

# **Instagram-based Recommendation System for Travel Destinations in Italy**

## **Creators**

Taize Kang  
Fanni Kovapohja  
Markku Ylitalo



**University of Helsinki**

Introduction to Data Science (DATA11001)  
18.10.2020

We have created a recommendation system for travel destinations in Italy, which is based on Instagram posts with the hashtag *italytravel*. We scraped 34 578 Instagram posts from the web, collected the locations and hashtags of them and analyzed this information to define the social media popularity and features of 839 different locations. With the analyzed data, we created a recommendation system, in which user can select different features and get location recommendations based on them. The recommendation system is implemented by exploiting k-means clustering and it includes 658 different locations. The final deliverable is a web page consisting of a feature selection window and a recommendation window with the popularity filter and the filtered recommendations on the map of Italy.

We gathered the data from Instagram with the web scraping program written in Python. When creating the web scraper we exploited the ideas of the blog post *Scraping Instagram Locations with Python* written by Ksenia Tikhomirova (<https://medium.com/@tik.dev/scrap-instagram-locations-with-python-d48ba6e56ebc>, read on 18.9.2020). The Instagram web page with the hashtag *italytravel* shows approximately 70 Instagram posts at a time in order of publication time where the newest posts are shown first. With the following URL address structure, it is possible to get at the previous loading page consisting of the next newest 70 posts:

When the parameter `<end_cursor>` is empty, the URL address (1) leads to the front page of Instagram's subgallery under the hashtag *italytravel*. From that page, it is possible to collect the shortcode of each post for the future processing and also fetch the `<end_cursor>` value for the next loading page of posts.

When this is done repeatedly, the shortcodes of posts can be collected from multiple loading pages.

The URL address (2) leads to the page where the location information and caption of a post can be gathered. The caption is the part of a post that includes text, user account links and hashtags.

On 18.9.2020, we managed to scrape 500 pages of Instagram posts, altogether 34 578 posts. The connection problems of Instagram complicated the scraping so we accepted the largest number of posts we succeeded to get as our data set instead of the planned 150 000 posts. We collected the location name and caption of each post and saved this information to JSON file.

[illegible]

We used regular expression methods to remove all the other parts of the caption strings but the hashtags. We only saved the strings starting with '#' and containing at least 3 alphabetical or '\_' characters. Then we removed all the '#' marks and created a list of distinct hashtag words for each post with the 'split()' method. We decided to use only hashtags and not the other parts of the captions, because the hashtags are more homogeneous than the other parts, which include more different languages and wordings, and typographical errors.

After that, with the help of programmed automatization and the self-collected 'stop words' of hashtags, for example, 'blue', 'usa', 'lol', 'sky', 'follow', 'follow4follow', 'woman', 'friday', 'morning', 'couple', 'brunette', 'adventure', 'september', 'settembre' and 'selfie', we managed to remove the most of the uninformative hashtags in the data frame (df2). We also removed all the hashtags that were as a part of the string of any location mentioned in the data frame (df2), because the location information in hashtags does not have any additional value when we already know the locations. For the comparison between the previous and current state, there is the cleaned version of the hashtags of *Monte Sant'Angelo* shown under:

```
['vieste', 'gargano', 'centrostorico', 'the stag society', 'hoscoc', 'hoscocmen', 'hoscocgreece', 'hman', 'landscape', 'bread', 'borghi', 'city', 'pane', 'villages', 'architettura', 'paesaggio', 'tour', 'foggia', 'gargano', 'splash_reflex', 'landscape', 'sea', 'nature', 'villages', 'village', 'architettura', 'gargano', 'forestaumbra', 'foggia', 'vieste', 'manfredonia', 'mareadriatico', 'adriaticsea', 'sea', 'collina', 'garganocast', 'splash_reflex', 'landscape', 'sea', 'nature', 'villages', 'village', 'architettura', 'castello', 'castle', 'templari', 'storia', 'paradisopugliese', 'terrasanta', 'unicatmondo', 'solocosebelle', 'peaceful', 'splash_reflex', 'landscape', 'sea', 'nature', 'villages', 'village', 'architettura']
```

## Duplicates of the Locations

We noticed that in the data frame (df2), there were locations that meant the same but had been written differently or in different languages, for example, in English, Italian, German, Spanish or Chinese. We tried to exploit programmed automatization when searching the duplicates as much as we could, but still, there was a lot of manual work to do in this process. When we combined these duplicates, we also combined all the hashtags and summed all the popularity values of them. After this process, there were 892 locations left.

	location	popularity
0	venice	1230
1	rome	867
2	positano amalfi coast	543
3	capri	350
4	amalfi coast	322
..	...	...
887	hotel puccini	4
888	giovinazzo bari puglia	4
889	villa tellus holiday villa	4
890	rosignano solvay toscana	4
891	sulmona	4
[892 rows x 2 columns]		

Later we realized that there still were many duplicates left. There were also unwanted locations outside the borders of Italy. We examined the locations again more thoroughly than the first time. We dealt with the duplicates in the same manner as above and removed the unwanted locations from the data frame (df2). Also, the rows without any hashtags were deleted at this point. We created a new cleaned data frame (df3), where there were 839 locations left.

## 4. Features of the Locations

We gathered all the hashtags (\*) from the data frame (df3) and counted the 200 most common hashtags of them. While analyzing the list of the most common hashtags, we decided to divide the hashtags into 7 different classes and 42 subclasses. For example, the class *NATURE* consists of subclasses *beach*, *lake*, *river*, *sea*, *waterfall*, *hills*, *mountain* and *volcanic*, and the class *ACTIVITIES* consists of subclasses *hiking*, *sailing*, *swimming*, *fashion*, *luxury*, *food*, *restaurant* and *wine*.

For each subclass, we formed a sublist of hashtags from the list of all the hashtags (\*) with the help of **fastText**. Each sublist included only hashtags that were related to the given class. For example, fastText suggested that the following hashtags describe the subclass *hiking* the best:

```
[('hikingtrails', 0.9842933416366577), ('hike', 0.9762042760848999), ('dolomitemountains', 0.970732881130981), ('wild', 0.97059889755249), ('dolomitemountain', 0.96478509029541), ('mountainscapes', 0.9641682505607605), ('dolomite', 0.9636821746826172), ('dolomitesunesco', 0.9633682370185852), ('mountainstones', 0.9630481004714966), ('dolomitibellunesi', 0.9619472026824951), ('dolomitici', 0.961671880613403), ('mountainscenery', 0.9614839467407227), ('dolomiten', 0.960288882255542), ('dolomitesyou', 0.960240595626831), ('dolomite', 0.9588807821273804), ('dolomitaleo', 0.958728031668091), ('dolomitunesco', 0.956875145433333), ('dolomitas', 0.956161320209502), ('mountainsarecalling', 0.9558297395706177), ('wilderness', 0.9553650617599487), ('dolomiti', 0.9549998641014099), ('cinqueteriorcinadinadamezzo', 0.9541064508088716), ('mountains', 0.9533324837684631), ('cortinadamezzo', 0.953248143196106), ('dolomitidasogno', 0.95317143201828), ('secedamountain', 0.9528548121452332), ('mountainscape', 0.9520748853683472), ('trekking', 0.951884388923645), ('natgeowild', 0.9514551758766174), ('dolomitisuperski', 0.9500757455825806), ('mountainbike', 0.9496070146560669), ('folkscenery', 0.9475473165512085), ('addictedtothemountains', 0.9478006121826172), ('secedam', 0.9468228816986084), ('trecedimilavaredo', 0.9458196759223938), ('dolomity', 0.945677638053894), ('letrecimilavaredo', 0.945537805557251), ('dolomitiemotions', 0.9453480243682861), ('wildernessnation', 0.9446677565574646), ('dolomitipatrimoniounesco', 0.9440531730651855), ('mountainview', 0.9422052502632141), ('visitdolomites', 0.9415383338928223), ('mountainside', 0.9413654208183289), ('mountaineering', 0.9398941993713379), ('mountaintop', 0.9398379921913147), ('cinquetorri', 0.9394592046737671), ('mountainviews', 0.9391295909881592), ('dreizinnen', 0.937866687746582), ('hiker', 0.937632679939277), ('mountainelopement', 0.9375296831130981)]
```

We added 42 empty columns to the data frame (df3) for each subclass. We created a program that went all the locations and features through one by one, so that the frequency of each feature in each location's

hashtag list was counted. For example, for the location *Venice* and subclass *sailing*, the program increased the value at the row *Venice*, in the column *sailing*, every time when some hashtag from the fastText processed hashtag list of *sailing* was found from the hashtag list of *Venice*. At this point, the hashtag lists of locations could have contained the same hashtag multiple times, so that the frequencies of features would be counted correctly when considering that the hashtag list of each location had been formed from several Instagram posts.

The frequencies were altered to relative frequencies along the row/location, so that the locations with different popularities would be comparable with each other. The most popular locations had the largest absolute frequencies of all of the features, so the less popular locations would never be recommended without relative frequencies. After that, the column of hashtag lists of locations was removed, and we had a final data frame (df4) for building the recommendation system:

	loc_id	location	popularity	hiddengem	unesco	renaissance	gothic	ancient	medieval	...	fashion	aesthetic	luxury	food	art	history	culture	archeology	architecture
0	0	venice	1230	0.000	1.515	1.732	0.649	0.000	0.000	...	2.386	0.386	2.492	1.030	7.595	3.331	6.162	4.420	7.934
1	1	rome	867	0.000	0.235	0.469	0.000	3.521	0.000	...	1.030	0.263	1.459	1.214	8.077	5.110	8.243	5.101	9.032
2	2	positano	543	0.000	1.064	0.000	0.000	1.064	0.532	...	1.942	0.161	2.604	1.058	2.641	1.528	2.048	1.491	2.678
3	3	capri	350	0.000	0.329	0.000	0.000	0.000	0.329	...	1.838	0.000	6.956	0.936	1.821	0.558	1.142	0.739	1.340
4	4	amalfi coast	322	0.704	2.113	0.000	0.000	0.000	0.000	...	2.252	0.081	2.948	1.265	1.497	0.847	1.625	0.754	1.776
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
834	834	marina di vasto	4	0.000	0.000	0.000	0.000	0.000	0.000	...	2.273	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
835	835	loro ciuffenna	4	0.000	0.000	0.000	0.000	0.000	0.000	...	0.962	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
836	836	mantua	4	0.000	0.000	0.000	0.000	0.000	50.000	...	0.000	0.000	0.000	0.000	11.538	7.692	9.615	7.692	11.538
837	837	abbazia di santantimo	4	0.000	0.000	0.000	0.000	0.000	0.000	...	0.000	0.000	0.000	0.000	0.000	44.444	16.667	0.000	5.556
838	838	sulmona	4	0.000	0.000	0.000	0.000	0.000	0.000	...	0.000	0.000	0.000	39.286	0.000	0.000	7.143	0.000	7.143

## 5. K-means Clustering for the Locations

The location data was not classified, so we tried to find a solution to divide the location data into classes based on different features. We planned to build a recommendation system, in which user can choose desired features and get location recommendations, with the help of these classes. We could have got better recommendation accuracy if we had exploited the exact relative frequencies of the features. However, within the time limit, we did not discover any clustering method, in which we could have used the relative frequencies. That is why, after experiments, we changed the feature data into binary form by setting values less than 2% to be 0 and the values greater than or equal to 2% to be 1. We interpreted that the relative frequencies less than 2% were not significant.

We decided to use k-means clustering method of Scikit-learn to get k classes of locations based on their feature similarity. After experimenting different k-values, we decided to set k = 15, because the sizes of the clusters were the most appropriate then. If we have more locations in the future, the value of k must be determined again. The first step was to set 15 centroids randomly and divide the data into 15 classes by them. Then the data was re-classified by calculating the distance until there was no change to the centroids. After training the data, we found that there were two very large clusters, which would make the recommendations inaccurate. We used k-means again for them with k = 4 to get smaller clusters. Unfortunately, this worked only with the one large cluster. The other large cluster was divided into one very large subcluster and three small subclusters, sizes of 2, 1 and 1. That is why only another one of the large clusters was divided into smaller clusters, which may affect negatively the recommendation results.

