

Project Report: Phase 1

Group 18, Fall 2018

Group Members

1. Aditya Kanchivakam Ananth
2. Aniket Dhole
3. Shatrujit Singh
4. Siddharth Pandey
5. Siddhesh Narvekar
6. Vasisht Sankaranarayanan

Abstract

Keywords

- TF - Term Frequency
- DF - Document Frequency
- IDF - Inverse Document Frequency
- TF-IDF - Term Frequency-Inverse Document Frequency
- CM - Global Color Moments on HSV Color Space
- CN - Global Color Naming Histogram
- HOG - Global Histogram of Oriented Gradients
- LBP - Global Locally Binary Patterns on Gray Scale
- CSD - Global Color Structure Descriptor
- GLRLM - Global Statistics on Gray Level Run Length Matrix
- Code 3x3 - Spatial Pyramid Representation (CN3x3, CM3x3, LBP3x3, GLRLM3x3)

ABSTRACT

Similarity measures play an important role in text related research and applications well as image processing applications. The purpose of this project is to compute similarity between users, images and locations. We have computed document similarity similarity, image similarity and location similarity with the help of visual and textual descriptors by utilising similarity measures such as Euclidean Distance and Cosine Similarity.

Introduction

- Terminology

- Term Frequency (TF): Term Frequency is computed using the weight of a given keyword k in a given document d . It is mathematically defined as

$$tf(k, d) = \frac{count(k, d)}{size(d)}$$

- Document Frequency (DF): Document Frequency is computed using the total number of documents in the database D containing the given keyword k . It is mathematically defined as

$$df(k, D) = \frac{number\ of\ documents\ containing(k, D)}{number\ of\ documents(D)}$$

- Inverse Document Frequency (IDF): Inverse Document Frequency is an inverse measure of DF. It is mathematically defined as

$$idf(k, D) = \log\left(\frac{number\ of\ documents(D)}{number\ of\ documents\ containing(k, D)}\right)$$

- Term Frequency-Inverse Document Frequency (TF-IDF): TF-IDF is the weight of the keyword k for document d in database D combines the concepts of TF and IDF:

$$tf_{idf(k, d, D)} = tf(k, d) * idf(k, D)$$

- Goal Description

- Task 1: Implement a program which given a user ID, a model (TF, DF, TF-IDF), and value “k”, returns the most similar k users based on textual descriptors. For each match, also list the overall matching score and three terms that have the highest similarity contribution.
- Task 2: Implement a program which given an image ID, a model (TF, DF, TF-IDF), and value “k”, returns the most similar k images based on textual descriptors. For each match, also list the overall matching score and the 3 terms that have the highest similarity contribution.
- Task 3: Implement a program which, given a location ID, a model (TF, DF, TF-IDF), and value “k”, returns the most similar k locations based on textual descriptors. For each match, also list the overall matching score and the 3 terms that have the highest similarity contribution. Note: In this phase, the

location IDs will always be specified as the “number” field (i.e., 1 to 30) in the devset topics.xml file.

- Task 4: Implement a program which, given a location ID, a model (CM, CM3x3, CN, CN3x3, CSD, GLRLM, GLRLM3x3, HOG, LBP, LBP3x3), and value “k”, returns the most similar k locations based on the corresponding visual descriptors of the images as specified in the “img” folder. For each match, also list the overall matching score as well as the 3 image pairs that have the highest similarity contribution.
- Task 5: Implement a program which, given a location ID and value “k”, returns the most similar k locations based on the corresponding visual descriptors of the images as specified in the “img” folder. For each match, also list the overall matching score and the individual contributions of the 10 visual models.
- Assumptions
 - All terms are keywords.
 - The term frequencies have been defined differently in the dataset.
 - Term frequency: Count of the term in a document.
 - Document frequency: Count of the documents in the database containing the term.
 - TF-IDF: This has been set to the value given by a division of term frequency by document frequency.

