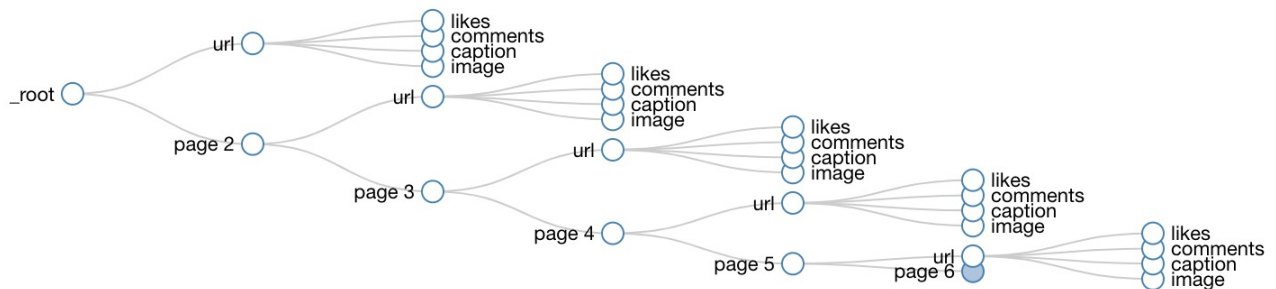


National Geographic - Instagram and image analytics

01. Scraping Instagram and collecting posts

I used a dedicated browser scraper extension for extracting image urls and corresponding likes and comments from the National Geographic's (NatGeo) instagram page. Since this scraper does not automatically scroll infinitely on the official instagram application, I used a secondary website 'picbon.com' which is an exact reflection of the instagram application and allows auto infinite scrolling while scraping.



02. Calculating Engagement

```
In [176]: 1 import pandas as pd
          2 pd.set_option('display.max_columns', None)
          3 pd.set_option('display.max_colwidth', 10000)
```

After scrapping posts from instagram I came up with a list of 634 posts, including information about the number of likes and comments realted to each post, the caption (description of the post) and the url of the post's image.

```
In [177]: 1 instagram_posts = pd.read_csv('instagram.csv')
```

```
In [178]: 1 instagram_posts.head()[ :1]
```

```
Out[178]:
```

	likes	comments	caption
0	444841	2107	<p>Photo By @BrianSkerry\nA West Indian Manatee calf nurses from its mom as they settle down on the sandy sea floor in the waters off the coast of Belize, an important country for these endangered animals. Manatees here live in mangroves and often sleep there overnight and feed on nearby sea grass beds during the day. Unlike Florida manatees, these animals are not nearly as acclimated to humans and are shyer. On this morning however, this very relaxed mom and calf allowed me into their world.\nFor more images and stories about ocean wildlife follow @BrianSkerry\n#manatees #belize #mesoamericanreef #endangeredspecies</p> <p>1.cdninstagram.com/vp/3a1c818472deb760c090c7f315/fr/e15/s1080x1080/39625475_248078652711752_6ig_cache_key=MTg1MzAwNjU3Nz</p>

Using the above collected urls and the Google Vision API, I came up with a list of labels for each post, while using the number of likes and comments I calculated engagement values for each post using the below formula:

$$EngagementValue = 0.4 \frac{likes}{maxLikes} + 0.6 \frac{comments}{maxComments}$$

Using this formula and the overall median engagement value of all the posts, I set the "Engagement" parameter as 1 (there was engagement related to a post) for all posts that had EngagementValue above or equal to the median and 0 for the rest of the posts, classifying all our posts into two discrete categories.

It was important to take into account the relevant engagement of each post, identifying the posts that were popular in relation to the rest of the National Geographic posts.

```
In [179]: 1 import statistics
          2 import csv
          3 from google.oauth2 import service_account
          4 from google.cloud import vision
```

```

5
6 SCRAPED_POSTS_FILE_PATH = 'instagram.csv'
7 MAX_POSTS = 634
8
9 class InstaPost:
10     def __init__(self, nlikes, ncomments, url='', caption='', labels=
11         self.url = url
12         self.caption = caption
13         self.nlikes = nlikes
14         self.ncomments = ncomments
15         self.labels = labels
16         self.engagement = engagement
17
18
19 def calc_engagement(instaposts):
20     max_likes = max(instaposts, key=lambda post: post.nlikes).nlikes
21     max_comments = max(instaposts, key=lambda post: post.ncomments).n
22     for i in range(0, len(instaposts)):
23         instaposts[i].nlikes = (instaposts[i].nlikes * 1.0) / max_lik
24         instaposts[i].ncomments = (instaposts[i].ncomments * 1.0) / m
25         instaposts[i].engagement = (4 * instaposts[i].nlikes) + (6 *
26         instaposts[i].engagement_dec = (4 * instaposts[i].nlikes) + (
27     median_engagement = statistics.median(map(lambda post: post.engag
28     for post in instaposts:
29         post.engagement = 1 if post.engagement >= median_engagement e
30
31
32 def get_posts_attrs(instaposts):
33     credentials = service_account.Credentials.from_service_account_fi
34     client = vision.ImageAnnotatorClient(credentials=credentials)
35     image = vision.types.Image()
36     for post in instaposts:
37         image.source.image_uri = post.url
38         response = client.label_detection(image=image)
39         post.labels = [label.description for label in response.label_
40
41
42 def import_scraped_posts(file_path):
43     post_counter = 0
44     instaposts = []
45     with open(file_path, 'r', encoding="utf8") as csv_file:
46         csv_reader = csv.reader(csv_file, quotechar='"', delimiter=',
47         next(csv_reader)
48         for row in csv_reader:
49             if post_counter > MAX_POSTS:
50                 break
51             post = InstaPost(int(row[0]), int(row[1]), row[3], row[2])
52             instaposts.append(post)
53             post_counter += 1
54     return instaposts

```

```

55
56 ###Testing only###
57 #post1 = InstaPost(20, 30, 'https://www.gettyimages.com/gi-resources/
58 #post2 = InstaPost(30, 40, 'https://il.wp.com/thefreshimages.com/wp-c
59 #instaposts = [post1,post2]
60
61 #calc_engagement(instaposts)
62 #get_posts_attrs(instaposts)
63
64 #print(instaposts[0].engagement)
65 #print(instaposts[1].engagement)
66 #print(instaposts[0].labels)
67 #print(instaposts[1].labels)
68
69 instaposts = import_scraped_posts(SCRAPED_POSTS_FILE_PATH)
70 calc_engagement(instaposts)
71 get_posts_attrs(instaposts)
72 for post in instaposts:
73     print(post.ncomments,post.nlikes,post.engagement,post.labels)
74
75 #####
76
77
, wood , animal source foods ]
0.19176983035304906 0.6878308130727165 1 ['surgeon', 'operating theate
r', 'service', 'medical', 'hospital']
0.012408298945437872 0.2946170420115282 1 ['sea', 'coastal and oceanic
landforms', 'body of water', 'ocean', 'shore', 'sky', 'horizon', 'coas
t', 'wave', 'water']
0.006189821182943604 0.15481301267366243 0 ['bird', 'seabird', 'albatr
oss', 'water', 'sea', 'wave', 'ocean', 'fauna', 'beak', 'wind wave']
0.025876891334250345 0.23144040515075376 1 ['climbing', 'rock climbing
', 'rock', 'sport climbing', 'adventure', 'tree', 'outdoor recreation'
, 'rock climbing equipment', 'outcrop', 'free climbing']
0.006476386978450252 0.13288870824711793 0 ['bazaar', 'public space',
'marketplace', 'market', 'shopkeeper', 'city']
0.015780223139232767 0.40485458080845405 1 ['wildlife', 'lion', 'mamma
l', 'fauna', 'terrestrial animal', 'masai lion', 'big cats', 'zoo', 'o
rganism', 'cat like mammal']
1.0 0.8895702085796631 1 ['rhinoceros', 'mammal', 'wildlife', 'terrest
rial animal', 'fauna', 'horn', 'eye', 'snout', 'organism', 'grass']
0.0069731010239951095 0.08569559654892107 0 ['fun', 'recreation', 'gir
l', 'leisure']

```

```

In [222]: 1 all_posts = [(post.ncomments, post.nlikes, post.labels) for post in
2 df_posts = pd.DataFrame(all_posts, columns=['ncomments', 'nlikes', '
3 df_posts['caption'] = instagram_posts.iloc[:, 'caption']
4 df_posts['engagement'] = [post.engagement for post in instaposts]
5 df_posts['engagement_dec'] = [post.engagement_dec for post in insta

