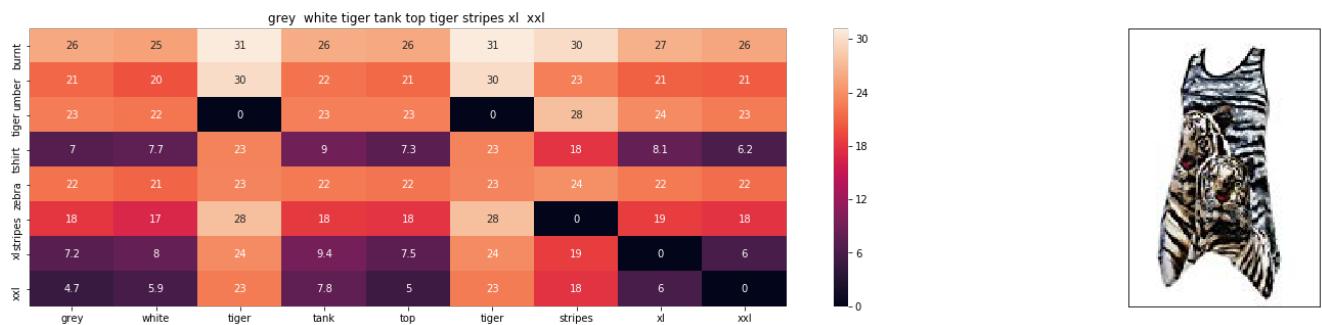
euclidean distance from input : 1.655093106685058

---



---

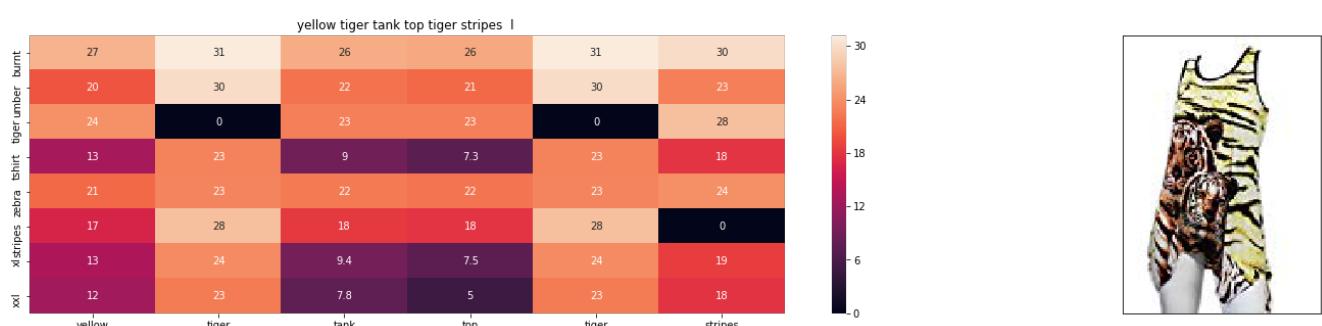




ASIN : B00JXQAFZ2

Brand : Si Row

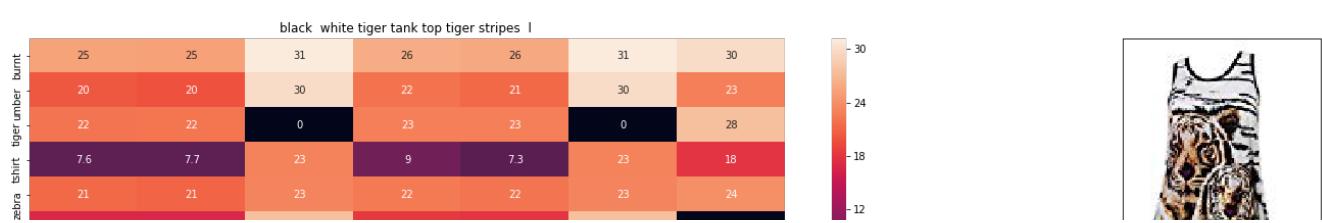
euclidean distance from input : 1.7729360410334245

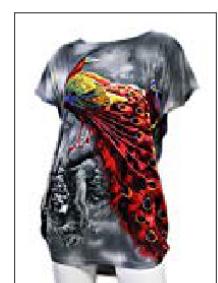
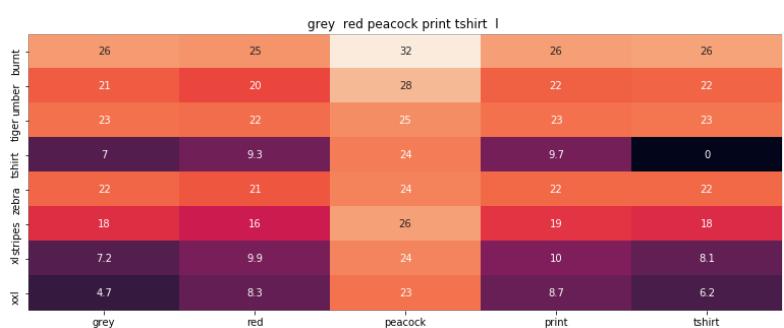
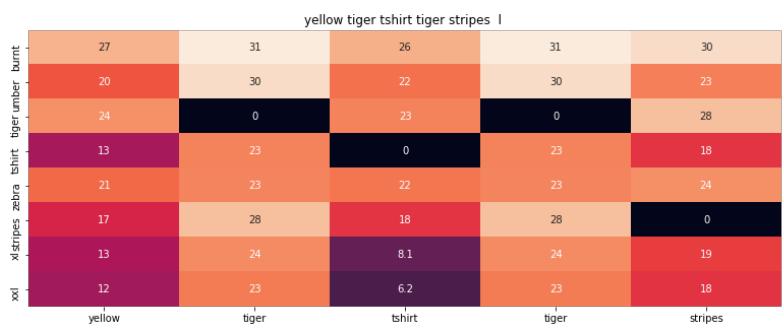
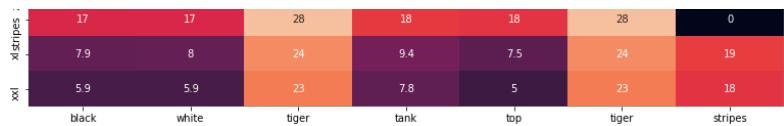