Description of the Proposed Solution/Implementation

Task 1:

Each user in the input file has a user id and information related to each term for all terms of a user: term value, term frequency (tf), document frequency (df) and tf-idf. Based on the input model requested at run time, one of the following program branches is executed:

- Model= TF: Vectors are prepared on basis of the term frequency of the given user and other user’s common terms. Euclidean distance is then calculated using these vectors and the users which have the least score are sent to the output. The highest contributing terms are those for which difference of term frequencies between the corresponding ones in the input user’s terms is the least, since they contribute to a lower Euclidean distance.

User 1 (Given user)

Term	A	B	C	D
Term Frequency	4	8	3	5

User 2

Term	A	C	D	E
Term Frequency	5	6	9	3

User 3

Term	A	B	D	F
Term Frequency	2	1	9	4

$$\begin{aligned}
 \text{Euclidean Distance}(\text{User1}, \text{User2}) &= \sqrt{(tf_A(\text{User1}) - tf_A(\text{User2}))^2 + (tf_B(\text{User1}) - tf_B(\text{User2}))^2 + (tf_C(\text{User1}) - tf_C(\text{User2}))^2 + (tf_D(\text{User1}) - tf_D(\text{User2}))^2} \\
 &= \sqrt{(4 - 5)^2 + (8 - 0)^2 + (3 - 6)^2 + (5 - 9)^2} \\
 &= \sqrt{(-1)^2 + (8)^2 + (-3)^2 + (-4)^2} \\
 &= \sqrt{1 + 64 + 9 + 16} = \sqrt{90} = 9.487
 \end{aligned}$$

$$\begin{aligned}
 \text{Euclidean Distance}(\text{User1}, \text{User3}) &= \sqrt{(tf_A(\text{User1}) - tf_A(\text{User3}))^2 + (tf_B(\text{User1}) - tf_B(\text{User3}))^2 + (tf_C(\text{User1}) - tf_C(\text{User3}))^2 + (tf_D(\text{User1}) - tf_D(\text{User3}))^2} \\
 &= \sqrt{(4 - 2)^2 + (8 - 1)^2 + (3 - 0)^2 + (5 - 9)^2} = \sqrt{(2)^2 + (7)^2 + (3)^2 + (-4)^2} \\
 &= \sqrt{4 + 49 + 9 + 16} = \sqrt{78} = 8.832
 \end{aligned}$$

Thus, User 3 is more similar to given user than User 2 due to lesser Euclidean distance. The best contributing term for user 3 is term “A” as it contributes the least to the distance.

- Model= DF: The similarity for DF is calculated by taking the intersection of the input terms and the other user’s terms. The users with highest number of common terms are classified as more similar. The highest contributed terms for each user are the ones with lowest DF value because they are rare, and thus important.

User 1 (Given user)

Term	A	B	C	D
DF	4	8	6	5

User 2

Term	A	B	D	G
DF	4	8	5	4

User 3

Term	A	C	E	F
DF	4	6	5	3

Number of common terms between User 1 (Given user) and User 2= 3 (Terms A, B and D)

Number of common terms between User 1 and User 3= 2 (Terms A and C)

Thus, User 2 is more similar to the given user than User 3 due to more common terms. The best contributing term for User 2 is term “A” as it has the least document frequency which means that it the most important term in the document.

- Model= TF-IDF: Similarity between given users and other users is computed by taking the cosine similarity. The users with the lowest cosine distance are classified as more similar to the given user. The highest contributing terms for each user are the terms which have the highest product when multiplied with the corresponding terms of the given users, thus contributing to a cosine score closer to 1.

User 1 (Given user)

