

Lab4 LSH

Wu Han, Class: SEIEE 2301

December 24, 2024

1 Experiment Overview

This experiment introduces the concept of Locality Sensitive Hashing (LSH). The basic idea of LSH is to first classify images based on their hash values, and then determine the final similarity by comparing feature vectors. This approach is much faster than directly comparing feature vectors one by one, significantly saving traversal time.

In this experiment, we need to complete such process: establish feature vectors and perform Hamming code analysis and hash value calculation.

2 Problem Restatement and Key Code Explanation

2.1 Creating Feature Vectors

According to the problem description, the feature vector consists of the color histograms of the four corners of the image: top-left, bottom-left, top-right, and bottom-right. We can use such code to achieve this, as shown below:

First we use this function to normalize a given vector to unit length.

```
1 def normalize_vector(v: List[float]) -> np.ndarray:  
2     """Normalize a vector to unit length."""  
3     return np.array(v) / np.linalg.norm(v)
```

It converts the input list `v` into a numpy array and then scales it by its Euclidean norm (magnitude), making it a unit vector.

Then we use this function to calculate the color histogram of an image by dividing the image into four quadrants and computing the sum of RGB values for each quadrant.

```
1 def compute_color_histogram(img: np.ndarray) -> np.ndarray:  
2     """Compute a color histogram for the image and return a feature  
3         vector."""  
4     # Split the image into 4 quadrants and calculate the sum of RGB  
5     values for each quadrant  
6     height, width, _ = img.shape
```

```

5     quadrants = [img[:height // 2, :width // 2], img[:height // 2,
6                   width // 2:], img[height // 2:, :width // 2],
7                   img[height // 2:, width // 2:]]
8
9     # Sum the RGB values for each quadrant
10    feature_vectors = [np.sum(q, axis=(0, 1)) for q in quadrants]
11
12    # Normalize each quadrant's feature vector
13    normalized_vectors = [normalize_vector(v) for v in
14                           feature_vectors]
15
16    # Concatenate all normalized vectors into a single feature
17    # vector
18    return np.concatenate(normalized_vectors)

```

The image is split into four quadrants (top-left, top-right, bottom-left, bottom-right), and the sum of RGB values in each quadrant is calculated. These values are normalized to form feature vectors, which are then concatenated to form a final feature vector for the image.

2.2 Getting Hamming Code

Although the PPT shows a method to obtain hash values without using Hamming code, Python's strong string manipulation capabilities allow us to directly use strings to get the Hamming code. The corresponding code is:

We use this function to normalize the histogram values into discrete values (0, 1, or 2) based on specific thresholds.

```

1 def normalize_histogram(v: np.ndarray) -> np.ndarray:
2     """Normalize the histogram into discrete values (0, 1, 2) based
3       on thresholds."""
4     v = v.copy()
5     for i in range(len(v)):
6         if v[i] < 0.3:
7             v[i] = 0
8         elif v[i] < 0.6:
9             v[i] = 1
10        else:
11            v[i] = 2
12    return v

```

This function loops through each element of the histogram and assigns a value of 0, 1, or 2 depending on which threshold the value falls into (0.3, 0.6).

Then we use this function to generate a Hamming code from the normalized histogram.

```

1 def generate_hamming_code(v: np.ndarray) -> str:
2     """Generate a Hamming code string from the normalized histogram
3       ."""
4     return ''.join(['00' if val == 0 else '10' if val == 1 else
5                     '11' for val in v])

```

It maps each value in the normalized histogram to a corresponding 2-bit string: '00' for 0, '10' for 1, and '11' for 2, and concatenates them to form the full Hamming code.

2.3 Using Locality Sensitive Hashing for Retrieval

Locality Sensitive Hashing involves directly taking certain bits from the Hamming code to obtain the hash function. The corresponding code is:

This function selects specific bits from the Hamming code based on a list of indices.

```
1 def select_bits(hamming_code: str, indices: List[int]) -> str:
2     """Select bits from the Hamming code based on the specified
3     indices."""
4     return ''.join(hamming_code[i] for i in indices)
```

It selects the bits from the Hamming code according to the provided list of indices.

This function computes the cosine similarity between two feature vectors.

```
1 def compute_similarity(vec1: np.ndarray, vec2: np.ndarray) -> float:
2     """Compute cosine similarity between two feature vectors."""
3     return np.dot(vec1, vec2)
```

It computes the cosine similarity between two vectors by taking their dot product, assuming both vectors are normalized.