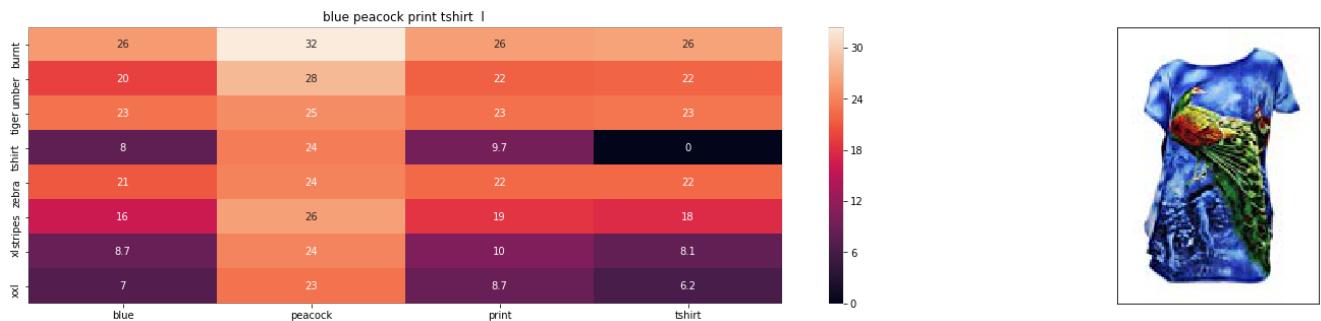
ASIN : B00JXQAUWA

Brand : Si Row

euclidean distance from input : 1.8028781200573059







ASIN : B00JXQC8L6

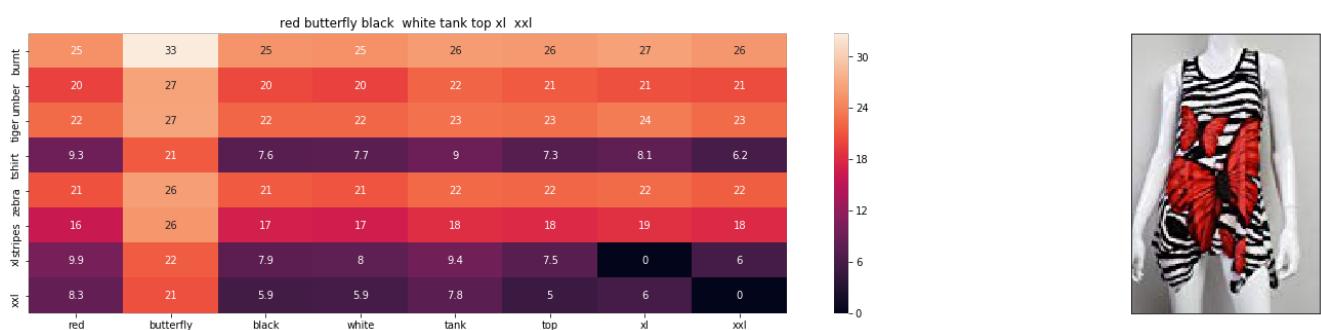
Brand : Si Row

euclidean distance from input : 1.9214293049716347

---



---



ASIN : B00JV63CW2

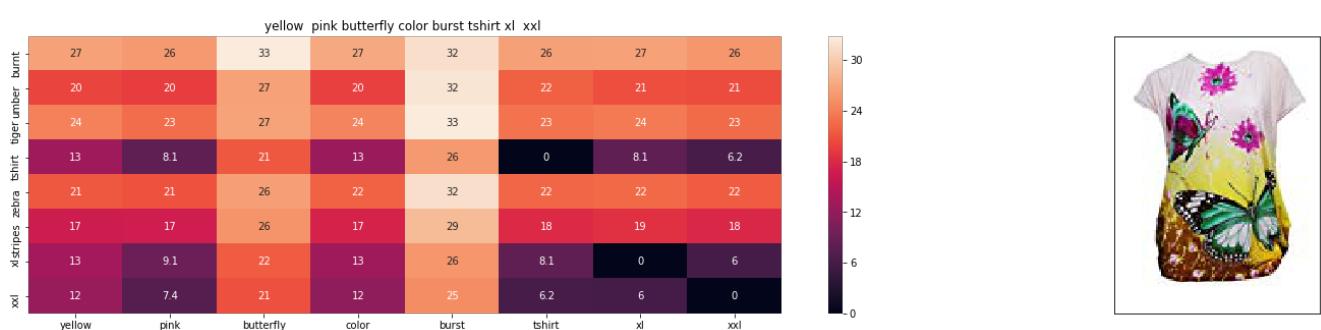
Brand : Si Row

euclidean distance from input : 1.9364632349698592

---



---

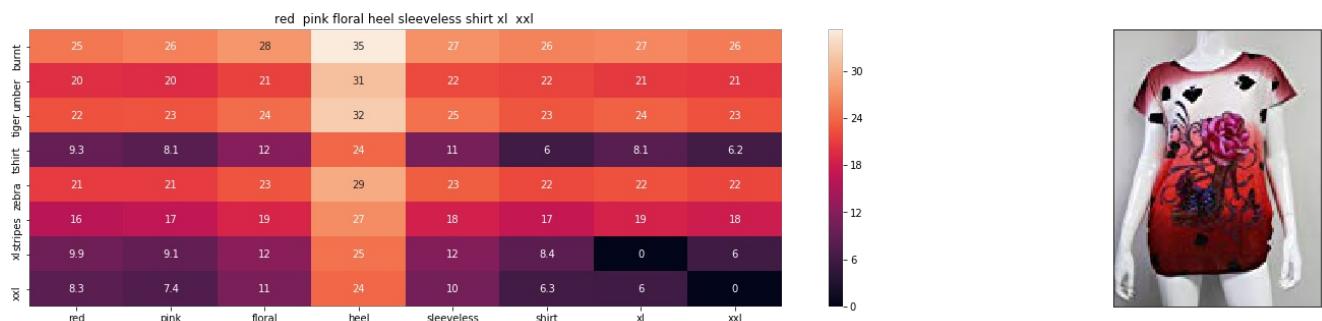


ASIN : B00JXQBBMI

Brand : Si Row

euclidean distance from input : 1.9567038799450012

=====

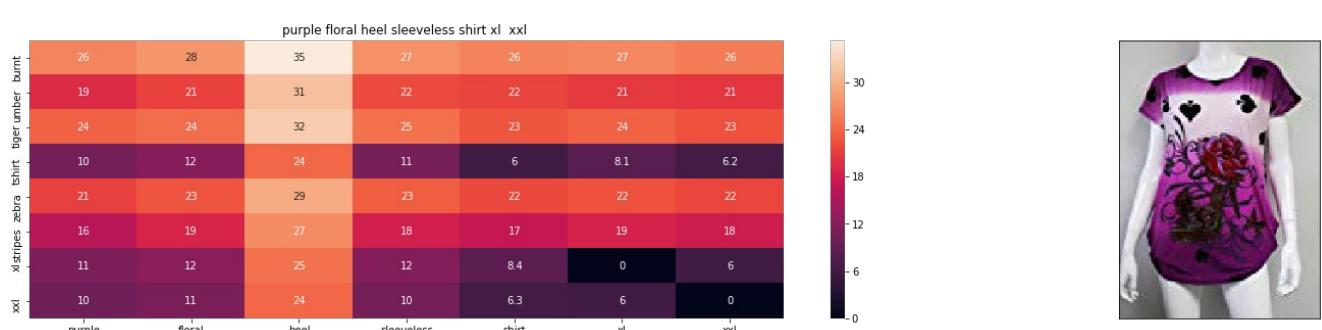


ASIN : B00JV63QQE

Brand : Si Row

euclidean distance from input : 1.9806065649370113

=====

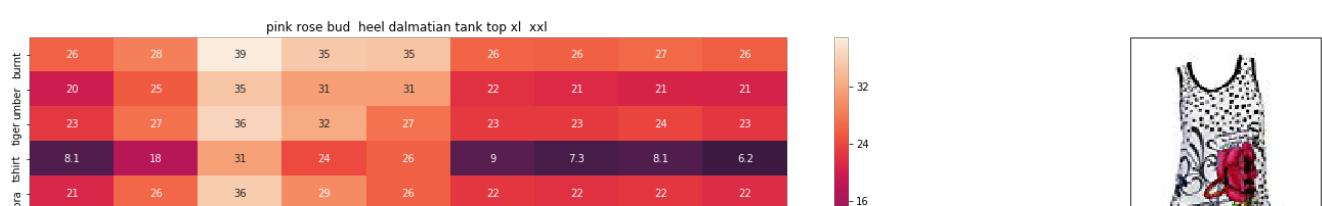


ASIN : B00JV63VC8

Brand : Si Row

euclidean distance from input : 2.0121855999157043

=====

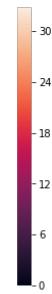
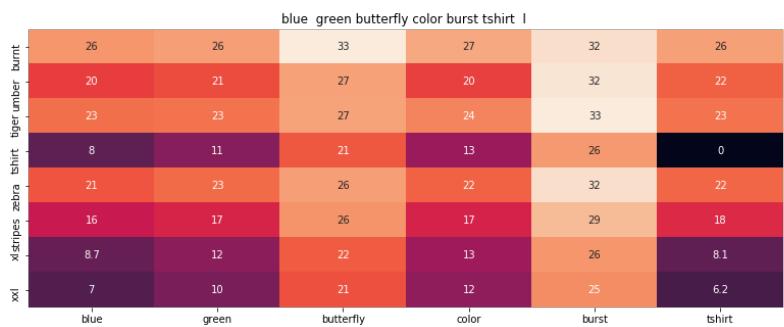
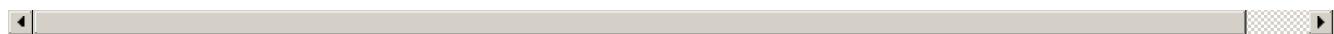




ASIN : B00JXQAX2C

Brand : Si Row

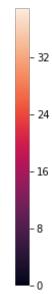
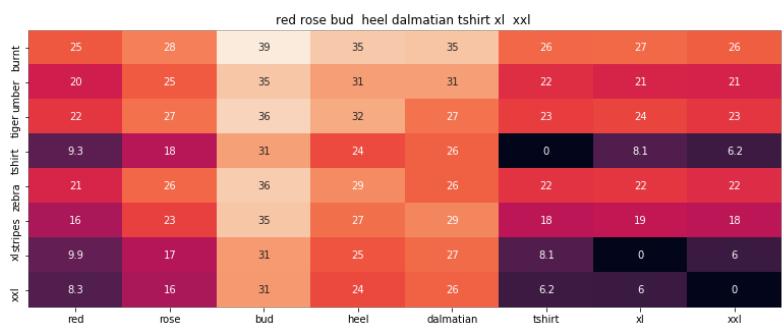
euclidean distance from input : 2.013351787548872



ASIN : B00JXQC0C8

Brand : Si Row

euclidean distance from input : 2.013883348273304

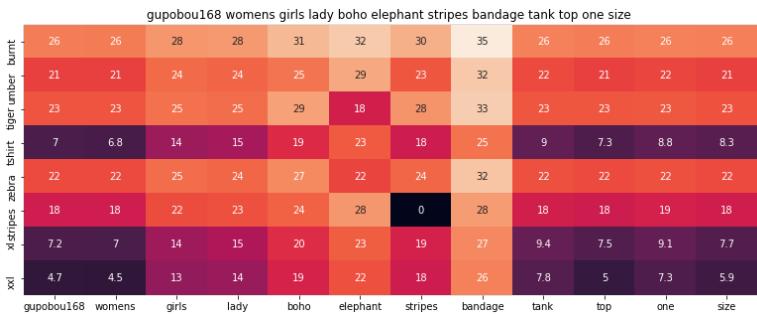


ASIN : B00JXQABB0

Brand : Si Row

euclidean distance from input : 2.0367258248579985





ASIN : B01ER18406

Brand : GuPoBoU168

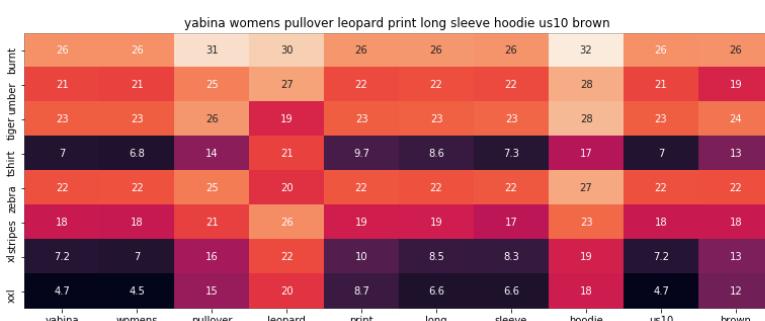
euclidean distance from input : 2.6562041677776853



ASIN : B01LZ7BQ4H

Brand : WAYF

euclidean distance from input : 2.684906782301833

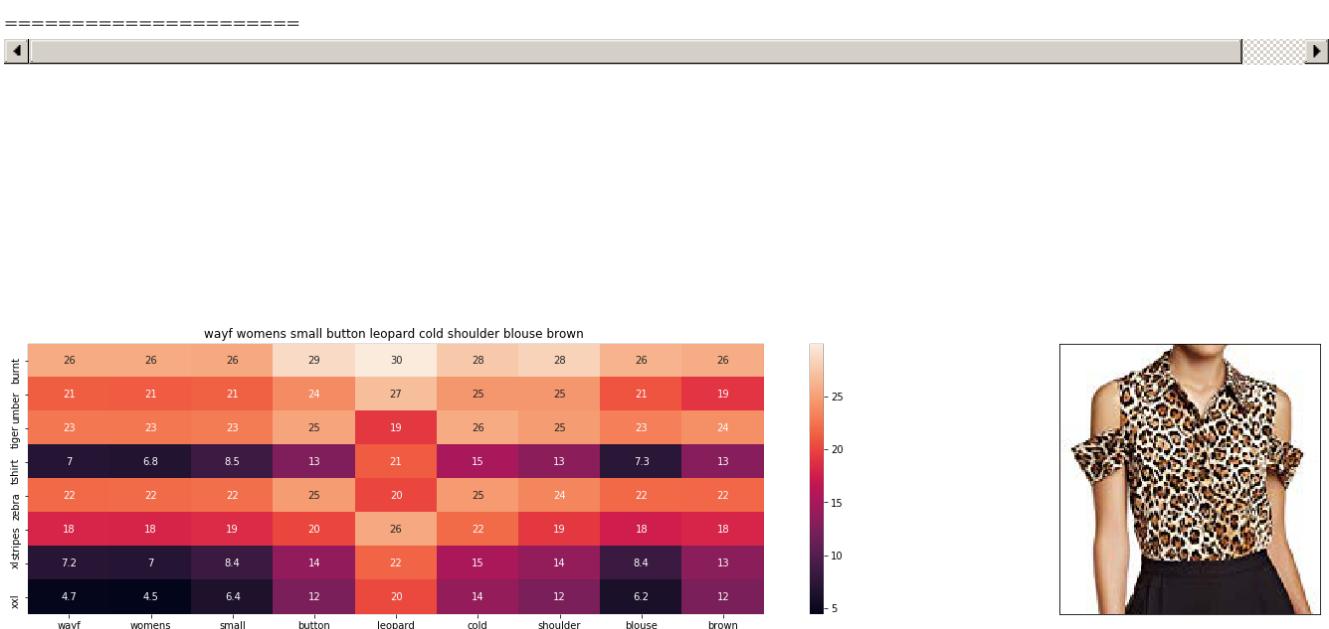


ASIN : B01KJUM6JI

Brand : YABINA

euclidean distance from input : 2.6858381926578345





ASIN : B01M06V4X1

Brand : WAYF

euclidean distance from input : 2.694761948650377

=====

=====

## [10.2] Keras and Tensorflow to extract features

In [115]:

```
import numpy as np
from keras.preprocessing.image import ImageDataGenerator
from keras.models import Sequential
from keras.layers import Dropout, Flatten, Dense
from keras import applications
from sklearn.metrics import pairwise_distances
import matplotlib.pyplot as plt
import requests
from PIL import Image
import pandas as pd
import pickle
```

Using TensorFlow backend.