## 6. Recommendation System

The purpose of a recommendation system for travel destinations is to find locations similar to the preferences of a user from the data. We implemented this as follows: User decides whether or not he wants to select each feature. If he selects the feature, the value of this feature is 1, otherwise the value of the feature is 0. After getting the input list from user, we use k-means to examine in which class the list should belong. If the class is the large class, which is divided into smaller classes, we use k-means model again to get smaller and more accurate classification. When the right class is found, the whole class is recommended for the user.



## 7. Final Deliverable

We decided to use a web page as a front end and our recommendation system as a back end. We exploited Flask to do that. The web page consists of two windows. In the first window, user can select desired features by ticking the checkboxes. Then we transfer the checkboxes to a list of zeros and ones, which we use as an input of the k-means model. When we got the result class, the locations of the class are shown on the map of Italy, which is the second window. The map of the second window is implemented with Plotly. We exploited GeoEarth's API to get the coordinates of the locations. However, some of the coordinates were not found, so we had to remove these locations from the recommendation system. In the final deliverable, there are 658 different locations left. The dots of different sizes on the map refer to the popularity of locations — the bigger dot means more popular location. There is also a filter beside the map for user to select four different levels of popularity. These levels are shown on the map with different colors of the dots.

## 8. Evaluation and Conclusions

In the project proposal, we discussed that we were going to scrape Instagram posts with the hashtag *italytravel* from the web, gather the locations of the posts, use the number of times a location appears in the posts as the popularity of it and sort the locations using hashtags. We also planned to build a recommendation system for travel destinations, in which user can filter the popularity and type of the destinations. All that is what we ended up to implement. Also the visualizations were carried out as planned, apart from showing the list of the recommended locations. We thought that it is clearer for the end-user to see the locations with their names only on the map. What we did not implement from the project proposal, was exploiting weather or average cost information of the locations. We decided to ignore these features, because the differences in them between different locations are not significant inside Italy. Because we were only able to collect the 34 578 latest Instagram posts, the most of them were published in September. That is why we could not exploit the season or dates of the posts as we planned. In the future work, it would be better to analyze more posts and distribute the gathering of posts around the year to make the data sample more diverse. In addition, we did not use the decision forest model we wrote about because we had only unlabeled data. The k-means clustering turned out to work better for our project.

Since the tool was designed to offer user optimized travel recommendations finding the both viral and hidden gems of Italy, it fulfills its purpose laudably. There are already many travel guides available, but because the most popular destinations gather all the hype, the lesser-known ones are often almost impossible to find. The added value of the tool is that users will find new and interesting destinations that they may not have found without it. That makes our idea to exploit Instagram posts unique and useful. The visual outcome of the system is functional and helps user to locate the destinations easily. There are also many expanding options, for example, a road trip organizer, which connects desired locations. In addition, there are a couple of things to improve: k-means do not perform accurately with all combinations of the features, there are still a few locations outside Italy included in, and program needs to be ran again for another search. Also the gathered data may not be so reliable due to the impacts of the COVID-19 situation on traveling in Italy.

Taize Kang majored his BSc degree in intelligent science and technology, and Fanni Kovapohja and Markku Ylitalo majored their BSc degrees in mathematics. Because the project hardly included any mathematics, and Kovapohja and Ylitalo had computer science only as a minor subject in their BSc degrees, nearly everything of the project, except from the basic knowledge of the Python programming, was new to them. Kang had more experience in programming but was not very familiar with the front end process and there were new libraries for him. None of us was very familiar with the machine learning techniques. Creating the idea of the project and making the whole recommendation system from scratch required huge efforts from the whole group. That is why the learning outcomes of this project and course were really massive. Each group member had equal workload and participation in every stage of the process. We had a lot of work to do in all stages of the data life cycle, so our project is suitable for this course and includes sufficient amount of work.