

# Amazon Apparel Recommendations

## [4.2] Data and Code:

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

## [4.3] Overview of the data

In [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")
```

In [2]:

```
# we have give a json file which consists of all information about
# the products
# loading the data using pandas' read_json file.
data = pd.read_json('tops_fashion.json')
```

In [3]:

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

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

## Terminology:

What is a dataset?

Rows and columns

Data-point

Feature/variable

In [4]:

```
# each product/item has 19 features in the raw dataset.
data.columns # prints column-names or feature-names.
```

Out[4]:

```
Index(['sku', 'asin', 'product_type_name', 'formatted_price', 'author',
       'color', 'brand', 'publisher', 'availability', 'reviews',
       'large_image_url', 'availability_type', 'small_image_url',
       'editorial_review', 'title', 'model', 'medium_image_url',
       'manufacturer', 'editorial_reivew'],
      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 apperal, 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 [5]:

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

In [6]:

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

Number of data points : 183138 Number of features: 7

Out[6]:

	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...	

## [5.1] Missing data for various features.

Basic stats for the feature: product\_type\_name

In [7]:

```
# 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
```

In [8]:

```
# 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']

In [9]:

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

Out[9]:

[('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 [10]:

```
# 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 [11]:

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

Out[11]:

```
[('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 [12]:

```
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 [13]:

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

Out[13]:

```
[(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 [14]:

```
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 [15]:

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

Out[15]:

```
[(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 [16]:

```
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
```

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

In [18]:

```
# 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

In [19]:

```
# consider products which have color information
# data['color'].isnull() => gives the information about the dataframe row's which have 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.

In [20]:

```
data.to_pickle('28k_apparel_data-1')
```

In [0]:

```
# 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[52]:

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

## [5.2] Remove near duplicate items

### [5.2.1] Understand about duplicates.

In [21]:

```
# read data from pickle file from previous stage
data = pd.read_pickle('28k_apparel_data-1')

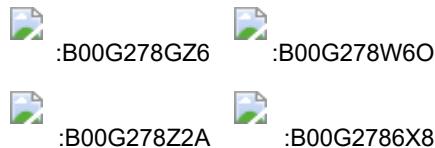
# find number of products that have duplicate titles.
print(sum(data.duplicated('title')))
# we have 2325 products which have same title but different color
```

2325

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



**These shirts exactly same except in color**



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 [22]:

```
# read data from pickle file from previous stage
data = pd.read_pickle('28k_apparel_data-1')
```

In [23]:

```
data.head()
```

Out[23]:

	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 [24]:

```
# 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 [25]:

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

Out[25]:

	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 Woman..

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 [26]:

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

In [27]:

```

import itertools
stage1_dedupe_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'
    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',
        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
        # example: a =['a', 'b', 'c', 'd']
        # b = ['a', 'b', 'd']
        # itertools.zip_longest(a,b): will give [('a', 'a'), ('b', 'b'), ('c', 'd'), ('d', Non
        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
        # if the number of words in which both strings differ are < 2 , we are considering
        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 in
            stage1_dedupe_asins.append(data_sorted['asin'].loc[indices[i]])

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

```

In [28]:

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

We removed the dupliactes which differ only at the end.

In [29]:

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

Number of data points : 17592

In [30]:

```
data.to_pickle('17k_apperal_data -1')
```

### [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 [31]:

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

In [32]:

```
data.shape
```

Out[32]:

(17592, 7)

In [33]:

```
# 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'
    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',

        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
        # 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
        if (length - count) < 3:
            indices.remove(j)
```

In [34]:

```
# 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 [35]:

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

Number of data points after stage two of dedupe: 16434

In [36]:

```
data.to_pickle('16k_apperal_data-1')
# 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 [37]:

```
data = pd.read_pickle('16k_apperal_data-1')

# 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 [38]:

```
# 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())
            # Conver all letters to lower-case
            word = word.lower()
            # stop-word removal
            if not word in stop_words:
                string += word + " "
        data[column][index] = string
```

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

In [39]:

```
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")
```

3.7303271 seconds

In [40]:

data.head()

Out[40]:

	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 [41]:

data.shape

Out[41]:

(16434, 7)

In [42]:

data.to\_pickle('16k\_apperal\_data\_preprocessed-1')

## Stemming

In [43]:

```
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

## [8] Text based product similarity

In [44]:

```
data = pd.read_pickle('16k_apperal_data_preprocessed-1')
data.head()
```

Out[44]:

	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 [45]:

```
# 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 occurrence of the word
    # 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 title
    ax = plt.subplot(gs[0])
    # it displays a cell in white color if the word is intersection(list of words of title, list of words of apparel)
    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 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 display_img based with parameter 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 apparel's vector, it is of a dict type {word:count}
    # vec2 : recommended apparel's vector, it is of a dict type {word:count}
    # url : apparel's image url
    # text: title of recommended 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 decision
    intersection = set(vec1.keys()) & set(vec2.keys())

```

```

# we set the values of non intersecting words to zero, this is just to show the difference
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(list of words of title1 and list of words of title2): val
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 word
        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
        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 words 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()' though
    return Counter(words) # Counter counts the occurrence of each word in list, it returns dict

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

    # vector1 = dict{word11:#count, word12:#count, etc.}
    vector1 = text_to_vector(text1)

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

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

```

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

In [46]:

```
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 occured in tha
```

Out[46]:

(16434, 12684)

In [47]:

```

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 apparel
    # the metric we used here is cosine, the coside distance is measured as K(X, Y) = <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])

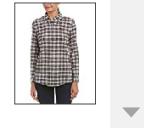
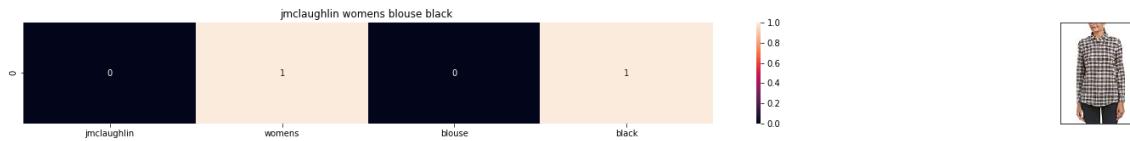
    for i in range(0,len(indices)):
        # we will pass 1. doc_id, 2. title1, 3. title2, url, model
        get_result(indices[i],data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
        print('ASIN : ',data['asin'].loc[df_indices[i]])
        print ('Brand: ', data['brand'].loc[df_indices[i]])
        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 : B074KN55WS

Brand: J. McLaughlin

Title: jmclaughlin womens blouse black

Euclidean similarity with the query image : 2.8284271247461903

---



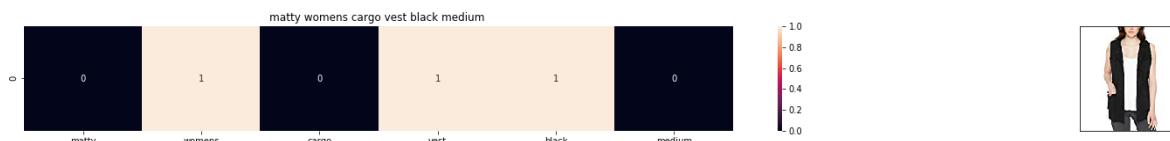
ASIN : B074TSFJHM

Brand: M Missoni

Title: missoni womens top black

Euclidean similarity with the query image : 2.8284271247461903

---



ASIN : B01LXQ0550

Brand: Matty M

Title: matty womens cargo vest black medium

Euclidean similarity with the query image : 2.8284271247461903

---



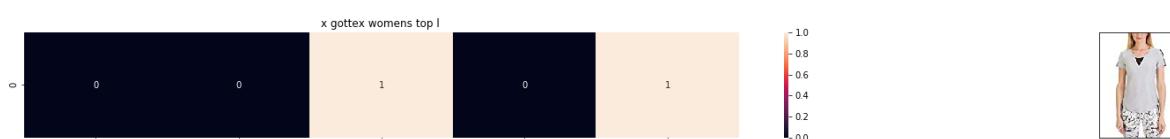
ASIN : B01EF0WPFI

Brand: ONEPICE

Title: onepice womens tegan sara waistcoat black

Euclidean similarity with the query image : 2.8284271247461903

---



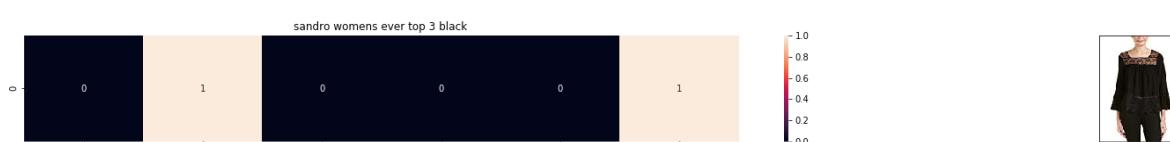
ASIN : B073M78XLP

Brand: X by Gottex

Title: x gottex womens top l

Euclidean similarity with the query image : 3.0

---



ASIN : B06XKX3BZL

Brand: Sandro

Title: sandro womens ever top 3 black

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

---



**ASIN : B072HLFGT5**

**Brand: kensie**

**Title: kensie womens solid shoulder casual top black**

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

---



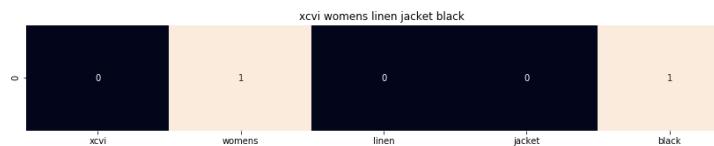
**ASIN : B013L1V8PK**

**Brand: J Brand Jeans**

**Title: j brand womens sheer shirt black**

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

---



**ASIN : B01M31Q4Z0**

**Brand: XCVI**

**Title: xcvi womens linen jacket black**

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

---



**ASIN : B074TWLMJ8**

**Brand: Badgley Mischka**

**Title: badgley mischka womens top 4 black**

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

---



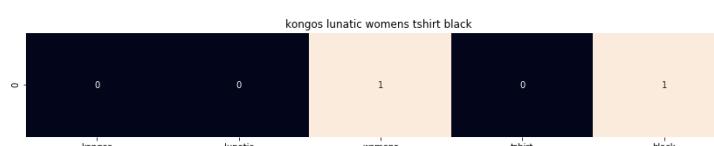
**ASIN : B075B1NXKV**

**Brand: Jay Godfrey**

**Title: jay godfrey womens blouse 6 black**

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

---



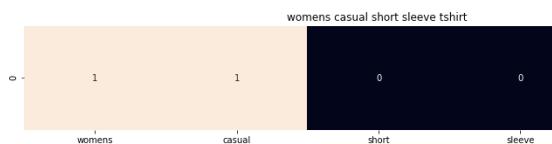
ASIN : B01IXZLJIA

Brand: Moonflow

Title: kongos lunatic womens tshirt black

Euclidean similarity with the query image : 3.0

---



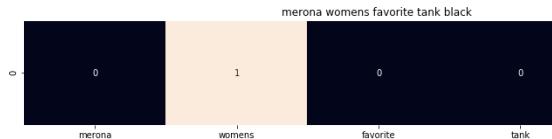
ASIN : B074T9KG9Q

Brand: Rain

Title: womens casual short sleeve tshirt

Euclidean similarity with the query image : 3.0

---



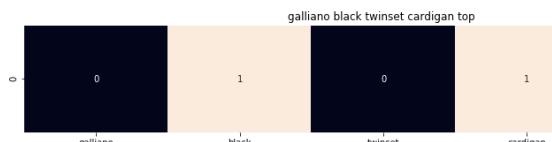
ASIN : B01KBGZE4Y

Brand: Merona

Title: merona womens favorite tank black

Euclidean similarity with the query image : 3.0

---



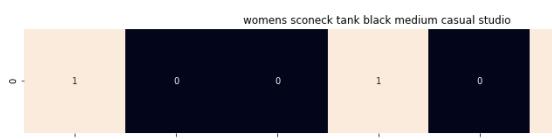
ASIN : B074G57HQJ

Brand: Galliano

Title: galliano black twinset cardigan top

Euclidean similarity with the query image : 3.0

---



ASIN : B06ZZCKDQW

Brand: Casual Studio

Title: womens sconeck tank black medium casual studio

Euclidean similarity with the query image : 3.0

---



ASIN : B01M34NPAM

Brand: WHAT ON EARTH

Title: womens faux vegan black leather vest purse

Euclidean similarity with the query image : 3.0

---

jill stuart womens wool shift dress 0 black



ASIN : B07591PX36

Brand: Jill Stuart

Title: jill stuart womens wool shift dress 0 black

Euclidean similarity with the query image : 3.0

=====

## [8.5] TF-IDF based product similarity

In [48]:

```
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
# tfidf_title_features[doc_id, index_of_word_in_corpus] = tfidf values of the word in given
```

In [49]:

```

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

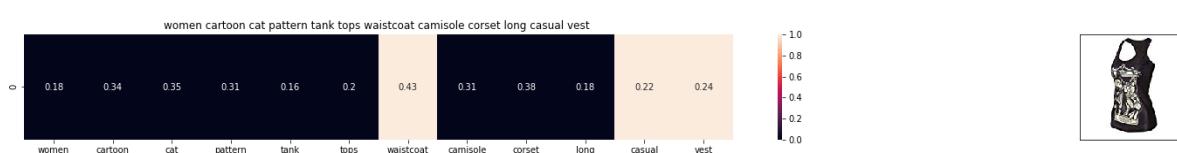
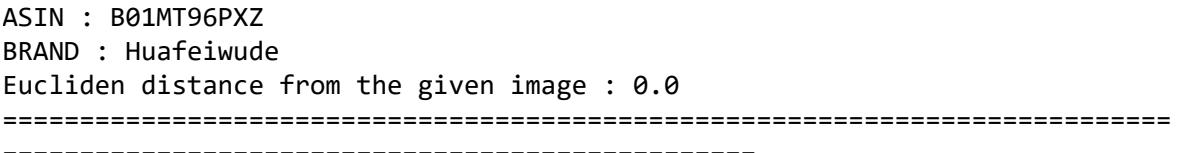
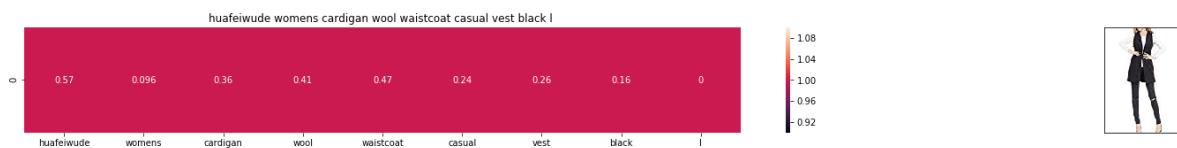
    # pairwise_dist will store the distance from given input apparel to all remaining apparel
    # the metric we used here is cosine, the coside distance is measured as K(X, Y) = <X, Y>
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
    pairwise_dist = pairwise_distances(tfidf_title_features,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

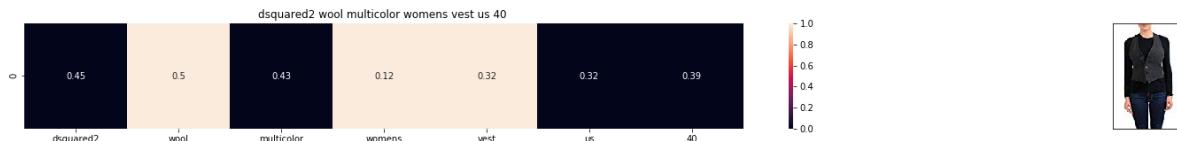
    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    #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]])
        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

```





ASIN : B01N2MU4DR

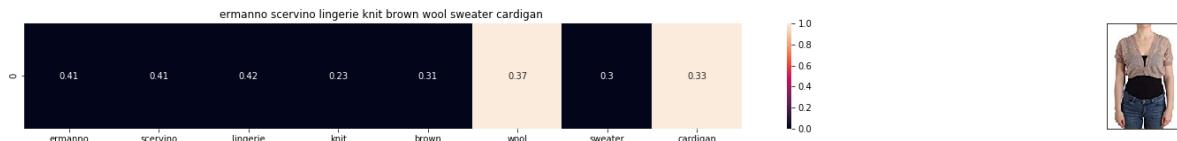
BRAND : DSQUARED2

Eucliden distance from the given image : 1.188600377321188

---



---



ASIN : B074G59T11

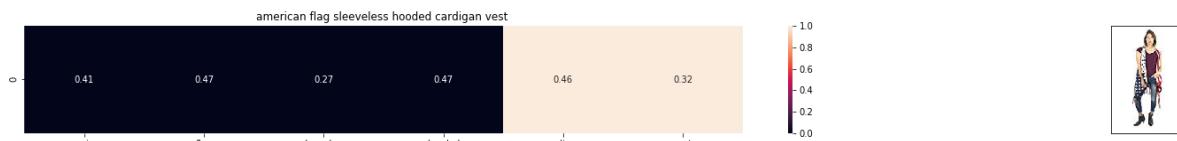
BRAND : ERMANNO SCERVINO

Eucliden distance from the given image : 1.206359485845001

---



---



ASIN : B01MF523X0

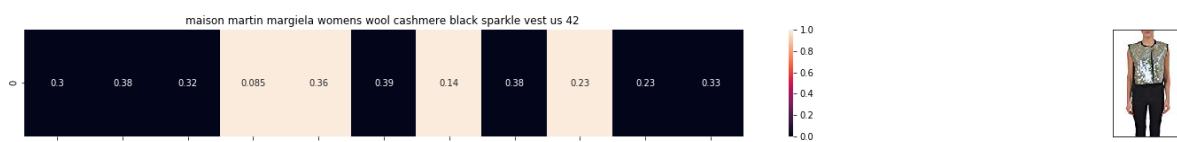
BRAND : A&O International

Eucliden distance from the given image : 1.225015799943407

---



---



ASIN : B0175G7C20

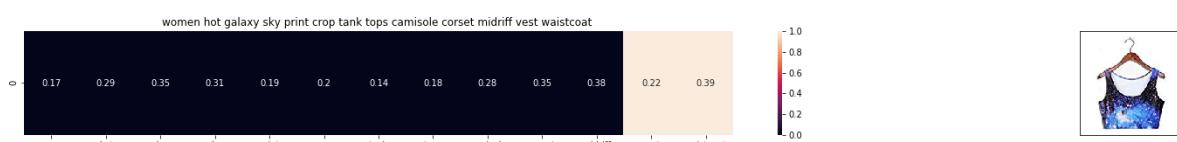
BRAND : Maison Martin Margiela

Eucliden distance from the given image : 1.2351958567475707

---



---



ASIN : B011R13I5Y

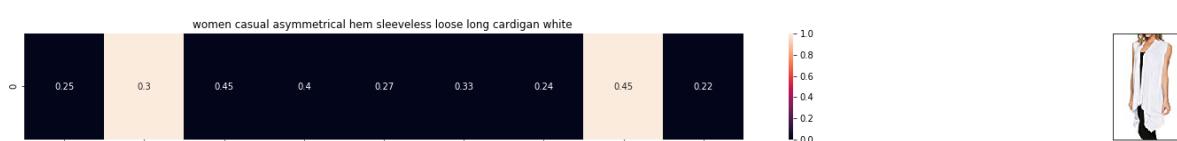
BRAND : Huayang

Eucliden distance from the given image : 1.236015538222882

---



---



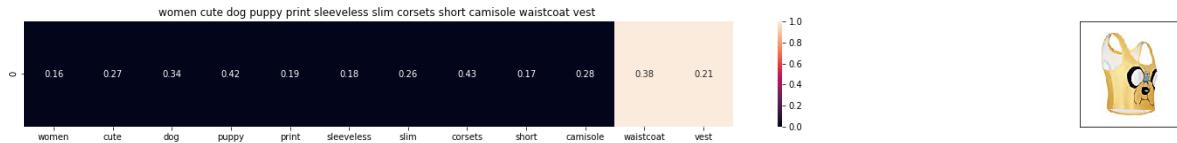
ASIN : B019V3MIP6

BRAND : KingField

Eucliden distance from the given image : 1.2364678792669093

=====

=====



ASIN : B011R13T2Q

BRAND : Huayang

Eucliden distance from the given image : 1.2395445383327484

=====

=====



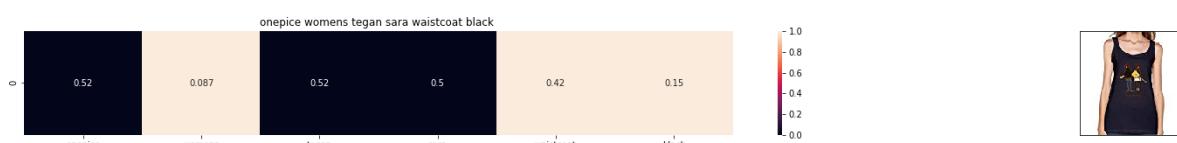
ASIN : B06XY144C5

BRAND : Missoni

Eucliden distance from the given image : 1.2415915469686623

=====

=====



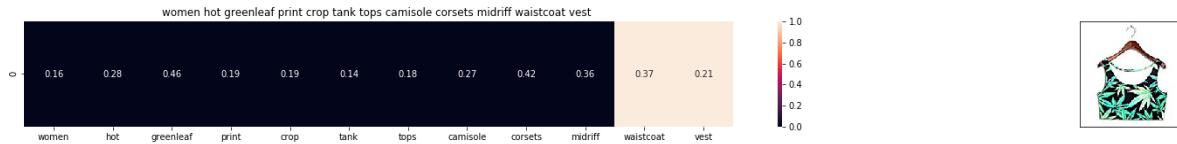
ASIN : B01EF0WPF1

BRAND : ONEPIECE

Eucliden distance from the given image : 1.2416544379264938

=====

=====



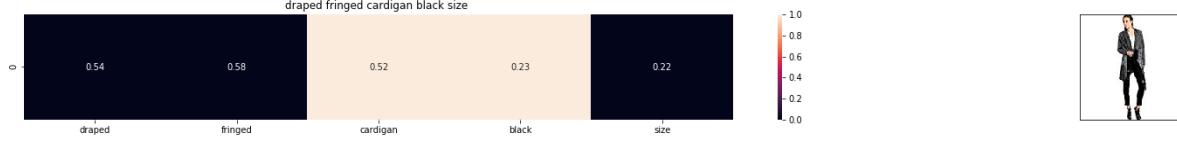
ASIN : B011R130VW

BRAND : Huayang

Eucliden distance from the given image : 1.2428237398407325

=====

=====



ASIN : B01CB33866

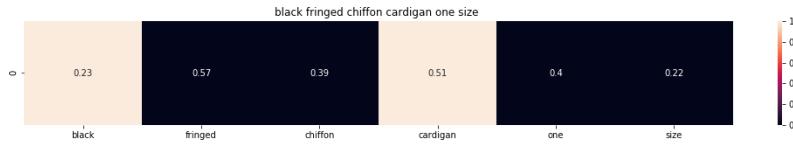
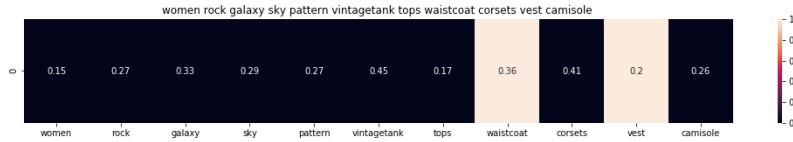
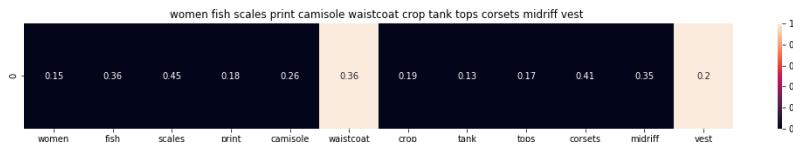
BRAND : Flying Tomato

Eucliden distance from the given image : 1.2439474778333672

=====

=====

galliano pink wool knit top

**ASIN : B074G4V6NJ****BRAND : Galliano****Eucliden distance from the given image : 1.246074624114252****ASIN : B010AMAMLY****BRAND : HP-LEISURE****Eucliden distance from the given image : 1.2463648318197242****ASIN : B011R13M10****BRAND : Huayang****Eucliden distance from the given image : 1.2479129830677684****ASIN : B011R13BDS****BRAND : Huayang****Eucliden distance from the given image : 1.247961302381963****ASIN : B01GIP26MA****BRAND : Luxury Divas****Eucliden distance from the given image : 1.2521865818430968****ASIN : B01AVX8IOU****BRAND : Sandistore****Eucliden distance from the given image : 1.2527438390346828**

## [8.5] IDF based product similarity

In [14]:

```
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
# idf_title_features[doc_id, index_of_word_in_corpus] = number of times the word occured in
```

In [51]:

```
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 [52]:

```
# 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 value
    # idf_title_features[:, idf_title_vectorizer.vocabulary_[i]].nonzero()[0] will return a
    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_courpus] = idf value of word
        idf_title_features[j,idf_title_vectorizer.vocabulary_[i]] = idf_val
```



In [53]:

```

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 apparel
    # the metric we used here is cosine, the coside distance is measured as K(X, Y) = <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]])
        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

```

In [0]:

```
# 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 [3]:

```
from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

# in this project we are using a pretrained model by google
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file which contains a dict ,
# and it contains all our corpus words as keys and model[word] as values
# To use this code-snippet, download "GoogleNews-vectors-negative300.bin"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNLNUTTLSS21pQmM/edit
# it's 1.9GB in size.

...
model = KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True
...
# if you do NOT have RAM >= 12GB, use the code below.
with open('word2vec_model', 'rb') as handle:
    model = pickle.load(handle)
```

In [4]:

```
# Utility functions

def get_word_vec(sentence, doc_id, m_name):
    # sentence : title of the apparel
    # doc_id: document id in our corpus
    # 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)
    vec = []
    for i in sentence.split():
        if i in vocab:
            if m_name == 'weighted' and i in idf_title_vectorizer.vocabulary_:
                vec.append(idf_title_features[doc_id, idf_title_vectorizer.vocabulary_[i]])
            elif m_name == 'avg':
                vec.append(model[i])
        else:
            # if the word in our corpus is not there in the google word2vec corpus, we are
            vec.append(np.zeros(shape=(300,)))
    # we will return a numpy array of shape (#number of words in title * 300 ) 300 = Len(w2
    # each row represents the word2vec representation of each word (weighted/avg) in given
    return np.array(vec)

def get_distance(vec1, vec2):
    # vec1 = np.array(#number_of_words_title1 * 300), each row is a vector of Length 300 co
    # vec2 = np.array(#number_of_words_title2 * 300), each row is a vector of Length 300 co

    final_dist = []
    # for each vector in vec1 we calculate the distance(euclidean) to all vectors in vec2
    for i in vec1:
        dist = []
        for j in vec2:
            # np.linalg.norm(i-j) will result the euclidean distance between vectors i, j
            dist.append(np.linalg.norm(i-j))
        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):
    # sentence1 : title1, input apparel
    # sentence2 : 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) o
    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) o
    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
```

```
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 [5]:

```
# 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):
    # sentence: its title of the apparel
    # num_features: the length 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 initialize a vector of size 300 with all zeros
    # we add each word2vec(wordi) to this featureVec
    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]])
            elif m_name == 'avg':
                featureVec = np.add(featureVec, model[word])
    if(nwords>0):
        featureVec = np.divide(featureVec, nwords)
    # returns the avg vector of given sentance, its of shape (1, 300)
    return featureVec
```

## [9.2] Average Word2Vec product similarity.

In [45]:

```
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 [46]:

```

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))

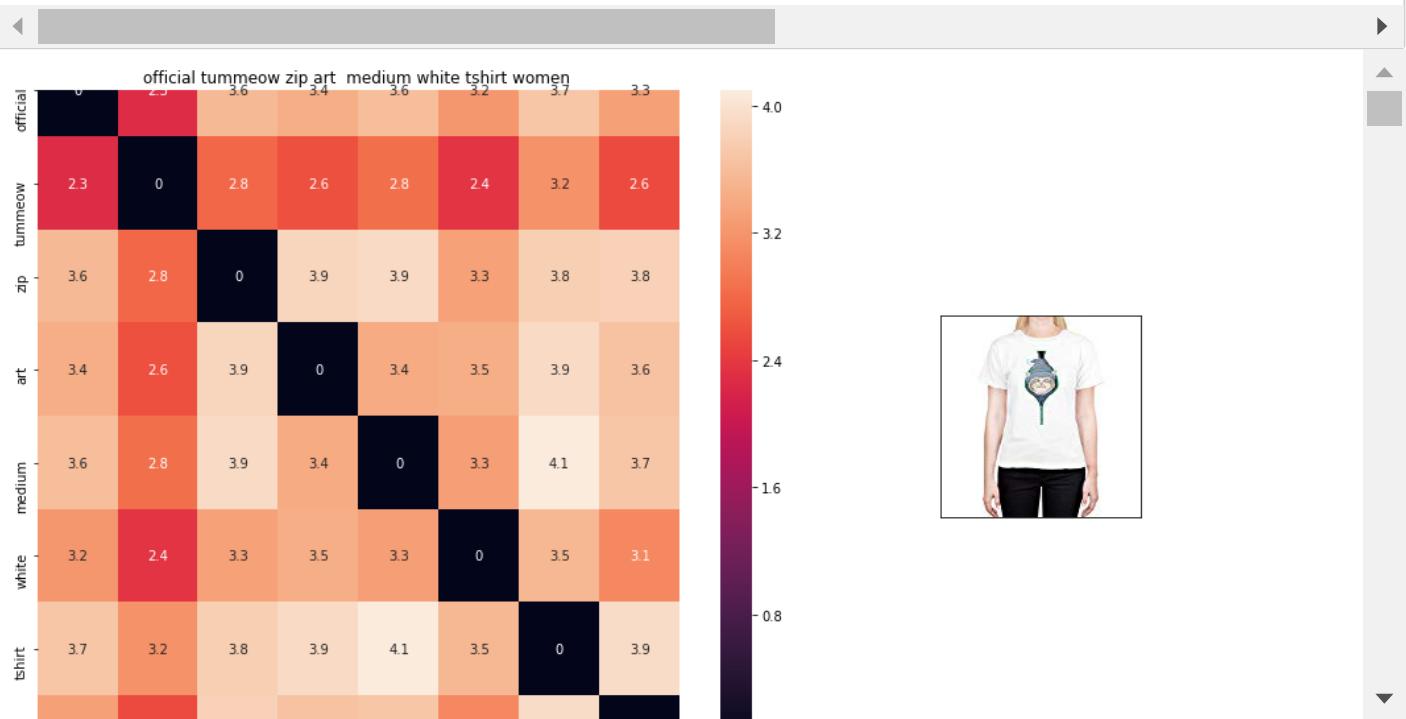
    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    #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)
        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*



## [9.4] IDF weighted Word2Vec for product similarity

In [57]:

```
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 [58]:

```

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 apparel
    # the metric we used here is cosine, the coside distance is measured as K(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))

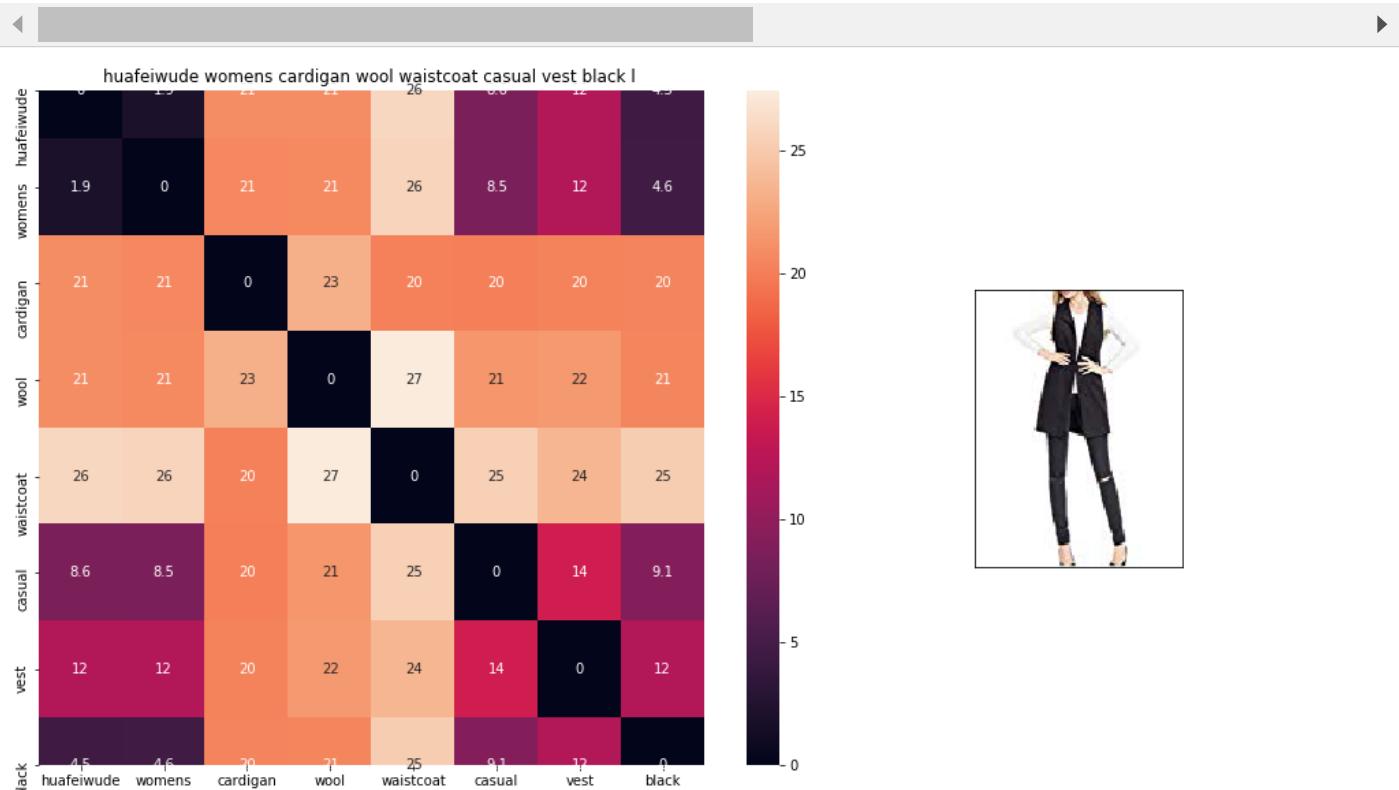
    # np.argsort will return indices of 9 smallest distances
    indices = np.argsort(pairwise_dist.flatten())[0:num_results]
    #pdists will store the 9 smallest distances
    pdists = np.sort(pairwise_dist.flatten())[0:num_results]

    #data frame indices of the 9 smallest distace'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['asin'].loc[df_indices[i]])
        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 : B01MT96PXZ

Brand : Huafeiwude

euclidean distance from input : 0.0

=====

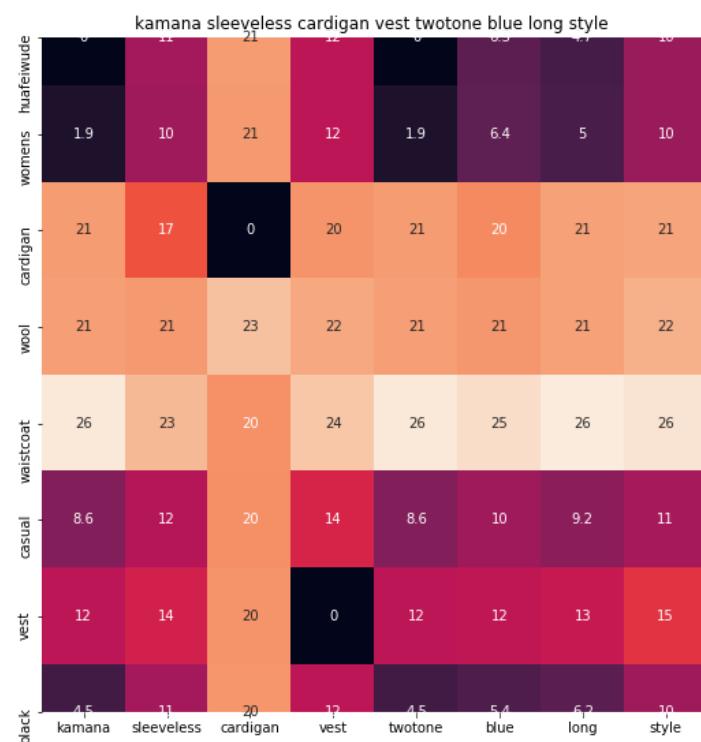
=====



ASIN : B07473KFK1

Brand : Sannysis

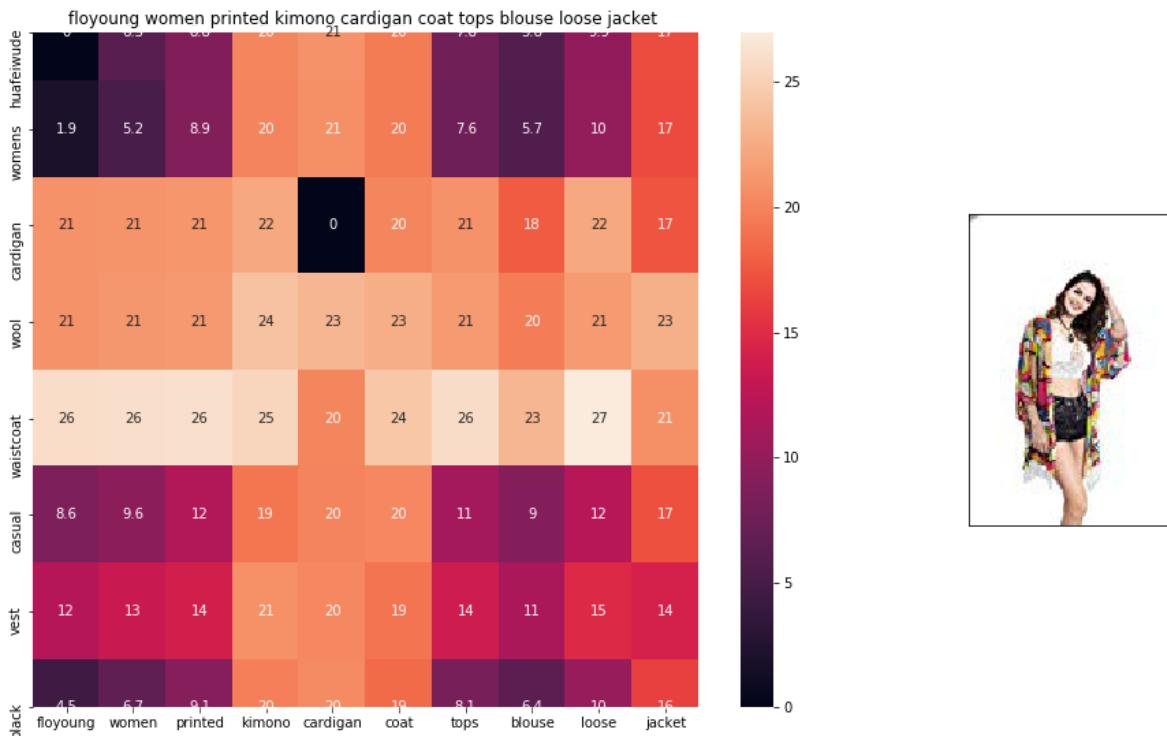
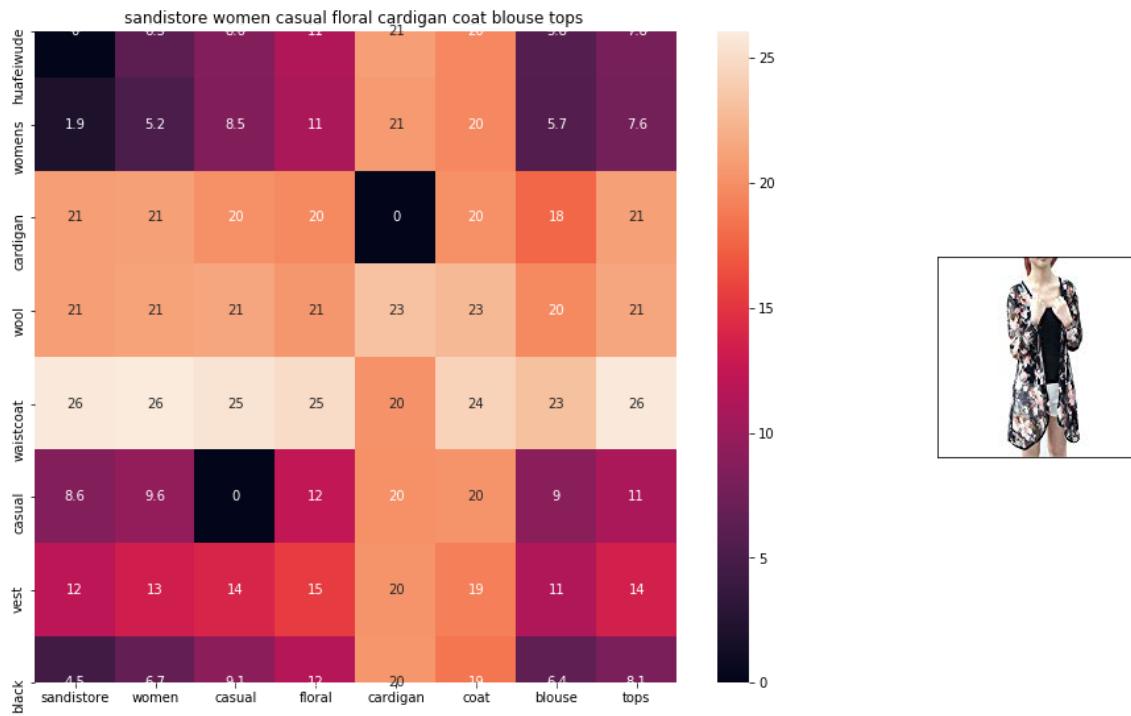
euclidean distance from input : 3.7969613

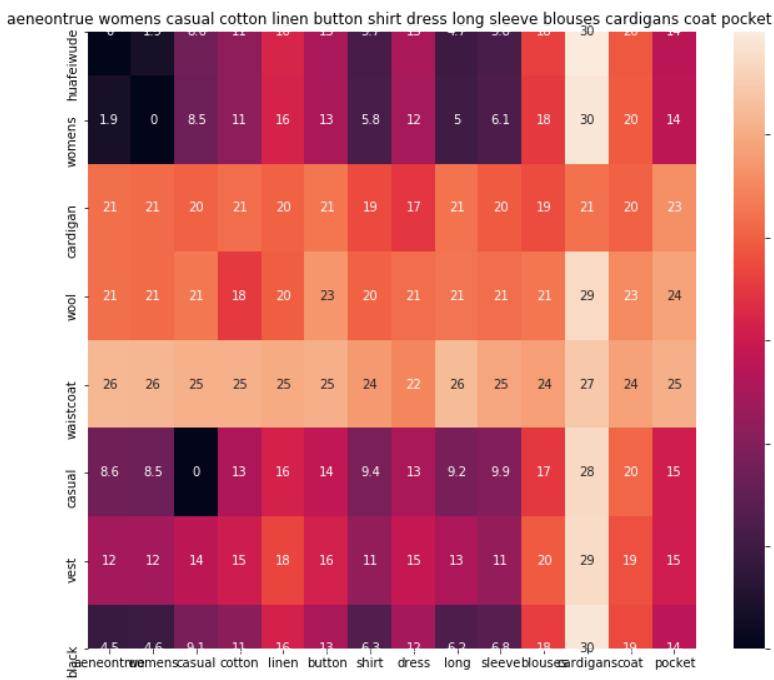


ASIN : B0751686K7

Brand : Kamana

euclidean distance from input : 3.9051597





ASIN : B074V1K5QJ

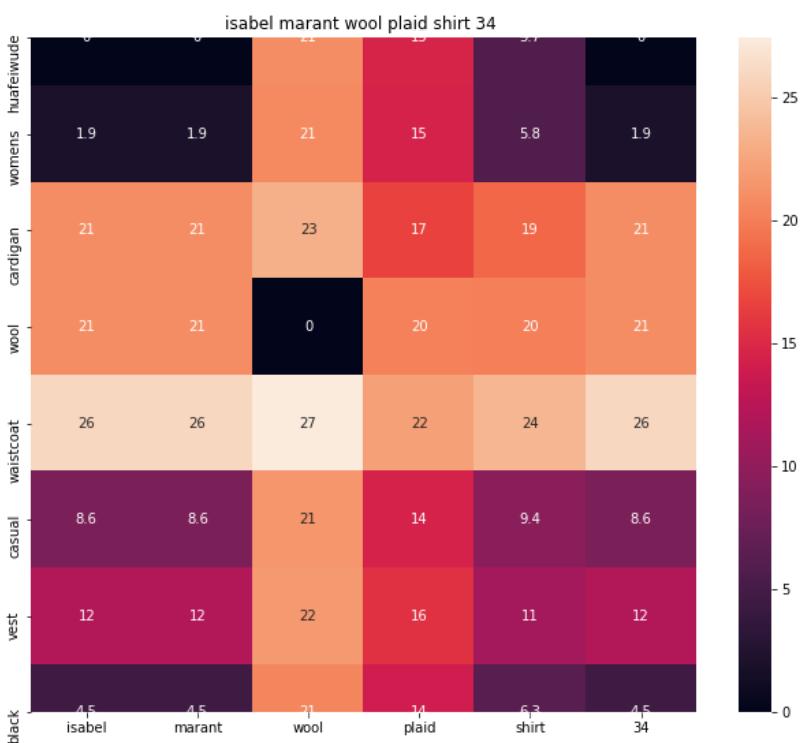
Brand : Aeneontrue

euclidean distance from input : 4.0558114

---



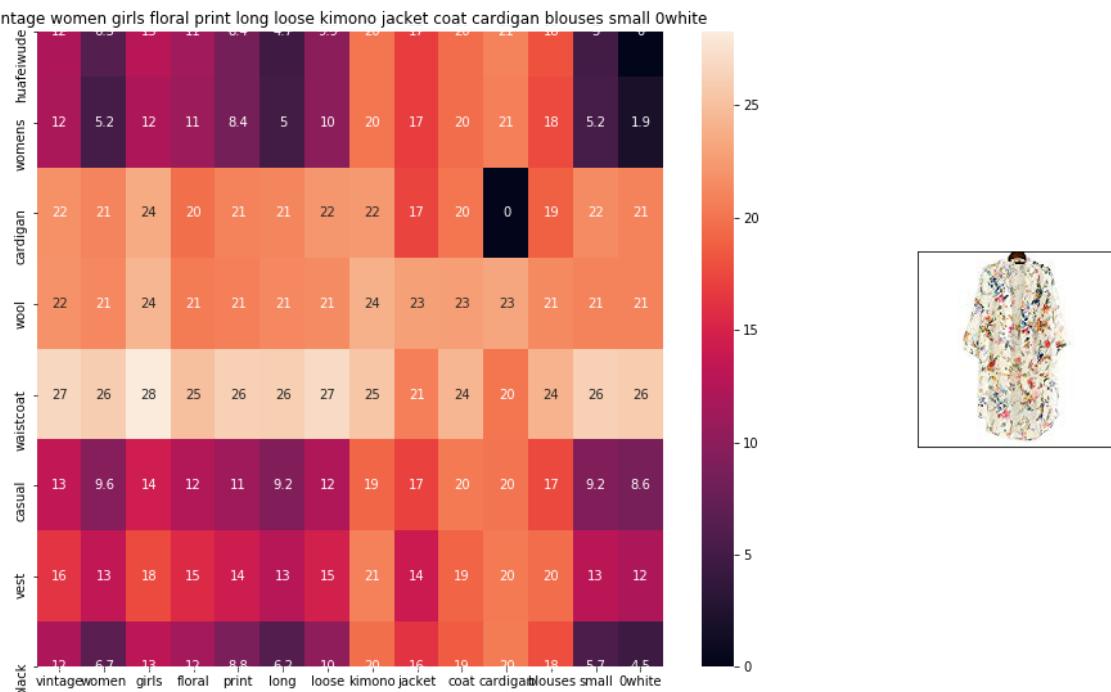
---



ASIN : B074NBDL2W

Brand : Isabel Marant

euclidean distance from input : 4.211965



ASIN : B07375JCKD

Brand : ACEFAST INC

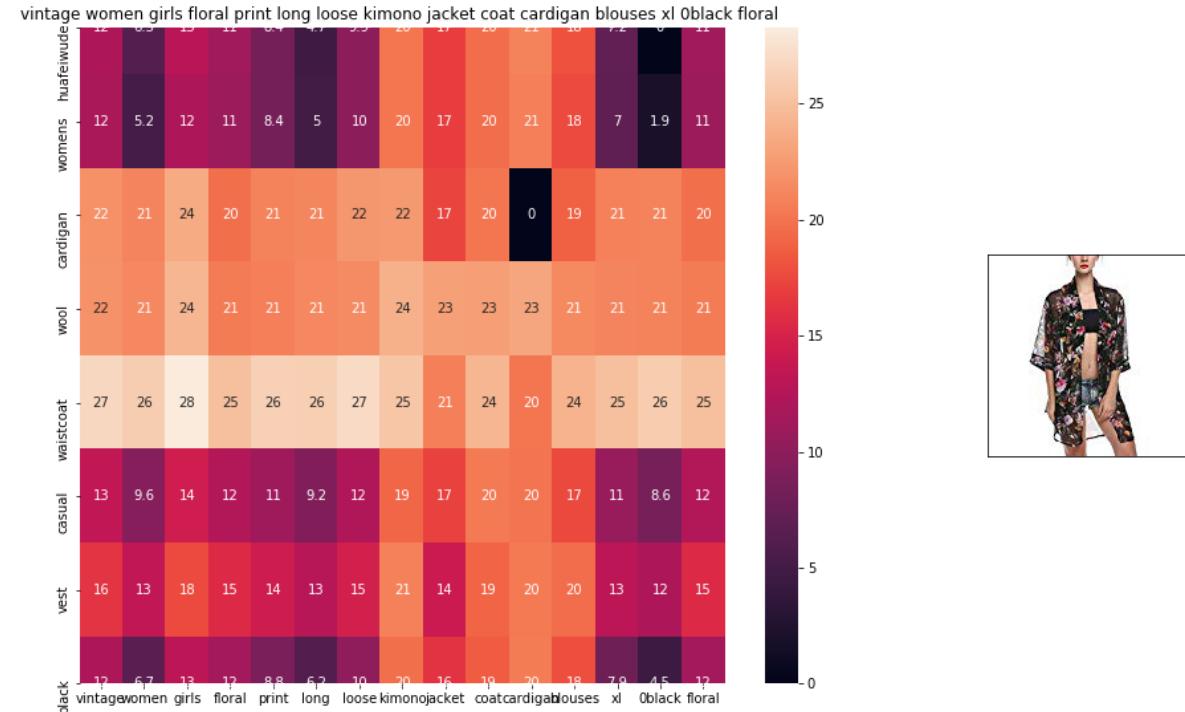
euclidean distance from input : 4.2146835



ASIN : B011R13YBM

Brand : Huayang

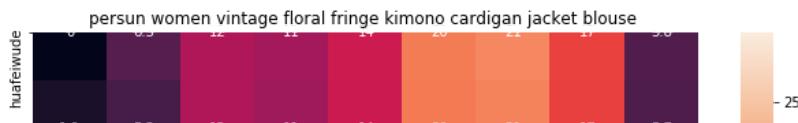
euclidean distance from input : 4.2270885



ASIN : B073PNYB2S

Brand : ACEFAST INC

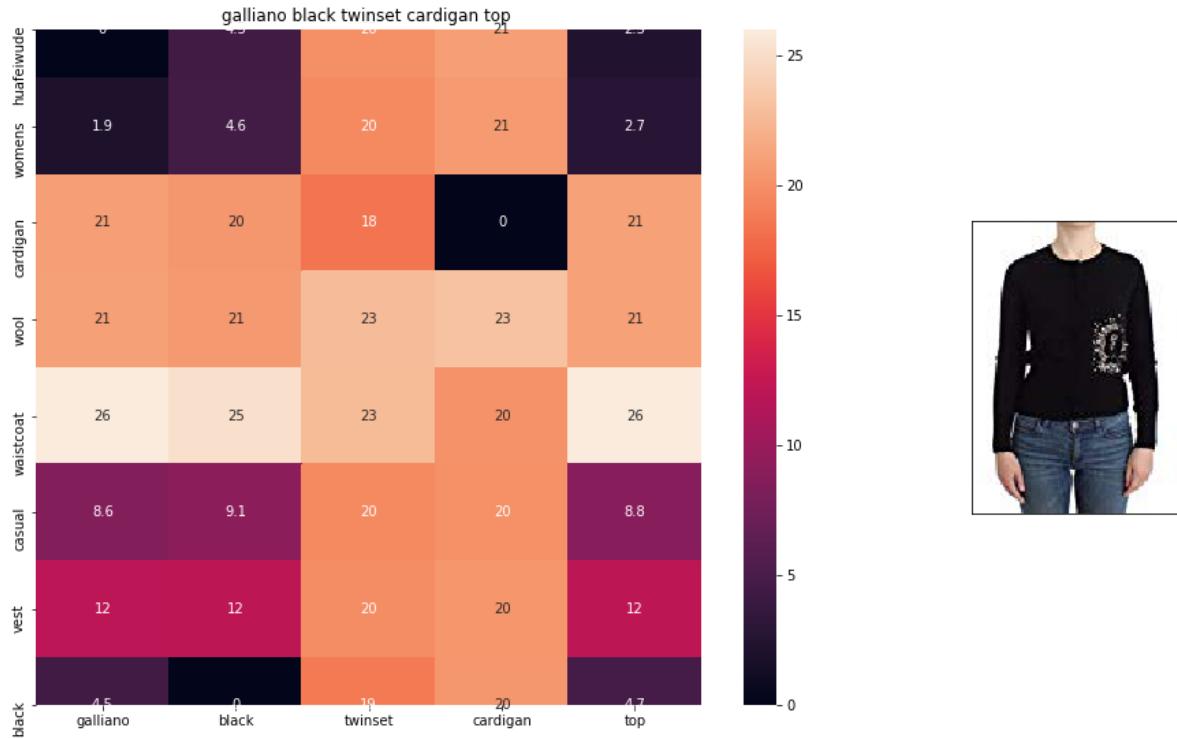
euclidean distance from input : 4.2848



ASIN : B00Y26RSRQ

Brand : Persun

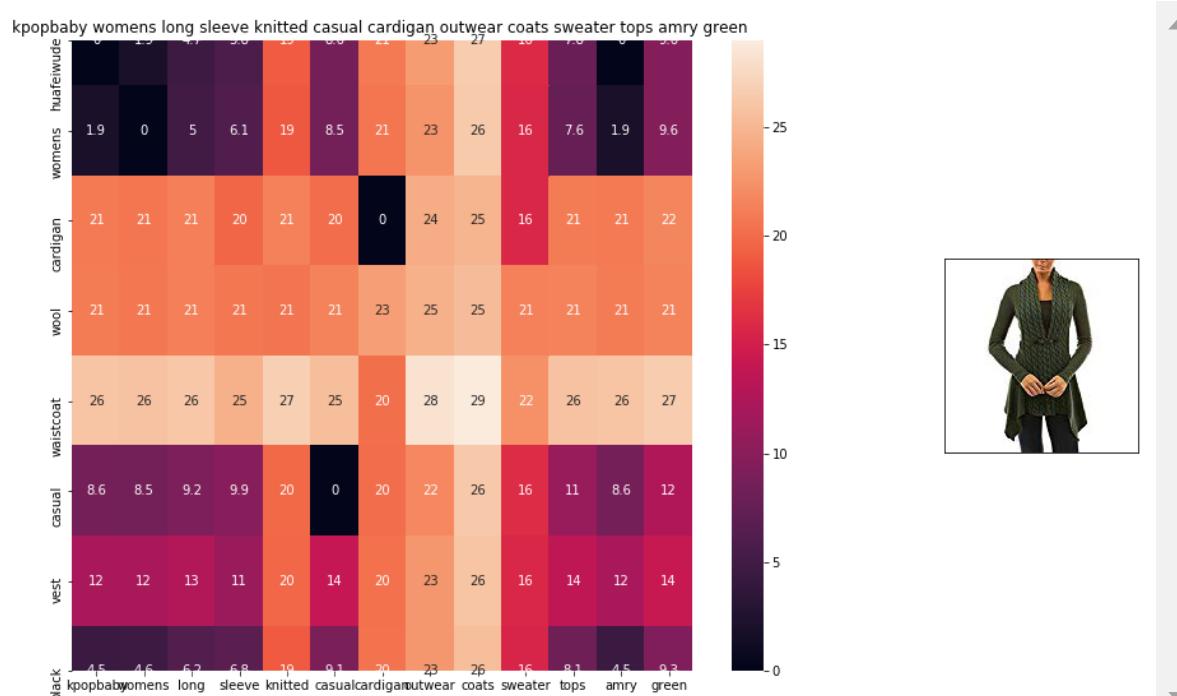
euclidean distance from input : 4.2976522



ASIN : B074G57HQJ

Brand : Galliano

euclidean distance from input : 4.3163776



ASIN : B074LD7G7K

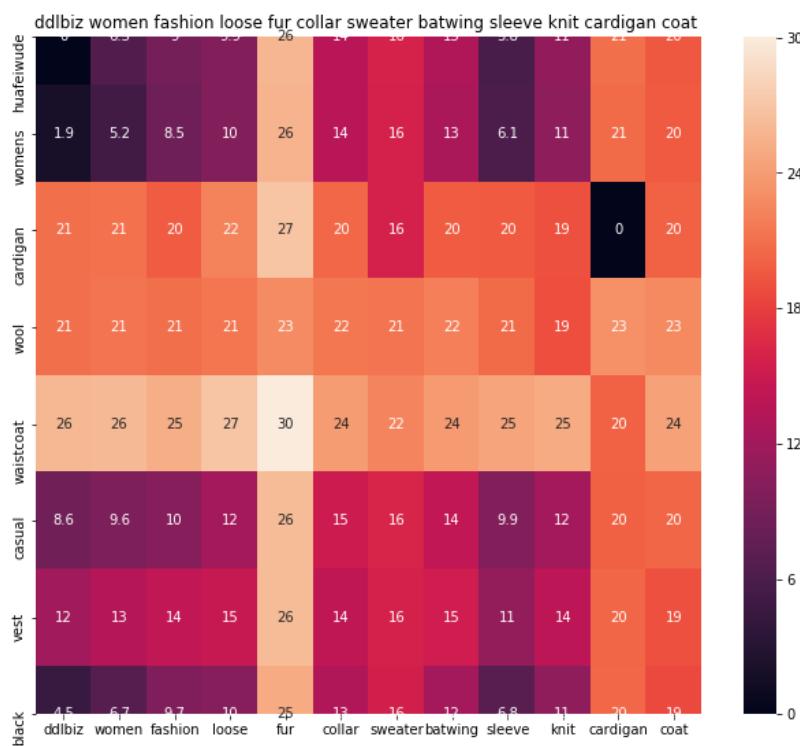
Brand : KpopBaby

euclidean distance from input : 4.35567

---



---



ASIN : B01MY76AQJ

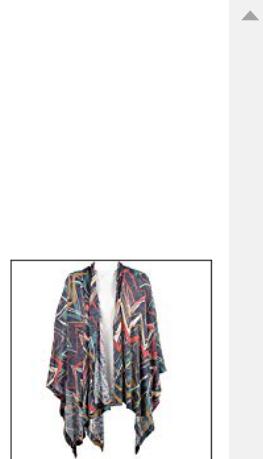
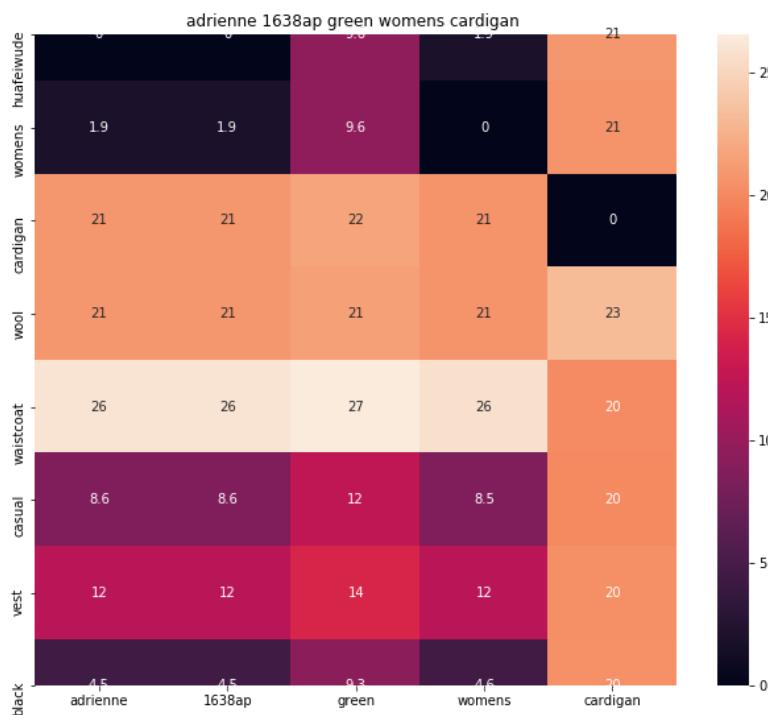
Brand : DDLBiz

euclidean distance from input : 4.408304

---



---



ASIN : B0001HW05W

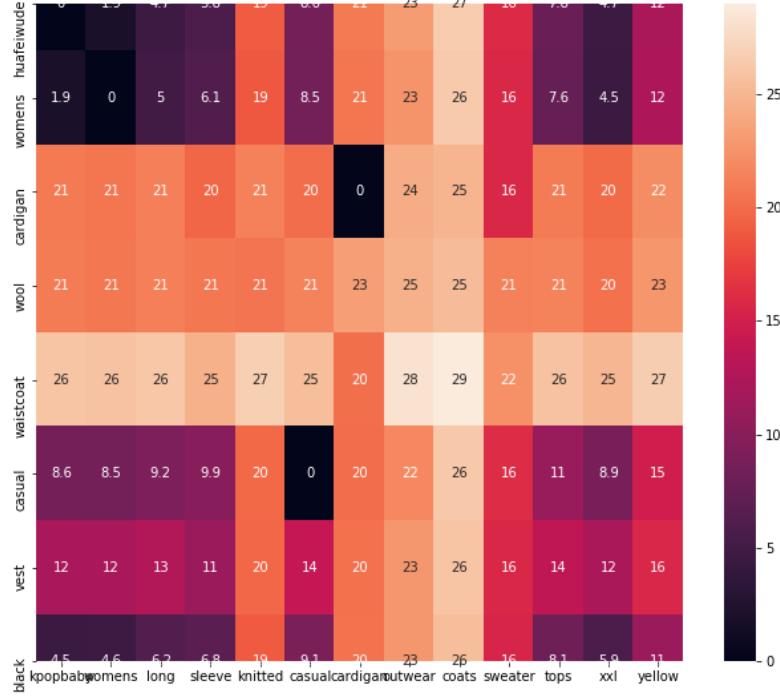
Brand : Liuqiuuh

euclidean distance from input : 4.422664

=====

=====

kpopbaby womens long sleeve knitted casual cardigan outwear coats sweater tops xxl yellow



ASIN : B074LCPJJZ

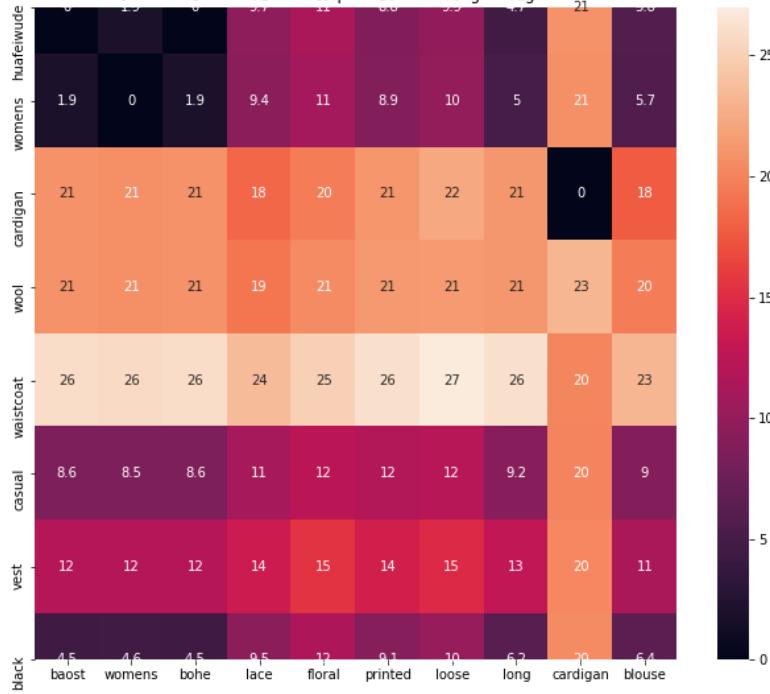
Brand : KpopBaby

euclidean distance from input : 4.499772

=====

=====

baost womens bohe lace floral printed loose long cardigan blouse



ASIN : B01GTA9352

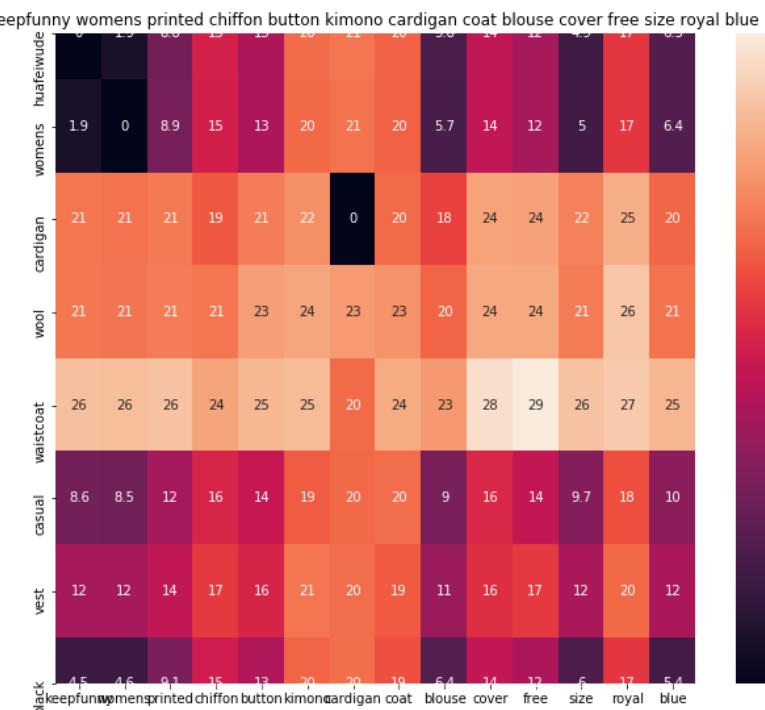
Brand : Baost

euclidean distance from input : 4.5200267

---



---



ASIN : B01IV2EEAK

Brand : KEEPFUNNY

euclidean distance from input : 4.5225873

---



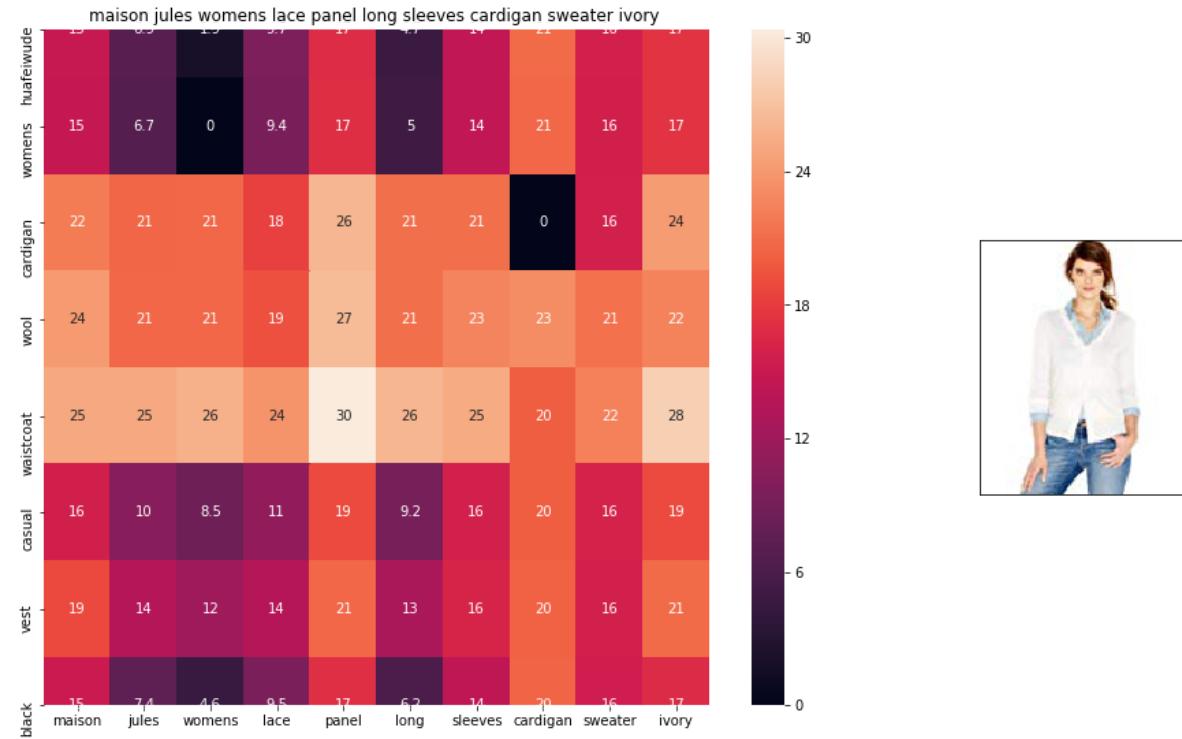
---



ASIN : B01AYVHQE2

Brand : Blansdi

euclidean distance from input : 4.534905



ASIN : B01DOKN4WS

Brand : Maison Jules

euclidean distance from input : 4.535721

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

In [59]:

```
# 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 [60]:

```

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) o
    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) o
    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]],
                   [data['asin'].loc[df_id2],brands[doc_id2], colors[doc_id2], types[doc_id2]]]

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

    # 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 lines 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 [61]:

```
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

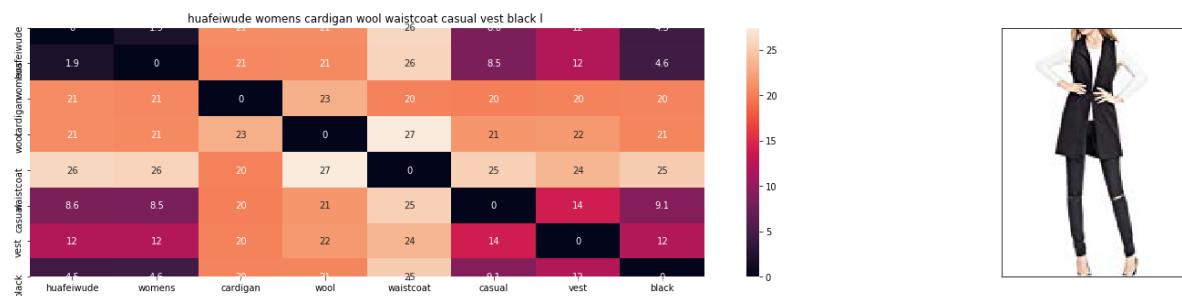
    # pairwise_dist will store the distance from given input apparel to all remaining apparel
    # the metric we used here is cosine, the coside distance is mesured as K(X, Y) = <X, Y>
    # http://scikit-learn.org/stable/modules/metrics.html#cosine-similarity
    idf_w2v_dist = pairwise_distances(w2v_title_weight, 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]

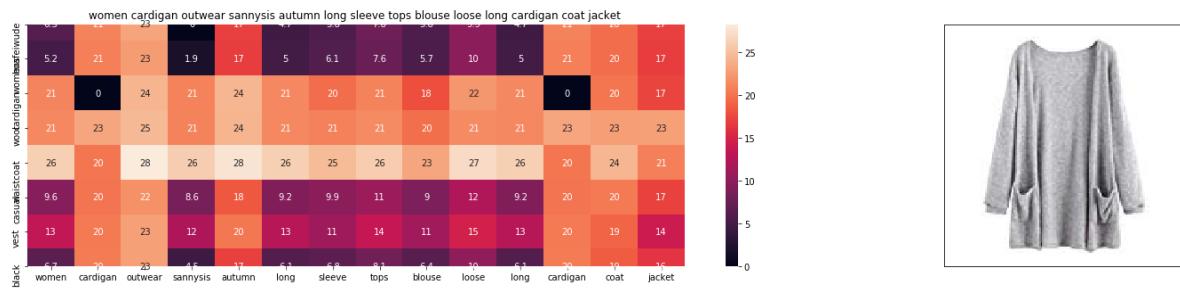
    #data frame indices of the 9 smallest distace's
    df_indices = list(data.index[indices])

    for i in range(0, len(indices)):
        heat_map_w2v_brand(data['title'].loc[df_indices[0]], data['title'].loc[df_indices[i]])
        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, 5, 5, 20)
# in the give heat map. each cell contains the euclidean distance between words i. i
```



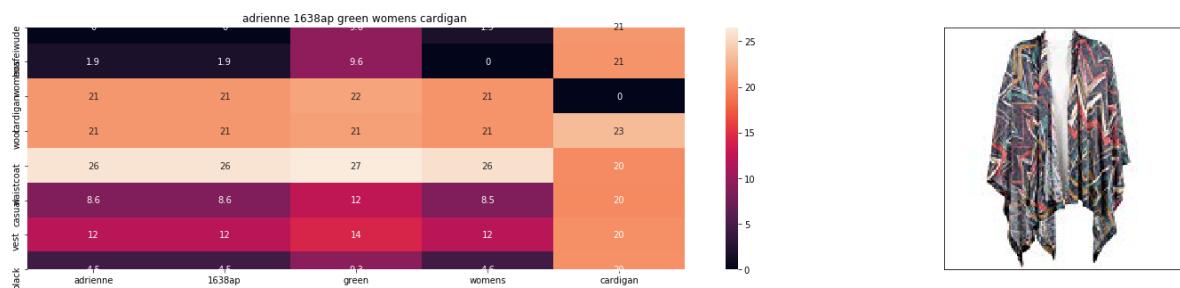
ASIN : B01MT96PXZ  
Brand : Huafeiwude  
euclidean distance from input : 0.6



ASIN : B07473KFK1

Brand : Sannysis

euclidean distance from input : 2.8984806060791017

=====  
=====


ASIN : B0001HW05W

Brand : Liuqiuuhu

euclidean distance from input : 2.9184389116186766

=====  
=====

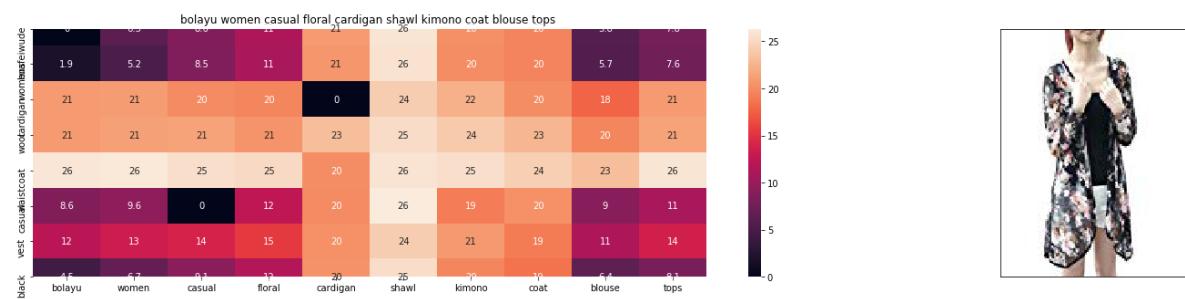

floyoung women printed kimono cardigan coat tops blouse loose jacket

**ASIN : B01D6EUG3W****Brand : FloYoung****euclidean distance from input : 2.9966798782348634**


---



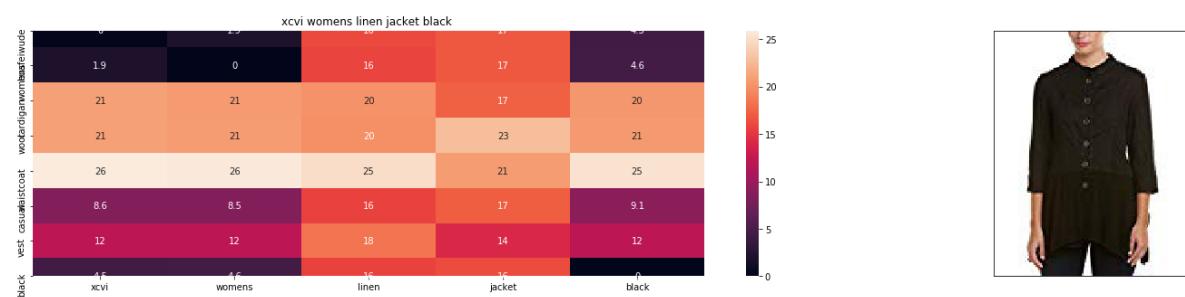
---

**ASIN : B01F6UHWEU****Brand : Bolayu****euclidean distance from input : 3.0006389619726805**


---



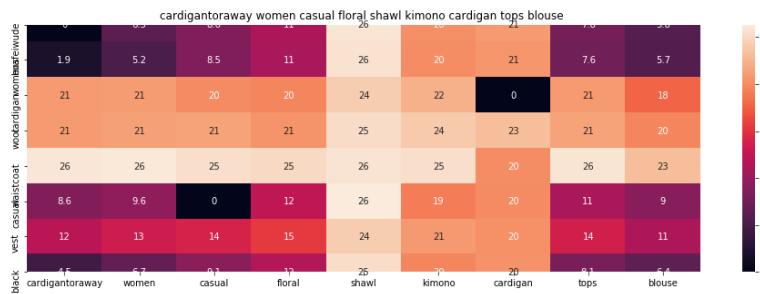
---

**ASIN : B01M31Q4Z0****Brand : XCVI****euclidean distance from input : 3.0484086992163326**


---



---



ASIN : B01CZMPCY4

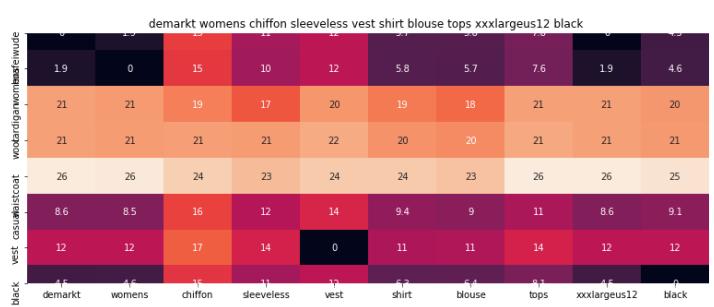
Brand : Toraway

euclidean distance from input : 3.050865936459985

---



---



ASIN : B00JKCQZJE

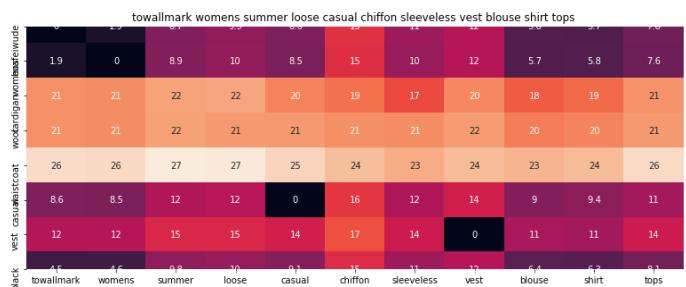
Brand : Demarkt

euclidean distance from input : 3.081585693540063

---



---



ASIN : B00QGEJ3MA

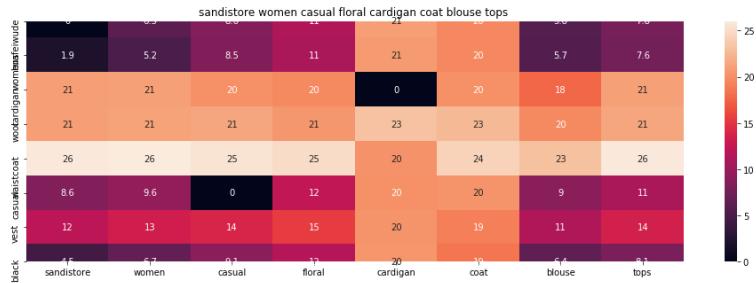
Brand : Towallmark

euclidean distance from input : 3.0982917787451414

---



---



ASIN : B01AVX8IOU

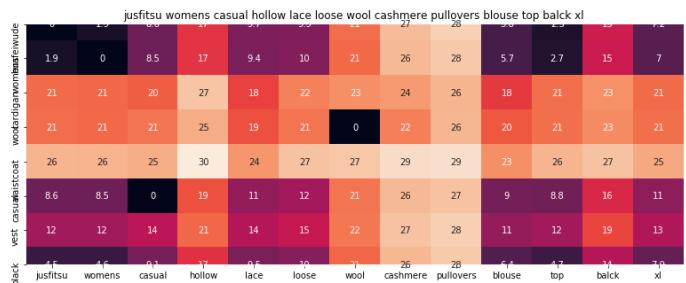
Brand : Sandistore

euclidean distance from input : 3.1044492795824143

---



---



ASIN : B01N96GX38

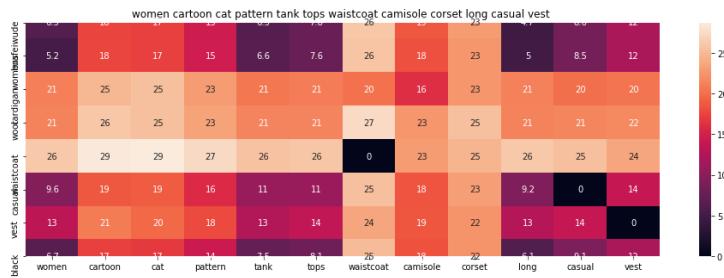
Brand : Jusfitsu

euclidean distance from input : 3.1124135972876217

---



---



ASIN : B011R13YBM

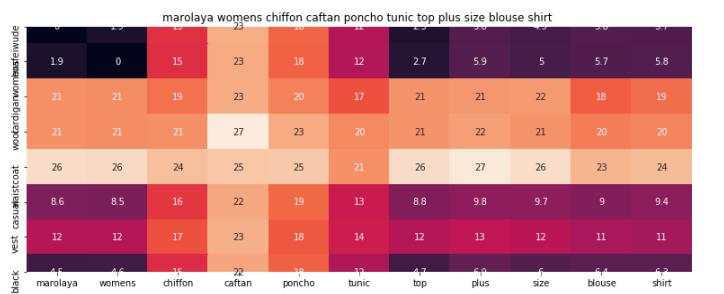
Brand : Huayang

euclidean distance from input : 3.1135442733764647

---



---



ASIN : B01CE40UW2

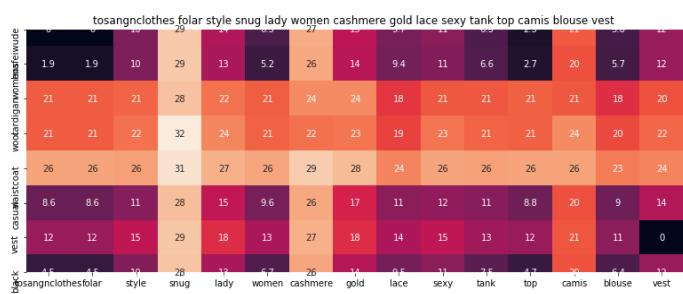
Brand : Marolaya

euclidean distance from input : 3.126744651975122

---



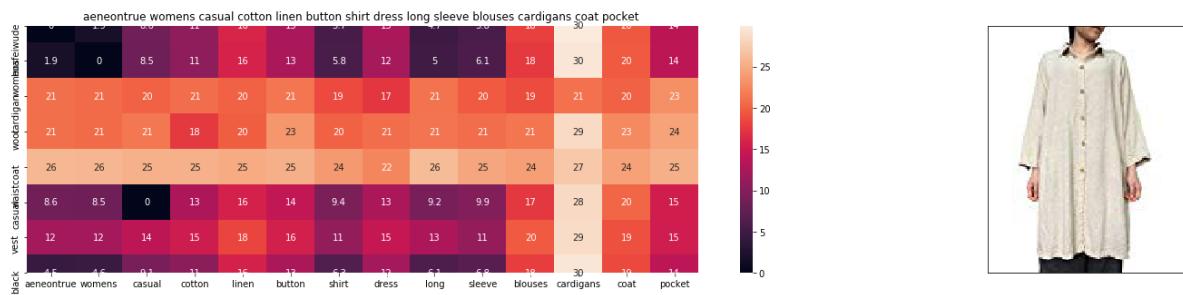
---



ASIN : B06XFRDH4J

Brand : Tosangn\_Clothes

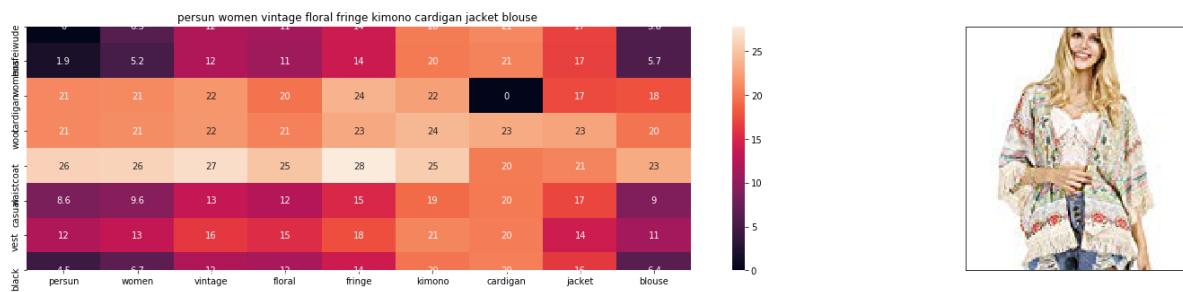
euclidean distance from input : 3.1292724611181884



ASIN : B074V1K5QJ

Brand : Aeneontrue

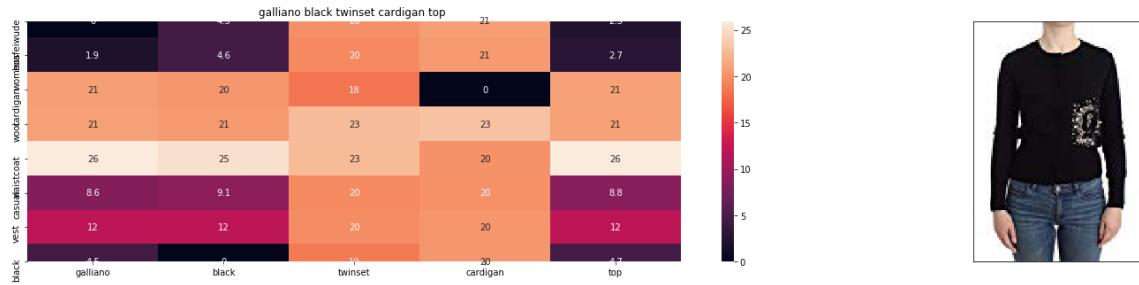
euclidean distance from input : 3.1459396436571216



ASIN : B00Y26RSRQ

Brand : Persun

euclidean distance from input : 3.1488262176513673



ASIN : B074G57HQJ

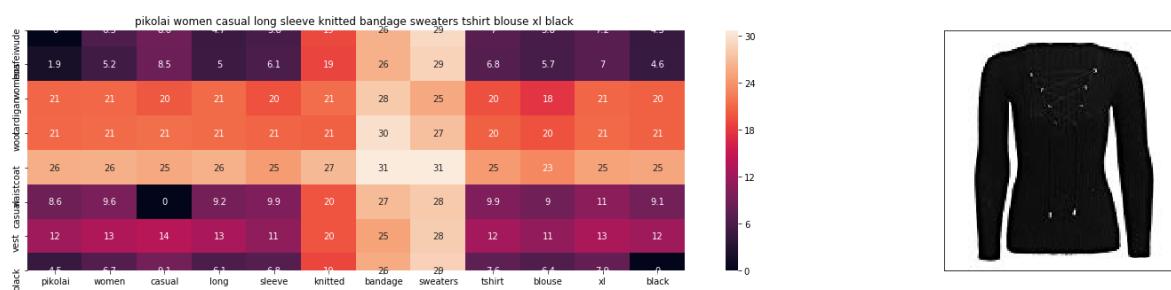
Brand : Galliano

euclidean distance from input : 3.158188819885254

---



---



ASIN : B01LZ2B1HT

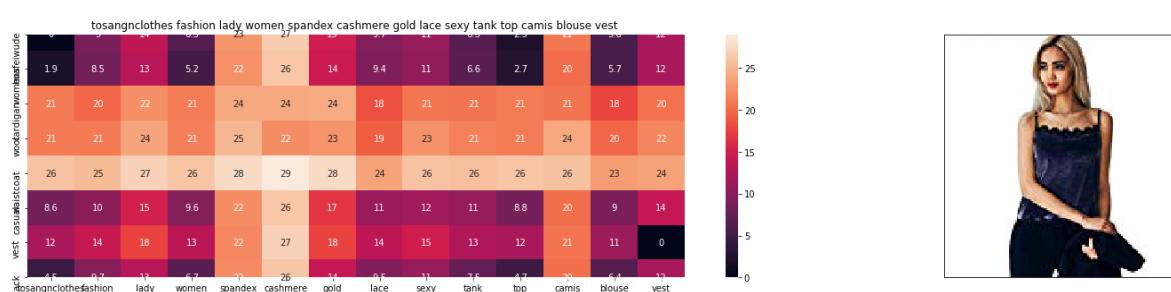
Brand : Pikolai

euclidean distance from input : 3.159334373654809

---



---



ASIN : B06XFFQLPL

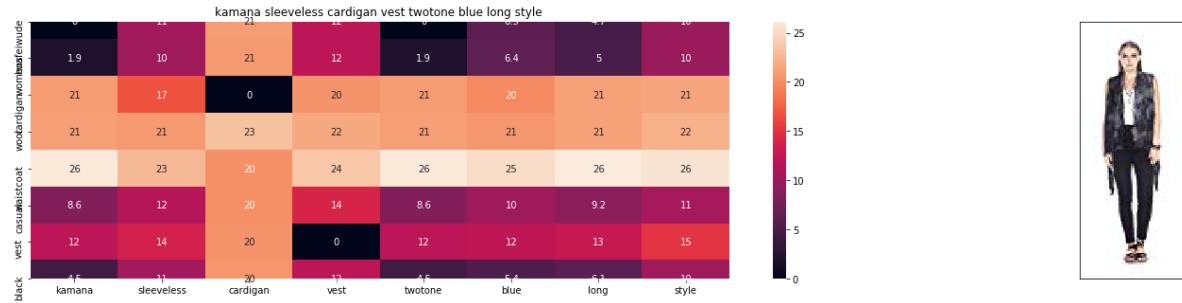
Brand : Tosangn\_Clothes

euclidean distance from input : 3.1688568117041256

---



---



ASIN : B0751686K7

Brand : Kamana

euclidean distance from input : 3.177324751152331

---

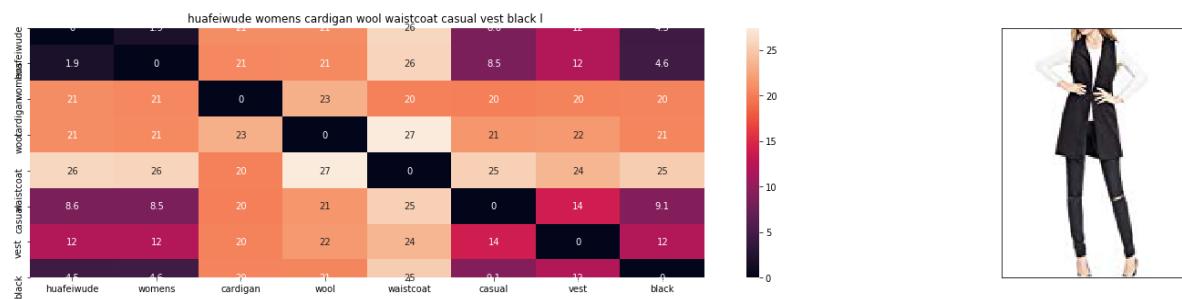


---

In [62]:

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

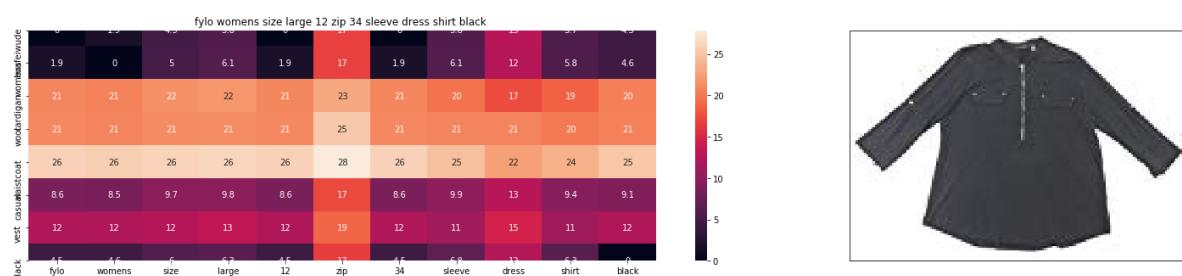
idf_w2v_brand(12566, 5, 50, 20)
```



ASIN : B01MT96PXZ

Brand : Huafeiwude

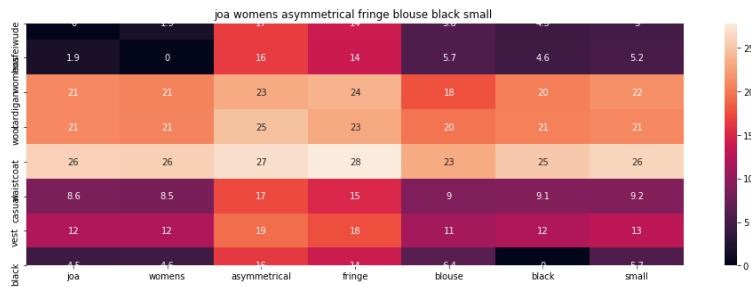
euclidean distance from input : 0.0



ASIN : B0718Y9J4M

Brand : f

euclidean distance from input : 1.4440810463645242



ASIN : B0721ZGF6S

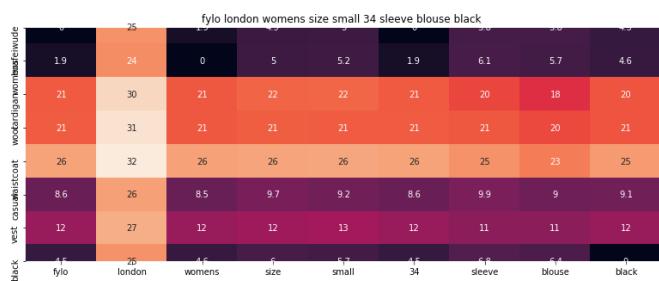
Brand : J.O.A.

euclidean distance from input : 1.4896895321932706

---



---



ASIN : B073HJRP1D

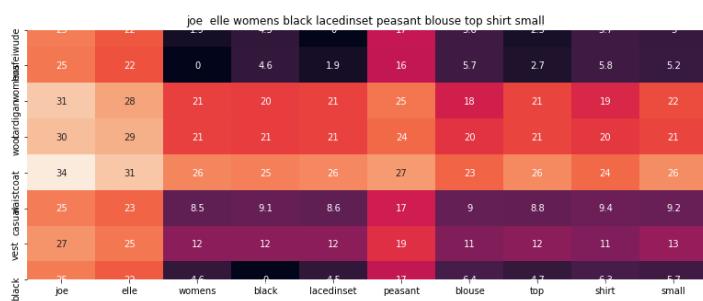
Brand : f

euclidean distance from input : 1.4966244784268465

---



---



ASIN : B073GDVFT2

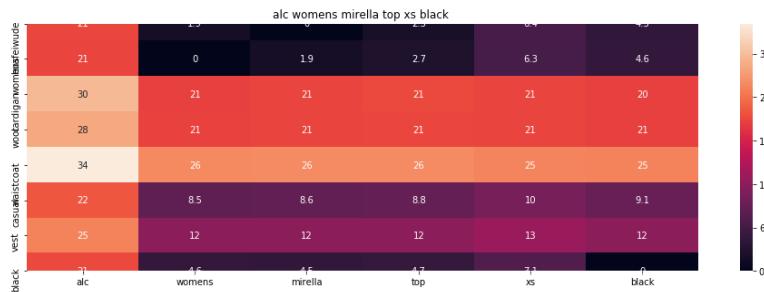
Brand : J&E

euclidean distance from input : 1.5201109452681107

---



---



ASIN : B072N8CS7N

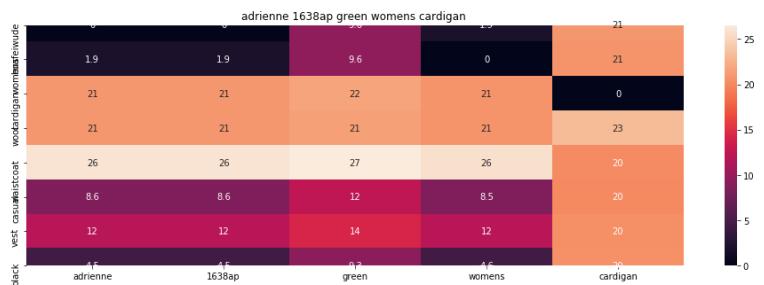
Brand : A.L.C.

euclidean distance from input : 1.6151734092018821

---



---



ASIN : B0001HW05W

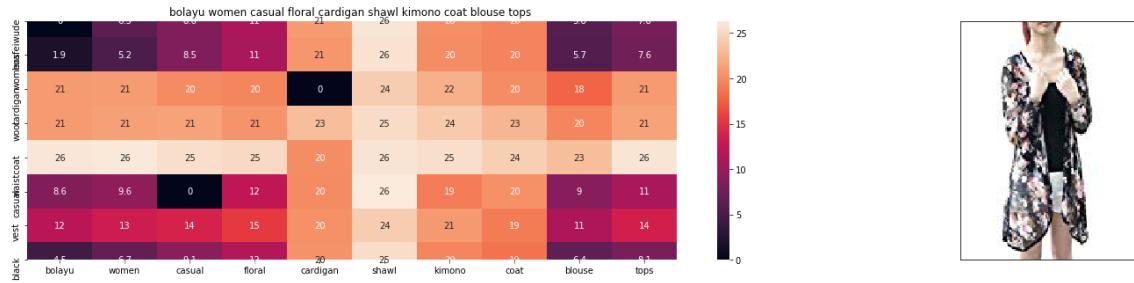
Brand : Liuqiuahu

euclidean distance from input : 1.6877090804177464

---



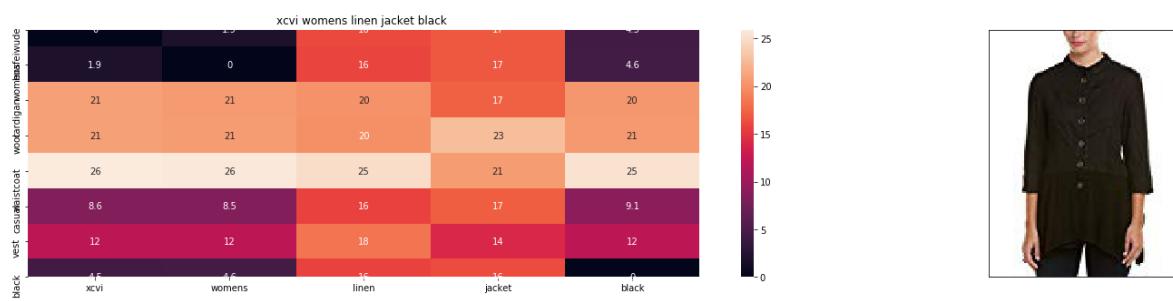
---



ASIN : B01F6UHWEU

Brand : Bolayu

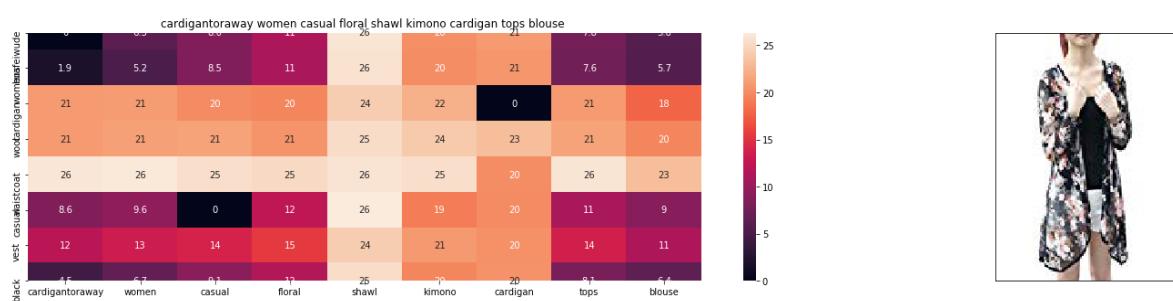
euclidean distance from input : 1.7026545441184742



ASIN : B01M31Q4Z0

Brand : XCVI

euclidean distance from input : 1.7113399508900473



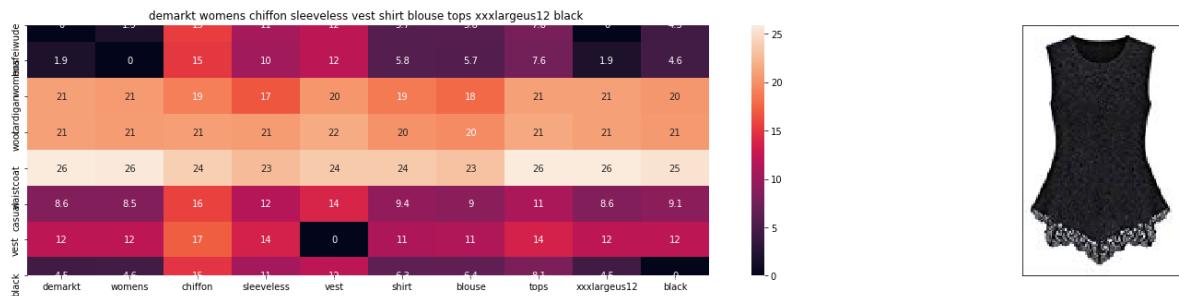
ASIN : B01CZMPCY4

Brand : Toraway

euclidean distance from input : 1.7117867212979843

=====

=====



ASIN : B00JKCQZJE

Brand : Demarkt

euclidean distance from input : 1.7173721316761803



ASIN : B00QGEJ3MA

Brand : Towallmark

euclidean distance from input : 1.7204096017134671



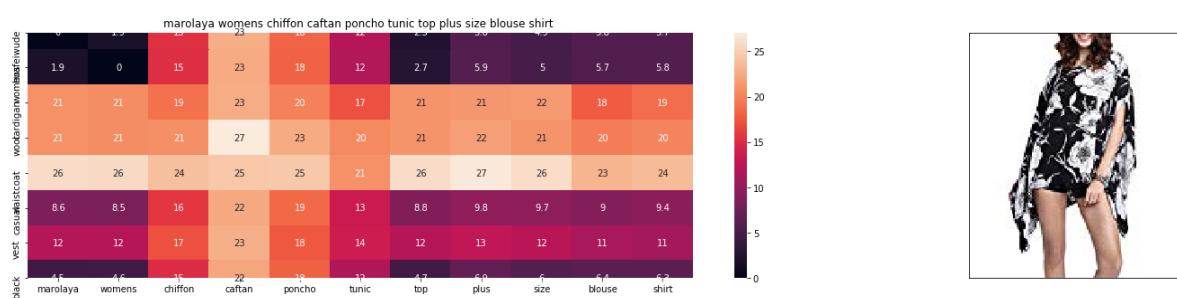


ASIN : B01N96GX38

Brand : Jusfitsu

euclidean distance from input : 1.7229772050848273

=====

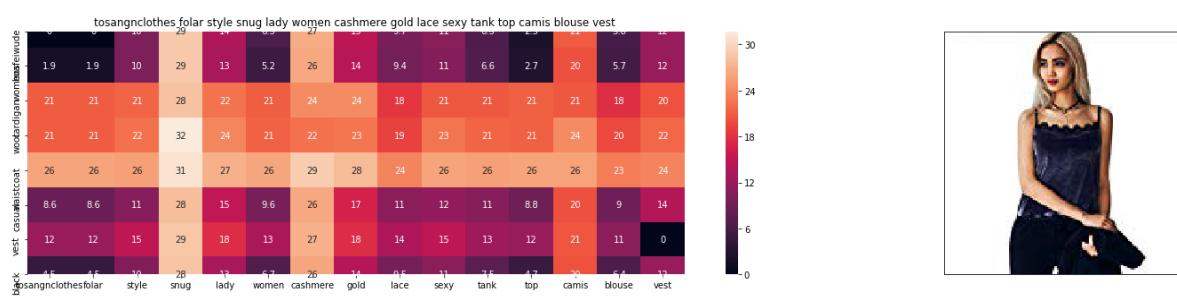


ASIN : B01CE40UW2

Brand : Marolaya

euclidean distance from input : 1.7255828513916454

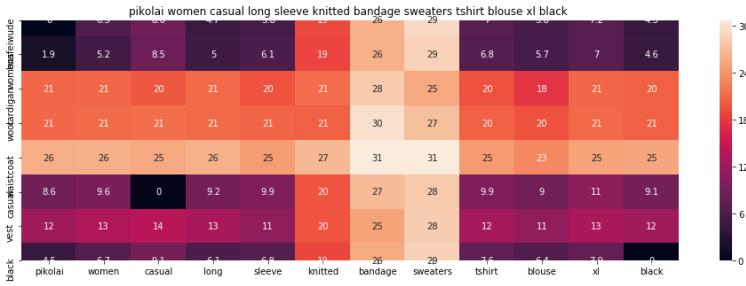
=====



ASIN : B06XFRDH4J

Brand : Tosangn\_Clothes

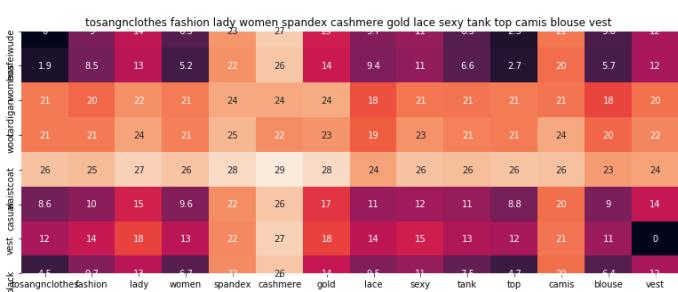
euclidean distance from input : 1.726042453054021



ASIN : B01LZ2B1HT

Brand : Pikolai

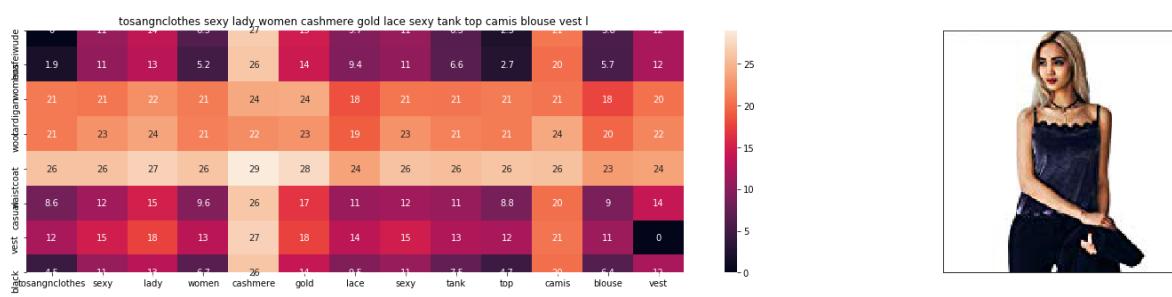
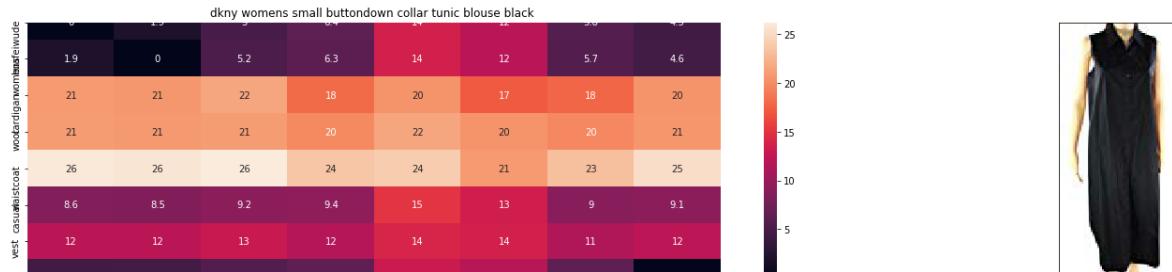
euclidean distance from input : 1.7315082553334067



ASIN : B06XFFQLPL

Brand : Tosangn\_Clothes

euclidean distance from input : 1.7332396077060097



## [10.2] Keras and Tensorflow to extract features

In [63]:

```
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 [0]:

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

# 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 codse 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 // batc
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 [ ]:

```
#Load the features and corresponding ASINS info.
bottleneck_features_train = np.load('16k_data_cnn_features.npy')
asins = np.load('16k_data_cnn_feature_asins.npy')
asins = list(asins)

# Load the original 16K dataset
data = pd.read_pickle('16k_apperal_data_preprocessed')
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

    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 idx, 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/' + asins[indices[i]]))

get_similar_products_cnn(931, 20)
```



Product Title: antthony maxine striped crossover top xl  
 Euclidean Distance from input image: 4.0460973e-06  
 Amazon Url: [www.amazon.com/dp/B01M73YDQQ](http://www.amazon.com/dp/B01M73YDQQ)



Product Title: ohconcept collection net cover tunic nude  
 Euclidean Distance from input image: 35.04788  
 Amazon Url: [www.amazon.com/dp/B01FGDBC6G](http://www.amazon.com/dp/B01FGDBC6G)



Product Title: mona b anchor fashion tee assorted sizes xlarge 16  
Euclidean Distance from input image: 35.384586  
Amazon Url: [www.amazon.com/dp/B06Y3PFHMB](http://www.amazon.com/dp/B06Y3PFHMB)



Product Title: perman womens long sleeve round neck splice shirt blouse tops pullover xl white  
Euclidean Distance from input image: 36.69707  
Amazon Url: [www.amazon.com/dp/B01LZEZIA6](http://www.amazon.com/dp/B01LZEZIA6)



Product Title: victorias secret pink campus henley longsleeve shirt large  
Euclidean Distance from input image: 36.9127  
Amazon Url: [www.amazon.com/dp/B0741W6L4W](http://www.amazon.com/dp/B0741W6L4W)



Product Title: sandistore womens girl slim tshirt blouses tops xl  
Euclidean Distance from input image: 37.293144  
Amazon Url: [www.amazon.com/dp/B00ZKZBEOQ](http://www.amazon.com/dp/B00ZKZBEOQ)

In [56]:

```
bottleneck_features_train.shape
```

Out[56]:

```
(16042, 25088)
```

In [ ]:

```
asins.shape
```

## ASSIGNMENT

In [6]:

```
#Load the features and corresponding ASINS info.
bottleneck_features_train = np.load('16k_data_cnn_features.npy')
asins = np.load('16k_data_cnn_feature_asins.npy')
asins = list(asins)

# Load the original 16K dataset
data = pd.read_pickle('16k_apperal_data_preprocessed-1')
df_asins = list(data['asin'])
```

In [7]:

```
bottleneck_features_train.shape
```

Out[7]:

```
(16042, 25088)
```

In [8]:

```
data.shape
```

Out[8]:

```
(16434, 7)
```

In [9]:

```
len(df_asins)
```

Out[9]:

```
16434
```

In [10]:

```
data = data[:16042]
data.shape
```

Out[10]:

```
(16042, 7)
```

In [11]:

```
df_asins = df_asins[:16042]
len(df_asins)
```

Out[11]:

16042

In [12]:

```
data.head()
```

Out[12]:

	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 [15]:

```
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 [16]:

```
len(w2v_title_weight)
```

Out[16]:

16042

In [17]:

```
# 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 [18]:

```
extra_features.shape
```

Out[18]:

```
(16042, 5740)
```

In [19]:

```
color_features.shape
```

Out[19]:

```
(16042, 1860)
```

In [20]:

```

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) o
    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) o
    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]],
                   [data['asin'].loc[df_id2],brands[doc_id2], colors[doc_id2], types[doc_id2]]]

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

    # 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 lines 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 [21]:

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

#Get similar products using CNN features (VGG-16)
def get_similar_products_cnn(doc_id, wt, wb, wc, wi, num_results):
    doc_id = asins.index(df_asins[doc_id])
    idf_w2v_dist = pairwise_distances(w2v_title_weight, w2v_title_weight[doc_id].reshape(1, -1))
    brand_dist = pairwise_distances(brand_features, brand_features[doc_id])
    color_dist = pairwise_distances(color_features, color_features[doc_id])
    bottleneck_features_dist = pairwise_distances(bottleneck_features_train, bottleneck_features_test)

    pairwise_dist = (wt * idf_w2v_dist + wb * brand_dist + wc * color_dist + wi * bottleneck_features_dist) / 4

    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 idx, 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/' + asins[indices[i]])
```

In [22]:

```
get_similar_products_cnn(12566, 1, 10, 5, 5, 20)
```



Product Title: huafeiwude womens cardigan wool waistcoat casual vest blac  
k 1

Euclidean Distance from input image: 1.3043908915105498e-06

Amazon Url: [www.amazon.com/dp/B01MT96PXZ](http://www.amazon.com/dp/B01MT96PXZ)



In [23]:

```
get_similar_products_cnn(12566,5,50,20,20,20)
```



Product Title: huafeiwude womens cardigan wool waistcoat casual vest black  
Euclidean Distance from input image: 1.153356156704065e-06

In [25]:

```
get_similar_products_cnn(12566,5,100,100,20,20)
```



Product Title: huafeiwude womens cardigan wool waistcoat casual vest black  
Euclidean Distance from input image: 4.86972599497272e-07  
Amazon Url: [www.amazon.com/dp/B01MT96PXZ](http://www.amazon.com/dp/B01MT96PXZ)



In [26]:

```
get_similar_products_cnn(12566,5,10,10,100,20)
```



Product Title: huafeiwude womens cardigan wool waistcoat casual vest black 1

Euclidean Distance from input image: 4.382753279060125e-06

Amazon Url: [www.amazon.com/dp/B01MT96PXZ](http://www.amazon.com/dp/B01MT96PXZ)



In [ ]: