

Configuration Manual

MSc Research Project Programme Name

Adelola Adebo Student ID: X19134711

School of Computing National College of Ireland

Supervisor: Dr Vladimir Milosavljevic

National College of Ireland



.....

MSc Project Submission Sheet

School of Computing

Adelola Adebo	jj
X19134711	

Student

Student TD

Name:

Student ID.	Data Analytics		2020
Programme:		Year:	
	Msc Research Project		
Module:			
	Dr Vladimir Milosavljevic		
Lecturer:	17 th August 2020		
Due Date	17 August 2020.		
	A Natural Language Processing Approach to a	a Skinca	are
Project Title:	Recommendation Engine		
	-		
			17
Word Count:	Page Count:		

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Signature:	
------------	------------	--

Date:

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple	
copies)	
Attach a Moodle submission receipt of the online project	
submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project,	
both for your own reference and in case a project is lost or mislaid. It is	
not sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Adelola Adebo Student ID: X19134711

1 Introduction

This configuration manual provides detailed documentation on the I.T solution deployed as part of the research thesis, A Natural Language Processing Approach to a Skincare Recommendation. The document covers all steps taken to implement the solution.

2 Hardware Setup

The deep learning models used in this research was implemented using Google Collaboratory cloud machine due to the large dataset that we used. The specifications for the host machine are:

- Operating system: MacOS 10.15.5
- Processor: 2.5 GHz Dual-Core Intel Core i5
- RAM: 8 GB
- Hard drive: 500 GB

3 Software environment and library dependencies

Python 3.7 is used for scripting. Jupyter notebook which is part of the Anaconda Python package is used for the overall implementation of this project. Other library dependencies used are:

- Numpy
- Pandas
- NLTK
- Genism
- Matplotlib
- Scikit-learn
- Tensorflow
- Keras

4 Datasets

A total of three datasets was used for this research. The first dataset used is skincare comments from Reddit which was sourced from Google BigQuery and downloaded as csv files. The second dataset was skincare products sourced from <u>here</u> and the third is a list of skincare ingredients sourced from <u>here</u>. All three datasets were ethically sourced.

Google Cloud Platform	My First Project 👻 🔍 Search products and resources
BigQuery () FEATURES & IN	IFO 🔤 SHORTCUT 🔀 BACK TO CLASSIC
y history	Unsaved query Edited
d queries	1 select subreddit, body, author, created_utc
history	3 where subreddit = 'SkincareAddiction'
sfers	
duled queries	
rvations	
igine	
ources + ADD DATA -	
earch for your tables and datasets 🛛 🕘	
ong-canto-257700	Run ▼ La Save query :::: Save view Schedule query ▼ And More
igquery-public-data	Query history C REFRESH
n-bigquery	4:52 PM Select subreddit,body,author,created_utc from `fh-bigquery.reddit_comments
roperati-data-public	4:48 PM Select subreddit,body,author,created_utc from `fh-bigquery.reddit_comments

A screenshot of downloading the Reddit comments from Google BigQuery.

Figure 1: Extracting the user review comments from Google BigQuery

The following sections show the project implementation.

5 Data cleaning and pre-processing

Loading and Cleaning the data: Removing punctuations, stopwords, HTML tags, special characters, and dropping irrelevant comments, etc was performed before modelling.



Load the data

In	[1:	1	<pre>csv_list = glob('reddit*.csv')</pre>
In	r	1:	1	csv list



Figure 2: Steps to clean the data.

6 Feature Engineering After cleaning the datasets, we used custom functions to match reviews with keywords from the dictionary created.

I	Matching
In []:	<pre>skincare_ingredients = pd.read_pickle('skincare_ingrednts.pkl') product_names = pd.read_pickle('product_names.pkl')</pre>
In []:	skincare_ingredients
In []:	1 product_names
In []:	<pre>1 skin_type = ['oily', 'dry', 'combination', 'normal', 'sensitive', 'acne', 2 'hyperpigmentation', 'anti-aging', 'discolouration', 'breakout']</pre>
In []:	<pre>1 extend_ingredients = ['ha', 'aha', 'bha', 'micellar water' 'vit c', 'pie', 'sunscreen', 'spf', 'snail cream',</pre>
	<pre>100 'LaNeige Water Bank Essence', 101 'Drunk Elephant C-Firma Day Serum', 102 'Boscia Makeup Breakup Cleansing Oil', 103 'Dr Jart Ceramidin Lotion', 104 'ole henriksen', 105 'shiseido sunscreen', 106 'Ole Henriksen Truth Serum', 107 'Ole Henriksen Banana Bright', 108 'Olay Regenerist Retinol24', 109 'aztec secret clay mask', 110 'la rocheposay anthelios antishine spf50' 111 112]</pre>
In []:	<pre>1 #extend lists 2 for item in extend_ingredients: 3 skincare_ingredients.append(item) 4 5 for item in extend_products: 6 product_names.append(item)</pre>
In []:	<pre>1 def skin_match(x, match): 2 splitreview = x.split(" ") 3 setreview = set(splitreview) 4 setmatch = set(match) 5 6 matches = setreview.intersection(setmatch) 7 matches = list(matches) 8 9 if len(matches) >= 1: 10 match_string = ','.join(matches) 11 return match_string 12 else: 13 return ''</pre>