In [ ]:

```
# https://gist.github.com/fchollet/f35fbc80e066a49d65f1688a7e99f069
# Code reference: https://blog.keras.io/building-powerful-image-classification-models-using-very-little-data.html

# This code takes 40 minutes to run on a modern GPU (graphics card)
# like Nvidia 1050.
# GPU (Nvidia 1050): 0.175 seconds per image

# This code takes 160 minutes to run on a high end i7 CPU
# CPU (i7): 0.615 seconds per image.

#Do NOT run this code unless you want to wait a few hours for it to generate output

# each image is converted into 25088 length dense-vector

'''

# dimensions of our images.
img_width, img_height = 224, 224
```

```

top_model_weights_path = 'bottleneck_fc_model.h5'
train_data_dir = 'images2/'
nb_train_samples = 16042
epochs = 50
batch_size = 1

def save_bottlebeck_features():
    #Function to compute VGG-16 CNN for image feature extraction.

    asins = []
    datagen = ImageDataGenerator(rescale=1. / 255)

    # build the VGG16 network
    model = applications.VGG16(include_top=False, weights='imagenet')
    generator = datagen.flow_from_directory(
        train_data_dir,
        target_size=(img_width, img_height),
        batch_size=batch_size,
        class_mode=None,
        shuffle=False)

    for i in generator.filenames:
        asins.append(i[2:-5])

    bottleneck_features_train = model.predict_generator(generator, nb_train_samples // batch_size)
    bottleneck_features_train = bottleneck_features_train.reshape((16042,25088))

    np.save(open('16k_data_cnn_features.npy', 'wb'), bottleneck_features_train)
    np.save(open('16k_data_cnn_feature_asins.npy', 'wb'), np.array(asins))

save_bottlebeck_features()
'''
```

## [10.3] Visual features based product similarity.

In [159]:

```

#load the features and corresponding ASINS info.
bottleneck_features_train =
np.load('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/16k_data_cnn_features.npy')
asins =
np.load('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/16k_data_cnn_feature_asins.npy')
asins = list(asins)

# load the original 16K dataset
data =
pd.read_pickle('/home/bhargav/AAIC/DataSets/Amzon_Fashion_Discovery/pickels/16k_apperial_data_prepesed')
df_asins = list(data['asin'])

from IPython.display import display, Image, SVG, Math, YouTubeVideo

#get similar products using CNN features (VGG-16)
def get_similar_products_cnn(doc_id, num_results):
    doc_id = asins.index(df_asins[doc_id])
    pairwise_dist = pairwise_distances(bottleneck_features_train, bottleneck_features_train[doc_id].reshape(1,-1))

    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    for i in range(len(indices)):
        rows = data[['medium_image_url','title']].loc[data['asin']==asins[indices[i]]]
        for indx, row in rows.iterrows():
            display(Image(url=row['medium_image_url'], embed=True))
            print('Product Title: ', row['title'])
            print('Euclidean Distance from input image:', pdists[i])
            print('Amazon URL: www.amazon.com/dp/1423103004/')
```

```
printf('Amazon Url: www.amazon.com/dp/' + asins['titles'][i])  
get_similar_products_cnn(12566, 20)
```



Product Title: burnt umber tiger tshirt zebra stripes xl xxl  
Euclidean Distance from input image: 0.0625  
Amazon Url: [www.amazon.com/dp/B00JXQB5FQ](http://www.amazon.com/dp/B00JXQB5FQ)



Product Title: pink tiger tshirt zebra stripes xl xxl  
Euclidean Distance from input image: 30.050072  
Amazon Url: [www.amazon.com/dp/B00JXQASS6](http://www.amazon.com/dp/B00JXQASS6)



Product Title: yellow tiger tshirt tiger stripes l  
Euclidean Distance from input image: 41.261078  
Amazon Url: [www.amazon.com/dp/B00JXQCUIC](http://www.amazon.com/dp/B00JXQCUIC)



Product Title: brown white tiger tshirt tiger stripes xl xxl  
Euclidean Distance from input image: 44.000187  
Amazon Url: [www.amazon.com/dp/B00JXQCWT0](http://www.amazon.com/dp/B00JXQCWT0)



Product Title: kawaii pastel tops tees pink flower design  
Euclidean Distance from input image: 47.38251  
Amazon Url: [www.amazon.com/dp/B071FCWD97](http://www.amazon.com/dp/B071FCWD97)



Product Title: womens thin style tops tees pastel watermelon print  
Euclidean Distance from input image: 47.718403  
Amazon Url: [www.amazon.com/dp/B01JUNHBRM](http://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](http://www.amazon.com/dp/B071SBCY9W)



Product Title: edv cheetah run purple multi xl  
Euclidean Distance from input image: 48.046516  
Amazon Url: [www.amazon.com/dp/B01CUPYBM0](http://www.amazon.com/dp/B01CUPYBM0)



Product Title: danskin womens vneck loose performance tee xsmall pink ombre  
Euclidean Distance from input image: 48.101917  
Amazon Url: [www.amazon.com/dp/B01F7PHXY8](http://www.amazon.com/dp/B01F7PHXY8)



Product Title: summer alpaca 3d pastel casual loose tops tee design  
Euclidean Distance from input image: 48.118877  
Amazon Url: [www.amazon.com/dp/B01I80A93G](http://www.amazon.com/dp/B01I80A93G)



Product Title: miss chievous juniors striped peplum tank top medium shadowpeach  
Euclidean Distance from input image: 48.13129  
Amazon Url: [www.amazon.com/dp/B0177DM70S](http://www.amazon.com/dp/B0177DM70S)



Product Title: red pink floral heel sleeveless shirt xl xxl  
Euclidean Distance from input image: 48.16947  
Amazon Url: [www.amazon.com/dp/B00JV63QOE](http://www.amazon.com/dp/B00JV63QOE)



Product Title: moana logo adults hot v neck shirt black xxl  
Euclidean Distance from input image: 48.25679  
Amazon Url: [www.amazon.com/dp/B01LX6H43D](http://www.amazon.com/dp/B01LX6H43D)



Product Title: abaday multicolor cartoon cat print short sleeve longline shirt large  
Euclidean Distance from input image: 48.265633  
Amazon Url: [www.amazon.com/dp/B01CR57YY0](http://www.amazon.com/dp/B01CR57YY0)



Product Title: kawaii cotton pastel tops tees peach pink cactus design  
Euclidean Distance from input image: 48.362583  
Amazon Url: [www.amazon.com/dp/B071WYLBZS](http://www.amazon.com/dp/B071WYLBZS)



Product Title: chicago chicago 18 shirt women pink  
Euclidean Distance from input image: 48.383617  
Amazon Url: [www.amazon.com/dp/B01GXAZTRY](http://www.amazon.com/dp/B01GXAZTRY)



Product Title: yichun womens tiger printed summer tshirts tops  
Euclidean Distance from input image: 48.449303  
Amazon Url: [www.amazon.com/dp/B010NN9RXO](http://www.amazon.com/dp/B010NN9RXO)



Product Title: nancy lopez whimsy short sleeve whiteblacklemon drop xs  
Euclidean Distance from input image: 48.478912  
Amazon Url: [www.amazon.com/dp/B01MPX6IDX](http://www.amazon.com/dp/B01MPX6IDX)



Product Title: womens tops tees pastel peach ice cream cone print  
Euclidean Distance from input image: 48.557964  
Amazon Url: [www.amazon.com/dp/B0734GRKZL](http://www.amazon.com/dp/B0734GRKZL)



Product Title: uswomens mary j blige without tshirts shirt  
 Euclidean Distance from input image: 48.614376  
 Amazon Url: [www.amazon.com/dp/B01M0XXFKK](http://www.amazon.com/dp/B01M0XXFKK)

## Assignment:

- Combining TextData vector Features,Brand,Color Features along with the CNN Features

In [157]:

```
#extra_features = hstack((w2v_title_weight,brand_features,
color_features,bottleneck_features_train)).tocsr()
print(w2v_title_weight.shape)
print(extra_features.shape)
print(bottleneck_features_train.shape)

(16042, 300)
(16042, 5735)
(16042, 25088)
```

In [163]:

```
def idf_w2v_brand(doc_id, w1, w2, w3, num_results):

    idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1,-1))
    ex_feat_dist = pairwise_distances(extra_features, extra_features[doc_id])
    bottleneck_features_train_dist = pairwise_distances(bottleneck_features_train,
bottleneck_features_train[doc_id].reshape(1,-1))

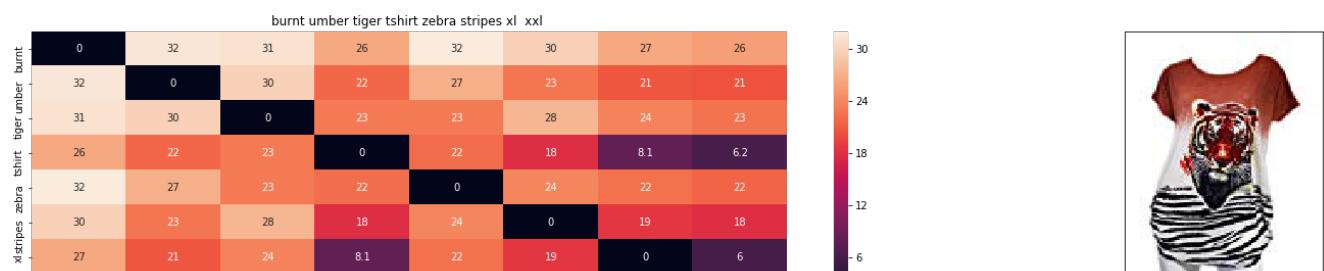
    pairwise_dist = (w1 * idf_w2v_dist + w2 * ex_feat_dist + w3 * bottleneck_features_train_dist
)/float(w1 + w2 + w3)

    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    df_indices = list(data.index[indices])

    for i in range(0, len(indices)):
        heat_map_w2v_brand(data['title'].loc[df_indices[0]],data['title'].loc[df_indices[i]], data[
'medium_image_url'].loc[df_indices[i]], indices[0], indices[i],df_indices[0], df_indices[i], 'weigh
ted')
        print('ASIN :',data['asin'].loc[df_indices[i]])
        print('Brand :',data['brand'].loc[df_indices[i]])
        print('euclidean distance from input :', pdists[i])
        print('='*125)

idf_w2v_brand(12566, 20, 5, 2, 20)
```





ASIN : B00JXQB5FQ

Brand : Si Row

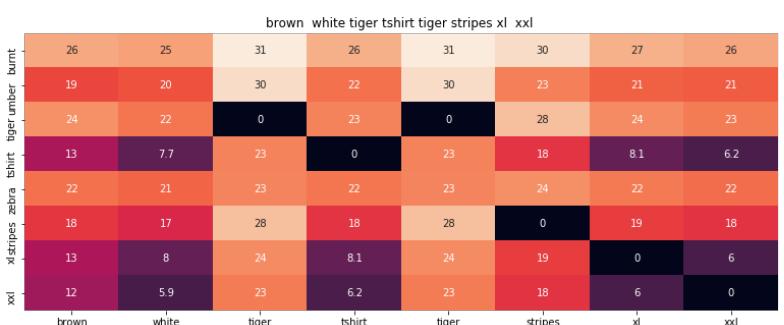
euclidean distance from input : 0.006167160967985789



ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from input : 6.865509315840215