```

```
In [223]: 1 df_posts[:1]
```

```
Out[223]:
```

	ncomments	nlikes	labels	caption	engagement	engagement_dec
0	0.020126	0.256822	[marine mammal, fauna, water, mammal, marine biology, underwater, manatee, organism, wildlife, dugong]	Photo By @BrianSkerry\nA West Indian Manatee calf nurses from its mom as they settle down on the sandy sea floor in the waters off the coast of Belize, an important country for these endangered animals. Manatees here live in mangroves and often sleep there overnight and feed on nearby sea grass beds during the day. Unlike Florida manatees, these animals are not nearly as acclimated to humans and are shyer. On this morning however, this very relaxed mom and calf allowed me into their world.\nFor more images and stories about ocean wildlife follow @BrianSkerry\n#manatees #belize #mesoamericanreef #endangeredspecies	1	1.148048

```
In [224]: 1 df_posts.to_csv('all_data.csv')
```

03. Logistic Regression using image labels and engagement

Having collected and preprocessed my data, I decided to run a logistic regression model to make predictions and identify the key features of a successful image post. But before I could do that, I first created a list of the features, which in this case were the unique labels from all the posts.

```
In [225]: 1 all_data = pd.read_csv('all_data.csv')
```

```
In [184]: 1 # Getting distinct labels to convert it into features
2 import ast
3
4 labels_list=[]
5
6 for post_labels in all_data.iloc[:,['labels']]:
7     labels_list += ast.literal_eval(post_labels)
8 label_set=set(labels_list)
9 label_set=sorted(list(label_set))
10
11 print('The total number of features is: '+ str(len(label_set)))
```

The total number of features is: 913

Having set the list of features, we created a dataframe indicating the features that appear in every post along with the engagement related to the post.

```
In [185]: 1 import pandas as pd
2 import numpy as np
3
4 # creating an empty dataframe
5 labels = [ast.literal_eval(post_labels) for post_labels in all_data.
6
7 colmns = label_set
8 df = pd.DataFrame(0.0, index = np.arange(len(labels)), columns = col
9 df['Engagement'] = all_data.iloc[:,['engagement']]
```

```
In [186]: 1 import math
2 from textblob import TextBlob as tb
3
4 def tf(word, blob):
5     return blob.words.count(word) / len(blob.words)
6
7 def n_containing(word, bloblist):
8     return sum(1 for blob in bloblist if word in blob.words)
9
10 def idf(word, bloblist):
11     return math.log(len(bloblist) / (1 + n_containing(word, bloblist)
12
13 def tfidf(word, blob, bloblist):
14     return tf(word, blob) * idf(word, bloblist)
```

```
In [187]: 1 # make label_sent which is a list that contains all the labels conca
2 label_sent = []
3
4 for i in range(len(labels)):
5     label_sent.append(" ".join(labels[i]))
```

```
In [188]: 1 # convert the post list to textblob format for passing into tf-idf f
2 labels_tb = []
3
4 for labels in label_sent:
5     labels_tb.append(tb(labels))
```

```
In [189]: 1 bloblist = labels_tb
2
3 tfidf_list = []
4
5 # list containing all the labels with tfidf score
6
7 for i, blob in enumerate(bloblist):
8     scores = {word: tfidf(word, blob, bloblist) for word in blob.words}
9     tfidf_list.append(scores)
10
```

```
In [190]: 1 #Assigning value of feature according to the label set of each post
2
3 data_ind = 0
4 for index, row in df.iterrows():
5     for column in df:
6         if column in tfidf_list[data_ind]:
7             # print(index,column,tfidf_list[data_ind][column])
8
9             df.at[index,column] = tfidf_list[data_ind][column]
10     data_ind=data_ind+1
```

```
In [191]: 1 columns=['Post']+label_set
2 #Divide data set into Features and Outcome
3 X=df[label_set]
4 y = df['Engagement']
```

```
In [192]: 1 import matplotlib
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score, confusion_matrix
7
8 #Divide into test and train data
9 X_train, X_test, y_train, y_test = train_test_split(
10     X, y, test_size=0.3, random_state=0)
11
12 #Logistic Regression
13 lr = LogisticRegression(C=2, random_state=42)
14 lr.fit(X_train, y_train)
```

```
Out[192]: LogisticRegression(C=2, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs
                             =1,
                             penalty='l2', random_state=42, solver='liblinear', tol=0.000
                             1,
                             verbose=0, warm_start=False)
```

```
In [193]: 1 #Predicting test outcome
2 y_pred = lr.predict(X_test)
3
4 #Calculate Accuracy
5 print('Accuracy: %.2f' % accuracy_score(y_test,y_pred))
6
7 #####Extra part to calculate probability of prediction
8 # y_pred = lr.predict_proba(X_test)
9 # from sklearn.metrics import roc_curve, auc, roc_auc_score
10 # false_positive_rate, true_positive_rate, thresholds = roc_curve(y_
11 # print (auc(false_positive_rate, true_positive_rate))
12 # print (roc_auc_score(y_train, lr.predict(X_train)))
```

Accuracy: 0.74

```
In [194]: 1 #Confusion Matrix
2 cnf_matrix = confusion_matrix(y_test, y_pred)
3 cnf_matrix
```

```
Out[194]: array([[64, 32],
                 [17, 78]], dtype=int64)
```


The small accuracy of the model was mostly attributed to the false positive values (based on the confusion matrix), while I saw that it improved drastically with the introduction of more data. Indicatively, with the use of 450 posts was only getting an accuracy slightly above 50%, while the above was a result of using 634.

04. Logistic Regression using captions and engagement

Not being satisfied with my results, my second attempt included the captions instead of the Google vision API image labels. In order to complete the regression this time, I first ran through the captions and made the relevant cleaning (removing stop words) and replacements in order to optimize my results.

```
In [195]: 1 from nltk import pos_tag
2 from nltk.tokenize import word_tokenize
3 from nltk.corpus import stopwords
4 from collections import Counter
5 import nltk
6 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\lovek\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

Out[195]: True

```
In [196]: 1 def count_words(dataframe, column):
2     captions = []
3     for i in range(dataframe.shape[0]):
4         captions.append(pd.DataFrame(dataframe.iloc[:, column]).loc[
5             wordcount = Counter(pos_tag(word_tokenize(''.join(captions))))
6             word_list = sorted(list(wordcount.items()), key = lambda w: -w[1]
7             stoplist = nltk.corpus.stopwords.words('english')
8             word_list = [word_list[i] for i in range(len(word_list))
9                 if len(word_list[i][0][0]) > 2 and word_list[i][0][0]
10                 and word_list[i][0][1] in ['NN', 'NNP', 'NNS', 'JJ']
11     return word_list
```

```
In [197]: 1 word_list = count_words(instagram_posts, 'caption')
2 len(word_list)
```

Out[197]: 9453

```
In [198]: 1 animal = ['species', 'wildlife', 'animals', 'elephants', 'lions', 'b
2             'bird', 'fish', 'insect', 'beast', 'mammal', 'habitat', 'fish'
```

```

2         'bird', 'lion', 'ivory', 'bear', 'pride', 'habitat', 'is
3         'jaguar', 'horses', 'snake', 'animal', 'dragon', 'lewa_wil
4         'whale', 'penguins', 'bull', 'rhinos', 'horse', 'flamingo'
5         'tortoises', 'rhinos', 'cat', 'tigers', 'insects', 'goat',
6
7     human = ['people', 'work', 'family', 'women', 'children', 'refugee',
8             'mother', 'camp', 'men', 'girls', 'rohingya', 'humans', 'ki
9             'son', 'baby', 'tourists', 'rangers', 'friends', 'father',
10            'families', 'woman', 'wife', 'residents', 'brother', 'poach
11            'americans', 'president', 'survivor', 'kid', 'chancellor',
12            'visitors']
13
14     photographer = ['timlaman', 'stephenwilkes', 'stevewinterphoto', 'mi
15                    'gabrielegalimbertipphoto', 'williamalbertallard', 'c
16                    'renaeffendipphoto', 'jimmy_chin', 'muhammedmuheisenp
17                    'mmuheisen', 'george', 'paleyphoto', 'muheisen', 'ca
18                    'hammond_robin', 'beverlyjoubert', 'brianskerry', 'p
19                    'christineeckstrom', 'salvarezphoto', 'florianschulz
20                    'katieorlinsky', 'davidalanharvey', 'ljohnphoto', 'p
21
22     place = ['myanmar', 'china', 'india', 'bangladesh', 'peru', 'borneo'
23            'california', 'matera', 'afghanistan', 'africa', 'capital',
24            'alaska', 'indonesia', 'madagascar', 'nation', 'greenland',
25            'ghana', 'france', 'north', 'south', 'east', 'west', 'greec
26            'kenya', 'japan', 'florida']
27
28     advertisement = ['africanparksnetwork', 'samsungmobileusa', 'withgala
29
30     nature = ['water', 'river', 'sea', 'ice', 'place', 'image', 'area',
31            'coast', 'mountain', 'mountains', 'ocean', 'tree', 'island
32            'rivers', 'cannabis', 'rock', 'waters', 'forest', 'dunes',
33            'environment', 'everglades', 'trees', 'garden', 'plants',
34            'field', 'sky', 'plains', 'rain', 'ecosystems', 'lakes', '
35            'glacier', 'socotra', 'flames', 'storm', 'fog', 'tide', 'o
36            'natureern', 'naturen', 'natures']
37
38     civilization = ['city', 'park', 'country', 'war', 'village', 'commun
39                    'border', 'crisis', 'school', 'culture', 'boat', 'ho
40                    'market', 'toys', 'education', 'cities', 'york', 'po
41                    'research', 'prayer', 'organizations', 'conflict', '
42                    'hotel', 'toy', 'poverty', 'cultures', 'music', 'art
43                    'schools', 'motel', 'bridge', 'law', 'street', 'room
44
45     time = ['years', 'time', 'day', 'night', 'days', 'today', 'year', 'c
46            'summer', 'months', 'future', 'hours', 'times', 'august', 'w
47            'monument', 'weeks', 'month', 'decade', 'moments', 'season',
48            'weekend', 'spring']
49
50     size = ['big', 'small', 'tiny', 'large', 'long', 'high', 'massive',
51            'narrow', 'thin', 'enormous', 'tall', 'heavy', 'distant', 'i
52            'hundreds', 'temperatures', 'half', 'size', 'distance', 'len

```

```

53
54 different = ['new', 'great', 'old', 'different', 'important', 'tradi
55             'perfect', 'popular', 'dangerous', 'famous', 'dramatic'
56             'private', 'successful', 'endangered', 'valuable', 'ori
57             'vulnerable', 'vital', 'spiritual', 'distinct', 'devast
58
59 nice = ['beautiful', 'incredible', 'amazing', 'magical', 'sweet', 'w
60         'stunning', 'excellent']

```

```

In [199]: 1 def less_caption_words(dataframe,column):
2           # creating list of captions with lower cases
3           captions = []
4           for i in range(dataframe.shape[0]):
5               captions.append(pd.DataFrame(dataframe.iloc[:,column]).loc[
6
7           # introducing replacements
8           word_lists = [animal, human, photographer, place, advertisement,
9                         different, nice]
10          replacement_words = ['animal', 'human', 'photographer', 'place',
11                               'time', 'size', 'different', 'nice']
12
13          # replacing words and creating new captions
14          simple_captions = []
15          for c in captions:
16              for i in range(len(replacement_words)):
17                  for word in word_lists[i]:
18                      c = c.replace(word, replacement_words[i])
19                  simple_captions.append(c)
20          token_captions = [word_tokenize(c) for c in simple_captions]
21          # counting words using the simplified captions with the replaced w
22          wordcount = Counter(pos_tag(word_tokenize(''.join(simple_caption
23          word_list = sorted(list(wordcount.items()), key = lambda w: -w[1
24          stoplist = nltk.corpus.stopwords.words('english')
25          word_list = [word_list[i] for i in range(len(word_list))
26                      if len(word_list[i][0][0]) > 2 and word_list[i][0][
27                      and word_list[i][0][1] in ['NN', 'NNP', 'NNS', 'JJ'
28          return [simple_captions, token_captions, word_list]

```

```

In [200]: 1 [simple_captions, token_captions, new_word_list] = less_caption_word

```

```

In [201]: 1 # creating an empty dataframe
2          wrds = [new_word_list[i][0][0] for i in range(len(new_word_list))]
3          engmnt = list(all_data['engagement'])
4          cllmns = wrds+['Engagement']
5          df_caption = pd.DataFrame(0.0, index = np.arange(len(token_captions)
6          df_caption['Engagement'] = ([engmnt[i] for i in range(len(token_capt

```

```
In [202]: 1 df_caption.shape
```

```
Out[202]: (634, 8894)
```

As seen through the extensive word replacements, I achieved the reduction of the captions' count by 500 words, reducing the complexity of the model and the features sacrificing only part of its accuracy of information.