2.4 Find Similar Images

This function processes a dataset of images and finds the most similar images to the target image based on the cosine similarity of their feature vectors.

```
1 def find_similar_images(target_img: np.ndarray, dataset_dir: str,
2     key_indices: List[int],
3     num_images: int = 40, top_n: int = 6) ->
4     List[Dict[str, float]]:
5     """Find the most similar images from the dataset to the target
6     image."""
7     vectors, dataset = [], []
8
9     # Process all images in the dataset
10    for i in range(1, num_images + 1):
11        img = cv2.imread(f"{dataset_dir}/{i}.jpg")
12        color_hist = compute_color_histogram(img)
13        normalized_hist = normalize_histogram(color_hist)
14        hamming_code = generate_hamming_code(normalized_hist)
15        selected_bits = select_bits(hamming_code, key_indices)
16        dataset.append(selected_bits)
17        vectors.append(color_hist)
18
19    # Process the target image
20    target_hist = compute_color_histogram(target_img)
21    target_normalized_hist = normalize_histogram(target_hist)
```

```

19     target_hamming_code = generate_hamming_code(
20         target_normalized_hist)
21     target_selected_bits = select_bits(target_hamming_code,
22         key_indices)
23
24     # Find the similarities between the target image and dataset
25     images
26     similarities = {}
27     for i, vector in enumerate(vectors):
28         if dataset[i] == target_selected_bits:
29             similarity = compute_similarity(vector, target_hist)
30             similarities[i] = similarity
31
32     # Sort the similarities in descending order and pick the top N
33     sorted_similarities = sorted(similarities.items(), key=lambda x
34         : x[1], reverse=True)
35     top_similar_images = sorted_similarities[:top_n]
36
37     return top_similar_images

```

This function processes all the images in the dataset and calculates their feature vectors, Hamming codes, and selected bits. It first compares whether the Hamming code between the target and the object images are the same. If so, it then compares the target image with each image in the dataset using cosine similarity and returns the top N most similar images.

3 Results Analysis and Comparison

3.1 Result Analysis

Here is the result of running the code above (LSE.py). We can see the most

```

C:\Users\WH\anaconda3\python.exe E:\MultiMedia\lab4-LSH-2\lab4-LSH\LSE.py
The most similar image is: 38 with similarity score 4.0000
The other top 5 similar images are:
Rank 2: Image 12 with similarity score 3.9952
Rank 3: Image 23 with similarity score 3.9727
Rank 4: Image 26 with similarity score 3.9692
Rank 5: Image 40 with similarity score 3.9667
Rank 6: Image 7 with similarity score 3.9656
Time taken: 0.0470 seconds

```

Figure 1: Result Of Running LSE.py

similar image is 38, and this is obviously right. And the second most similar image is image12, which is also a picture of a similar tree. Other images, such as 26 and 40, have a dominant yellow tone, which is consistent with the target. Image 7, on the other hand, is a tree and shares similar shapes with the target.

3.2 Comparison of Different Projection Sets

First, we use the full Hamming code analysis, i.e., `key = list(range(0,24))`, and obtain the following result:

```
The most similar image is: 38 with similarity score 4.0000
The other top 5 similar images are:
Rank 2: Image 12 with similarity score 3.9952
Time taken: 0.0450 seconds
```

We thus found that when using this encoding method, only two images and the target image were in the same set. Upon observation, these two images were both of similar types—large trees—which indicates that our hash matching is highly effective

By comparing their distributions, we find that the differences in the 3th, 7th, 11th, 19th bits are quite significant. Therefore, we change the key and get the following results:

```
The most similar image is: 38 with similarity score 4.0000
The other top 5 similar images are:
Rank 2: Image 12 with similarity score 3.9952
Rank 3: Image 26 with similarity score 3.9692
Rank 4: Image 40 with similarity score 3.9667
Time taken: 0.0470 seconds
```

We found that when using this method, only 4 images and the target image were in the same set. Upon observation, Among these four images, two are the same tree images as mentioned above, and the other two are images of leopards. It is worth noting that yellow is the dominant color in all these images.

Considering that our Hamming code is based on the BGR color order, and the target image has blue on top and yellow at the bottom, we try comparing the 2nd, 8th, 18th, and 23th bits, resulting in:

```
The most similar image is: 38 with similarity score 4.0000
The other top 5 similar images are:
Rank 2: Image 12 with similarity score 3.9952
Rank 3: Image 40 with similarity score 3.9667
Rank 4: Image 25 with similarity score 3.9629
Rank 5: Image 8 with similarity score 3.9625
Rank 6: Image 21 with similarity score 3.9590
Time taken: 0.0490 seconds
```

In this case, Images with more blue and yellow colors are included . It is clear that our feature vectors are somewhat crude, but they generally meet our needs.