Term	A	B	C	D
TF-IDF	0.5	0.91	1.33	2.14

User 2

Term	A	B	D	G
TF-IDF	0.33	0.25	0.5	1

User 3

Term	A	C	E	F
------	---	---	---	---

TF-IDF	0.89	0.44	0.9	0.86
--------	------	------	-----	------

$$\begin{aligned}
\text{Cosine similarity}(\text{User1}, \text{User2}) &= \cos\theta = \frac{\sum (tfidf_{\text{User1}} * tfidf_{\text{User2}})}{\sqrt{\sum (tfidf_{\text{User1}})^2} * \sqrt{\sum (tfidf_{\text{User2}})^2}} \\
&= \frac{(0.5 * 0.33) + (0.91 * 0.25) + (1.33 * 0.5) + (2.14 * 1)}{\sqrt{(0.5)^2 + (0.91)^2 + (1.33)^2 + (2.14)^2} * \sqrt{(0.33)^2 + (0.25)^2 + (0.5)^2 + (1)^2}} \\
&= \frac{0.165 + 0.2275 + 0.665 + 2.14}{\sqrt{7.4266} * \sqrt{1.4214}} = \frac{3.1975}{3.24275} = 0.986
\end{aligned}$$

$$\begin{aligned}
\text{Cosine similarity}(\text{User1}, \text{User3}) &= \cos\theta = \frac{\sum (tfidf_{\text{User1}} * tfidf_{\text{User3}})}{\sqrt{\sum (tfidf_{\text{User1}})^2} * \sqrt{\sum (tfidf_{\text{User3}})^2}} \\
&= \frac{(0.5 * 0.89) + (0.91 * 0.44) + (1.33 * 0.9) + (2.14 * 0.86)}{\sqrt{(0.5)^2 + (0.91)^2 + (1.33)^2 + (2.14)^2} * \sqrt{(0.89)^2 + (0.44)^2 + (0.9)^2 + (0.86)^2}} \\
&= \frac{0.445 + 0.4004 + 1.197 + 1.8404}{\sqrt{7.4266} * \sqrt{2.5353}} = \frac{3.8828}{4.3382} = 0.895
\end{aligned}$$

Since, cosine similarity between user 1 and user 2 is closer to one than the cosine similarity between user 1 and user 3, User 2 is more similar to User 1. The highest contributing term for user 2 is the term "A" as the difference between the idf values of respective terms of user 1 and 2 is the least for term A, i.e. $|0.5 - 0.33| = 0.17$.

Task 2:

Each image in the input file has an image id and information related to each term for all terms of an image: term value, term frequency (tf), document frequency (df) and tf-idf. Based on the input model requested at run time, one of the following program branches is executed:

- Model= TF: Vectors are prepared on basis of the term frequency of the given image and other image's common terms. Euclidean distance is then calculated using these vectors and the images which have the least score are sent to the output. The highest contributing terms are those for which difference of term frequencies between the

corresponding ones in the input image's terms is the least, since they contribute to a lower Euclidean distance.

Image 1 (Given image)

Term	A	B	C	D
Term Frequency	4	8	3	5

Image 2

Term	A	C	D	E
Term Frequency	5	6	9	3

Image 3

Term	A	B	D	F
Term Frequency	2	1	9	4

$$\begin{aligned}
 \text{Euclidean Distance}(\text{Image1}, \text{Image2}) &= \sqrt{\left(t_{f_A(\text{Image1})} - t_{f_A(\text{Image2})}\right)^2 + \left(t_{f_B(\text{Image1})} - t_{f_B(\text{Image2})}\right)^2 + \left(t_{f_C(\text{Image1})} - t_{f_C(\text{Image2})}\right)^2 + \left(t_{f_D(\text{Image1})} - t_{f_D(\text{Image2})}\right)^2} \\
 &= \sqrt{(4 - 5)^2 + (8 - 0)^2 + (3 - 6)^2 + (5 - 9)^2} \\
 &= \sqrt{(-1)^2 + (8)^2 + (-3)^2 + (-4)^2} \\
 &= \sqrt{1 + 64 + 9 + 16} \qquad \qquad \qquad = \sqrt{90} \qquad \qquad \qquad = 9.487
 \end{aligned}$$

$$\begin{aligned}
 \text{Euclidean Distance}(\text{Image1}, \text{Image3}) &= \sqrt{\left(t_{f_A(\text{Image1})} - t_{f_A(\text{Image3})}\right)^2 + \left(t_{f_B(\text{Image1})} - t_{f_B(\text{Image3})}\right)^2 + \left(t_{f_C(\text{Image1})} - t_{f_C(\text{Image3})}\right)^2 + \left(t_{f_D(\text{Image1})} - t_{f_D(\text{Image3})}\right)^2} \\
 &= \sqrt{(4 - 2)^2 + (8 - 1)^2 + (3 - 0)^2 + (5 - 9)^2} \qquad \qquad \qquad = \sqrt{(2)^2 + (7)^2 + (3)^2 + (-4)^2} \\
 &= \sqrt{4 + 49 + 9 + 16} \qquad \qquad \qquad = \sqrt{78} \qquad \qquad \qquad = 8.832
 \end{aligned}$$