```
In [203]: 1 # convert the post list to textblob format for passing into tf-idf f
2 captions_tb = []
3
4 for caption in list(simple_captions):
5     captions_tb.append(tb(caption))
```

```
In [204]: 1 # import time
2 # t0 = time.time()
3
4 caption_tfidf_list = []
5
6 # list containing all the captions with tfidf score
7
8 for i, blob in enumerate(captions_tb):
9     scores = {word: tfidf(word, blob, captions_tb) for word in blob.
10     caption_tfidf_list.append(scores)
11
12 # t1 = time.time()
13 # print(t1-t0)
```

```
In [205]: 1 #Assigning value of feature according to the caption set of each pos
2
3 caption_data_ind = 0
4 for index, row in df_caption.iterrows():
5     for column in df_caption:
6         if column in caption_tfidf_list[caption_data_ind]:
7             df_caption.loc[index, column] = caption_tfidf_list[capti
8             caption_data_ind=caption_data_ind+1
```

```
In [206]: 1 #Divide data set into Features and Outcome
2 XX = df_caption[wrds]
3 yy = df_caption['Engagement']
```

```
In [207]: 1 import matplotlib
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score, confusion_matrix
7
8 #Divide into test and train data
9 XX_train, XX_test, yy_train, yy_test = train_test_split(XX, yy, test
10
11 #Logistic Regression
12 caption_lr = LogisticRegression(C=3, random_state=0)
13 caption_lr.fit(XX_train, yy_train)
```

```
Out[207]: LogisticRegression(C=3, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs
=1,
                             penalty='l2', random_state=0, solver='liblinear', tol=0.0001
,
                             verbose=0, warm_start=False)
```

```
In [208]: 1 #Predicting test outcome
2 yy_pred = caption_lr.predict(XX_test)
3
4 #Calculate Accuracy
5 print('Accuracy: %.2f' % accuracy_score(yy_test, yy_pred))
```

Accuracy: 0.73

Based on the above, it appeared that the overall accuracy did not improve over the original model. Interestingly enough, it appears that the information collected from the image recognition labeling was of the same value as the thorough description of the image by the photographer.

05. Logistic Regression using (captions, labels) and engagement

In this step, I combined all the information together to a third model making engagement predictions using both image labels and captions.

```
In [209]: 1 # tfidf values of labels and tfidf values of captions
          2 all_tfidf_table = pd.concat([df.loc[:, : 'zoo'], df_caption], axis=1)
          3 all_tfidf_table.shape
```

Out[209]: (634, 9807)

```
In [210]: 1 #Divide data set into Features and Outcome
          2 X_all = all_tfidf_table.loc[:, : 'benefits']
          3 y_all = all_tfidf_table['Engagement']
```

```
In [211]: 1 #Divide into test and train data
          2 X_all_train, X_all_test, y_all_train, y_all_test = train_test_split(
          3
          4 #Logistic Regression
          5 all_lr = LogisticRegression(C=2, random_state=42)
          6 all_lr.fit(X_all_train, y_all_train)
```

Out[211]: LogisticRegression(C=2, class_weight=None, dual=False, fit_intercept=True,
 intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs
=1,
 penalty='l2', random_state=42, solver='liblinear', tol=0.000
1,
 verbose=0, warm_start=False)

```
In [212]: 1 #Predicting test outcome
          2 y_all_pred = all_lr.predict(X_all_test)
          3
          4 #Calculate Accuracy
          5 print('Accuracy: %.2f' % accuracy_score(y_all_test, y_all_pred))
```

Accuracy: 0.75

While it appeared to be a slight increase in the overall accuracy on our results, the added complexity of the model is not justified. It seems that the information of the image labels and description are somewhat correlated essentially describing the same thing and carrying the same amount of information. The image caption is describing what the google vision API is summarising in its labels.

06. Clustering images - LDA on image labels

Aiming to identify what makes an image popular and attracts engagement, I used LDA (Latent Dirichlet allocation) on image labels. Running the clustering algorithm multiple times, in order to find the right number of topics that best separated / described my images, I ensured that the topics were clearly distinctive from each other while also their themes were intuitive.

```
In [232]: 1 import gensim
2 from gensim.utils import simple_preprocess
3 from gensim.parsing.preprocessing import STOPWORDS
4 from nltk.stem import WordNetLemmatizer, SnowballStemmer
5 from nltk.stem.porter import *
6 import numpy as np
7
8 import nltk
9 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\lovek\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Out[232]: True

```
In [233]: 1 post_labels = [ast.literal_eval(post) for post in all_data.iloc[:, 1]
2
3 dict_labels = {}
4 for i in range(len(post_labels)):
5     for j in range(len(post_labels[i])):
6         if post_labels[i][j] in dict_labels:
7             a = dict_labels[post_labels[i][j]] + 1
8             dict_labels[post_labels[i][j]] = a
9         else:
10            dict_labels[post_labels[i][j]] = 1
```

```
In [3]: 1 import lda
2 from sklearn.feature_extraction.text import CountVectorizer
```

```
In [235]: 1 stopwords_nltk=set(stopwords.words('english'))
```

```
In [236]: 1 words_freq_vec = CountVectorizer(input=post_labels, stop_words=stopwo
```

```
In [237]: 1 labels_under_pro = words_freq_vec.fit_transform([' '.join(p) for p i
```

```
In [238]: 1 ntopics = 4
2 corpus_arr = np.array(corpus)
3 model = lda.LDA(n_topics=int(ntopics), n_iter=500, random_state=1)
4 model.fit(labels_under_pro)
```

```
INFO:lda:n_documents: 634
INFO:lda:vocab_size: 935
INFO:lda:n_words: 6691
INFO:lda:n_topics: 4
INFO:lda:n_iter: 500
C:\Users\lovek\Anaconda3\lib\site-packages\lda\utils.py:55: FutureWarning: Conversion of the second argument of issubdtype from `int` to `np.signedinteger` is deprecated. In future, it will be treated as `np.int32 == np.dtype(int).type`.
    if sparse and not np.issubdtype(doc_word.dtype, int):
INFO:lda:<0> log likelihood: -59701
INFO:lda:<10> log likelihood: -42999
INFO:lda:<20> log likelihood: -41478
INFO:lda:<30> log likelihood: -40933
INFO:lda:<40> log likelihood: -40750
INFO:lda:<50> log likelihood: -40545
INFO:lda:<60> log likelihood: -40366
INFO:lda:<70> log likelihood: -40233
INFO:lda:<80> log likelihood: -40263
INFO:lda:<90> log likelihood: -40169
INFO:lda:<100> log likelihood: -40192
INFO:lda:<110> log likelihood: -40129
INFO:lda:<120> log likelihood: -40141
INFO:lda:<130> log likelihood: -40126
INFO:lda:<140> log likelihood: -40102
INFO:lda:<150> log likelihood: -40085
INFO:lda:<160> log likelihood: -40067
INFO:lda:<170> log likelihood: -39967
INFO:lda:<180> log likelihood: -39998
INFO:lda:<190> log likelihood: -39919
INFO:lda:<200> log likelihood: -39802
INFO:lda:<210> log likelihood: -39783
INFO:lda:<220> log likelihood: -39737
INFO:lda:<230> log likelihood: -39771
INFO:lda:<240> log likelihood: -39702
INFO:lda:<250> log likelihood: -39770
INFO:lda:<260> log likelihood: -39707
INFO:lda:<270> log likelihood: -39636
INFO:lda:<280> log likelihood: -39684
INFO:lda:<290> log likelihood: -39566
INFO:lda:<300> log likelihood: -39618
INFO:lda:<310> log likelihood: -39610
INFO:lda:<320> log likelihood: -39597
INFO:lda:<330> log likelihood: -39545
INFO:lda:<340> log likelihood: -39603
INFO:lda:<350> log likelihood: -39561
INFO:lda:<360> log likelihood: -39566
INFO:lda:<370> log likelihood: -39551
INFO:lda:<380> log likelihood: -39590
INFO:lda:<390> log likelihood: -39487
```



```
INFO:lda:<400> log likelihood: -39519
INFO:lda:<410> log likelihood: -39559
INFO:lda:<420> log likelihood: -39534
INFO:lda:<430> log likelihood: -39526
INFO:lda:<440> log likelihood: -39584
INFO:lda:<450> log likelihood: -39575
INFO:lda:<460> log likelihood: -39623
INFO:lda:<470> log likelihood: -39598
INFO:lda:<480> log likelihood: -39552
INFO:lda:<490> log likelihood: -39612
INFO:lda:<499> log likelihood: -39534
```

Out[238]: <lda.lda.LDA at 0x27187329358>

```
In [239]: 1 type(labels_under_pro)
          2 # dataframe_label = pd.DataFrame(total_features_words)
          3 dataframe_label = pd.DataFrame({'post_no': range(len(post_labels)),
```

```
In [240]: 1 topic_label = model.topic_word_
          2 topic_post=model.doc_topic_
          3 topic_post=pd.DataFrame(topic_post)
          4 dataframe_label=dataframe_label.join(topic_post)
          5 Insta_Posts=all_data.iloc[:,[3, 5, 6]]
```

In [241]: 1 all_data[:1]

Out[241]:

	Unnamed: 0	ncomments	nlikes	labels	caption	engagement	engag
0	0	0.020126	0.256822	['marine mammal', 'fauna', 'water', 'mammal', 'marine biology', 'underwater', 'manatee', 'organism', 'wildlife', 'dugong']	Photo By @BrianSkerry\r\nA West Indian Manatee calf nurses from its mom as they settle down on the sandy sea floor in the waters off the coast of Belize, an important country for these endangered animals. Manatees here live in mangroves and often sleep there overnight and feed on nearby sea grass beds during the day. Unlike Florida manatees, these animals are not nearly as acclimated to humans and are shy. On this morning however, this very relaxed mom and calf allowed me into their world.\r\nFor more images and stories about ocean wildlife follow @BrianSkerry\r\n#manatees #belize #mesoamericanreef #endangeredspecies		1

In [242]: 1 for i in range(int(ntopics)):
2 topic="topic_"+str(i)
3 Insta_Posts[topic]=dataframe_label.groupby(['post_no'])[i].mean()

C:\Users\lovek\Anaconda3\lib\site-packages\ipykernel_launcher.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: <http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>
(<http://pandas.pydata.org/pandas-docs/stable/indexing.html#indexing-view-versus-copy>)

This is separate from the ipykernel package so we can avoid doing imports until

```
In [243]: 1 Insta_Posts=Insta_Posts.reset_index()
2 topics=pd.DataFrame(topic_label)
3 topics.columns=words_freq_vec.get_feature_names()
4 topics1=topics.transpose()
5 print ("Topics label distribution written in file topic_label_dist.xls")
6 topics1.to_csv("topic_label_dist.csv")
7 Insta_Posts.to_csv("Insta_posts_topic_dist.csv",index=False)
8 print ("Restaurant topic distribution written in file Insta_posts_to
```

Topics label distribution written in file topic_label_dist.xlsx
Restaurant topic distribution written in file Insta_posts_topic_dist.xlsx

Analysing influential topics

Based on my analysis and the labels related to each category, I identified that sorting my labels into 4 topics was the most effective way of splitting them, as the labels were split to the below self explanatory topics:

1. Topic I: plants
2. Topic II: human factor
3. Topic III: Animal
4. Topic IV: Scenery

NOTE: The actual words per each topic are listed in the "topic_word_dist.xlsx" file.

In [244]:

```

1  import csv
2  posts = []
3  with open('Insta_posts_topic_dist.csv','r', encoding="utf8") as csv_
4      csv_reader = csv.reader(csv_file, quotechar='"', delimiter=',',
5      next(csv_reader)
6      for row in csv_reader:
7          posts.append((float(row[3]),float(row[4]),float(row[5]),floa
8
9  posts.sort(key=lambda x: x[4], reverse=True)
10 quartile_amount = int(len(posts)/4)
11
12 top_quartile_topic1_avg = sum(map(lambda x: x[0], posts[:quartile_am
13 top_quartile_topic2_avg = sum(map(lambda x: x[1], posts[:quartile_am
14 top_quartile_topic3_avg = sum(map(lambda x: x[2], posts[:quartile_am
15 top_quartile_topic4_avg = sum(map(lambda x: x[3], posts[:quartile_am
16
17 bottom_quartile_topic1_avg = sum(map(lambda x: x[0], posts[-quartile
18 bottom_quartile_topic2_avg = sum(map(lambda x: x[1], posts[-quartile
19 bottom_quartile_topic3_avg = sum(map(lambda x: x[2], posts[-quartile
20 bottom_quartile_topic4_avg = sum(map(lambda x: x[3], posts[-quartile
21
22 print('Top quartile topic avgs: ',top_quartile_topic1_avg,top_quarti
23 print('\n')
24 print('Bottom quartile topic avgs: ',bottom_quartile_topic1_avg,bott

```

Top quartile topic avgs: 1.0242003584503878 0.06629849408088037 0.0511766861781931 0.023865919085372403

Bottom quartile topic avgs: 1.1922729380081662 0.2750263838998251 0.19092801061018513 0.5266541625872856

Based on my analysis of the collected data, it appeared that the first topic (related to nature and plants) was prominent in all the images irrespective of their engagement, which was expected given that the images were collected from national geographic's page and thus it did not appear to be the differentiating factor.

On the other hand, the rest of the topics (human factor, animal and scenery) played a crucial role in separating images from mediocre to great and it appeared that more unsuccessful images had twice the amount of "human factor", more than 5 times "animal" related themes, and twice the "scenery" aspect.

07. Advice for National Geographic

Based on my analysis, I would advice the people responsible to keep focusing on natural themes, but rather adding more context to their images enhancing the human factor and somewhat avoiding the plain classic beautiful scenery themes.

In []:

1