Figure 3: Matching reviews with keywords from the skincare dictionary created.

All reclassified comments are joined into one final dataframe. The final output of keyword matching.

	Joinir	ng all the reclas	sified comments into one dataframe			
In [200]:	1 2 3 4 5 6 7	all_matches list_ = [] for file in df = pd lista final_match	<pre>= glob('newmatching*.csv') all_matches: .read_csv(file, lineterminator='\n ppend(df) = pd.concat(list_, ignore_index=T</pre>	') 'rue)		
In [201]:	1	final_match				
Out[201]:		Unnamed: 0	review	skin_type	products	ingredients
		0 0	moisturize every day things aveeno drugstorebu	acne	Vitamin C Serum, The Ordinary	niacinamide,salicylic acid,vitamin C
		1 1	similar routines skin typeconcerns going 28 ye	dry,acne	Salicylic Acid Cleanser	benzoyl peroxide, retinol, salicylic acid, vitami
	:	2 2	say anything else primary care appointment ja	dry,sensitive	cetaphil	hyaluronic acid,soap
		3 3	drugstore cleansers general tend harsh irritat	acne,oily	cerave	fragrance,menthol,salicylic acid
		4 4	acne pretty bad growing face back back somewha	acne,sensitive	Cerave Foaming Cleanser, cerave	fragrance
		.				
	2311	8 1151	dry skin burn lotion applied even without alco	dry	cerave,Cerave healing ointment	alcohol,coconut,coconut oil
	2311	9 1152	thayers hydrating toner nothing acne cerave cl	acne	cerave	salicylic acid
	2312	0 1153	sunscreenfor much products approach multiple h	dry	cerave	amino acids,squalane,urea
	2312	1 1154	ive used glow tonic impressed formulated prope	acne	Glow Tonic	ascorbic acid, fragrance, glycolic acid, salicyli
	2312	2 1155	cerave cleanser moisturizer treat acne acne sp	acne	cerave	salicylic acid,retinoid
	23123	3 rows × 5 colu	mns			

Figure 4: Final output of reclassified comments.

Sentiment Analysis

We perform sentiment analysis to infer ratings. The sentiment scores were derived using the VADER library in NLTK. The VADER sentiment lexicon was first updated with word sentiment scores from the SocialSent sentiment lexicon.

	Senti	ment Ar	nalys	is
In [228]:	1 dat	a = pd.re	ad_cs	v('final_match2.csv', lineterminator='\n')
In [229]:	1 ski	.ncare_sen	timent	<pre>t = pd.read_csv('sentiment lexicon.csv', header=None)</pre>
In [205]:	1 ski 2 ski	.ncare_sen .ncare_sen	timen timen	<pre>t = skincare_sentiment.drop(skincare_sentiment.columns[2], axis=1) t</pre>
Out[205]:		0	1	
	0	severe	-5.83	
	1	suffered	-5.68	
	2	developed	-5.57	
	3	painful	-5.54	
	4	massive	-5.48	
	4975	congrats	2.87	
altale An average	4976	beautiful	2.89	
click to expan	4977 co	ngratulations	2.90	
	4978	gorgeous	2.92	
	4979	lovely	2.95	
	4980 row	rs × 2 colum	าร	
In [206]:	1 new	_dict = d	ict(z	<pre>ip(skincare_sentiment[0], skincare_sentiment[1]))</pre>

Figure 5: Updating VADER sentiment lexicon with SocialSent sentiment lexicon

232]:	1	data[<mark>'sent</mark>	<pre>iment'] = data['review'].apply(la</pre>	ambda x: ser	ntiment(x))		
33]:	1	data					
33]:		Unnamed: 0	review	skin_type	products	ingredients	sentiment
		o 0	moisturize every day things aveen drugstorebu	acne	Vitamin C Serum, The Ordinary	niacinamide,salicylic acid,vitamin C	0.9798
		1 1	similar routines skin typeconcerns going 28 ye	dry,acne	Salicylic Acid Cleanser	benzoyl peroxide,retinol,salicylic acid,vitami	-0.7380
		2 2	say anything else primary care appointment ja	dry,sensitive	cetaphil	hyaluronic acid,soap	-0.9952
		3 3	drugstore cleansers general tend harsh irritat	acne,oily	cerave	fragrance,menthol,salicylic acid	0.9569
		4 4	acne pretty bad growing face back back somewha	acne,sensitive	Cerave Foaming Cleanser,cerave	fragrance	0.6286
							