Thus, Image 3 is more similar to given image than Image 2 due to lesser Euclidean distance. The best contributing term for image 3 is term “A” as it contributes the least to the distance.

- Model= DF: The similarity for DF is calculated by taking the intersection of the input terms and the other image's terms. The images with highest number of common terms are classified as more similar. The highest contributed terms for each image are the ones with lowest DF value because they are rare, and thus important.

Image 1 (Given image)

Term	A	B	C	D
DF	4	8	6	5

Image 2

Term	A	B	D	G
DF	4	8	5	4

Image 3

Term	A	C	E	F
DF	4	6	5	3

Number of common terms between Image 1 (Given image) and Image 2= 3 (Terms A, B and D)

Number of common terms between Image 1 and Image 3= 2 (Terms A and C)

Thus, Image 2 is more similar to the given image than Image 3 due to more common terms. The best contributing term for Image 2 is term "A" as it has the least document frequency which means that it is the most important term in the document.

- Model= TF-IDF: Similarity between given image and other images is computed by taking the cosine similarity. The images with the lowest cosine distance are classified as more similar to the given image. The highest contributing terms for each image are the terms which have the highest product when multiplied with the corresponding terms of the given users, thus contributing to a cosine score closer to 1.

Image 1 (Given image)

Term	A	B	C	D
------	---	---	---	---

TF-IDF	0.5	0.91	1.33	2.14
--------	-----	------	------	------

Image 2

Term	A	B	D	G
TF-IDF	0.33	0.25	0.5	1

Image 3

Term	A	C	E	F
TF-IDF	0.89	0.44	0.9	0.86

$$\begin{aligned}
 \text{Cosine similarity}(\text{Image1}, \text{Image2}) &= \cos\theta = \frac{\sum (tfidf_{\text{Image1}} * tfidf_{\text{Image2}})}{\sqrt{\sum (tfidf_{\text{Image1}})^2} * \sqrt{\sum (tfidf_{\text{Image2}})^2}} \\
 &= \frac{(0.5 * 0.33) + (0.91 * 0.25) + (1.33 * 0.5) + (2.14 * 1)}{\sqrt{(0.5)^2 + (0.91)^2 + (1.33)^2 + (2.14)^2} * \sqrt{(0.33)^2 + (0.25)^2 + (0.5)^2 + (1)^2}} \\
 &= \frac{0.165 + 0.2275 + 0.665 + 2.14}{\sqrt{7.4266} * \sqrt{1.4214}} = \frac{3.1975}{3.24275} = 0.986
 \end{aligned}$$

$$\begin{aligned}
 \text{Cosine similarity}(\text{Image1}, \text{Image3}) &= \cos\theta = \frac{\sum (tfidf_{\text{Image1}} * tfidf_{\text{Image3}})}{\sqrt{\sum (tfidf_{\text{Image1}})^2} * \sqrt{\sum (tfidf_{\text{Image3}})^2}} \\
 &= \frac{(0.5 * 0.89) + (0.91 * 0.44) + (1.33 * 0.9) + (2.14 * 0.86)}{\sqrt{(0.5)^2 + (0.91)^2 + (1.33)^2 + (2.14)^2} * \sqrt{(0.89)^2 + (0.44)^2 + (0.9)^2 + (0.86)^2}} \\
 &= \frac{0.445 + 0.4004 + 1.197 + 1.8404}{\sqrt{7.4266} * \sqrt{2.5353}} = \frac{3.8828}{4.3382} = 0.895
 \end{aligned}$$

Since, cosine similarity between image 1 and image 2 is closer to one than the cosine similarity between image 1 and image 3, Image 2 is more similar to Image 1. The highest contributing term for image 2 is the term “A” as the difference between the idf values of respective terms of image 1 and 2 is the least for term A, i.e. $|0.5 - 0.33| = 0.17$.