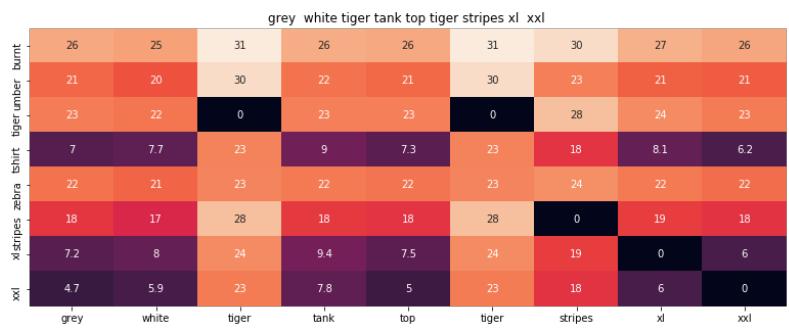


ASIN : B00JXQCWTO

Brand : Si Row

euclidean distance from input : 7.648690117730035





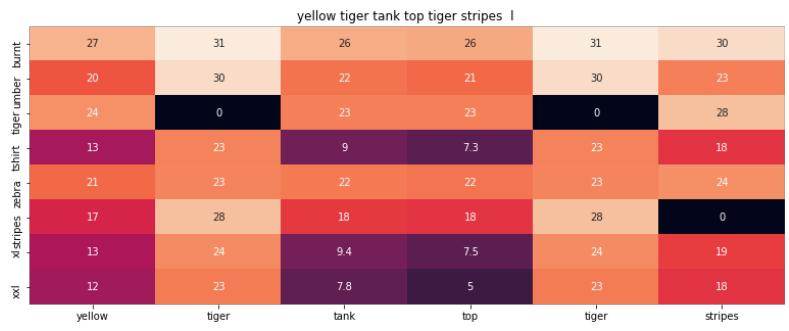
ASIN : B00JXQAFZ2

Brand : Si Row

euclidean distance from input : 8.22694453493382

=====

=====



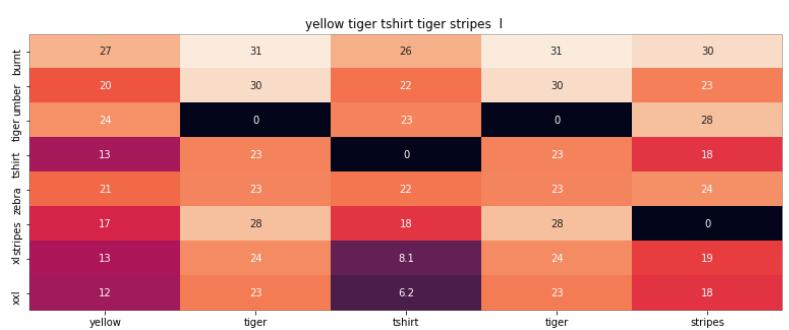
ASIN : B00JXQAUWA

Brand : Si Row

euclidean distance from input : 8.437752335226788

=====

=====



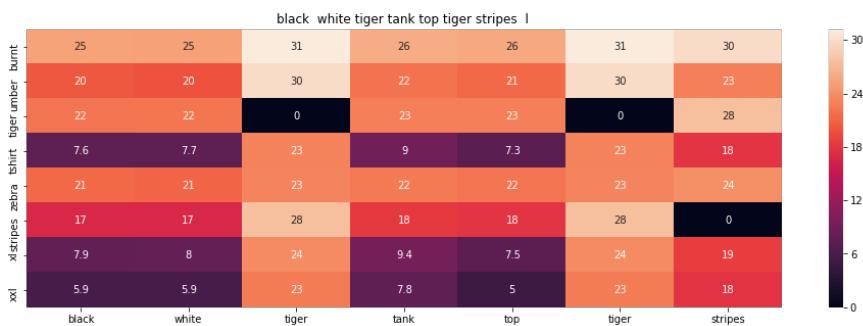
ASIN : B00JXQCUIC

Brand : Si Row

euclidean distance from input : 8.492368910114665

=====

=====



ASIN : B00JXQAO94

Brand : Si Row

euclidean distance from input : 8.570007889426007

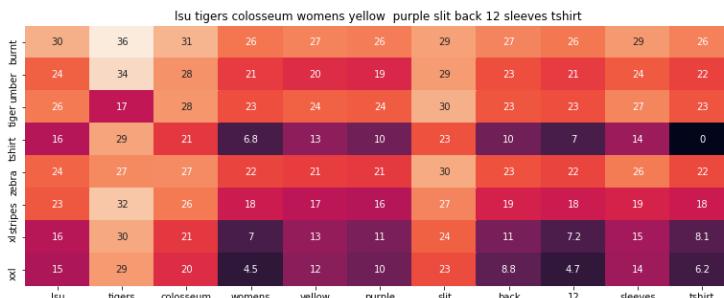
---



---



---



ASIN : B073R5Q8HD

Brand : Colosseum

euclidean distance from input : 8.792907030803338

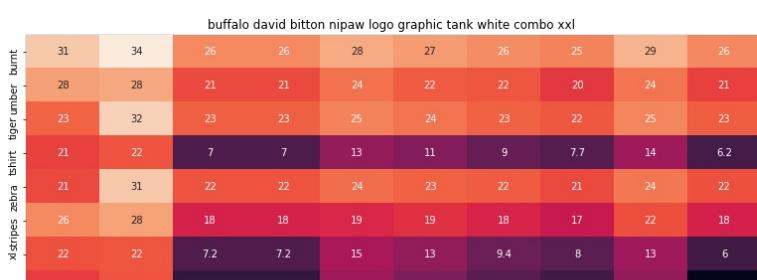
---



---



---





ASIN : B018H5AZXQ

Brand : Buffalo

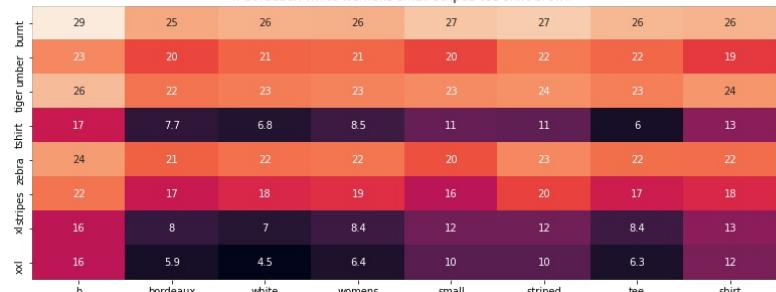
euclidean distance from input : 8.951488918862141

=====

=====



h bordeaux white womens small striped tee shirt brown



ASIN : B072BVB47Z

Brand : H By Bordeaux

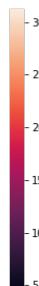
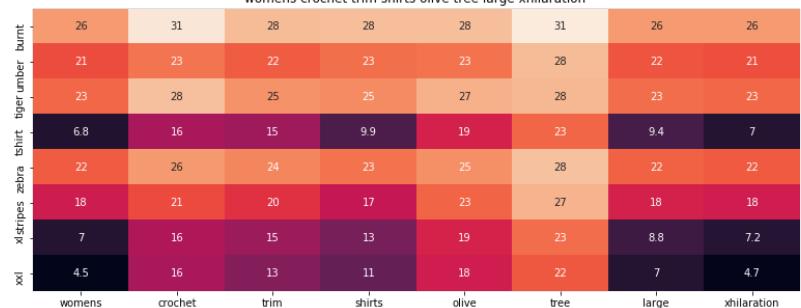
euclidean distance from input : 9.122048440955941

=====

=====



womens crochet trim shirts olive tree large xhilaration



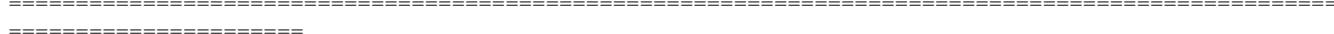
ASIN : B06XBHNM7J

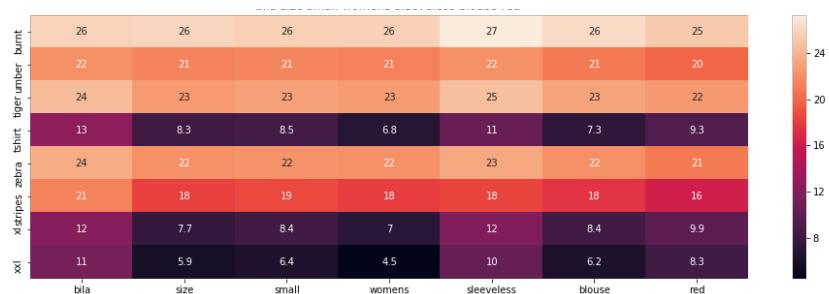
Brand : Xhilaration

euclidean distance from input : 9.153653462861708

=====

=====

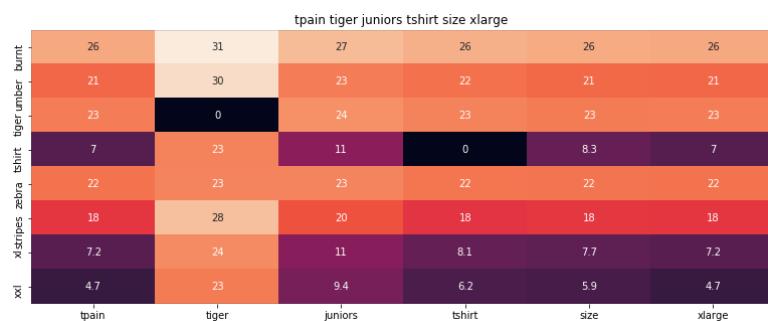




ASIN : B01L7ROZNC

Brand : Bila

euclidean distance from input : 9.160928960508922



ASIN : B01K0H02OG

Brand : Tultex

euclidean distance from input : 9.168128565814913



ASIN : B073ZHRBV8

Brand : Exotic India

euclidean distance from input : 9.212424101786766





ASIN : B0722DJVQP

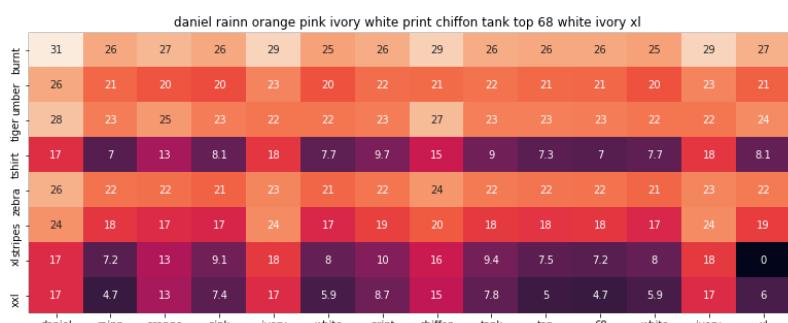
Brand : Kasper

euclidean distance from input : 9.241456937216098

---



---



ASIN : B01IPV1SFQ

Brand : Daniel Rainn

euclidean distance from input : 9.264452051233363

---



---





ASIN : B01DNNI1RO

Brand : Usstore

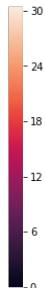
euclidean distance from input : 9.271516588132975

=====



kirkland signature womens long sleeve crew neck striped sweater size xxl color whitegrey