	231	18 23118	dry skin burn lotion applied even without alco	dry	cerave,Cerave healing ointment	alcohol,coconut,coconut oil	-0.3356
	231	19 23119	thayers hydrating toner nothing acne cerave cl	acne	cerave	salicylic acid	0.0155
	231	20 23120	sunscreenfor much products approach multiple h	dry	cerave	amino acids,squalane,urea	0.9870
	231	21 23121	ive used glow tonic impressed formulated prope	acne	Glow Tonic	ascorbic acid, fragrance, glycolic acid, salicyli	0.9872
	231	22 23122	cerave cleanser moisturizer treat acne acne sp	acne	cerave	salicylic acid, retinoid	-0.8552
	2312	23 rows × 6 co	lumns				

Figure 6: Performing sentiment analysis using the VADER library in NLTK.

After deriving sentiment scores, we created a user-item matrix. The code below shows creating the final structured dataset used for modelling.



	37 38 39 40 41 42 43 44		item_i p_id = produc produc else: produc produc	<pre>ndex = prod product_id t_name.appe t_id.append t_name.appe t_id.append</pre>	uct_name [item_in nd(ing) (p_id) nd(ing) (len(pro	.index] dex] duct_io	(ing) d)+1)	
In [236]:	1 d	= pd.D	ataFrame({	'user_id':u	ser_id,	'produc	ct_id':product_id,'product_name':p	<pre>product_name,'skin_type':s</pre>
In [237]:	1 d							
Out[237]:		user_id	product_id	product_name	skin_type	rating	review	
		0 1	1	Vitamin C Serum	acne	0.9798	moisturize every day things aveeno drugstorebu	
		1 1	2	The Ordinary	acne	0.9798	moisturize every day things aveeno drugstorebu	
		2 1	3	niacinamide	acne	0.9798	moisturize every day things aveeno drugstorebu	
		3 1	4	salicylic acid	acne	0.9798	moisturize every day things aveeno drugstorebu	
		4 1	5	vitamin C	acne	0.9798	moisturize every day things aveeno drugstorebu	
	-	•						
	11524	3 23122	4	salicylic acid	acne	0.9872	ive used glow tonic impressed formulated prope	
	11524	4 23122	5	vitamin C	acne	0.9872	ive used glow tonic impressed formulated prope	
	11524	5 23123	15	cerave	acne	-0.8552	cerave cleanser moisturizer treat acne acne sp	
	11524	6 23123	4	salicylic acid	acne	-0.8552	cerave cleanser moisturizer treat acne acne sp	

Figure 7: Steps to create the final user-item matrix used for modeling

7 Modelling The four machine learning models implemented and evaluated are MF, HFT, NCF and NCF + reviews.

7.1 Modelling Matrix Factorisation

Matrix factorisation

```
1 from __future__ import print_function, division
2 from builtins import range, input
n [6]:
n [7]:
         1 import pickle
         2 import numpy as np
         3 import pandas as pd
         4 import matplotlib.pyplot as plt
         5 from sklearn.utils import shuffle
         6 from datetime import datetime
         1 # load in the data
n [8]:
         2 import os
         3 if not os.path.exists('user2product.json') or \
               not os.path.exists('product2user.json') or \
         4
               not os.path.exists('userproduct2rating.json') or \
not os.path.exists('userproduct2rating_test.json'):
         5
         6
               import preprocess2dict
         8
         9
        10 with open('user2product.json', 'rb') as f:
                user2product = pickle.load(f)
        11
        12
        13 with open('product2user.json', 'rb') as f:
        14
               product2user = pickle.load(f)
        15
        16 with open('userproduct2rating.json', 'rb') as f:
        17
                userproduct2rating = pickle.load(f)
        18
        19 with open('userproduct2rating_test.json', 'rb') as f:userproduct2rating_test = pickle.load(f)
        20
   21
   22 N = np.max(list(user2product.keys())) + 1
   23 # the test set may contain productss the train set doesn't have data on
   24 m1 = np.max(list(product2user.keys()))
   25 m2 = np.max([m for (u, m), r in userproduct2rating_test.items()])
   26 M = max(m1, m2) + 1
27 print("N:", N, "M:", M)
   28
   29
   30 # initialize variables
   31 K = 10 # latent dimensionality
   32 W = np.random.randn(N, K)
   33 b = np.zeros(N)
   34 U = np.random.randn(M, K)
   35 c = np.zeros(M)
   36 mu = np.mean(list(userproduct2rating.values()))
   37
   38 # prediction[i,j] = W[i].dot(U[j]) + b[i] + c.T[j] + mu
   39
   40 def get_loss(d):
   41
           N = float(len(d))
   42
           sse = 0
           for k, r in d.items():
   43
   44
               i, j = k
p = W[i].dot(U[j]) + b[i] + c[j] + mu
   45
   46
                sse += (p - r)*(p - r)
   47
            return sse / N
   48
   49
   50 # train the parameters
   51 epochs = 10
   52 reg = 0.0001 # regularization penalty
   53 train_losses = []
   54 test_losses = []
   55 for epoch in range(epochs):
   56
          print("epoch:", epoch)
opach_start = datatime_now()
```



Figure 8: Implementation of matrix factorization model.

MSE: 0.39

7.2 Modelling Hidden Factors as Topics (HFT)



D	<pre>r = model.fit(x=[df_train.userId.values, df_train.product_idx.values, tp_train], y=df_train.rating.values - mu, epochs=epochs, batch_size=128, validation_data=([df_test.userId.values, df_test.product_idx.values, tp_test], df_test.rating.values - mu))</pre>
	<pre># plot losses plt.plot(r.history['loss'], label="train loss") plt.plot(r.history['val_loss'], label="test loss") plt.legend() # plot mse plt.plot(r.history['mean_squared_error'], label="train mse") plt.plot(r.history['val_mean_squared_error'], label="test mse") plt.legend() plt.show()</pre>
⊳	Epoch 1/10 721/721 [====================================

Figure 9: Implementation of Hidden Factors as Topics model.

MSE: 0.66

7.3 Modelling Neural Collaborative Filtering (NCF)

```
# keras model
D
   u = Input(shape=(1,), name='Users')
   m = Input(shape=(1,), name='Items')
   u_embedding = Embedding(N, K)(u) # (N, 1, K)
   m_embedding = Embedding(M, K)(m) # (N, 1, K)
    u_embedding = Flatten()(u_embedding) # (N, K)
   m_embedding = Flatten()(m_embedding) # (N, K)
    x = Concatenate()([u_embedding, m_embedding]) # (N, 2K)
   x = Dense(400)(x)
   x = BatchNormalization()(x)
   x = Activation('relu')(x)
   x = Dropout(0.5)(x)
   x = Dense(100)(x)
   x = Activation('relu')(x)
   x = Dense(1)(x)
   model = Model(inputs=[u, m], outputs=x)
   model.compile(
     loss='mse',
      # optimizer='adam',
     # optimizer=Adam(lr=0.01),
     optimizer=SGD(lr=0.08, momentum=0.9),
     metrics=['mean_squared_error'],
    r = model.fit(
      x=[df_train.userId.values, df_train.product_idx.values],
      y=df_train.rating.values - mu,
      epochs=epochs,
```



Figure 10: Implementation of Neural Collaborative Filtering model.

MSE: 0.32

7.4 Modelling Proposed Methodology, Neural Collaborative Filtering + reviews.

```
N = df.userId.max() + 1 # number of users
M = df.product idx.max() + 1 # number of product
# split into train and test
df = shuffle(df)
cutoff = int(0.8*len(df))
df_train = df.iloc[:cutoff]
df_test = df.iloc[cutoff:]
# tokenizing word and converting into sequences.
vocab_size = 1000
embedding dim = 128
max_length = 400
trunc_type='post'
padding_type='post'
oov_tok = "<00V>"
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(df_train['review'])
word_index = tokenizer.word_index
sequences = tokenizer.texts_to_sequences(df_train['review'])
review_padded = pad_sequences(sequences,maxlen=max_length, padding=padding_type,
                       truncating=trunc_type)
testing_sequences = tokenizer.texts_to_sequences(df_test['review'])
review testing padded = pad sequences(testing sequences,maxlen=max length,
```



Figure 11: Implementation of Neural Collaborative Filtering + reviews model.

MSE: 0.30

References

Hamilton, W.L., Clark, K., Leskovec, J. and Jurafsky, D. (2016) 'Inducing domain-specific sentiment lexicons from unlabeled corpora', In *Proceedings of the Conference on Empirical Methods in Natural Language Processing. Conference on Empirical Methods in Natural Language Processing*, 2016, p. 595, ArXiv preprint arxiv:1606.02820.