Task 3:

Each location in the first input file has a location id and location title. The second input file has the location title and information related to each term for all terms of a location: term value, term frequency (tf), document frequency (df) and tf-idf. The location title is fetched on basis of the location id given as an input and this location title is used to get the terms and their frequencies. Based on the input model requested at run time, one of the following program branches is executed:

- Model= TF: Vectors are prepared on basis of the term frequency of the given location and other location’s common terms. Euclidean distance is then calculated using these vectors and the locations which have the least score are sent to the output. The highest contributing terms are those for which difference of term frequencies between the corresponding ones in the input location’s terms is the least, since they contribute to a lower Euclidean distance.

Location 1 (Given location)

Term	A	B	C	D
Term Frequency	4	8	3	5

Location 2

Term	A	C	D	E
Term Frequency	5	6	9	3

Location 3

Term	A	B	D	F
Term Frequency	2	1	9	4

$$Euclidean\ Distance(Location1, Location2) = \sqrt{(tf_{A(Location1)} - tf_{A(Location2)})^2 + (tf_{B(Location1)} - tf_{B(Location2)})^2 + (tf_{C(Location1)} - tf_{C(Location2)})^2 + (tf_{D(Location1)} - tf_{D(Location2)})^2}$$

$$\begin{aligned}
&= \sqrt{(4-5)^2 + (8-0)^2 + (3-6)^2 + (5-9)^2} \\
&= \sqrt{(-1)^2 + (8)^2 + (-3)^2 + (-4)^2} \\
&= \sqrt{1 + 64 + 9 + 16} \qquad \qquad \qquad = \sqrt{90} \qquad \qquad \qquad = 9.487
\end{aligned}$$

$$\begin{aligned}
\text{Euclidean Distance}(\text{Location1}, \text{Location3}) &= \sqrt{(tf_{A(\text{Location1})} - tf_{A(\text{Location3})})^2 + (tf_{B(\text{Location1})} - tf_{B(\text{Location3})})^2 + (tf_{C(\text{Location1})} - tf_{C(\text{Location3})})^2 + (tf_{D(\text{Location1})} - tf_{D(\text{Location3})})^2} \\
&= \sqrt{(4-2)^2 + (8-1)^2 + (3-0)^2 + (5-9)^2} \qquad \qquad \qquad = \sqrt{(2)^2 + (7)^2 + (3)^2 + (-4)^2} \\
&= \sqrt{4 + 49 + 9 + 16} \qquad \qquad \qquad = \sqrt{78} \qquad \qquad \qquad = 8.832
\end{aligned}$$

Thus, Location 3 is more similar to given location than Location 2 due to lesser Euclidean distance. The best contributing term for location 3 is term “A” as it contributes the least to the distance.

- Model= DF: The similarity for DF is calculated by taking the intersection of the input terms and the other location’s terms. The locations with the highest number of common terms are classified as more similar. The highest contributed terms for each location are the ones with lowest DF value because they are rare, and thus important.

Location 1 (Given location)

Term	A	B	C	D
DF	4	8	6	5

Location 2

Term	A	B	D	G
DF	4	8	5	4

Location 3

Term	A	C	E	F
DF	4	6	5	3

Number of common terms between Location 1 (Given location) and Location 2= 3 (Terms A, B and D)

Number of common terms between Location 1 and Location 3= 2 (Terms A and C)

Thus, Location 2 is more similar to the given location than Location 3 due to more common terms. The best contributing term for Location 2 is term “A” as it has the least document frequency which means that it is the most important term in the document.

- Model= TF-IDF: Similarity between given location and other locations is computed by taking the cosine similarity. The location with the lowest cosine distance are classified as more similar to the given location. The highest contributing terms for each user are the terms which have the highest product when multiplied with the corresponding terms of the given locations, thus contributing to a cosine score closer to 1.