	26	30	26	26	26	28	26	27	30	26	26	27	26
26	26	21	22	22	25	22	20	25	21	21	20	21	21
21	25	21	22	22	25	22	20	25	21	21	20	21	21
23	26	23	23	23	25	23	23	27	23	23	24	23	23
7	17	6.8	8.6	7.3	15	9.8	11	14	8.3	6.2	13	7	7
22	25	22	22	22	25	22	20	25	22	22	22	22	22
18	21	18	19	17	22	18	16	20	18	18	17	18	18
7.2	17	7	8.5	8.3	16	10	12	17	7.7	6	13	7.2	7.2
4.7	16	4.5	6.6	6.6	14	8.7	10	15	5.9	0	12	4.7	4.7
xxl													
	kirkland	signature	womens	long	sleeve	crew	neck	striped	sweater	size	xxl	color	whitegrey



ASIN : B06XTPC3FP

Brand : Kirkland Signature

euclidean distance from input : 9.276525709497992

=====



bobeau peach womens small tribal lace tank blouse orange

	26	29	26	26	31	27	26	26	27
26	26	21	21	28	21	22	21	20	20
21	23	21	21	28	21	22	21	20	20
23	26	23	23	27	24	23	23	25	25
7	17	6.8	8.5	21	10	9	7.3	13	13
22	24	22	22	28	23	22	22	22	22
18	22	18	19	26	17	18	18	17	17
7.2	18	7	8.4	22	11	9.4	8.4	13	13
4.7	17	4.5	6.4	21	9.3	7.8	6.2	13	13
xxl									
	bobeau	peach	womens	small	tribal	lace	tank	blouse	orange



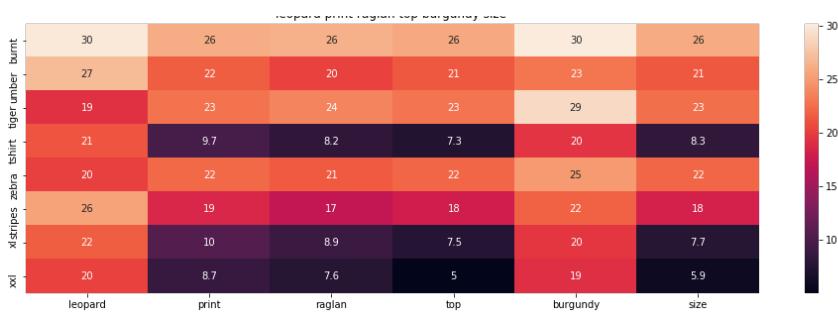
ASIN : B072JTHCX6

Brand : Bobeau

euclidean distance from input : 9.292742857535833

=====





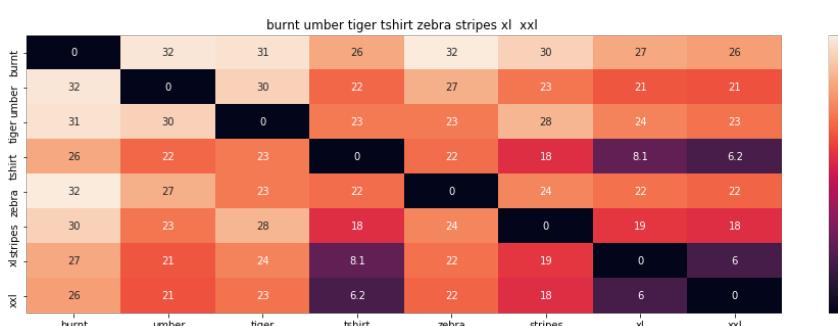
ASIN : B01C60RLDQ

Brand : 1 Mad Fit

euclidean distance from input : 9.308053193402795

In [164]:

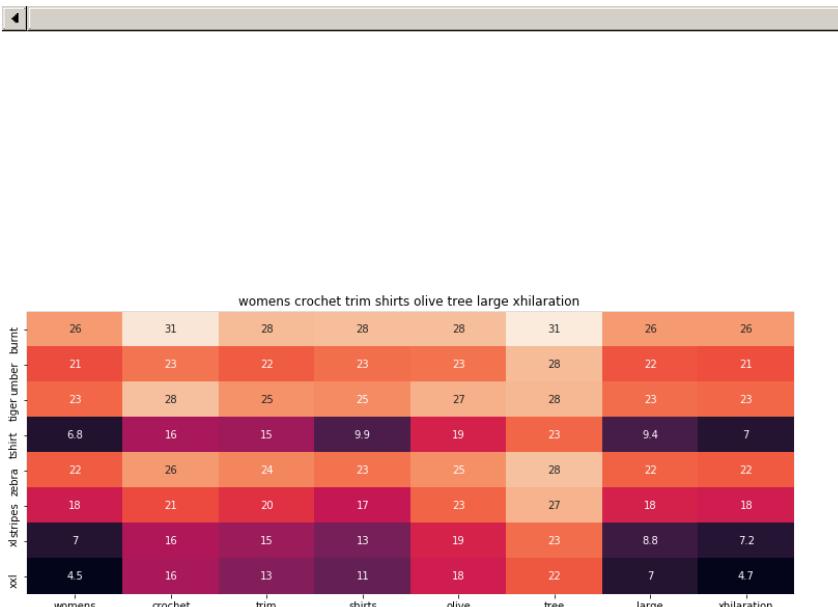
```
idf_w2v_brand(12566, 2, 5, 30, 20)
```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from input : 0.0360442625509726



ASIN : B06XBHNM7J

Brand : Xhilaration

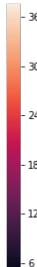
euclidean distance from input : 30.96323634492486

=====



breast cancer awareness juniors vneck shirt fight cancer

	breast	cancer	awareness	juniors	vneck	shirt	fight	cancer
burnt	34	33	38	27	26	26	34	33
tiger	30	31	35	23	21	22	32	31
umber	31	31	35	24	23	23	31	31
zebra	22	24	28	11	7	6	24	24
stripes	30	32	35	23	22	22	32	32
xxl	28	30	33	20	18	17	28	30
burnt	23	24	29	11	7.2	8.4	24	24
tiger	22	24	29	9.4	4.7	6.3	24	24



ASIN : B016CU40IY

Brand : Juiceclouds

euclidean distance from input : 32.08005932559849

=====



completely liz lange long flyaway vest 249682 turquoise

	completely	liz	lange	long	flyaway	vest	249682	turquoise
burnt	32	30	31	26	36	28	26	29
tiger	33	24	23	22	33	24	21	21
umber	32	27	28	23	34	25	23	27
zebra	25	15	17	8.6	26	13	7	18
stripes	32	26	26	22	32	24	22	23
xxl	30	23	24	19	32	19	18	21
burnt	25	16	17	8.5	26	13	7.2	18
tiger	25	15	16	6.6	26	12	4.7	17



ASIN : B074LTBWSW

Brand : Liz Lange

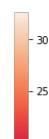
euclidean distance from input : 32.76823337667799

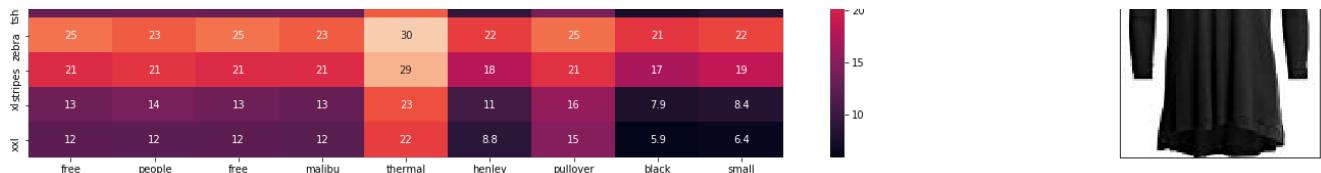
=====



free people free malibu thermal henley pullover black small

	28	27	28	27	33	27	31	25	26
burnt	28	27	28	27	30	21	25	20	21
tiger	24	24	24	22	31	25	26	22	23
umber	26	25	26	26	31	25	26	22	23
zebra	13	13	13	13	23	9.3	14	7.6	8.5
stripes	13	13	13	13	23	9.3	14	7.6	8.5

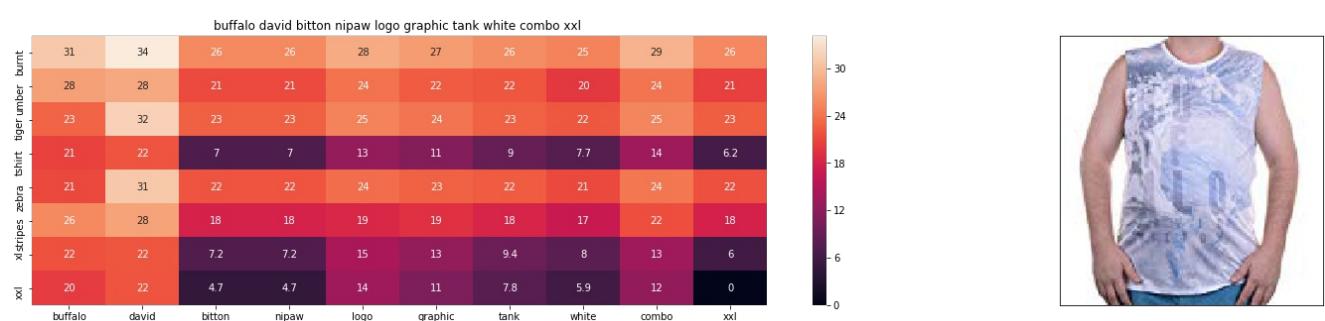




ASIN : B074MXY984

Brand : We The Free

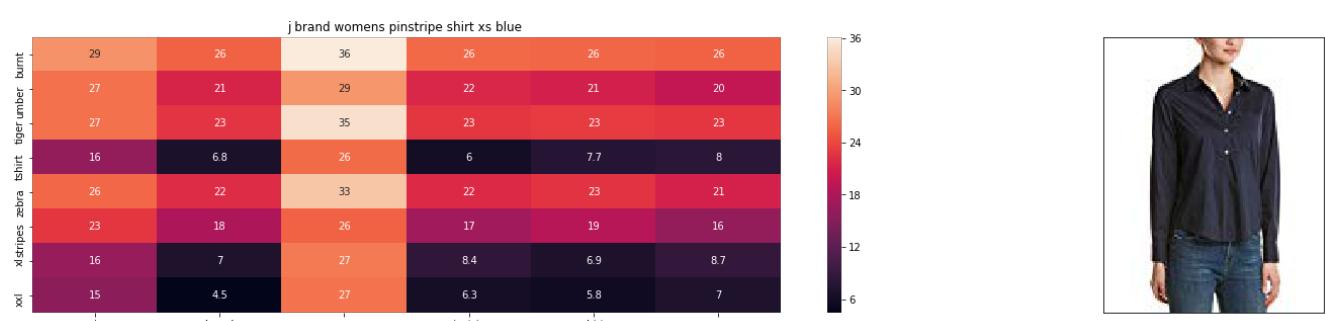
euclidean distance from input : 32.77801740491712



