

# Amazon Apparel Recommendations

## Loading all required libraries

In [0]:

```

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

```

In [0]:

```
1 data = pd.read_json('tops_fashion.json')
```

In [0]:

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

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

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 e.g: red and black stripes )
4. product\_type\_name (type of the apparel, ex: SHIRT/TSHIRT )
5. medium\_image\_url ( url of the image )
6. title (title of the product.)
7. formatted\_price (price of the product)

In [0]:

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

In [0]:

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

Number of data points : 183138 Number of features: 7

Out[37]:

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

Basic stats for the feature: product\_type\_name

In [0]:

```

1 print(data['product_type_name'].describe())
2
3 # 91.62% (167794/183138) of the products are shirts.

```

```

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

```

In [0]:

```

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

```

```

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

```

In [0]:

```

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

```

Out[40]:

```

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

```

**Basic stats for the feature: brand**

In [0]:

```

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

```

```

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

```

In [0]:

```

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

```

Out[42]:

```

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

```

**Basic stats for the feature: color**

In [0]:

```
1 print(data['color'].describe())
```

```

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

```

In [ ]:

```

1 # we have 7380 unique colors
2 # 7.2% of products are black in color
3 # 64956 of 183138 products have brand information. That's approx 35.4%.

```

In [0]:

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

Out[44]:

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

**Basic stats for the feature: formatted\_price**

In [0]:

```
1 print(data['formatted_price'].describe())
```

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

In [ ]:

```
1 # Only 28,395 (15.5% of whole data) products with price information
```

In [0]:

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

Out[46]:

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

**Basic stats for the feature: title**

In [0]:

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

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

In [ ]:

```
1 # All of the products have a title.
```

In [0]:

```
1 # consider products which have price information
2 data = data.loc[~data['formatted_price'].isnull()]
3 print('Number of data points After eliminating price=NULL :', data.shape[0])
```

Number of data points After eliminating price=NULL : 28395

In [0]:

```
1 # consider products which have color information
2 data = data.loc[~data['color'].isnull()]
3 print('Number of data points After eliminating color=NULL :', data.shape[0])
```

Number of data points After eliminating color=NULL : 28385

## Remove near duplicate items

**Understand about duplicates.**

In [0]:

```
1 data = pd.read_pickle('pickels/28k_apparel_data')
2
3 # find number of products that have duplicate titles.
4 print(sum(data.duplicated('title')))
```

2325

In [ ]:

```
1 # we have 2325 products which have same title but different color
```

### Removing duplicates - 1

In [0]:

```
1 data = pd.read_pickle('pickels/28k_apparel_data')
```

In [0]:

```
1 data.head()
```

Out[103]:

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

In [0]:

```
1 # Remove ALL products with very few words in title
2
3 data_sorted = data[data['title'].apply(lambda x: len(x.split())>4)]
4 print("After removal of products with short description:", data_sorted.shape[0])
```

After removal of products with short description: 27949

In [0]:

```

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

```

Out[105]:

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

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 [ ]:

```
1 # Code to remove the duplicates which differ only at the end.
```

In [0]:

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

In [0]:

```
1 import itertools
2 stage1_dedupe_asins = []
3 i = 0
4 j = 0
5 num_data_points = data_sorted.shape[0]
6 while i < num_data_points and j < num_data_points:
7
8     previous_i = i
9
10    a = data['title'].loc[indices[i]].split()
11
12
13    j = i+1
14    while j < num_data_points:
15
16        b = data['title'].loc[indices[j]].split()
17        length = max(len(a), len(b))
18        count = 0
19        for k in itertools.zip_longest(a, b):
20            if (k[0] == k[1]):
21                count += 1
22
23            them
24            if (length - count) > 2:
25                stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[i]])
26                i = j
27                break
28            else:
29                j += 1
30        if previous_i == i:
31            break
```

In [0]:

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

In [0]:

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

Number of data points : 17593

In [0]:

```
1 data.to_pickle('pickels/17k_apperal_data')
```

**Remove duplicates - 2**

In [ ]:

```
1 # Removing products with similar Title
```

In [0]:

```
1 data = pd.read_pickle('pickels/17k_apperal_data')
```

In [0]:

```
1 indices = []
2 for i, row in data.iterrows():
3     indices.append(i)
4
5 stage2_dedupe_asins = []
6
7 while len(indices) != 0:
8
9     i = indices.pop()
10
11     stage2_dedupe_asins.append(data['asin'].loc[i])
12
13     a = data['title'].loc[i].split()
14
15     for j in indices:
16
17         b = data['title'].loc[j].split()
18
19         length = max(len(a), len(b))
20         count = 0
21         for k in itertools.zip_longest(a, b):
22             if (k[0] == k[1]):
23                 count += 1
24         if (length - count) < 3:
25             indices.remove(j)
```

In [0]:

```
1 data = data.loc[data['asin'].isin(stage2_dedupe_asins)]
```

In [0]:

```
1 print('Number of data points after stage two of dedupe: ',data.shape[0])
```

Number of data points after stage two of dedupe: 16042

## Text pre-processing

In [0]:

```
1 # we use the list of stop words that are downloaded from nltk lib.
2 stop_words = set(stopwords.words('english'))
3 print ('list of stop words:', stop_words)
4
5 def nlp_preprocessing(total_text, index, column):
6     if type(total_text) is not int:
7         string = ""
8         for words in total_text.split():
9
10            word = ("").join(e for e in words if e.isalnum())
11            word = word.lower()
12            if not word in stop_words:
13                string += word + " "
14        data[column][index] = string
```

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

In [0]:

```
1 start_time = time.clock()
2
3 for index, row in data.iterrows():
4     nlp_preprocessing(row['title'], index, 'title')
5
6 print(time.clock() - start_time, "seconds")
```

3.572722000000006 seconds

In [0]:

```
1 data.head()
```

Out[6]:

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

## Stemming

In [0]:

```
1 from nltk.stem.porter import *
2 stemmer = PorterStemmer()
3 print(stemmer.stem('arguing'))
4 print(stemmer.stem('fishing'))
5
6
7 # Stemming on titles did not work very well.
8
```

argu  
fish

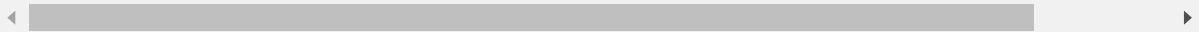
## Text based product similarity

In [0]:

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

Out[10]:

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



In [0]:

```
1 def display_img(url,ax,fig):
2
3     response = requests.get(url)
4     img = Image.open(BytesIO(response.content))
5
6     plt.imshow(img)
7
8
9 def plot_heatmap(keys, values, labels, url, text):
10
11
12     gs = gridspec.GridSpec(2, 2, width_ratios=[4,1], height_ratios=[4,1])
13     fig = plt.figure(figsize=(25,3))
14
15
16     ax = plt.subplot(gs[0])
17
18     ax = sns.heatmap(np.array([values]), annot=np.array([labels]))
19     ax.set_xticklabels(keys)
20     ax.set_title(text)
21     ax = plt.subplot(gs[1])
22
23     ax.grid(False)
24     ax.set_xticks([])
25     ax.set_yticks([])
26
27     display_img(url, ax, fig)
28
29     plt.show()
30
31 def plot_heatmap_image(doc_id, vec1, vec2, url, text, model):
32
33     intersection = set(vec1.keys()) & set(vec2.keys())
34
35     for i in vec2:
36         if i not in intersection:
37             vec2[i]=0
38
39
40     keys = list(vec2.keys())
41
42     values = [vec2[x] for x in vec2.keys()]
43
44     if model == 'bag_of_words':
45         labels = values
46     elif model == 'tfidf':
47         labels = []
48         for x in vec2.keys():
49
50             if x in tfidf_title_vectorizer.vocabulary_:
51                 labels.append(tfidf_title_features[doc_id, tfidf_title_vectorizer.vocabulary_[x]])
52             else:
53                 labels.append(0)
54     elif model == 'idf':
55         labels = []
56         for x in vec2.keys():
57
58             if x in idf_title_vectorizer.vocabulary_:
59                 labels.append(idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[x]])
```

```
60     else:
61         labels.append(0)
62
63     plot_heatmap(keys, values, labels, url, text)
64
65 def text_to_vector(text):
66     word = re.compile(r'\w+')
67     words = word.findall(text)
68
69     return Counter(words)
70
71
72
73 def get_result(doc_id, content_a, content_b, url, model):
74     text1 = content_a
75     text2 = content_b
76
77
78     vector1 = text_to_vector(text1)
79
80     vector2 = text_to_vector(text2)
81
82     plot_heatmap_image(doc_id, vector1, vector2, url, text2, model)
```

## Bag of Words (BoW) on product titles.

In [0]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 title_vectorizer = CountVectorizer()
3 title_features = title_vectorizer.fit_transform(data['title'])
4 title_features.get_shape()
```

Out[17]:

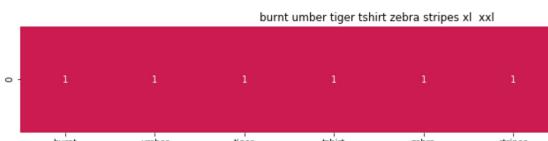
(16042, 12609)

In [0]:

```

1 def bag_of_words_model(doc_id, num_results):
2
3     pairwise_dist = pairwise_distances(title_features,title_features[doc_id])
4
5     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
6
7     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
8     df_indices = list(data.index[indices])
9
10    for i in range(0,len(indices)):
11        get_result(indices[i],data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
12        print('ASIN :',data['asin'].loc[df_indices[i]])
13        print ('Brand:', data['brand'].loc[df_indices[i]])
14        print ('Title:', data['title'].loc[df_indices[i]])
15        print ('Euclidean similarity with the query image :', pdists[i])
16        print('='*60)
17
18
19 bag_of_words_model(12566, 20)

```



ASIN : B00JXQB5FQ

Brand: Si Row

Title: burnt umber tiger tshirt zebra stripes xl xxl

Euclidean similarity with the query image : 0.0

---



ASIN : B00JXQASS6

Brand: Si Row

Title: pink tiger tshirt zebra stripes xl xxl

Euclidean similarity with the query image : 1.73205080757

---



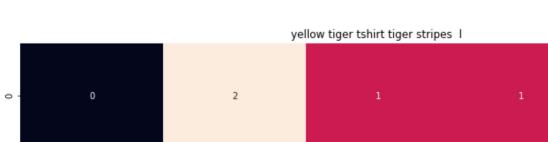
ASIN : B00JXQCWTO

Brand: Si Row

Title: brown white tiger tshirt tiger stripes xl xxl

Euclidean similarity with the query image : 2.44948974278

---



ASIN : B00JXQCUIC

**Brand:** Si Row

**Title:** yellow tiger tshirt tiger stripes 1

**Euclidean similarity with the query image :** 2.64575131106



**ASIN :** B07568NZX4

**Brand:** Rustic Grace

**Title:** believed could tshirt

**Euclidean similarity with the query image :** 3.0



**ASIN :** B01NB0NKRO

**Brand:** Ideology

**Title:** ideology graphic tshirt xl white

**Euclidean similarity with the query image :** 3.0

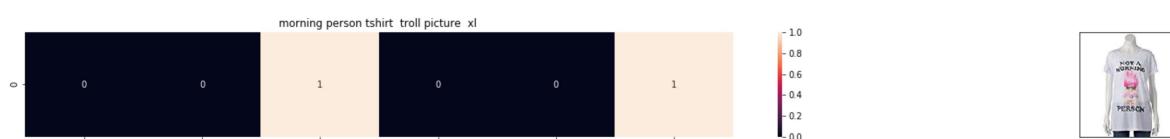


**ASIN :** B00JXQAFZ2

**Brand:** Si Row

**Title:** grey white tiger tank top tiger stripes xl xxl

**Euclidean similarity with the query image :** 3.0



**ASIN :** B01CLS8LMW

**Brand:** Awake

**Title:** morning person tshirt troll picture xl

**Euclidean similarity with the query image :** 3.16227766017

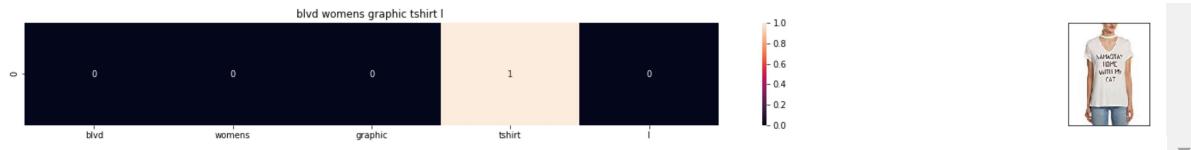


**ASIN :** B01KVZUB6G

**Brand:** Merona

**Title:** merona green gold stripes

**Euclidean similarity with the query image :** 3.16227766017



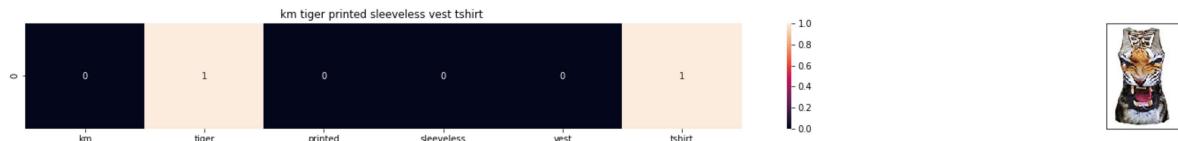
**ASIN : B0733R2CJK**

**Brand: BLVD**

**Title: blvd womens graphic tshirt l**

**Euclidean similarity with the query image : 3.16227766017**

---



**ASIN : B012VQLT6Y**

**Brand: KM T-shirt**

**Title: km tiger printed sleeveless vest tshirt**

**Euclidean similarity with the query image : 3.16227766017**

---



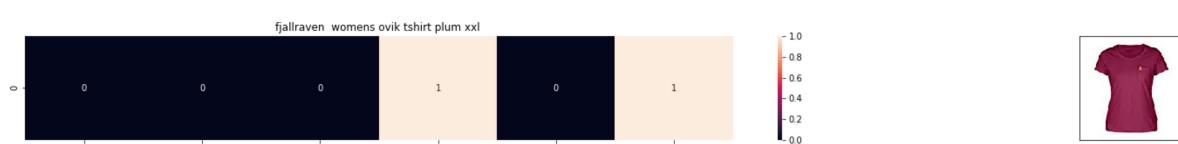
**ASIN : B00JXQC8L6**

**Brand: Si Row**

**Title: blue peacock print tshirt l**

**Euclidean similarity with the query image : 3.16227766017**

---



**ASIN : B06XC3CF6**

**Brand: Fjallraven**

**Title: fjallraven womens ovik tshirt plum xxl**

**Euclidean similarity with the query image : 3.16227766017**

---



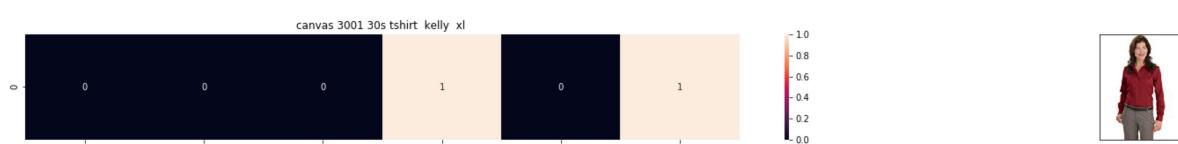
**ASIN : B005IT80BA**

**Brand: Hetalia**

**Title: hetalia us girl tshirt**

**Euclidean similarity with the query image : 3.16227766017**

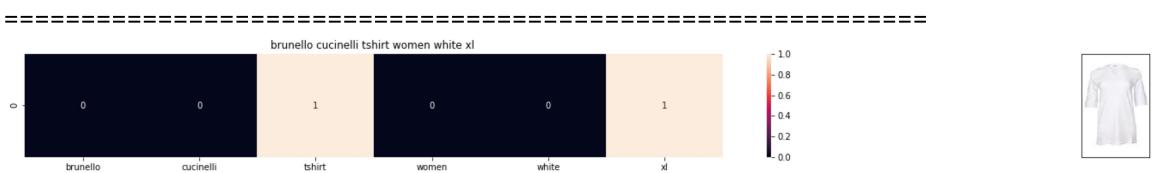
---



**ASIN : B0088PN0LA**

**Brand: Red House**

Title: canvas 3001 30s tshirt kelly xl  
 Euclidean similarity with the query image : 3.16227766017



ASIN : B06X99V6WC

Brand: Brunello Cucinelli

Title: brunello cucinelli tshirt women white xl

Euclidean similarity with the query image : 3.16227766017



ASIN : B06Y1JPW1Q

Brand: Xhilaration

Title: xhilaration womens lace tshirt salmon xxl

Euclidean similarity with the query image : 3.16227766017

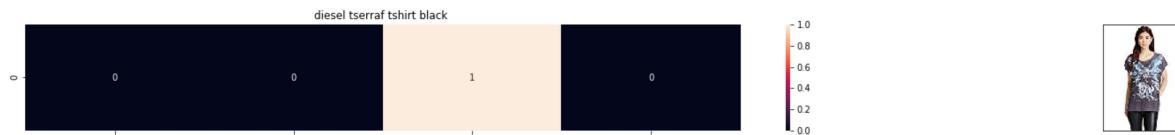


ASIN : B06X6GX6WG

Brand: Animal

Title: animal oceania tshirt yellow

Euclidean similarity with the query image : 3.16227766017

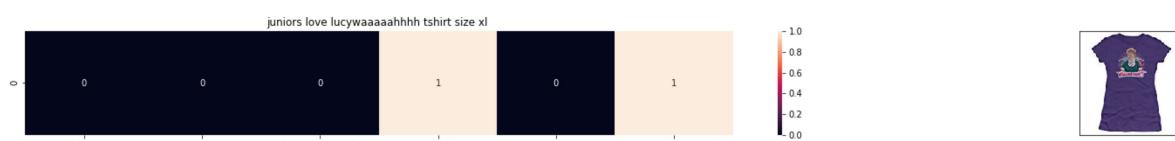


ASIN : B017X8PW9U

Brand: Diesel

Title: diesel tserraf tshirt black

Euclidean similarity with the query image : 3.16227766017



ASIN : B00IAA4JIQ

Brand: I Love Lucy

Title: juniors love lucyaaaaahhhh tshirt size xl

Euclidean similarity with the query image : 3.16227766017

## TF-IDF based product similarity

In [0]:

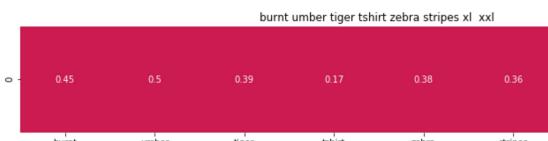
```
1 tfidf_title_vectorizer = TfidfVectorizer(min_df = 0)
2 tfidf_title_features = tfidf_title_vectorizer.fit_transform(data['title'])
```

In [0]:

```

1 def tfidf_model(doc_id, num_results):
2
3     pairwise_dist = pairwise_distances(tfidf_title_features, tfidf_title_features[doc_id])
4
5     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
6
7     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
8
9
10    df_indices = list(data.index[indices])
11
12    for i in range(0, len(indices)):
13
14        get_result(indices[i], data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
15        print('ASIN : ', data['asin'].loc[df_indices[i]])
16        print('BRAND : ', data['brand'].loc[df_indices[i]])
17        print('Euclidean distance from the given image : ', pdists[i])
18        print('*125')
19 tfidf_model(12566, 20)

```



ASIN : B00JXQB5FQ

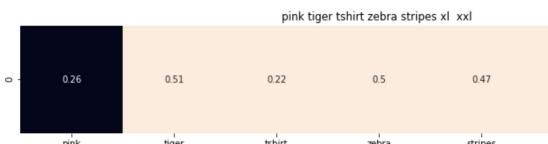
BRAND : Si Row

Euclidean distance from the given image : 0.0

---



---



ASIN : B00JXQASS6

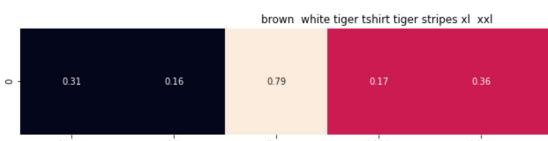
BRAND : Si Row

Euclidean distance from the given image : 0.753633191245

---



---



ASIN : B00JXQCWTO

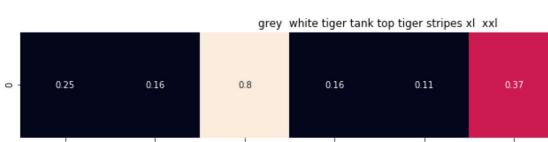
BRAND : Si Row

Euclidean distance from the given image : 0.935764394377

---



---



ASIN : B00JXQAFZ2

**BRAND : Si Row**

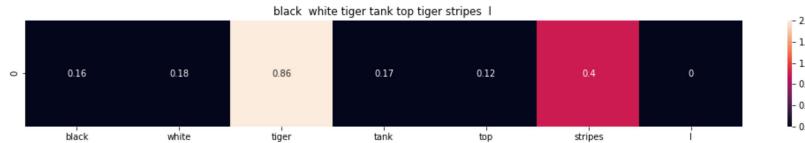
**Eucliden distance from the given image : 0.95861535242**



**ASIN : B00JXQCUIC**

**BRAND : Si Row**

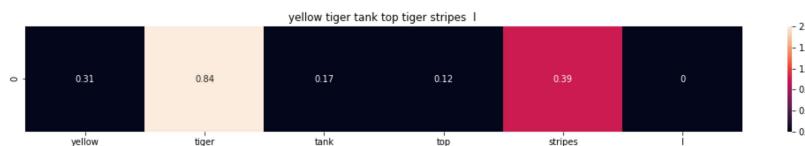
**Eucliden distance from the given image : 1.00007496145**



**ASIN : B00JXQA094**

**BRAND : Si Row**

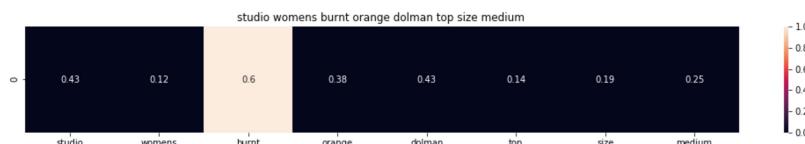
**Eucliden distance from the given image : 1.02321555246**



**ASIN : B00JXQAUWA**

**BRAND : Si Row**

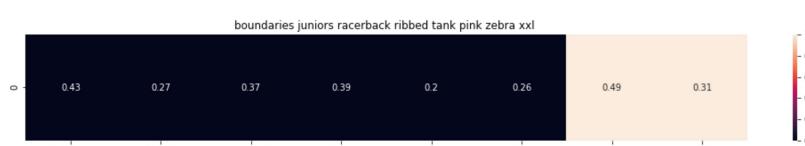
**Eucliden distance from the given image : 1.0319918463**



**ASIN : B06XSCVFT5**

**BRAND : Studio M**

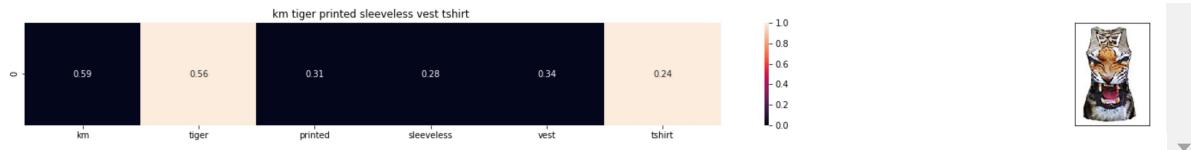
**Eucliden distance from the given image : 1.21068436704**



**ASIN : B06Y2GTYPM**

**BRAND : No Boundaries**

**Eucliden distance from the given image : 1.21216838107**



ASIN : B012VQLT6Y

BRAND : KM T-shirt

Eucliden distance from the given image : 1.21979064028

---



---



ASIN : B06Y1VN8WQ

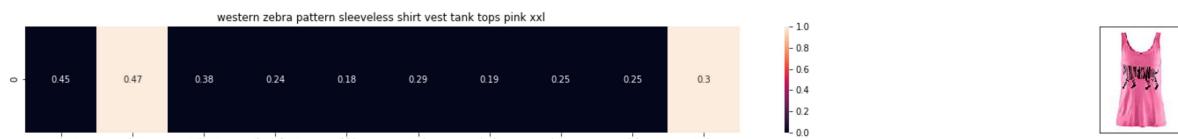
BRAND : Black Swan

Eucliden distance from the given image : 1.220684966

---



---



ASIN : B00Z6HEXWI

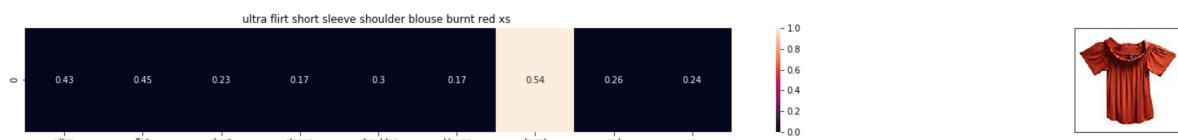
BRAND : Black Temptation

Eucliden distance from the given image : 1.22128139212

---



---



ASIN : B074TR12BH

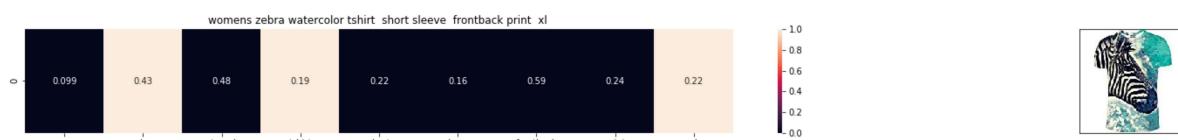
BRAND : Ultra Flirt

Eucliden distance from the given image : 1.23133640946

---



---



ASIN : B072R2JXKW

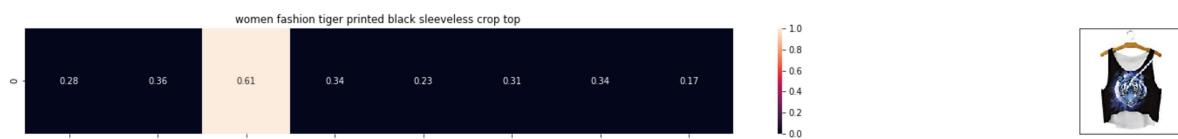
BRAND : WHAT ON EARTH

Eucliden distance from the given image : 1.23184519726

---



---

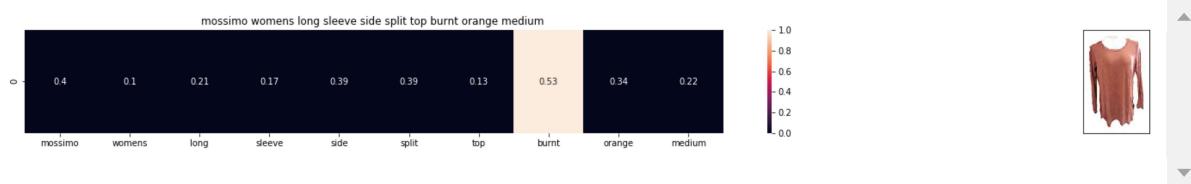


ASIN : B074T8ZYGX

BRAND : MKP Crop Top

Eucliden distance from the given image : 1.23406074574

---



ASIN : B071ZDF6T2

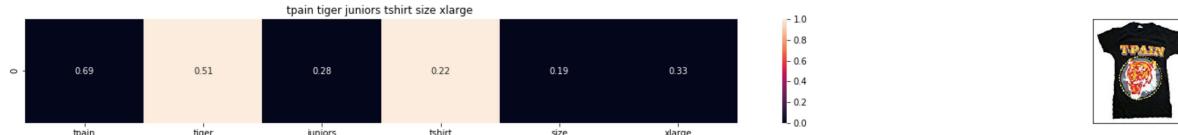
BRAND : Mossimo

Eucliden distance from the given image : 1.23527855777

---



---



ASIN : B01K0H020G

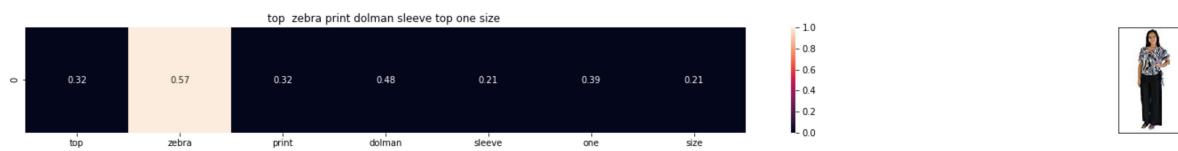
BRAND : Tultex

Eucliden distance from the given image : 1.23645729881

---



---



ASIN : B00H8A6ZLI

BRAND : Vivian's Fashions

Eucliden distance from the given image : 1.24996155053

---



---



ASIN : B010NN9RX0

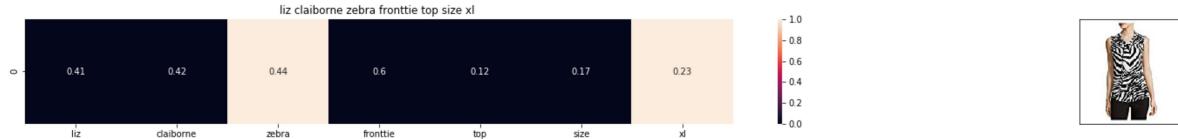
BRAND : YICHUN

Eucliden distance from the given image : 1.25354614209

---



---



ASIN : B06XBY5QXL

BRAND : Liz Claiborne

Eucliden distance from the given image : 1.25388329384

---



---

## IDF based product similarity

In [0]:

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

In [0]:

```
1 def nContaining(word):  
2     return sum(1 for blob in data['title'] if word in blob.split())  
3  
4 def idf(word):  
5     return math.log(data.shape[0] / (nContaining(word)))  
6  
7
```

In [0]:

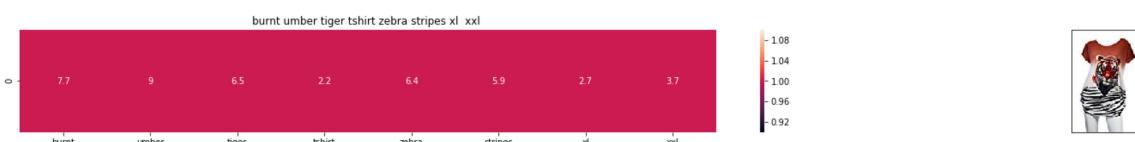
```
1 idf_title_features = idf_title_features.astype(np.float)  
2  
3 for i in idf_title_vectorizer.vocabulary_.keys():  
4     idf_val = idf(i)  
5  
6     for j in idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0]:  
7         idf_title_features[j,idf_title_vectorizer.vocabulary_[i]] = idf_val  
8  
9  
10  
11
```

In [0]:

```

1 def idf_model(doc_id, num_results):
2     pairwise_dist = pairwise_distances(idf_title_features,idf_title_features[doc_id])
3
4     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
5
6     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
7
8
9     df_indices = list(data.index[indices])
10
11    for i in range(0,len(indices)):
12        get_result(indices[i],data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
13        print('ASIN :',data['asin'].loc[df_indices[i]])
14        print('Brand :',data['brand'].loc[df_indices[i]])
15        print ('euclidean distance from the given image :', pdists[i])
16        print('*'*125)
17
18
19
20 idf_model(12566,20)

```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from the given image : 0.0



ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from the given image : 12.2050713112

## Text Semantics based product similarity

In [0]:

```

1 from gensim.models import Word2Vec
2 from gensim.models import KeyedVectors
3 import pickle
4
5 with open('word2vec_model', 'rb') as handle:
6     model = pickle.load(handle)
7

```

In [0]:

```
1 def get_word_vec(sentence, doc_id, m_name):
2
3     vec = []
4     for i in sentence.split():
5         if i in vocab:
6             if m_name == 'weighted' and i in idf_title_vectorizer.vocabulary_:
7                 vec.append(idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[i]])
8             elif m_name == 'avg':
9                 vec.append(model[i])
10            else:
11
12                vec.append(np.zeros(shape=(300,)))
13    return np.array(vec)
14
15 def get_distance(vec1, vec2):
16
17     final_dist = []
18     for i in vec1:
19         dist = []
20         for j in vec2:
21             dist.append(np.linalg.norm(i-j))
22         final_dist.append(np.array(dist))
23    return np.array(final_dist)
24
25
26 def heat_map_w2v(sentence1, sentence2, url, doc_id1, doc_id2, model):
27
28     s1_vec = get_word_vec(sentence1, doc_id1, model)
29     s2_vec = get_word_vec(sentence2, doc_id2, model)
30
31     s1_s2_dist = get_distance(s1_vec, s2_vec)
32
33
34
35
36     gs = gridspec.GridSpec(2, 2, width_ratios=[4,1],height_ratios=[2,1])
37     fig = plt.figure(figsize=(15,15))
38
39     ax = plt.subplot(gs[0])
40
41     ax = sns.heatmap(np.round(s1_s2_dist,4), annot=True)
42
43     ax.set_xticklabels(sentence2.split())
44
45     ax.set_yticklabels(sentence1.split())
46
47     ax.set_title(sentence2)
48
49     ax = plt.subplot(gs[1])
50
51     ax.grid(False)
52     ax.set_xticks([])
53     ax.set_yticks([])
54     display_img(url, ax, fig)
55
56     plt.show()
```

In [0]:

```

1 vocab = model.keys()
2
3 def build_avg_vec(sentence, num_features, doc_id, m_name):
4
5     featureVec = np.zeros((num_features,), dtype="float32")
6
7     nwords = 0
8
9     for word in sentence.split():
10        nwords += 1
11        if word in vocab:
12            if m_name == 'weighted' and word in idf_title_vectorizer.vocabulary_:
13                featureVec = np.add(featureVec, idf_title_features[doc_id, idf_title_ve
14                elif m_name == 'avg':
15                    featureVec = np.add(featureVec, model[word])
16            if(nwords>0):
17                featureVec = np.divide(featureVec, nwords)
18
19    return featureVec

```

## Average Word2Vec product similarity.

In [0]:

```

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

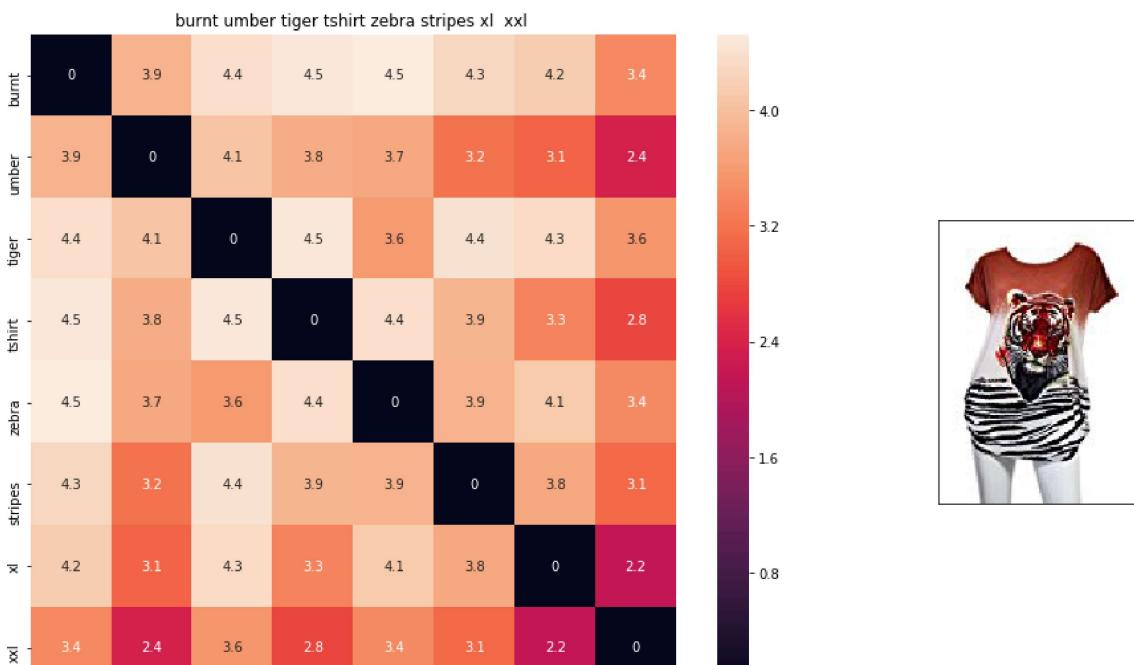
```

In [0]:

```

1 def avg_w2v_model(doc_id, num_results):
2
3     pairwise_dist = pairwise_distances(w2v_title, w2v_title[doc_id].reshape(1,-1))
4
5     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
6     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
7
8     df_indices = list(data.index[indices])
9
10    for i in range(0, len(indices)):
11        heat_map_w2v(data['title'].loc[df_indices[0]],data['title'].loc[df_indices[i]])
12        print('ASIN :',data['asin'].loc[df_indices[i]])
13        print('BRAND :',data['brand'].loc[df_indices[i]])
14        print ('euclidean distance from given input image :', pdists[i])
15        print('*125')
16
17
18 avg_w2v_model(12566, 20)
19

```



## IDF weighted Word2Vec for product similarity

In [0]:

```

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

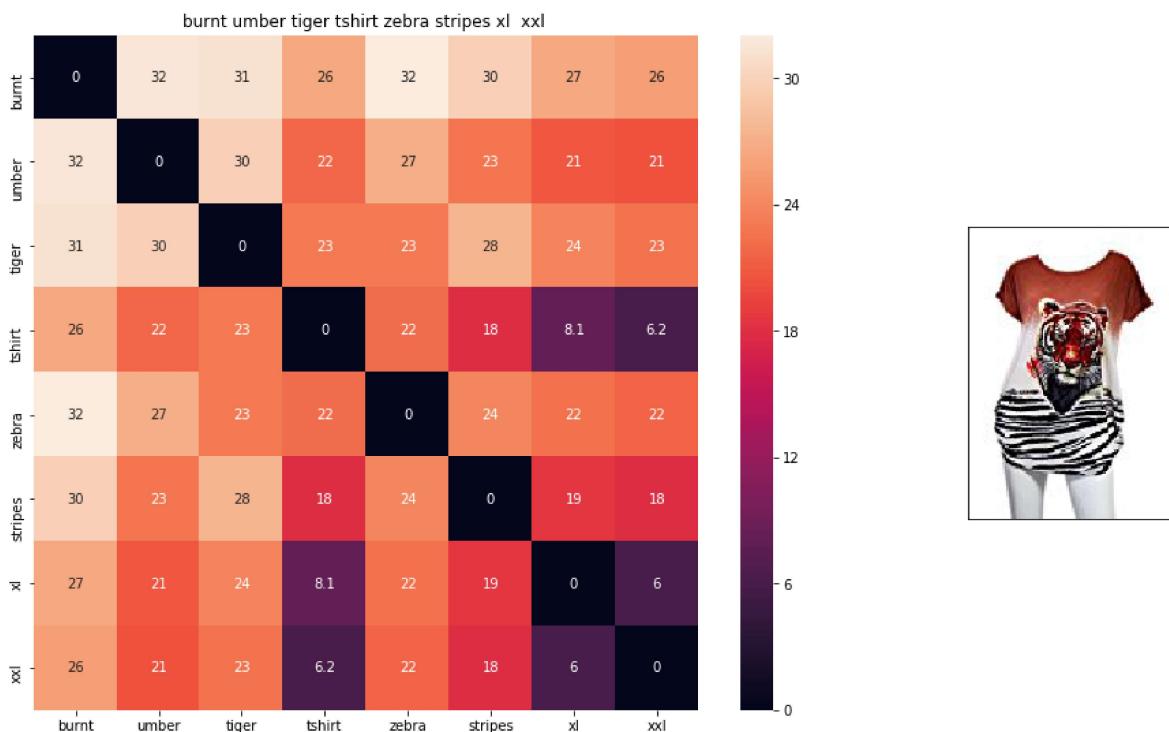
```

In [0]:

```

1 def weighted_w2v_model(doc_id, num_results):
2
3     pairwise_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1, -1))
4
5     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
6     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
7
8     df_indices = list(data.index[indices])
9
10    for i in range(0, len(indices)):
11        heat_map_w2v(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
12        print('ASIN :', data['asin'].loc[df_indices[i]])
13        print('Brand :', data['brand'].loc[df_indices[i]])
14        print('euclidean distance from input :', pdists[i])
15        print('*'*125)
16
17 weighted_w2v_model(12566, 20)

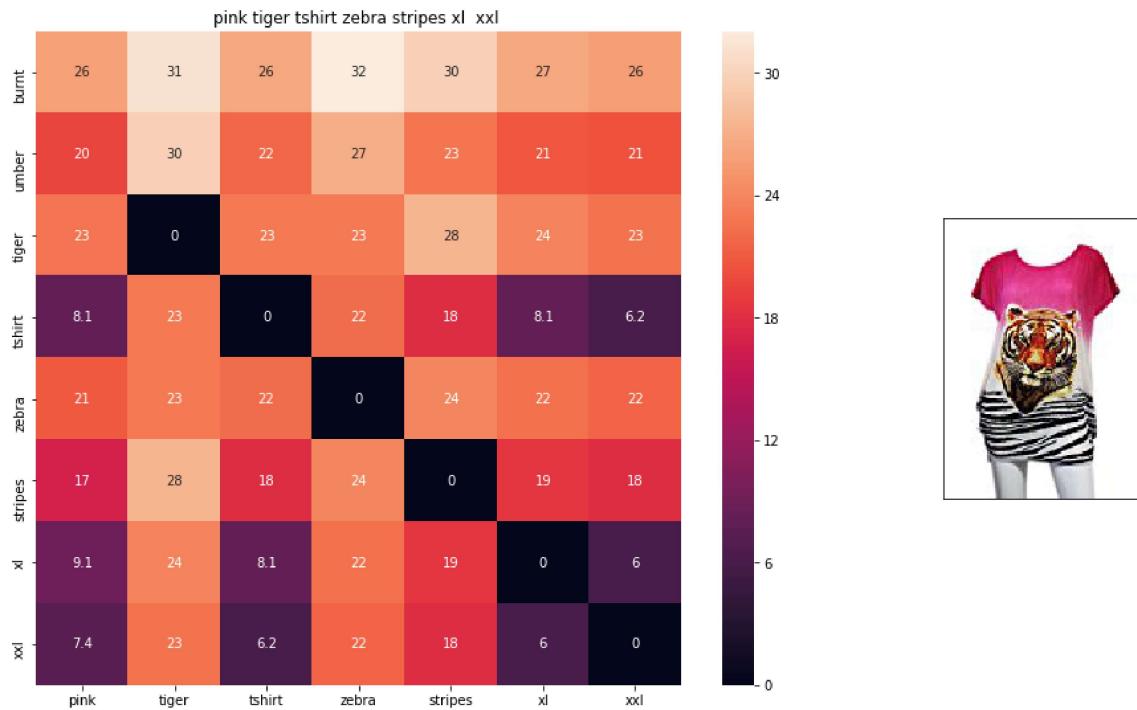
```



ASIN : B00JXQB5FQ

Brand : Si Row

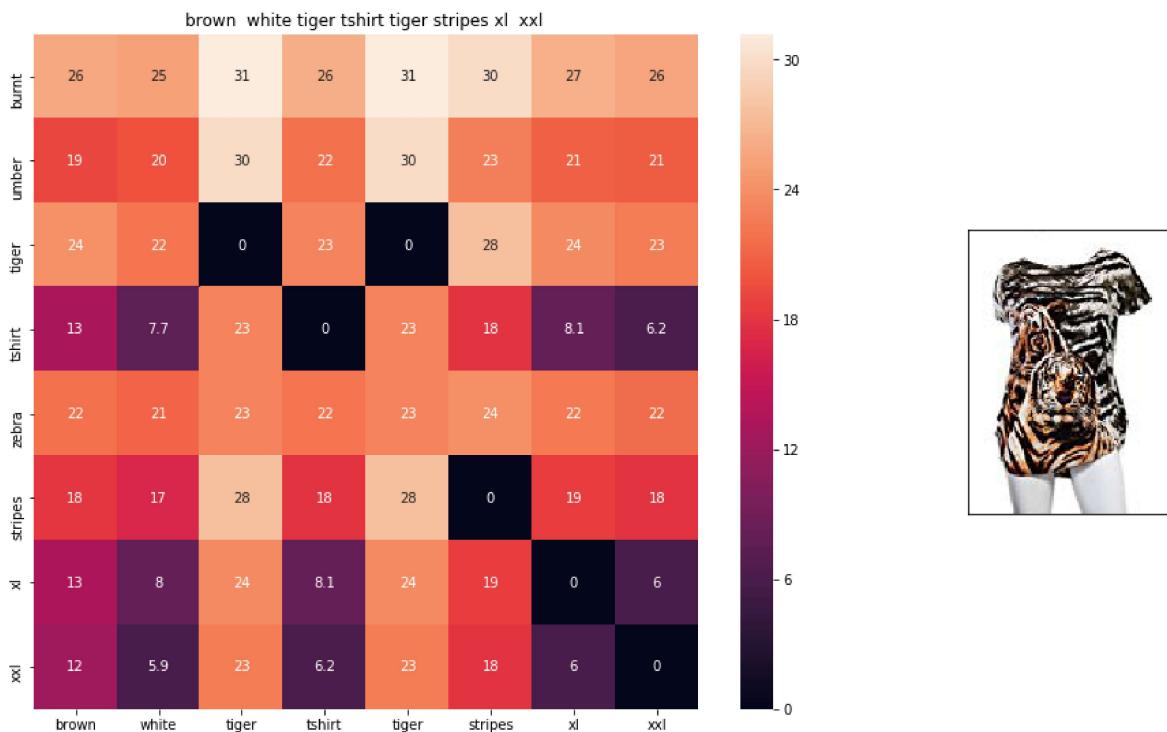
euclidean distance from input : 0.00390625



ASIN : B00JXQASS6

Brand : Si Row

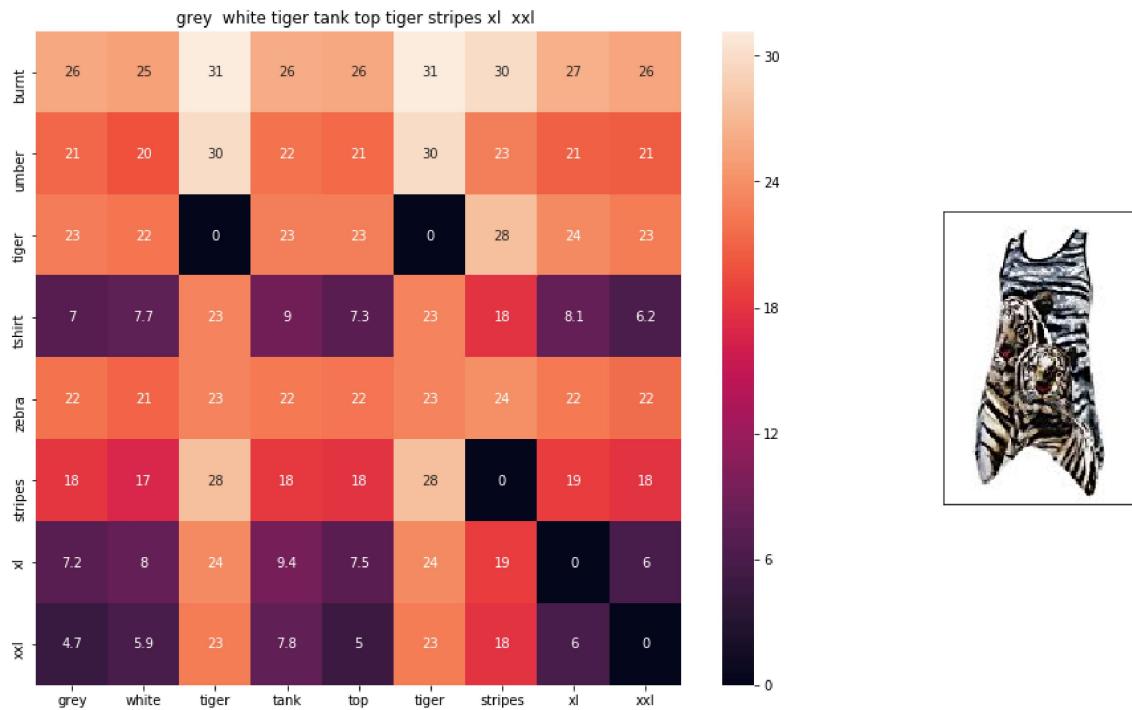
euclidean distance from input : 4.06389



ASIN : B00JXQCWT0

Brand : Si Row

euclidean distance from input : 4.77094



ASIN : B00JXQAFZ2

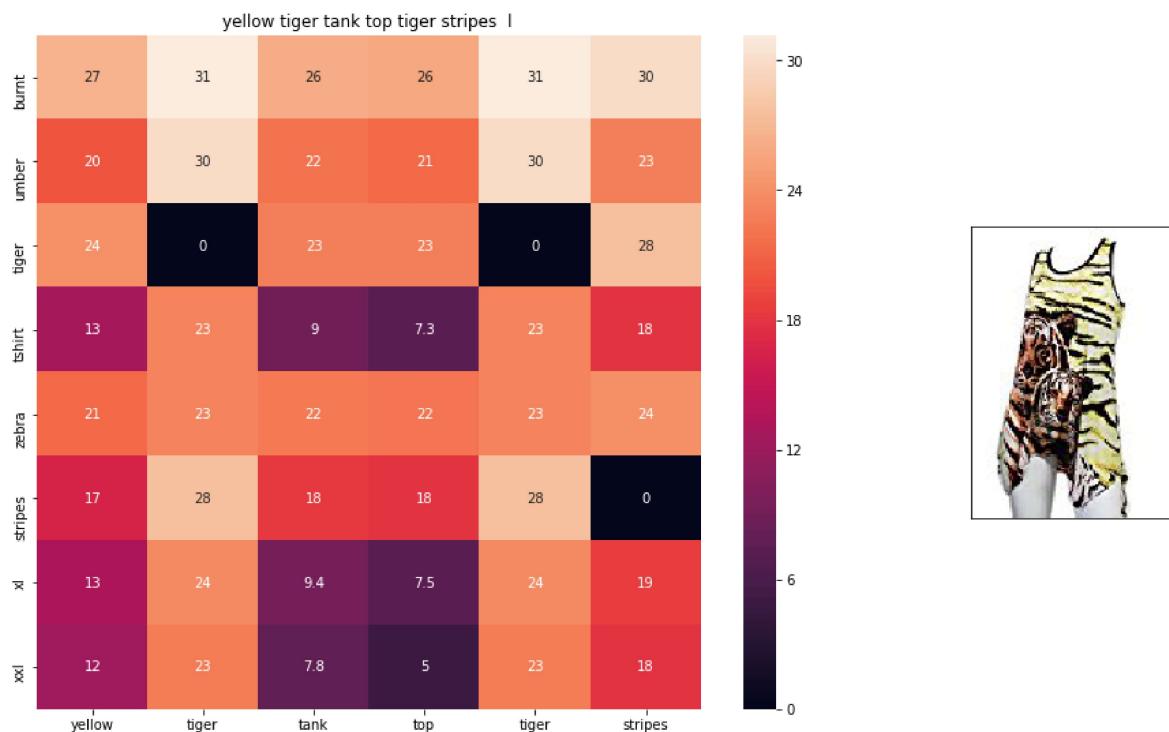
Brand : Si Row

euclidean distance from input : 5.36016

---



---



ASIN : B00JXQAUWA

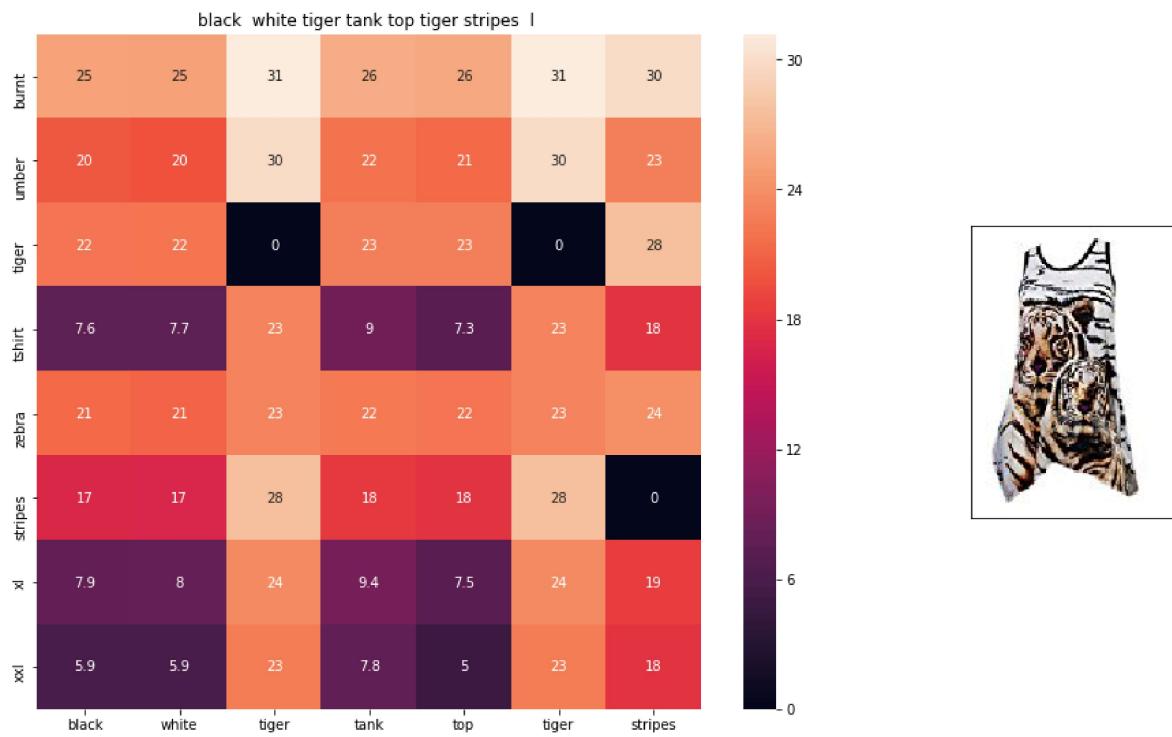
Brand : Si Row

euclidean distance from input : 5.68952

---



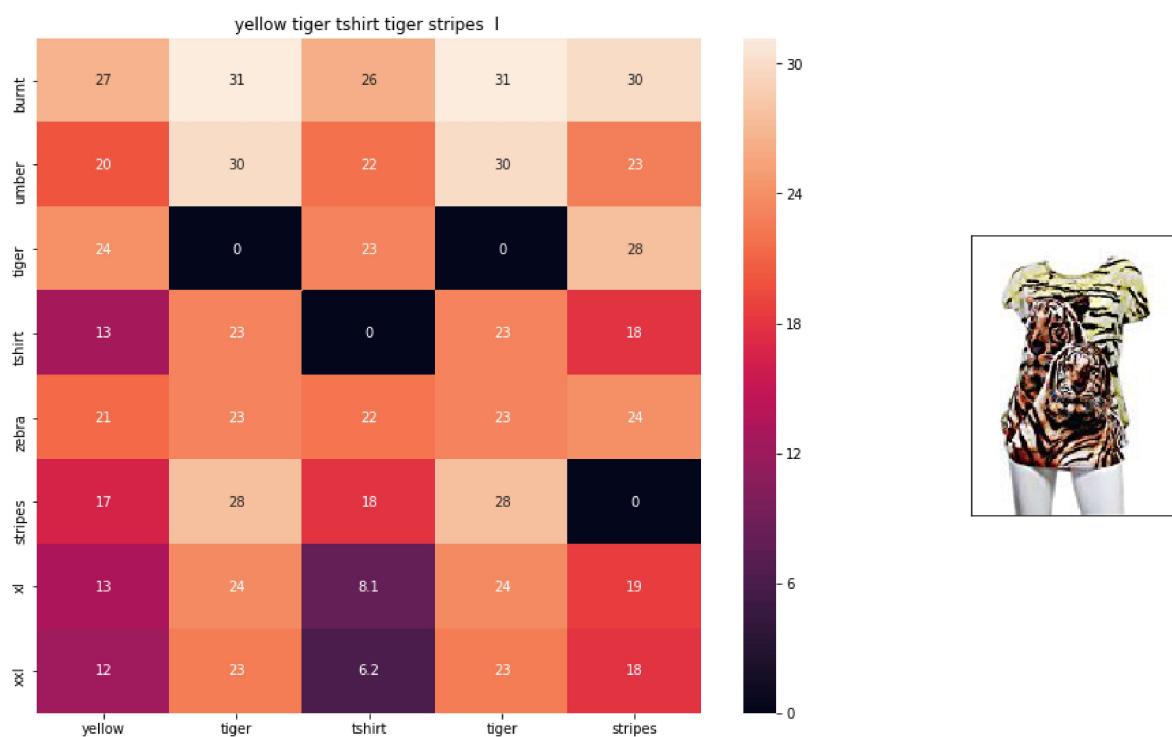
---



ASIN : B00JXQA094

Brand : Si Row

euclidean distance from input : 5.69302



ASIN : B00JXQCUIC

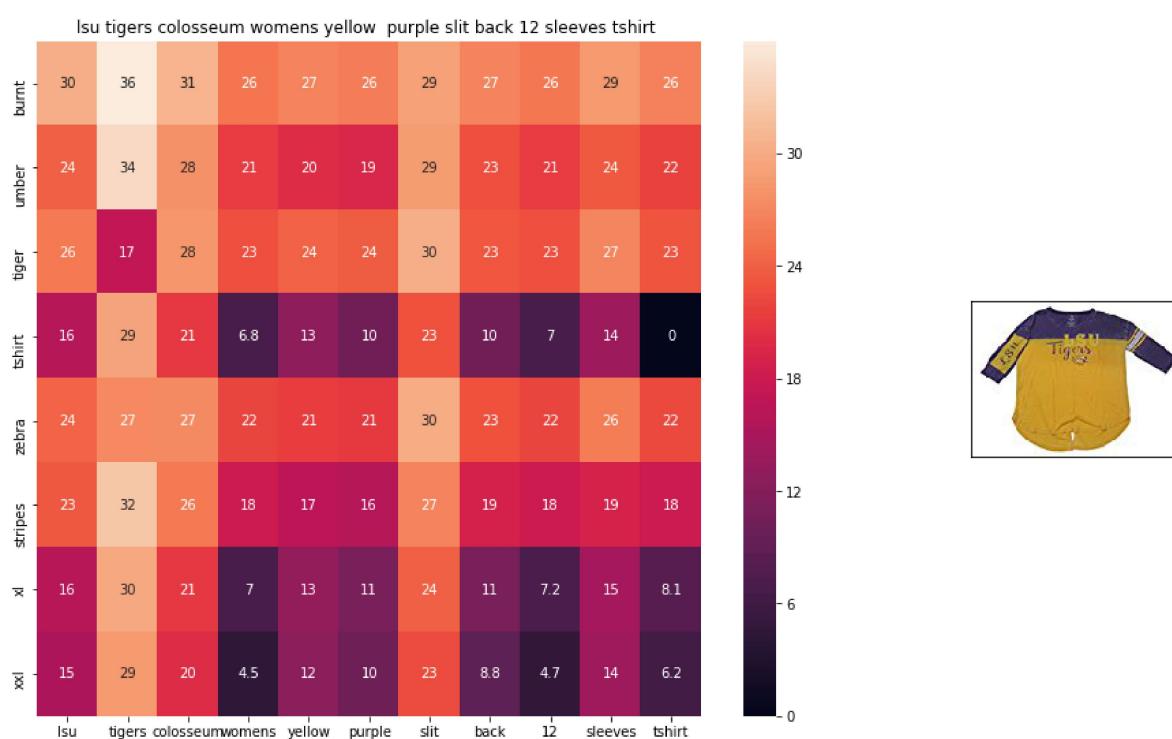
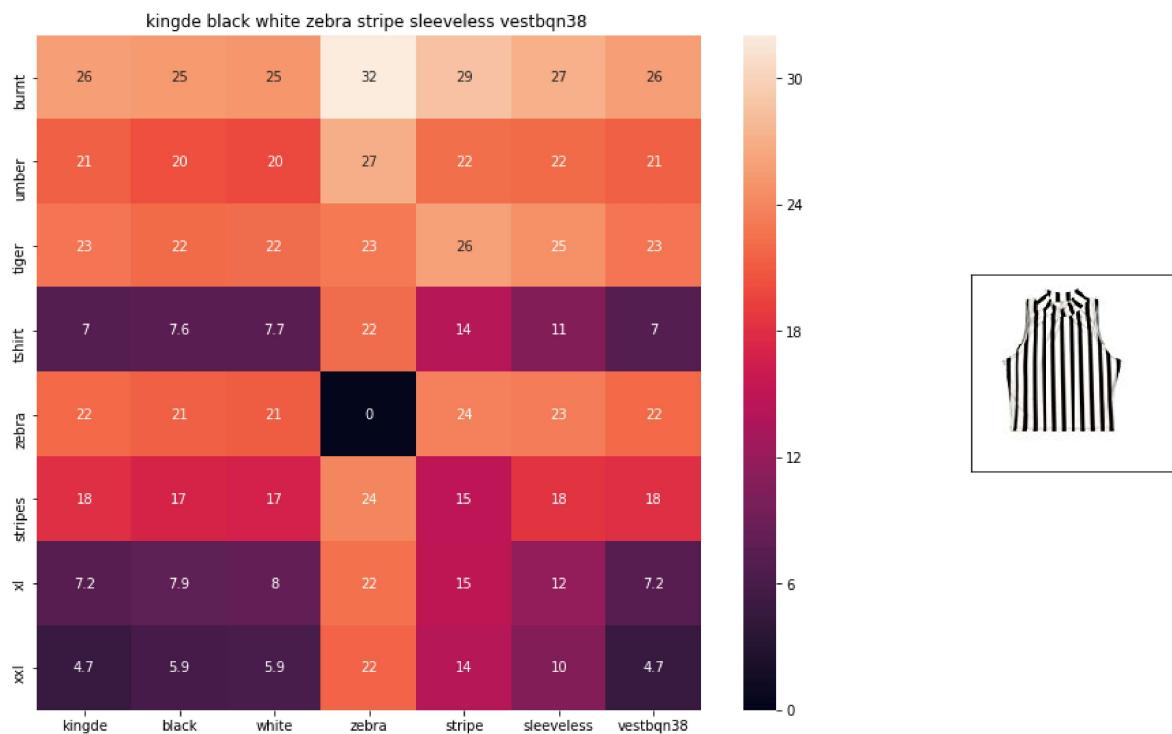
Brand : Si Row

euclidean distance from input : 5.89344

---



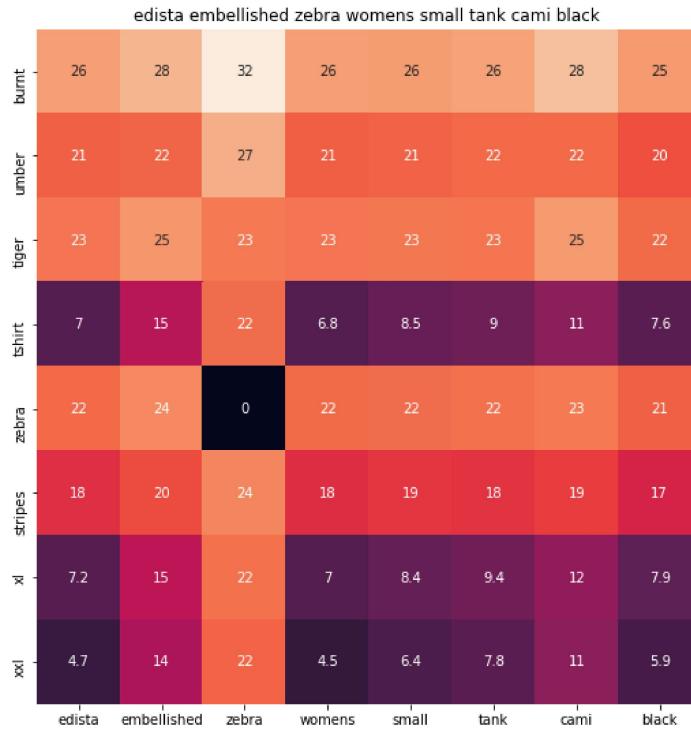
---



ASIN : B073R5Q8HD

Brand : Colosseum

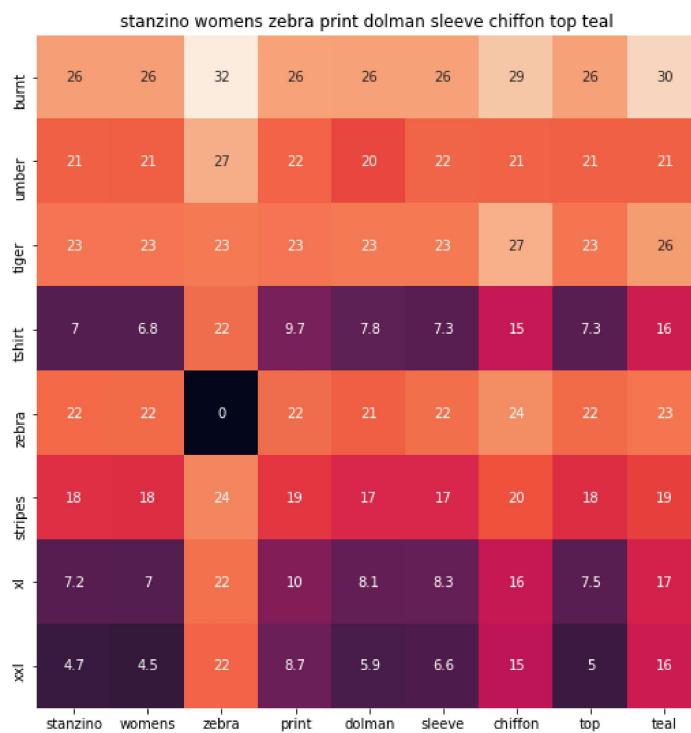
euclidean distance from input : 6.25671



ASIN : B074P8MD22

Brand : Edista

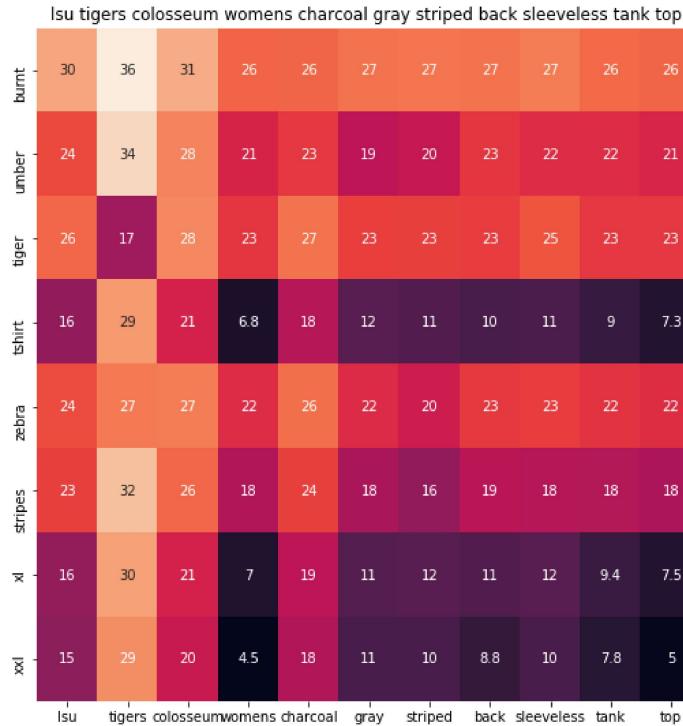
euclidean distance from input : 6.3922



ASIN : B00C0I3U3E

Brand : Stanzino

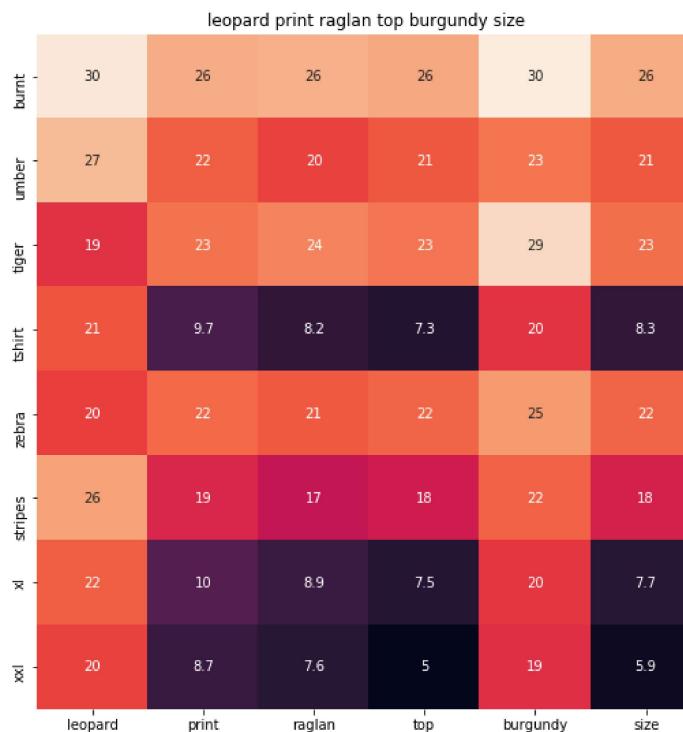
euclidean distance from input : 6.4149



ASIN : B073R4ZM7Y

Brand : Colosseum

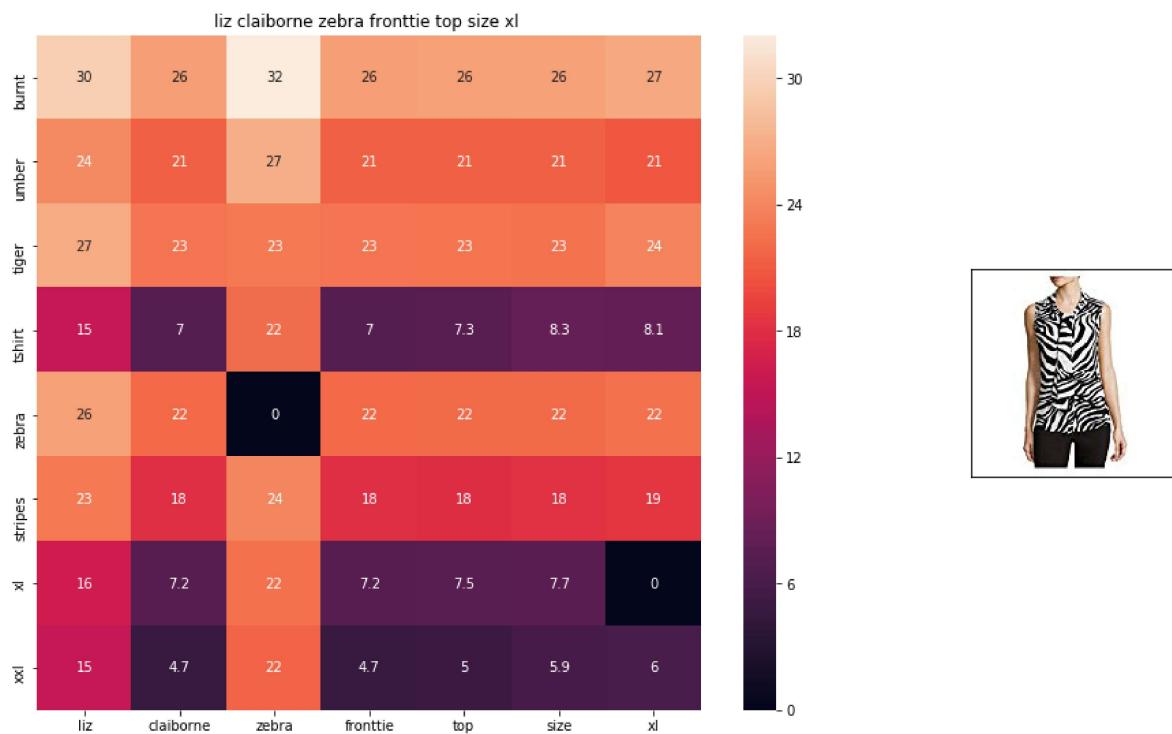
euclidean distance from input : 6.45096



ASIN : B01C60RLDQ

Brand : 1 Mad Fit

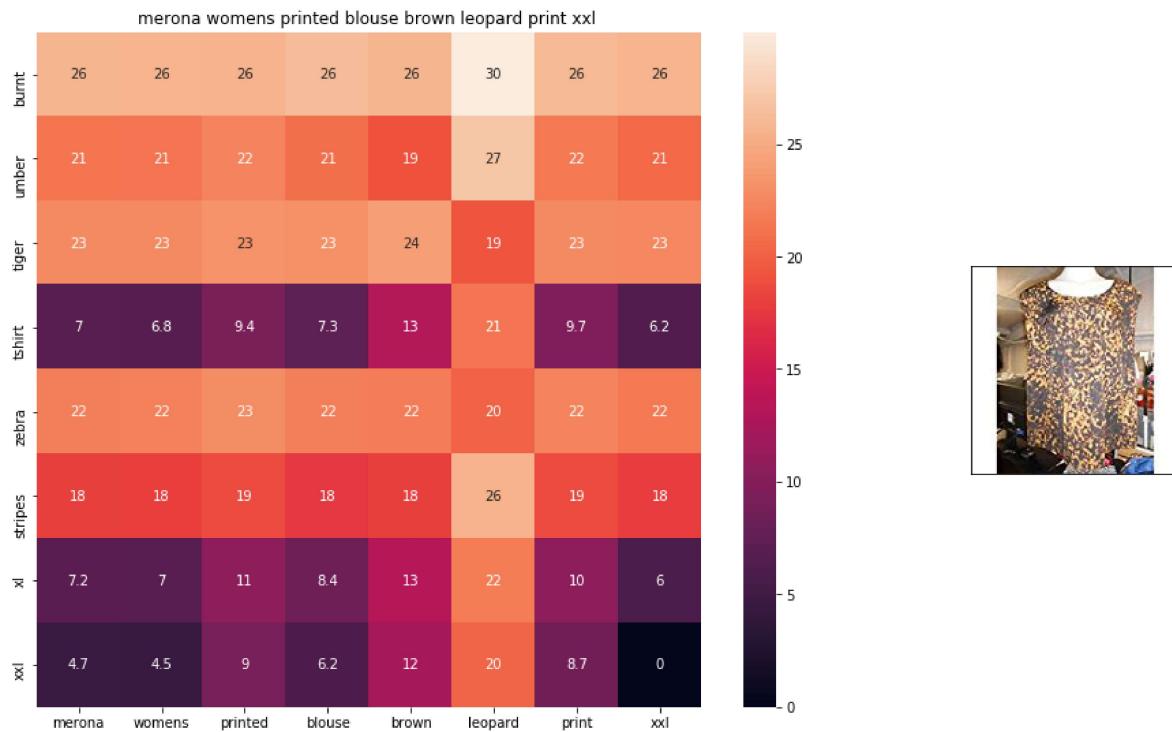
euclidean distance from input : 6.46341



ASIN : B06XBY5QXL

Brand : Liz Claiborne

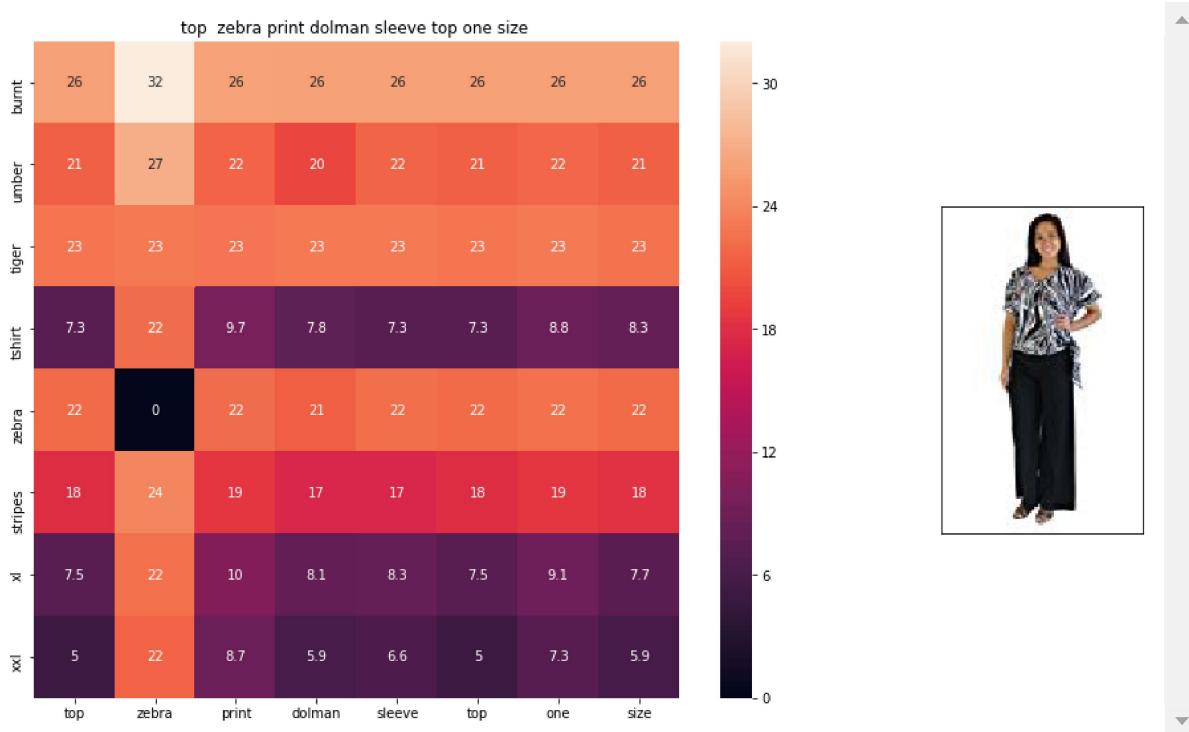
euclidean distance from input : 6.53922



ASIN : B071YF3WDD

Brand : Merona

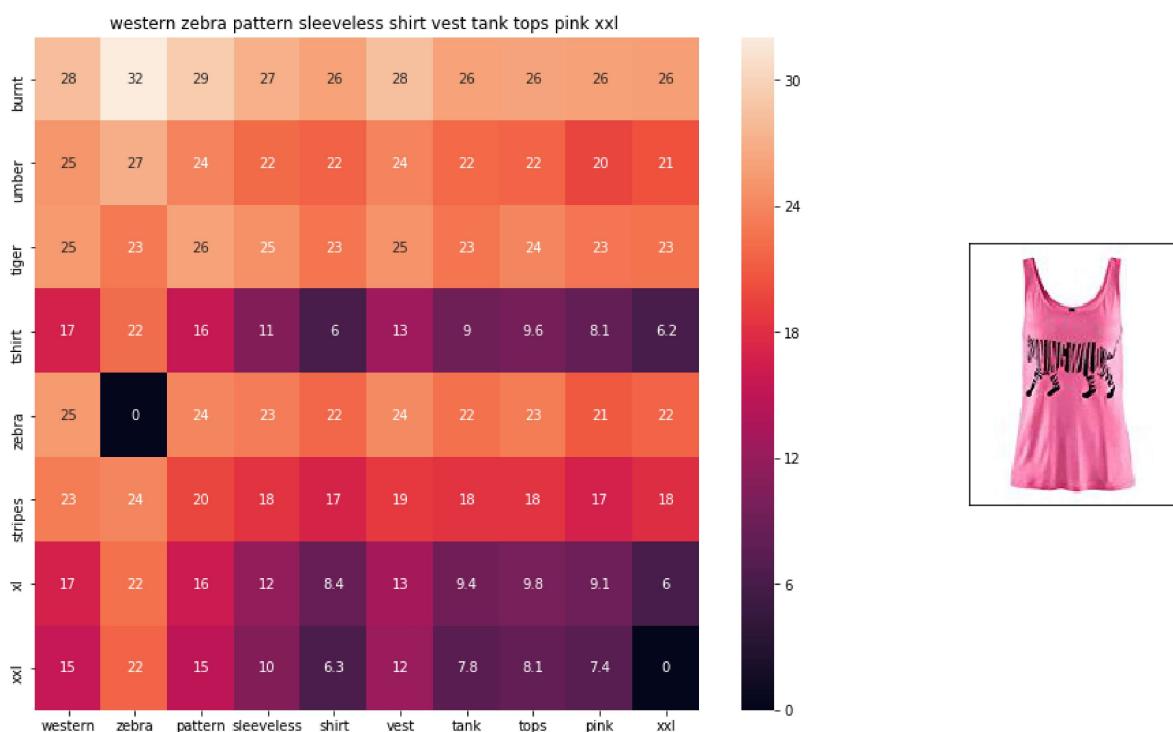
euclidean distance from input : 6.5755



ASIN : B00H8A6ZLI

Brand : Vivian's Fashions

euclidean distance from input : 6.63821



ASIN : B00Z6HEXWI

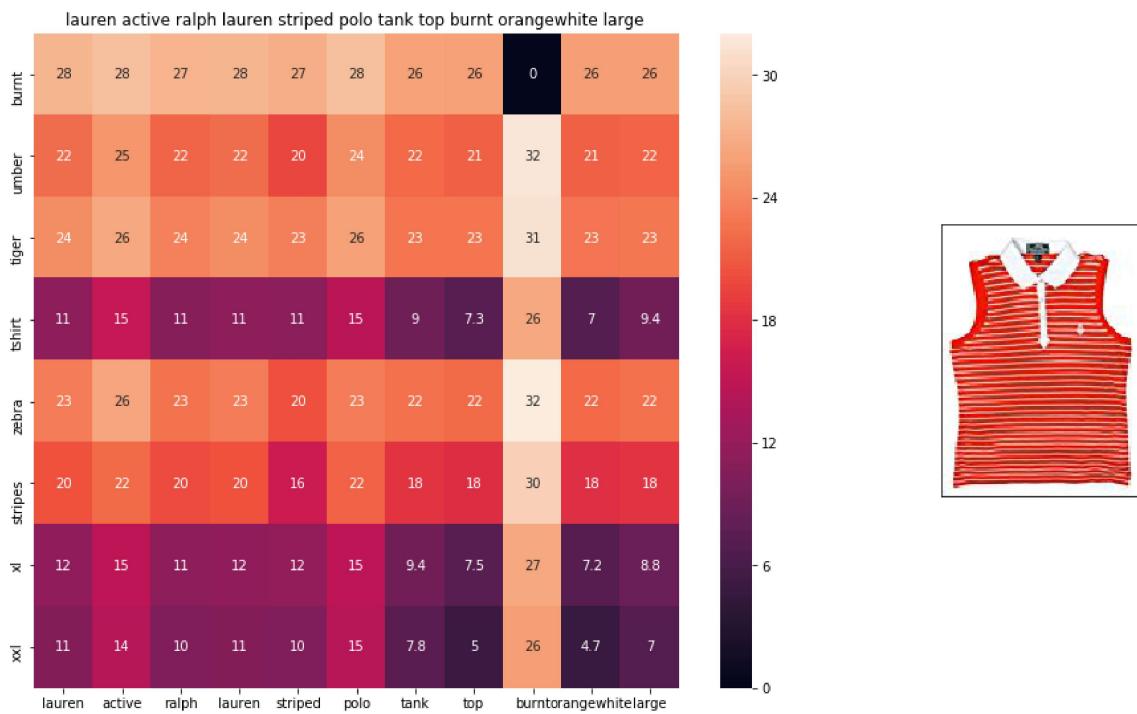
Brand : Black Temptation

euclidean distance from input : 6.66074

---



---



ASIN : B00ILGH50Y

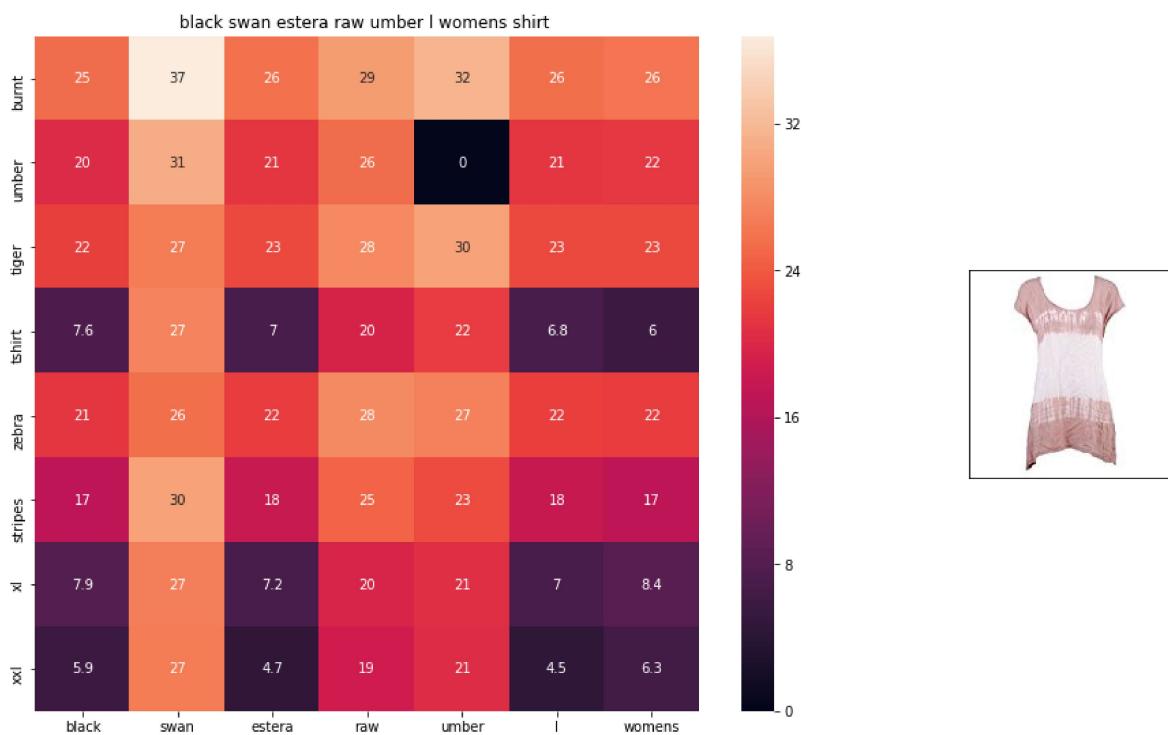
Brand : Ralph Lauren Active

euclidean distance from input : 6.68391

---



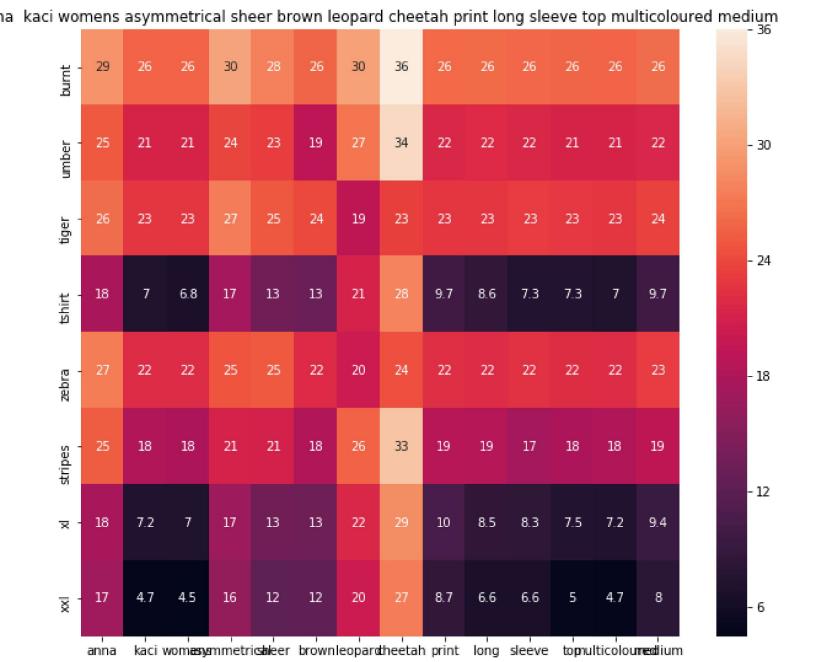
---



ASIN : B06Y1VN8WQ

Brand : Black Swan

euclidean distance from input : 6.70576



ASIN : B00KSNTY7Y

Brand : Anna-Kaci

euclidean distance from input : 6.70612

## Weighted similarity using brand and color.

In [0]:

```

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

```

In [0]:

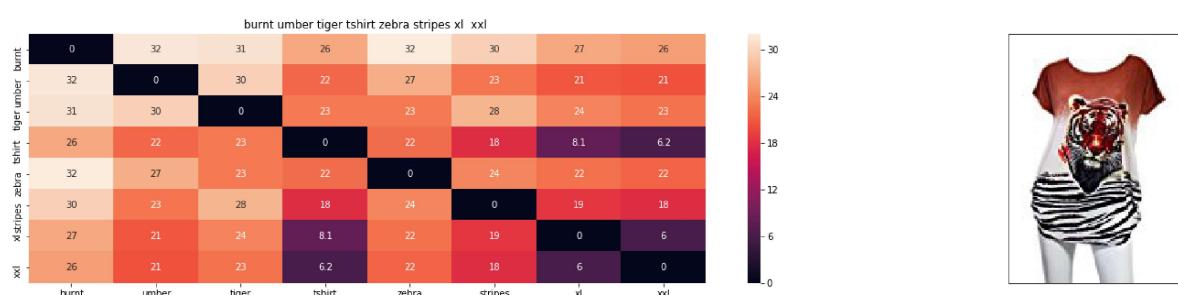
```
1 def heat_map_w2v_brand(sentance1, sentance2, url, doc_id1, doc_id2, df_id1, df_id2, mod
2
3     s1_vec = get_word_vec(sentance1, doc_id1, model)
4     s2_vec = get_word_vec(sentance2, doc_id2, model)
5
6     s1_s2_dist = get_distance(s1_vec, s2_vec)
7
8     data_matrix = [['Asin', 'Brand', 'Color', 'Product type'],
9                     [data['asin'].loc[df_id1], brands[doc_id1], colors[doc_id1], types[doc_id1]],
10                    [data['asin'].loc[df_id2], brands[doc_id2], colors[doc_id2], types[doc_id2]]]
11
12    colorscale = [[0, '#1d004d'], [.5, '#f2e5ff'], [1, '#f2e5d1']] # to color the heading
13
14    table = ff.create_table(data_matrix, index=True, colorscale=colorscale)
15    plotly.offline.iplot(table, filename='simple_table')
16
17    gs = gridspec.GridSpec(25, 15)
18    fig = plt.figure(figsize=(25,5))
19
20    ax1 = plt.subplot(gs[:, :-5])
21    # plotting the heap map based on the pairwise distances
22    ax1 = sns.heatmap(np.round(s1_s2_dist,6), annot=True)
23    # set the x axis labels as recommended apparels title
24    ax1.set_xticklabels(sentance2.split())
25    # set the y axis labels as input apparels title
26    ax1.set_yticklabels(sentance1.split())
27    # set title as recommended apparels title
28    ax1.set_title(sentance2)
29
30    # in last 25 * 10:15 grids we display image
31    ax2 = plt.subplot(gs[:, 10:16])
32    # we dont display grid lines and axis labels to images
33    ax2.grid(False)
34    ax2.set_xticks([])
35    ax2.set_yticks([])
36
37    # pass the url it display it
38    display_img(url, ax2, fig)
39
40    plt.show()
```

In [0]:

```

1 def idf_w2v_brand(doc_id, w1, w2, num_results):
2     idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1, -1))
3     ex_feat_dist = pairwise_distances(extra_features, extra_features[doc_id])
4     pairwise_dist = (w1 * idf_w2v_dist + w2 * ex_feat_dist)/float(w1 + w2)
5
6     indices = np.argsort(pairwise_dist.flatten())[0:num_results]
7     pdists = np.sort(pairwise_dist.flatten())[0:num_results]
8
9     df_indices = list(data.index[indices])
10
11
12 for i in range(0, len(indices)):
13     heat_map_w2v_brand(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
14     print('ASIN :', data['asin'].loc[df_indices[i]])
15     print('Brand :', data['brand'].loc[df_indices[i]])
16     print('euclidean distance from input :', pdists[i])
17     print('*'*125)
18
19 idf_w2v_brand(12566, 5, 5, 20)

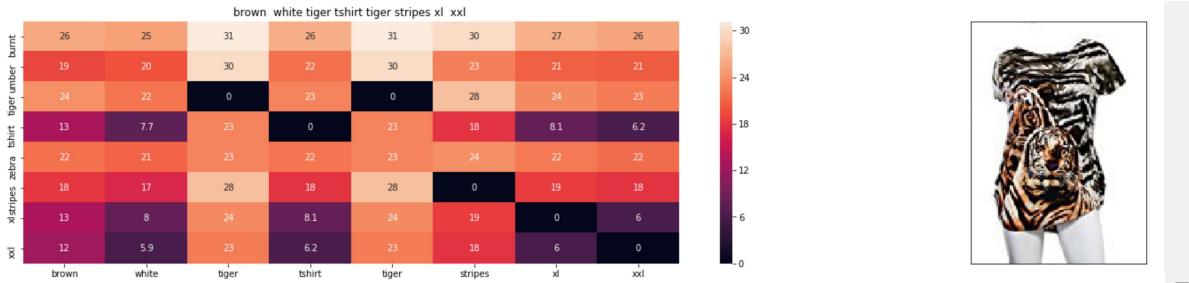
```



ASIN : B00JXQB5FQ

Brand : Si Row

euclidean distance from input : 0.001953125



ASIN : B00JXQCWTO

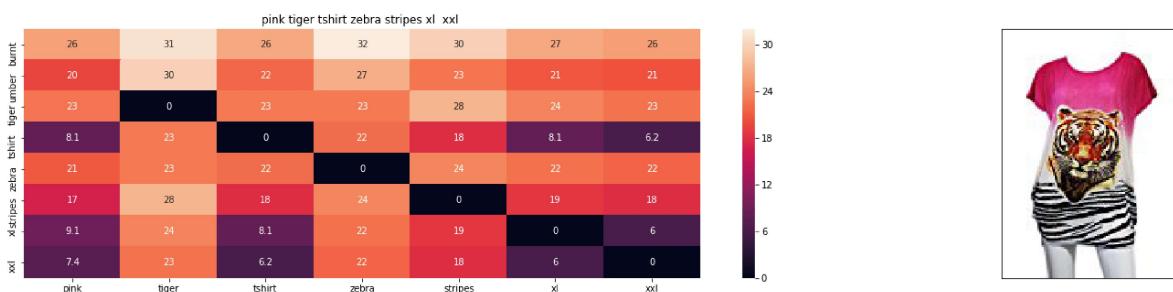
Brand : Si Row

euclidean distance from input : 2.38547115326

---



---



ASIN : B00JXQASS6

Brand : Si Row

euclidean distance from input : 2.73905105609

---



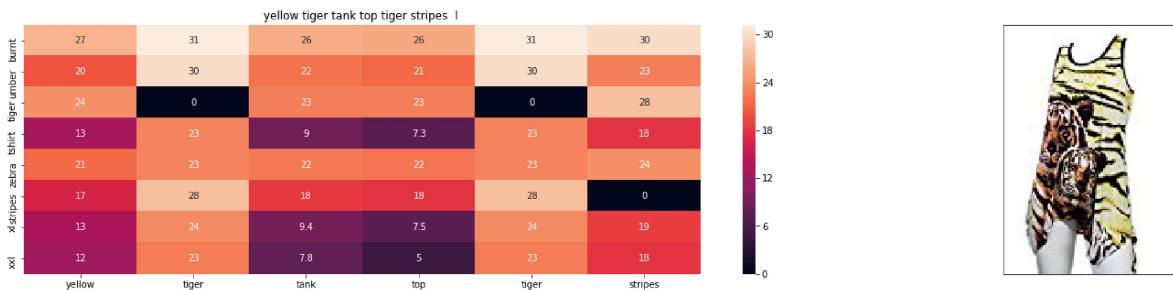
---



ASIN : B00JXQAFZ2

Brand : Si Row

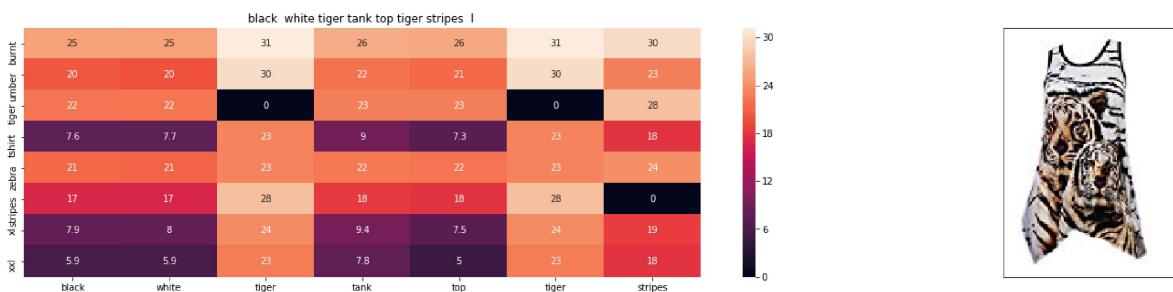
euclidean distance from input : 3.387187195



ASIN : B00JXQAUWA

Brand : Si Row

euclidean distance from input : 3.5518684389



ASIN : B00JXQA094

Brand : Si Row

euclidean distance from input : 3.5536174776



ASIN : B00JXQCUIC

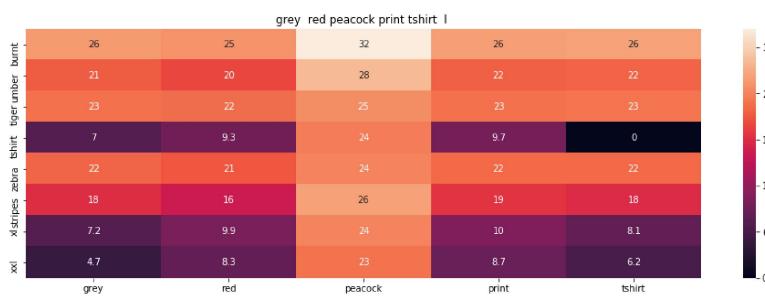
Brand : Si Row

euclidean distance from input : 3.65382804889

---



---



ASIN : B00JXQCFRS

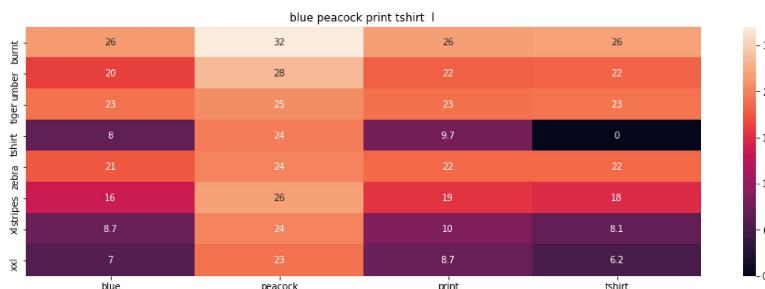
Brand : Si Row

euclidean distance from input : 4.12881164569

---



---



ASIN : B00JXQC8L6

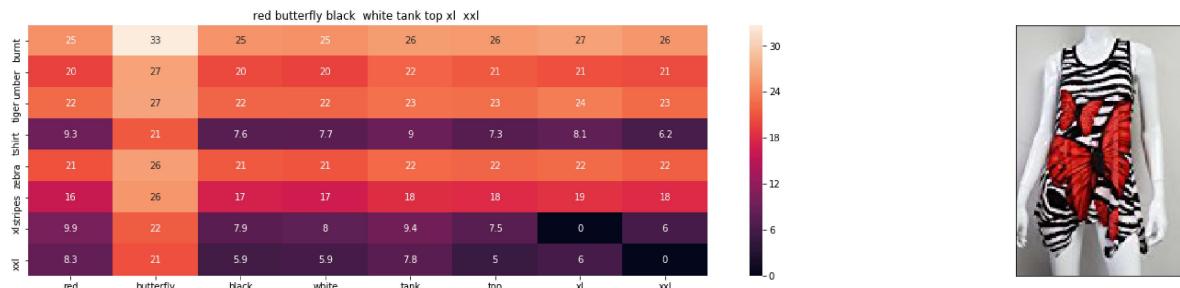
Brand : Si Row

euclidean distance from input : 4.20390052813

---



---



ASIN : B00JV63CW2

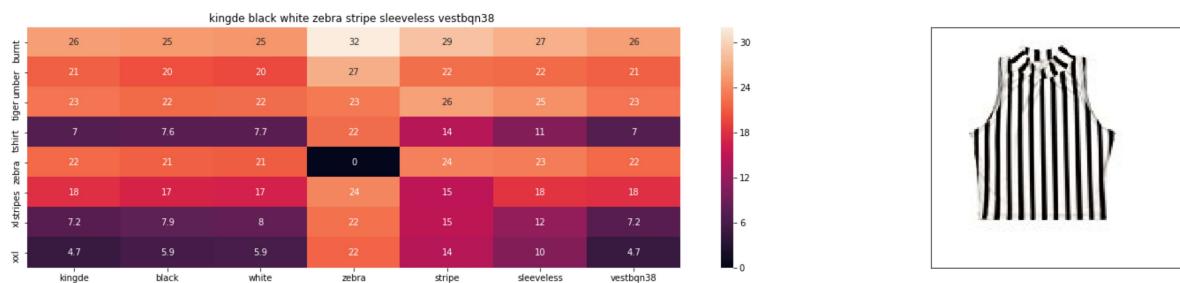
Brand : Si Row

euclidean distance from input : 4.28658676166

---



---



ASIN : B015H41F6G

Brand : KINGDE

euclidean distance from input : 4.38937078798

---



---



**ASIN : B00JXQBBMI**

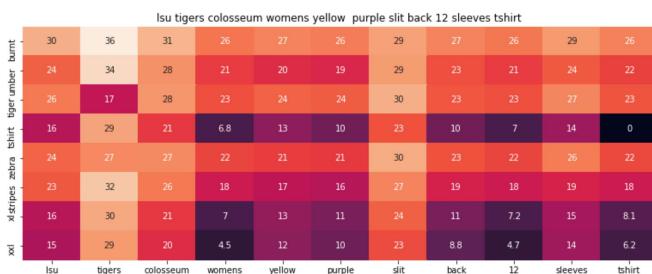
**Brand : Si Row**

**euclidean distance from input : 4.39790992755**

---



---



**ASIN : B073R5Q8HD**

**Brand : Colosseum**

**euclidean distance from input : 4.45122858369**

---



---



**ASIN : B074P8MD22**

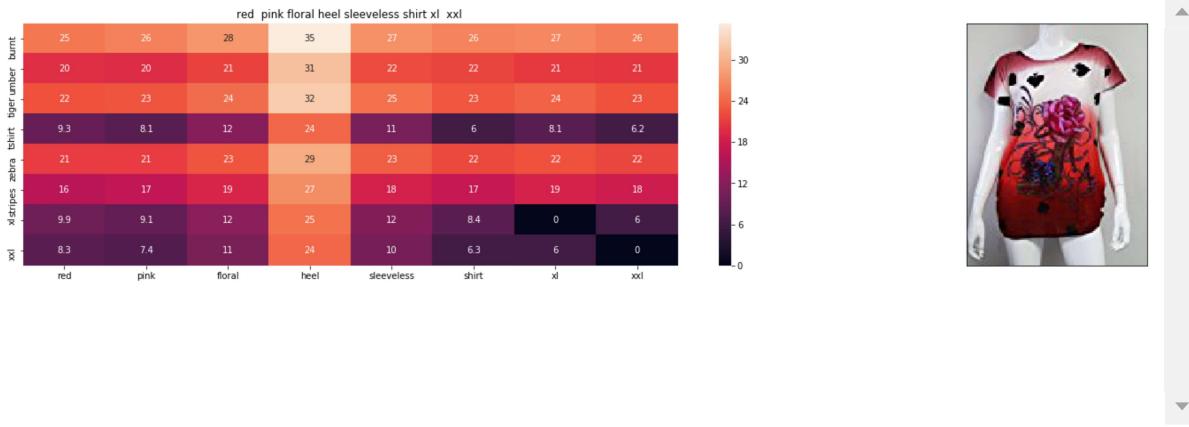
**Brand : Edista**

**euclidean distance from input : 4.51897779787**

---



---



ASIN : B00JV63QQE

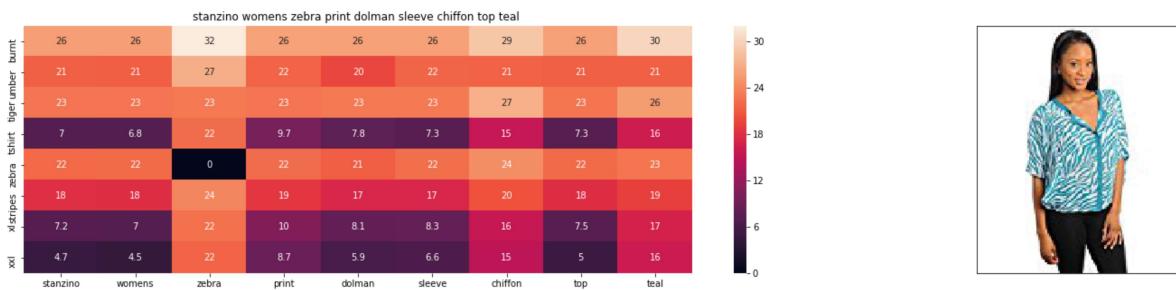
Brand : Si Row

euclidean distance from input : 4.52937545794

---



---



ASIN : B00C0I3U3E

Brand : Stanzino

euclidean distance from input : 4.53032614076

---



---

gupobou168 womens girls lady boho elephant stripes bandage tank top one size



**ASIN : B01ER18406**

**Brand : GuPoBoU168**

**euclidean distance from input : 4.54681702403**

---



---



**ASIN : B073R4ZM7Y**

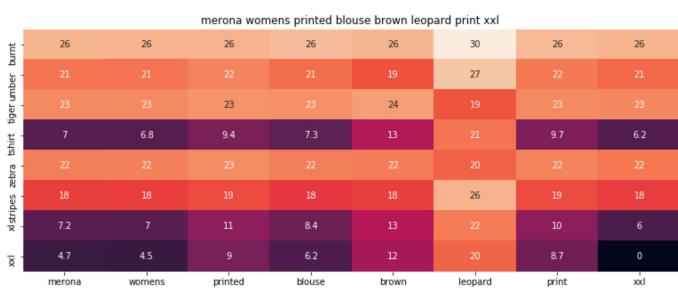
**Brand : Colosseum**

**euclidean distance from input : 4.54835554445**

---



---



**ASIN : B071YF3WDD**

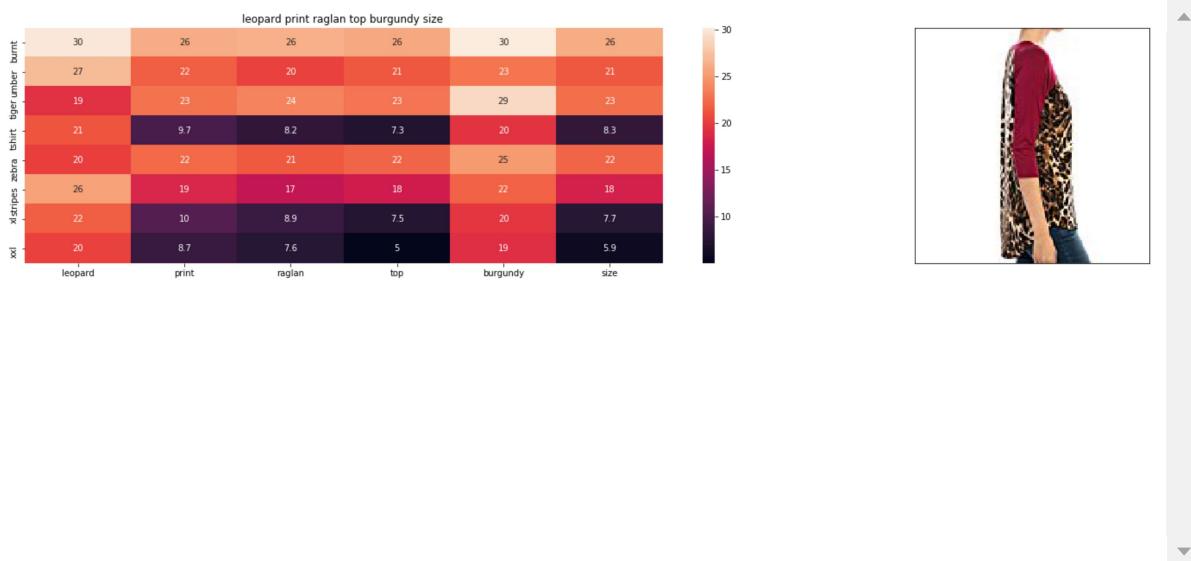
**Brand : Merona**

**euclidean distance from input : 4.61062742555**

---



---



ASIN : B01C60RLDQ

Brand : 1 Mad Fit

euclidean distance from input : 4.64591789282

---



---

In [0]:

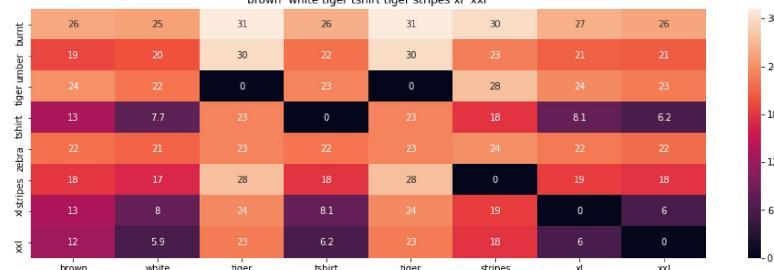
1 idf\_w2v\_brand(12566, 5, 50, 20)



ASIN : B00JXQB5FQ

Brand : Si Row

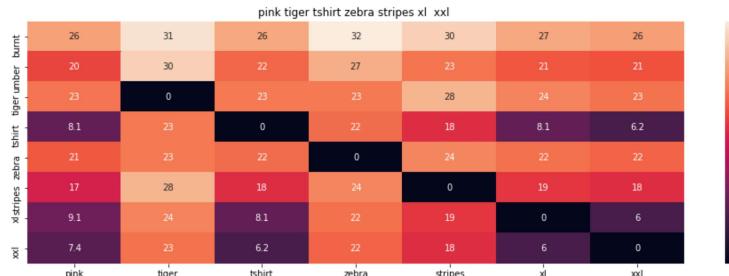
euclidean distance from input : 0.000355113636364



ASIN : B00JXQCWTO

Brand : Si Row

euclidean distance from input : 0.433722027865



ASIN : B00JXQASS6

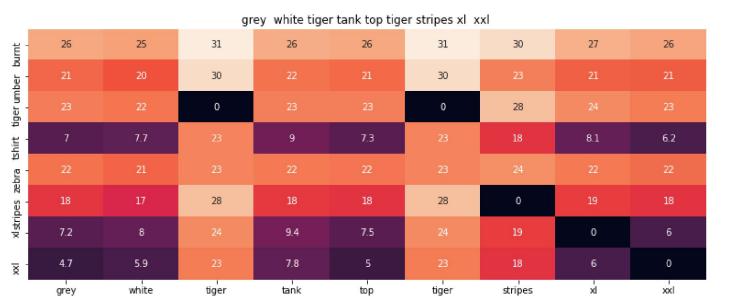
Brand : Si Row

euclidean distance from input : 1.65509310669

---



---



ASIN : B00JXQAFZ2

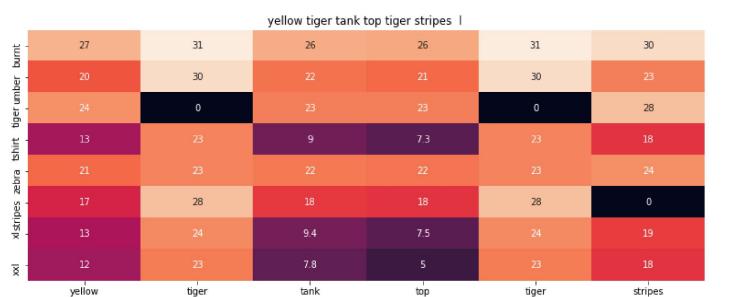
Brand : Si Row

euclidean distance from input : 1.77293604103

---



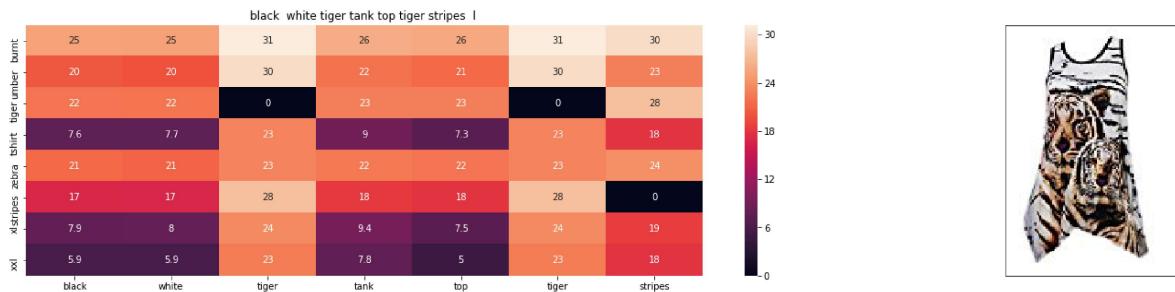
---



ASIN : B00JXQAUWA

Brand : Si Row

euclidean distance from input : 1.80287808538



ASIN : B00JXQA094

Brand : Si Row

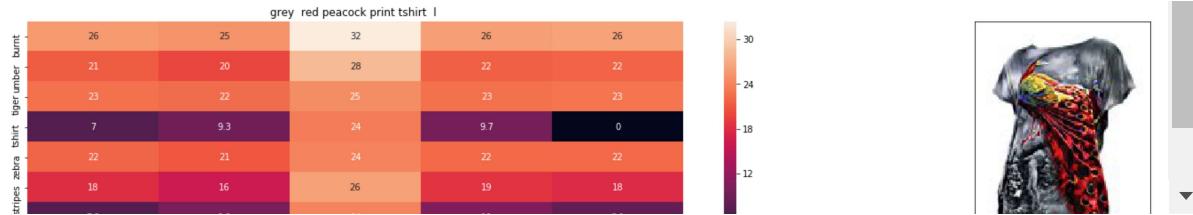
euclidean distance from input : 1.80319609241



ASIN : B00JXQCUIC

Brand : Si Row

euclidean distance from input : 1.82141619628



ASIN : B00JXQCFRS

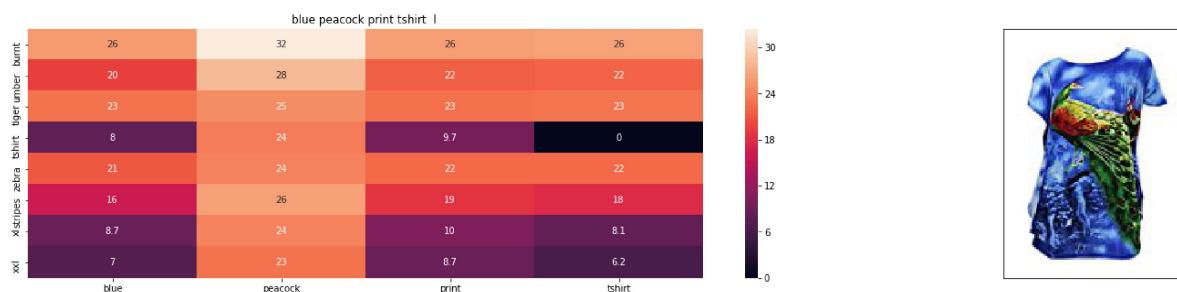
Brand : Si Row

euclidean distance from input : 1.90777685025

---



---



ASIN : B00JXQC8L6

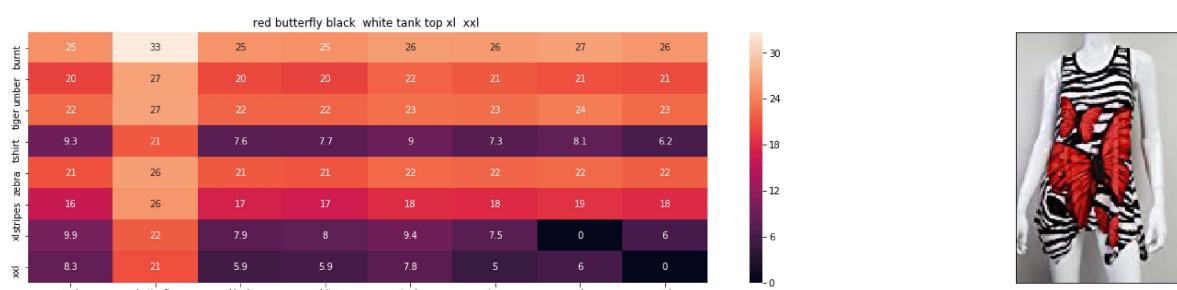
Brand : Si Row

euclidean distance from input : 1.92142937433

---



---



ASIN : B00JV63CW2

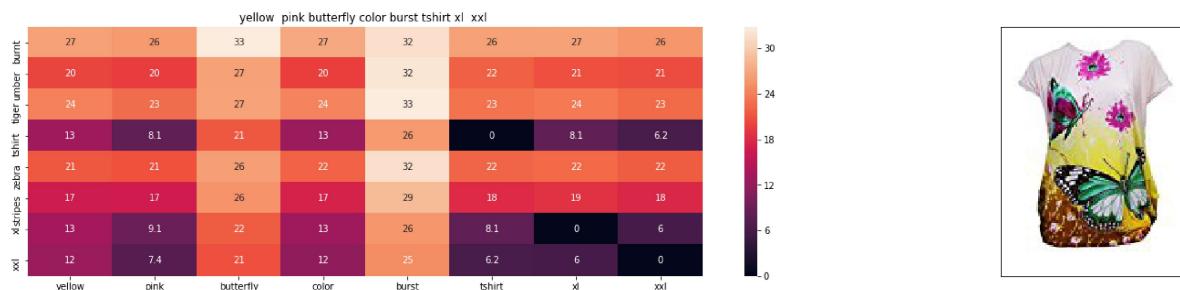
Brand : Si Row

euclidean distance from input : 1.93646323497

---



---



ASIN : B00JXQBBMI

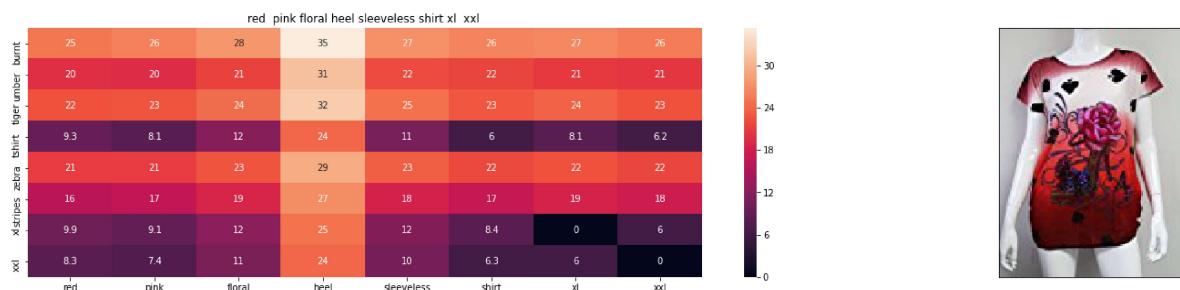
Brand : Si Row

euclidean distance from input : 1.95670381059

---



---



ASIN : B00JV63QQE

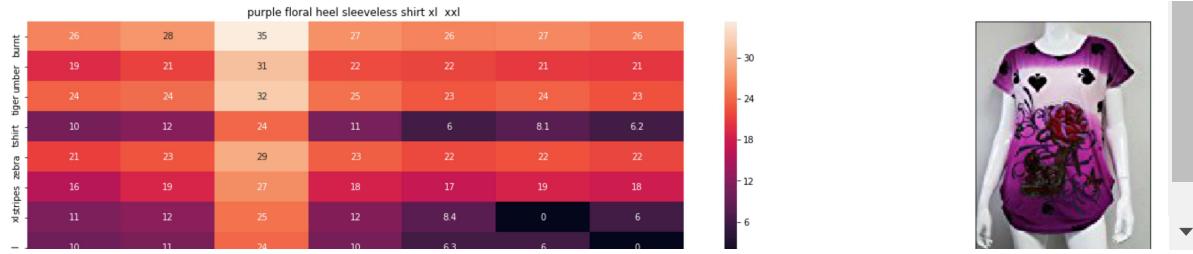
Brand : Si Row

euclidean distance from input : 1.9806066343

---



---



ASIN : B00JV63VC8

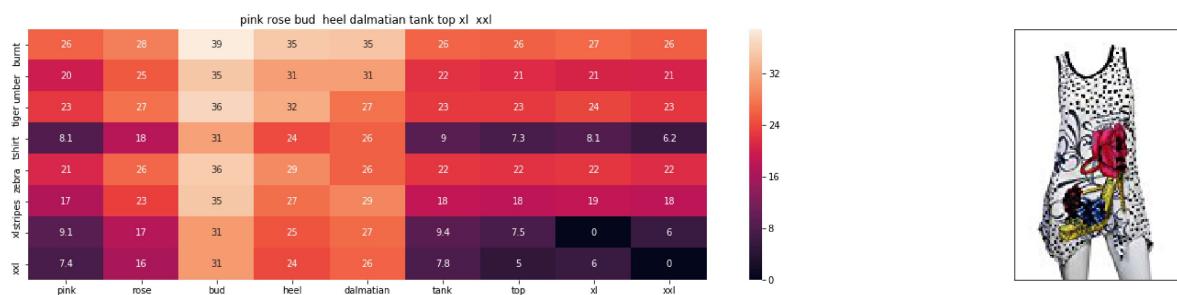
Brand : Si Row

euclidean distance from input : 2.01218559992

---



---



ASIN : B00JXQAX2C

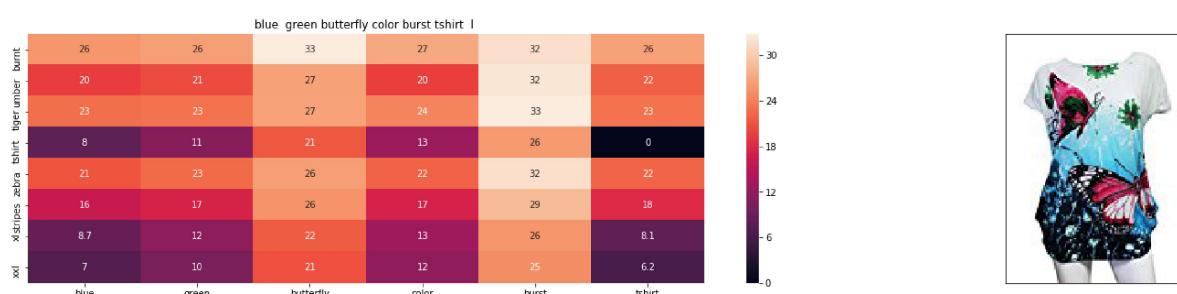
Brand : Si Row

euclidean distance from input : 2.01335178755

---



---



ASIN : B00JXQC0C8

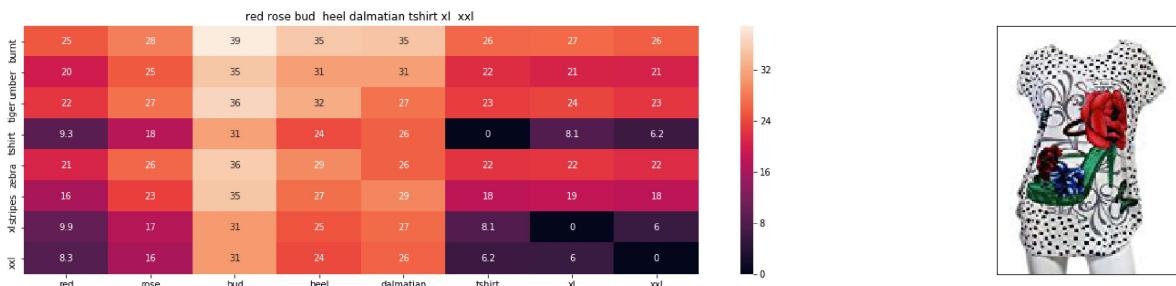
Brand : Si Row

euclidean distance from input : 2.01388334827

---



---



ASIN : B00JXQABB0

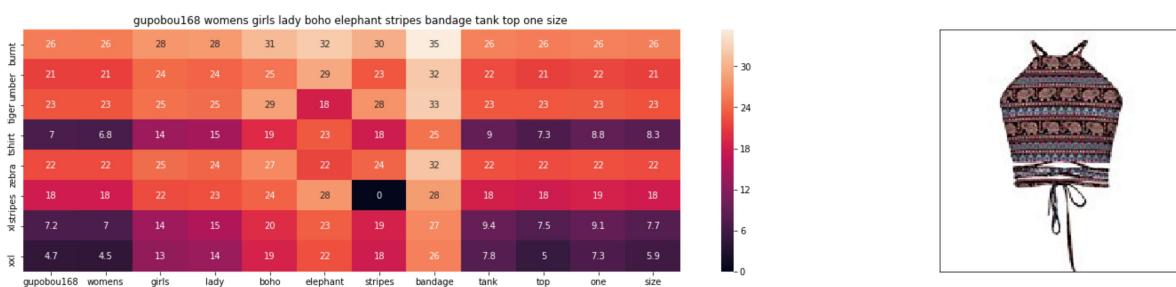
Brand : Si Row

euclidean distance from input : 2.0367257555

---



---



ASIN : B01ER18406

Brand : GuPoBoU168

euclidean distance from input : 2.65620416778

---



---



ASIN : B01LZ7BQ4H

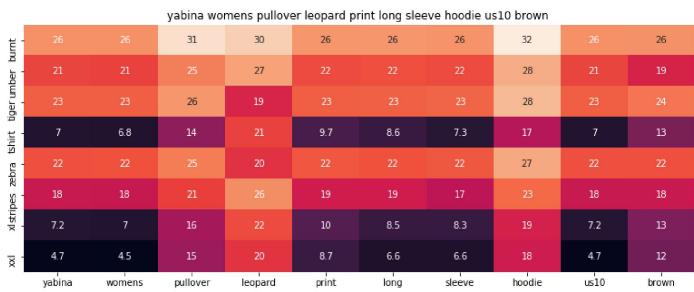
Brand : WAYF

euclidean distance from input : 2.6849067823

---



---



ASIN : B01KJUM6JI

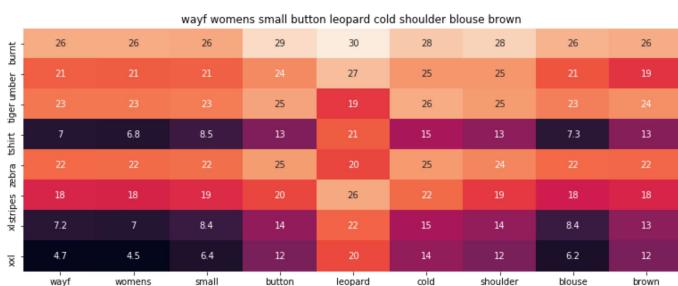
Brand : YABINA

euclidean distance from input : 2.68583819266

---



---



ASIN : B01M06V4X1

Brand : WAYF

euclidean distance from input : 2.69476194865

---



---

## Keras and Tensorflow to extract features

In [0]:

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

Using TensorFlow backend.

## Visual features based product similarity.

In [0]:

```

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

```



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



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



Product Title: yellow tiger tshirt tiger stripes 1

Euclidean Distance from input image: 41.2611

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: 44.0002

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: 47.3825

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: 47.7184

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: 48.0465

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: 48.1019

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: 48.1189

Amazon Url: [www.amazon.com/dp/B01I80A93G](http://www.amazon.com/dp/B01I80A93G)



Product Title: miss chievous juniors striped peplum tank top medium shadow peach

Euclidean Distance from input image: 48.1313

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: 48.1695

Amazon Url: [www.amazon.com/dp/B00JV63QQE](http://www.amazon.com/dp/B00JV63QQE)



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

Euclidean Distance from input image: 48.2568

Amazon Url: [www.amazon.com/dp/B01LX6H43D](http://www.amazon.com/dp/B01LX6H43D)



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

Euclidean Distance from input image: 48.2657

Amazon Url: [www.amazon.com/dp/B01CR57YY0](http://www.amazon.com/dp/B01CR57YY0)

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



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



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



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



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



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

In [0]:

1	
---	--