Location 1 (Given location)

Term	A	B	C	D
TF-IDF	0.5	0.91	1.33	2.14

Location 2

Term	A	B	D	G
TF-IDF	0.33	0.25	0.5	1

Location 3

Term	A	C	E	F
TF-IDF	0.89	0.44	0.9	0.86

$$\begin{aligned}
 \text{Cosine similarity}(\text{Location1}, \text{Location2}) &= \cos\theta = \frac{\sum (tfidf_{\text{Location1}} * tfidf_{\text{Location2}})}{\sqrt{\sum (tfidf_{\text{Location1}})^2} * \sqrt{\sum (tfidf_{\text{Location2}})^2}} \\
 &= \frac{(0.5*0.33) + (0.91*0.25) + (1.33*0.5) + (2.14*1)}{\sqrt{(0.5)^2 + (0.91)^2 + (1.33)^2 + (2.14)^2} * \sqrt{(0.33)^2 + (0.25)^2 + (0.5)^2 + (1)^2}}
 \end{aligned}$$

$$= \frac{0.165 + 0.2275 + 0.665 + 2.14}{\sqrt{7.4266} * \sqrt{1.4214}} = \frac{3.1975}{3.24275} = 0.986$$

$$\begin{aligned} \text{Cosine similarity}(\text{Location1}, \text{Location3}) &= \cos\theta = \frac{\sum (tfidf_{\text{Location1}} * tfidf_{\text{Location3}})}{\sqrt{\sum (tfidf_{\text{Location1}})^2} * \sqrt{\sum (tfidf_{\text{Location3}})^2}} \\ &= \frac{(0.5 * 0.89) + (0.91 * 0.44) + (1.33 * 0.9) + (2.14 * 0.86)}{\sqrt{(0.5)^2 + (0.91)^2 + (1.33)^2 + (2.14)^2} * \sqrt{(0.89)^2 + (0.44)^2 + (0.9)^2 + (0.86)^2}} \\ &= \frac{0.445 + 0.4004 + 1.197 + 1.8404}{\sqrt{7.4266} * \sqrt{2.5353}} = \frac{3.8828}{4.3382} = 0.895 \end{aligned}$$

Since, cosine similarity between location 1 and location 2 is closer to one than the cosine similarity between location 1 and location 3, Location 2 is more similar to Location 1. The highest contributing term for location 2 is the term “A” as the difference between the idf values of respective terms of location 1 and 2 is the least for term A, i.e.

$$|0.5 - 0.33| = 0.17.$$

Task 4:

There are three inputs to this task: the location id, a colour model (CM/CM3x3/CN/CN3x3/CSD/GLRLM/GLRLM3x3/HOG/LBP/LBP3x3) and k (number of results to be returned). For given location and a model, csv files corresponding to that model of other locations need to be parsed. Similarity is calculated in the following manner:

For each image in the given location, a best match based on cosine similarity is calculated between the other images for other locations. For example, cosine similarity of Image 1 in given location 1 is computed with all images in location 2 and location 3.

Location 1 (Given location)

Image ID	Mean H	STD H	Skew H	Mean S	STD S	Skew S	Mean V	STD V	Skew V
Image 1	4	8	6	5	4	3	1	5	4
Image 2	1	8	3	6	7	9	2	2	1

Location 2

Image ID	Mean H	STD H	S k e w H	Mean S	S T D S	Skew S	Mean V	STD V	S k e w V
Image 1	1	2	3	4	5	6	7	8	9
Image 2	9	8	7	6	5	4	3	2	1

Location 3

Image ID	Mean H	STD H	S k e w H	Mean S	S T D S	Skew S	Mean V	STD V	S k e w V
Image 1	2	4	6	8	0	1	3	5	7
Image 2	1	3	5	7	9	0	2	4	6

$$\text{Cosine similarity}(\text{image1}(\text{location 1}), \text{image1}(\text{location 2})) = \cos\theta = \frac{\sum (cmf_{\text{image1}(\text{location 1})} * cmf_{\text{image1}(\text{location 2})})}{\sqrt{\sum (cmf_{\text{image1}(\text{location 1})})^2} * \sqrt{\sum (cmf_{\text{image1}(\text{location 2})})^2}}$$

= 0.735

$$\text{Cosine similarity}(\text{image1}(\text{location 1}), \text{image2}(\text{location 2})) = \cos\theta = \frac{\sum (cmf_{\text{image1}(\text{location 1})} * cmf_{\text{image2}(\text{location 2})})}{\sqrt{\sum (cmf_{\text{image1}(\text{location 1})})^2} * \sqrt{\sum (cmf_{\text{image2}(\text{location 2})})^2}}$$

= 0.907

$$\text{Cosine similarity}(\text{image1}(\text{location 1}), \text{image1}(\text{location 3})) = \cos\theta = \frac{\sum (cmf_{\text{image1}(\text{location 1})} * cmf_{\text{image1}(\text{location 3})})}{\sqrt{\sum (cmf_{\text{image1}(\text{location 1})})^2} * \sqrt{\sum (cmf_{\text{image1}(\text{location 3})})^2}}$$

= 0.84

$$\text{Cosine similarity}(\text{image1}(\text{location 1}), \text{image2}(\text{location 3})) = \cos\theta = \frac{\sum (cmf_{\text{image1}(\text{location 1})} * cmf_{\text{image2}(\text{location 3})})}{\sqrt{\sum (cmf_{\text{image1}(\text{location 1})})^2} * \sqrt{\sum (cmf_{\text{image2}(\text{location 3})})^2}}$$

= 0.81

Thus, for image 1 the best matches in location 2 and 3 are as follows:

Location 1 (Given Location)	Location 2	Location 3
Image 1	Image 2	Image 1

In this manner best match for each image is calculated across each location. Then the scores for each location are averaged and the one with the highest score is classified as the most similar to the given location. The most contributing images for a location are the ones with the highest cosine similarities.

Task 5:

There are two inputs to this task: location id and k(number of results to be returned). In order to return the most similar location to the given location we compare the given location's colour model csv files to the corresponding colour model files of other locations. Similarity is calculated in the following manner:

We take the average of each column which represents a feature in the colour model file.

Location 1

Image ID	Mean H	STD H	S k e w H	Mean S	S T D S	Skew S	Mean V	STD V	S k e w V
Image 1	4	8	6	5	4	3	1	5	4
Image 2	1	8	3	6	7	9	2	2	1

=

Image ID	Mean H	STD H	S k e w H	Mean S	S T D S	Skew S	Mean V	STD V	S k e w V
Image 1	2.5	8	4.5	5.5	5.5	4.5	1.5	3.5	2.5

Location 2

Image ID	Mean H	STD H	S k e w H	Mean S	S T D S	Skew S	Mean V	STD V	S k e w V
Image 1	1	2	2	3	8	9	3	8	9
Image 2	9	8	6	5	6	4	5	2	1

=

Image ID	Mean H	STD H	S k e w H	M e a n S	S T D S	Skew S	Mean V	STD V	S k e w V
Image 1	5	8	4	4	7	6.5	4	5	5

Once we have taken the average we calculate the similarity between given location and each location in a colour model. We do the same for all models until we have a computation which resembles the following:

Location	CM	CM3x3	CN	CN3x3	CSD	GLRL M	GLRL M3x3	HOG	LBP	LBP3x3	Total
1	0.90	0.70	0.65	0.34	0.90	0.31	0.88	0.31	0.88	0.77	6.64
2	0.76	0.33	0.09	0.31	0.66	0.31	0.78	0.98	0.99	0.12	5.33
3	0.56	0.98	0.43	0.89	0.90	0.96	0.43	0.52	0.49	0.68	6.84

Once we get the scores for all the models, we calculate the total for each location. Then we rank the locations in descending order. In this case, Location 3 is the most similar to the given location. The contributing scores of each visual model are computed during this process.

Bibliography

[1] K. Selçuk Candan, Data Management for Multimedia Retrieval