ASIN : B018H5AZXQ

Brand : Buffalo

euclidean distance from input : 32.860679291487756

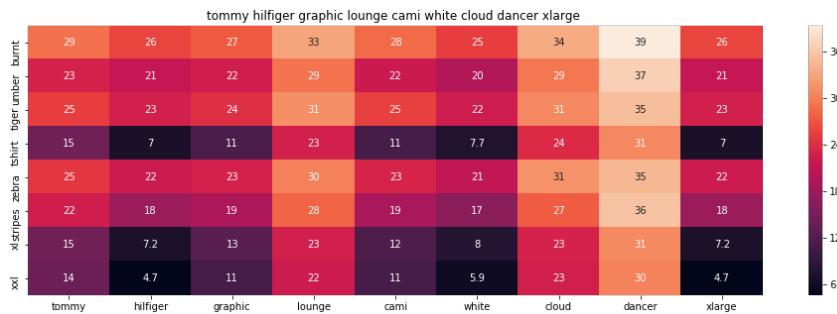


ASIN : B06XYP1X1F

Brand : J Brand Jeans

euclidean distance from input : 32.966362832016394

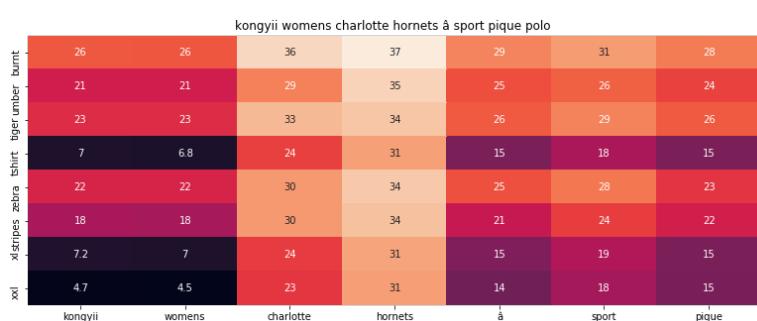




ASIN : B01BMSFYW2

Brand : igertommy hilf

euclidean distance from input : 33.12740295768283



ASIN : B01FJVZST2

Brand : KONGYII

euclidean distance from input : 33.8816550134647

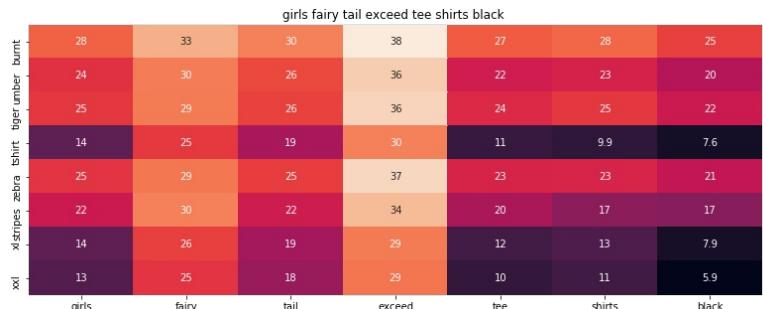


ASIN : B01L7ROZNC

Brand : Bila

euclidean distance from input : 34.041289088114

=====



ASIN : B01L9F153U

Brand : ATYPEMX

euclidean distance from input : 34.55468230978874

=====

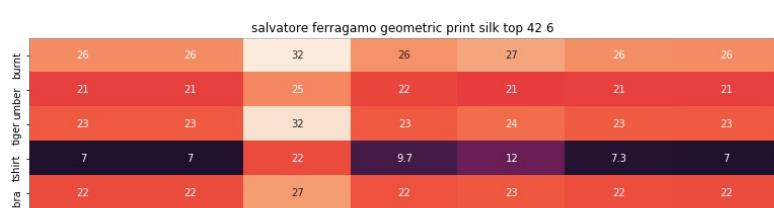


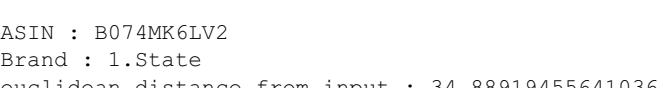
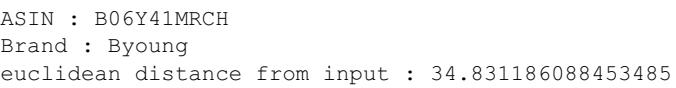
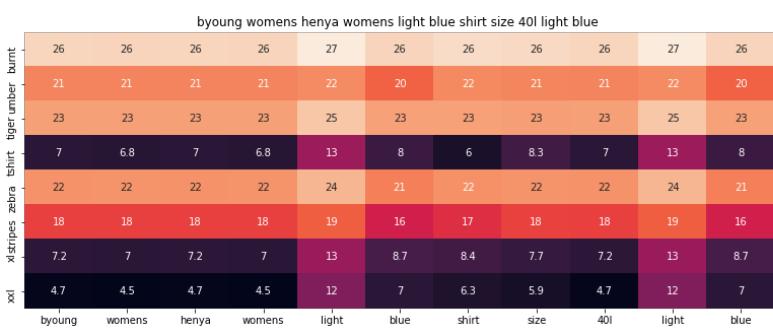
ASIN : B01EXXFS4M

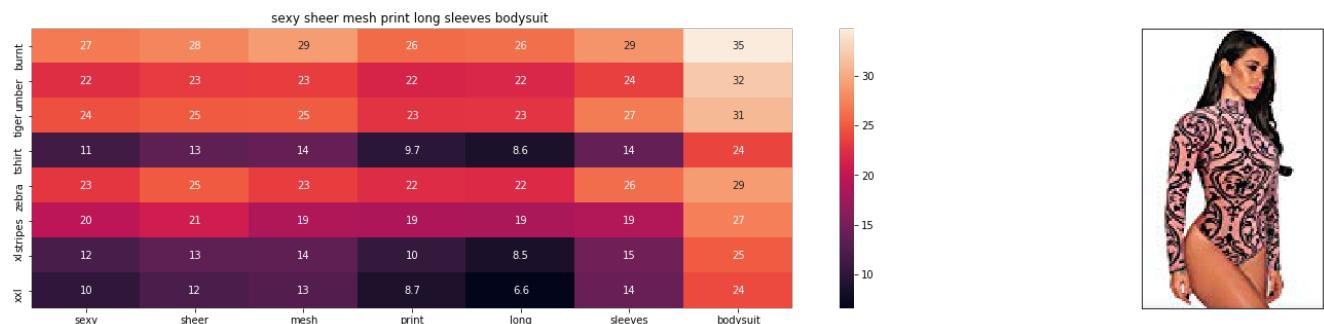
Brand : No Boundaries

euclidean distance from input : 34.62691956607173

=====







ASIN : B074Z5C98D

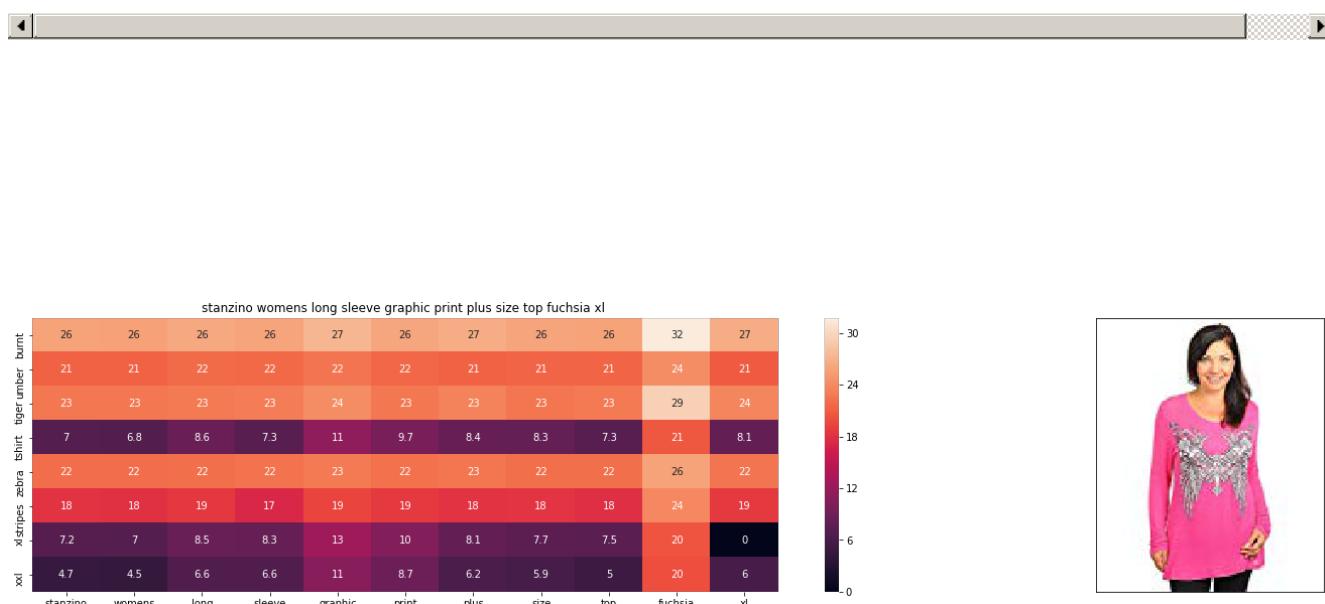
Brand : Ariella's closet

euclidean distance from input : 35.01680579780523

---



---



ASIN : B00DP4VH1I

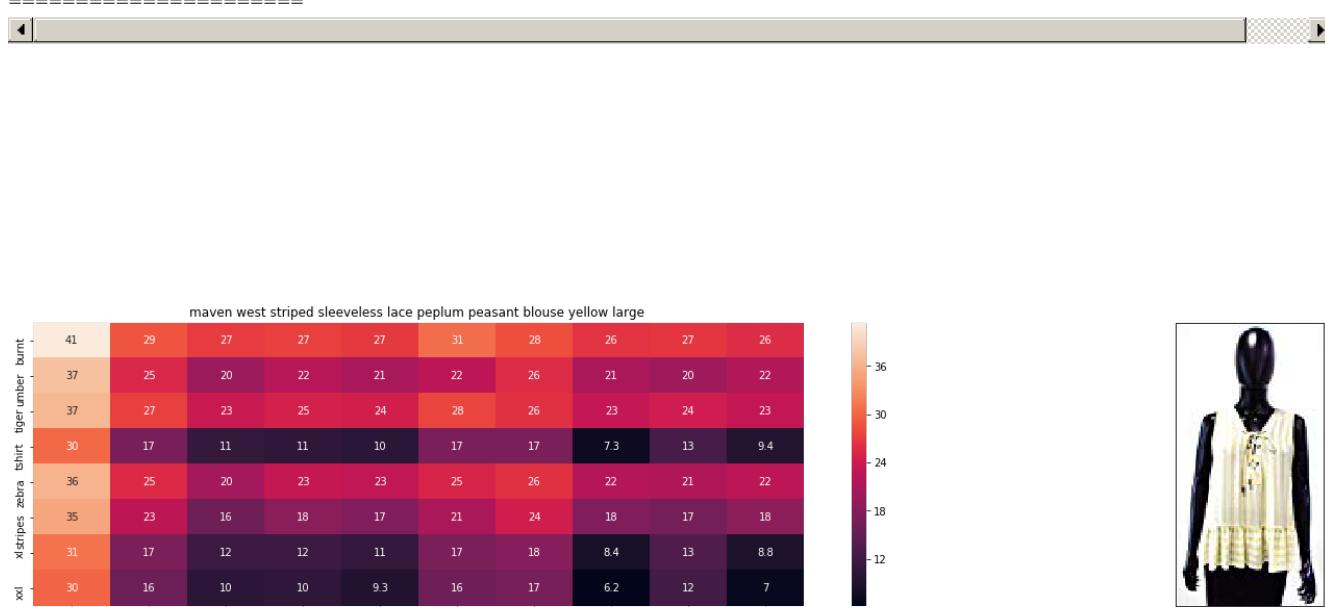
Brand : Stanzino

euclidean distance from input : 35.134891758319604

---



---



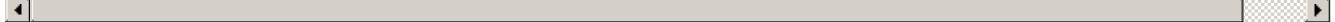
ASIN : B01M8GB3AL

Brand : Maven West

euclidean distance from input : 35.15784165675677

---

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



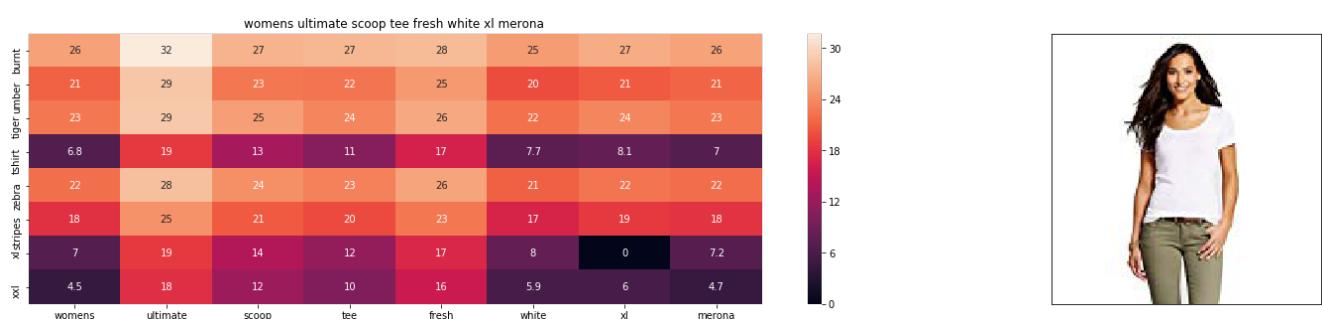
ASIN : B00JMAASRO

Brand : Wotefusi

euclidean distance from input : 35.176515389358784

=====

=====



ASIN : B01G7XE50E

Brand : Merona

euclidean distance from input : 35.19410238095243

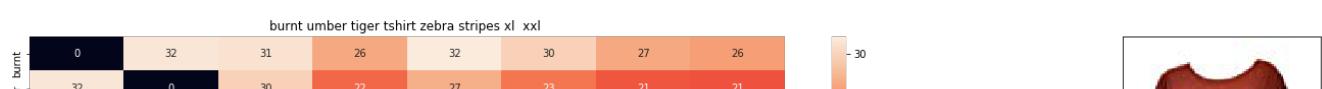
=====

=====



In [165]:

```
idf_w2v_brand(12566, 20, 1, 1, 20)
```

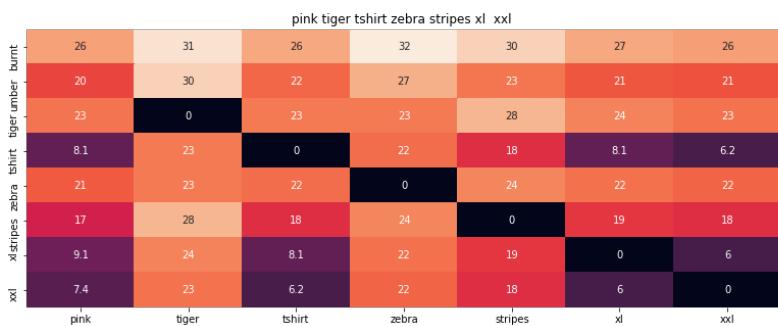




ASIN : B00JXQB5FQ

Brand : Si Row

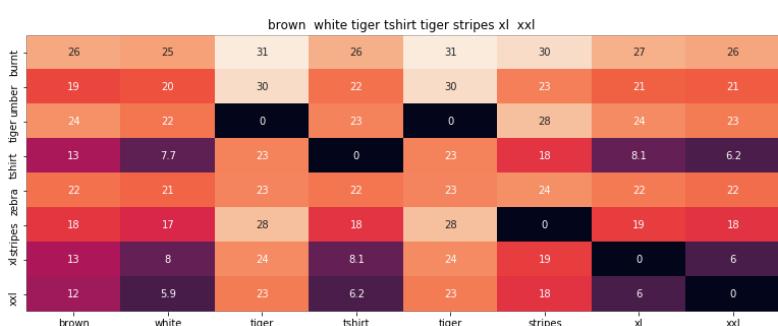
euclidean distance from input : 0.005559962412173098



ASIN : B00JXQASS6

Brand : Si Row

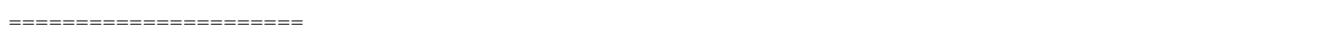
euclidean distance from input : 5.9637245871968245

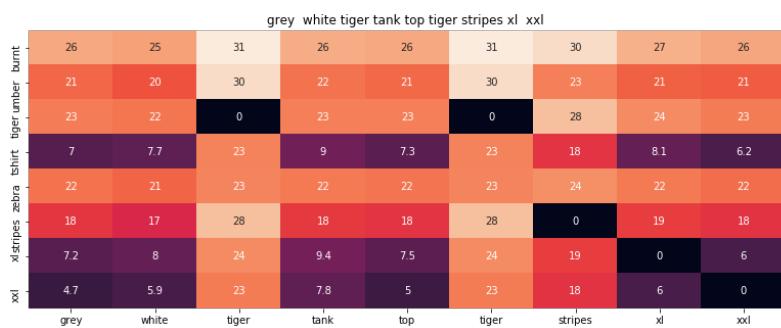


ASIN : B00JXQCWTO

Brand : Si Row

euclidean distance from input : 6.862124529751864





ASIN : B00JXQAFZ2

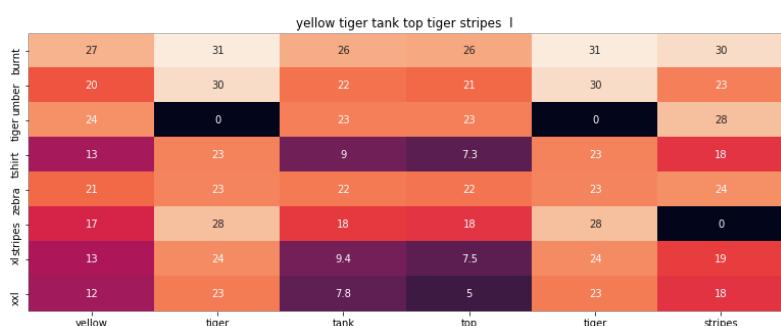
Brand : Si Row

euclidean distance from input : 7.388365416109733

---



---



ASIN : B00JXQAUWA

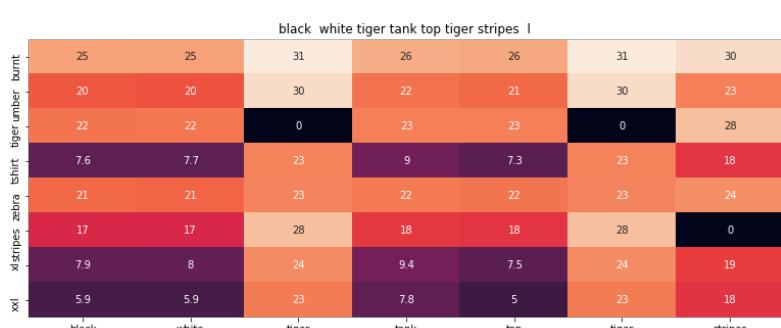
Brand : Si Row

euclidean distance from input : 7.667435143227098

---



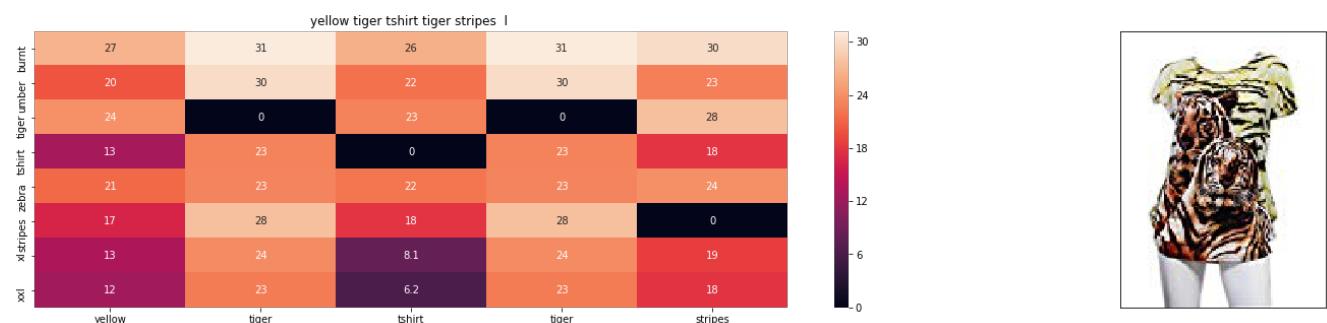
---



ASIN : B00JXQAO94

Brand : Si Row

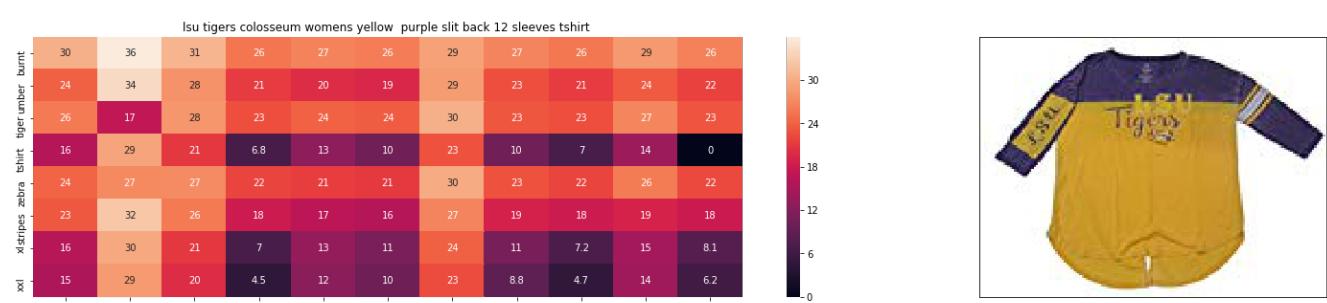
euclidean distance from input : 7.750181822359733



ASIN : B00JXQCUIC

Brand : Si Row

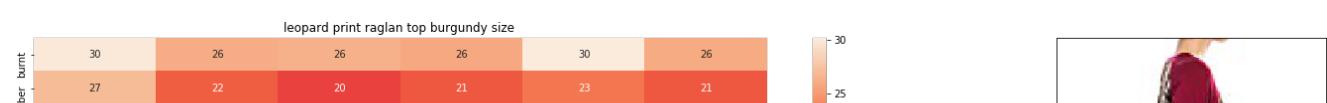
euclidean distance from input : 7.793640240772375