3.3 Comparison with NN

We replace the method of LSH searching, instead, we use NN search. This is how we implement it. For each feature vector in the dataset:

1. Compute the Euclidean distance between the target image's feature vector and the dataset image's feature vector using:

$$\text{distance} = \sqrt{\sum_{i=1}^n (\text{vec1}_i - \text{vec2}_i)^2}$$

2. Store the image index and its corresponding distance as a tuple in the list `distances`.

Just like this:

```
1 def compute_euclidean_distance(vec1: np.ndarray, vec2: np.ndarray)
2   -> float:
3     """Compute the Euclidean distance between two feature vectors
4     ."""
5     return np.linalg.norm(vec1 - vec2)
6
7 # Main function to process the dataset and find the most similar
8 # image using Nearest Neighbor
9 def find_nn_similar_images(target_img: np.ndarray, dataset_dir: str
10    , num_images: int = 40, top_n: int = 6) -> List[int]:
11     """Find the most similar images from the dataset to the target
12     image using Nearest Neighbor (NN)."""
13     vectors = []
14
15     # Process all images in the dataset
16     for i in range(1, num_images + 1):
17         img = cv2.imread(f"{dataset_dir}/{i}.jpg")
18         color_hist = compute_color_histogram(img)
19         vectors.append(color_hist)
20
21     # Process the target image
22     target_hist = compute_color_histogram(target_img)
23
24     # Calculate the Euclidean distances between the target image
25     # and all dataset images
26     distances = []
27     for i, vector in enumerate(vectors):
28         distance = compute_euclidean_distance(vector, target_hist)
29         distances.append((i, distance))
30
31     # Sort the distances in ascending order and pick the top N
32     # closest images
33     sorted_distances = sorted(distances, key=lambda x: x[1])
34
35     # Extract the indices of the top N most similar images
36     top_similar_images = sorted_distances[:top_n]
37
38     return top_similar_images
```

The details are in NN.py. And the result of running NN.py is:

```
The most similar image is: 38 with similarity score (distance) 0.0000
The other top 5 similar images are:
Rank 2: Image 12 with similarity score (distance) 0.0978
Rank 3: Image 23 with similarity score (distance) 0.2338
Rank 4: Image 26 with similarity score (distance) 0.2482
Rank 5: Image 40 with similarity score (distance) 0.2582
Rank 6: Image 7 with similarity score (distance) 0.2623
Time taken: 0.0456 seconds
```

After comparison, we found that the runtime of the two processing methods is essentially the same. This is partly because our database is very small, and our device performance is excellent, making it difficult to demonstrate the advantages of the algorithm on a small database. On the other hand, we have optimized the implementation of the NN algorithm, achieving high efficiency. Regarding the output image results, both algorithms produce almost identical results. The NN algorithm also tends to output images with more blue and yellow colors.

4 Further Considerations

4.1 Use of Color Histogram Features

As mentioned earlier, the retrieval effect meets our expectations. Although the feature information is somewhat crude, it effectively filters out irrelevant images. Specifically, when exploring the full set, we can unexpectedly find quite good results.

The similarity of the retrieved images is primarily reflected in the color. Since our target image has cool tones, the retrieved images are mostly blue and yellow, reflecting the main color scheme.

4.2 Can Other Features Be Designed?

Shape features describe the geometric properties of objects within an image, capturing information about their contours, boundaries, or regions. These features are particularly useful for tasks such as object recognition, image classification, and retrieval, especially when the shape plays a key role in defining the object. We can use methods like edge detector to implement it.

4.3 The Explanation Of the Content on the 10th Page of the PPT

On the 10th page of the ppt, we show another approach to implement the hashing method. The X1, X2, X3... just refers to the number in the projection sets. For example, there is a projection set [1, 3, 5, 9], then X1 refers to 1, X2 refers to 3...

4.4 Reflection on the Experiment

This experiment significantly deepened my understanding of the role of mathematics and algorithms in solving real-world problems. During the feature extraction process, I realized the importance of fundamental mathematical tools, such as vector normalization and distance computation. These seemingly basic concepts are crucial for the proper functioning of algorithms. By implementing Locality Sensitive Hashing (LSH) and Nearest Neighbor (NN) algorithms, I also gained insights into algorithm optimization, particularly in selecting feature dimensions or indexing strategies to enhance efficiency.