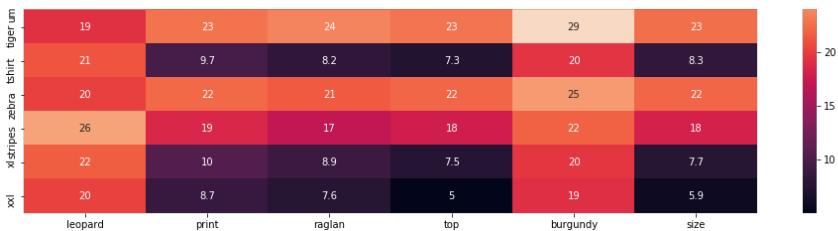


ASIN : B073R5Q8HD

Brand : Colosseum

euclidean distance from input : 8.05921291553978

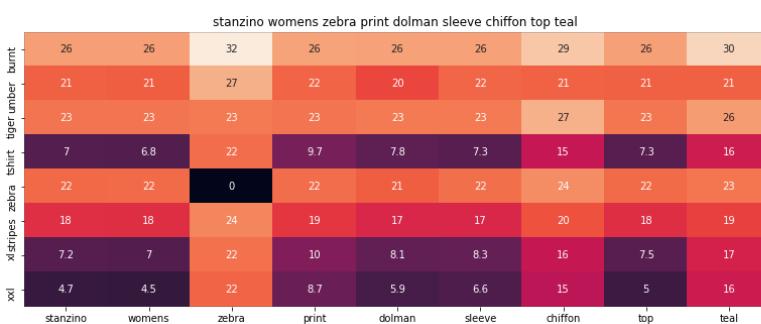




ASIN : B01C60RLDQ

Brand : 1 Mad Fit

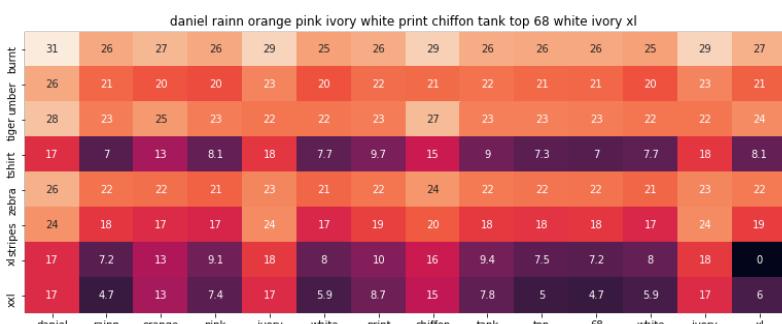
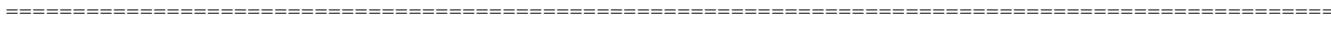
euclidean distance from input : 8.456825984640895



ASIN : B00C0I3U3E

Brand : Stanzino

euclidean distance from input : 8.51727534756881

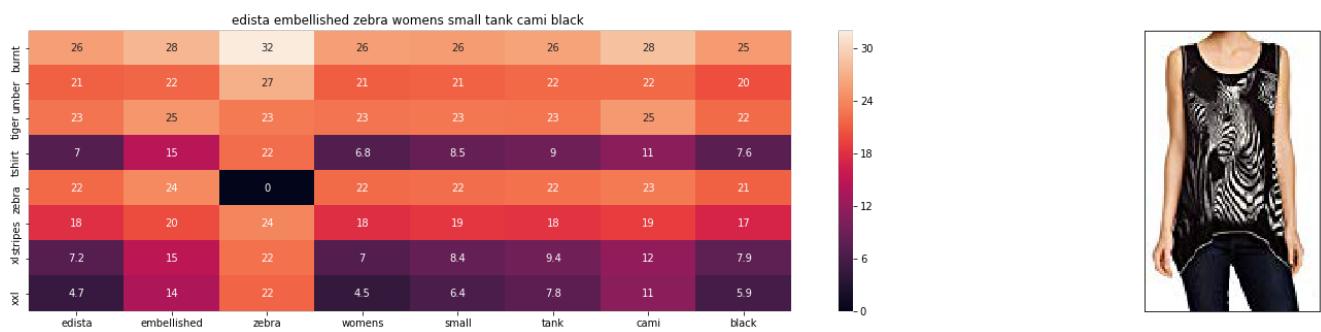


ASIN : B01IPV1SFQ

Brand : Daniel Rainn

euclidean distance from input : 8.535964445634322





ASIN : B074P8MD22

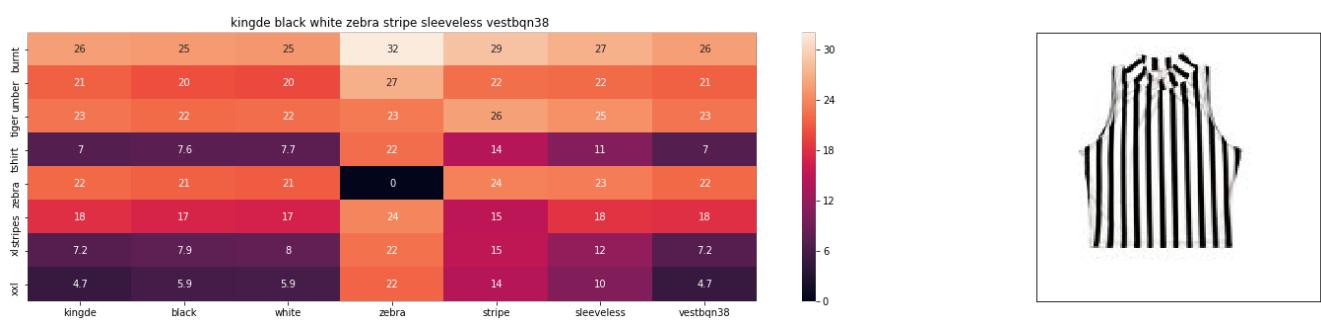
Brand : Edista

euclidean distance from input : 8.550863063411317

---



---



ASIN : B015H41F6G

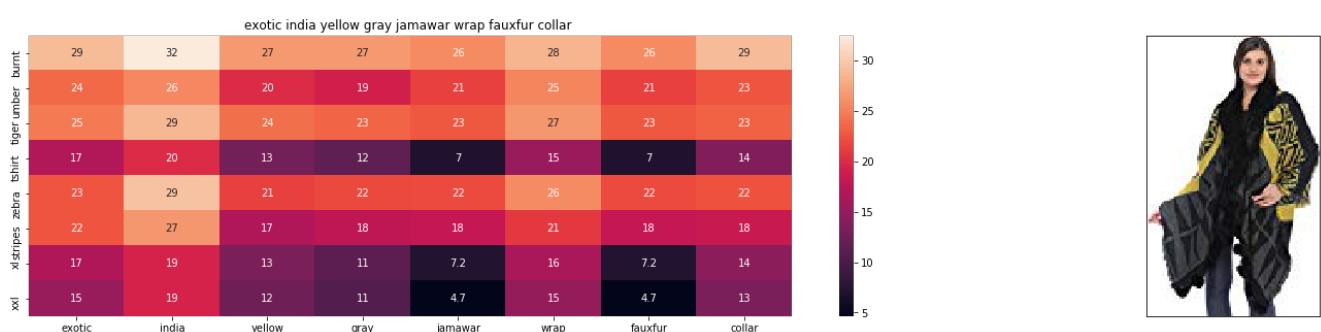
Brand : KINGDE

euclidean distance from input : 8.55733400374113

---



---

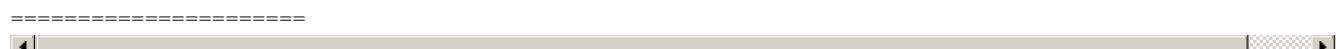


ASIN : B073ZHRBV8

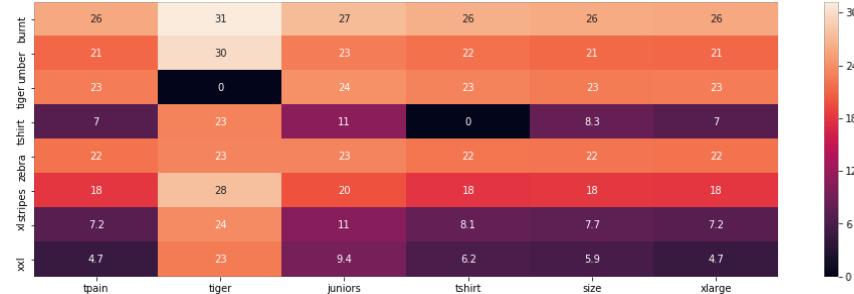
Brand : Exotic India

euclidean distance from input : 8.571198584849524

=====



tpain tiger juniors tshirt size xlarge



ASIN : B01K0H02OG

Brand : Tultex

euclidean distance from input : 8.591895681153682

=====



studio womens burnt orange dolman top size medium



ASIN : B06XSCVFT5

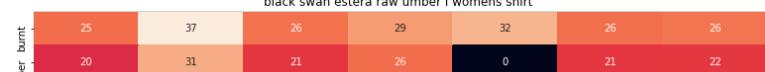
Brand : Studio M

euclidean distance from input : 8.61866035464711

=====



black swan estera raw umber l womens shirt



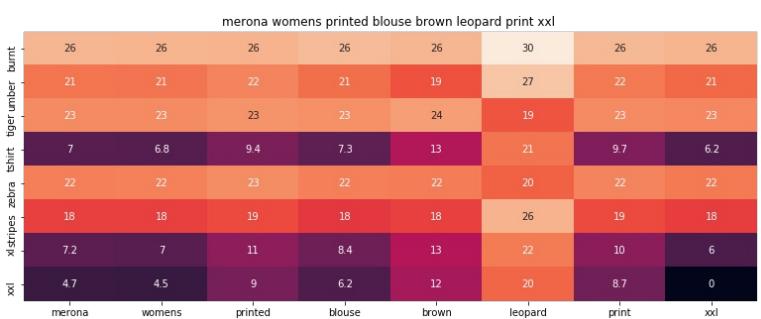


ASIN : B06Y1VN8WQ

Brand : Black Swan

euclidean distance from input : 8.619180505925959

-----  
-----  
-----

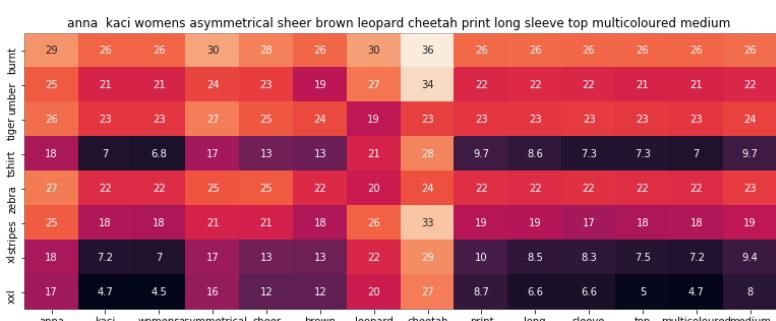


ASIN : B071YF3WDD

Brand : Merona

euclidean distance from input : 8.625227812279393

=====  
=====

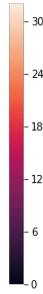


ASTN : BOOKSNTY7Y

Brand : Anna-Kaci

euclidean distance from input : 8.640583766623317

=====



ASIN : B00H8A6ZLI

Brand : Vivian's Fashions

euclidean distance from input : 8.642264487559485

---



---

In [ ]: