

# CSE 515 Multimedia and Web Databases

## Phase #3

(Due Nov 25th 2018, midnight)

**Description:** In this project, you will experiment with

- graph analysis
- clustering
- indexing
- classification

### Tasks:

- **Task 1:** Implement a program which, given a value  $k$ , creates an image-image similarity graph, such that from each image, there are  $k$  outgoing edges to  $k$  most similar/related images to it.
- **Task 2:** Given the image-image graph, identify  $c$  clusters (for a user supplied  $c$ ) using two distinct algorithms. You can use the graph partitioning/clustering algorithms of your choice for this task. Visualize the resulting image clusters.
- **Task 3:** Given an image-image graph, identify and visualize  $K$  most dominant images using Page Rank (PR) for a user supplied  $K$ .

See

- “S. Brin and L. Page. ”The anatomy of a large-scale hypertextual Web search engine”. Computer Networks and ISDN Systems 30: 107117, 1998”

- **Task 4:** Given an image-image graph and 3 user specified imageids identify and visualize  $K$  most relevant images using personalized PageRank (PPR) for a user supplied  $K$ .

See

- “J.-Y. Pan, H.-J. Yang, C. Faloutsos, and P. Duygulu. Automatic multimedia cross-modal correlation discovery. In KDD, pages 653658, 2004”

for a personalized PageRank formulation based on RandomWalks with re-start.

- **Task 5:**

- **5a:** Implement a Locality Sensitive Hashing (LSH) tool, for a and similarity/distance function of your choice, which takes as input (a) the number of layers,  $L$ , (b) the number of hashes per layer,  $k$ , and (c) a set of vectors as input and creates an in-memory index structure containing the given set of vectors. See

”Near-Optimal Hashing Algorithms for Approximate Nearest Neighbor in High Dimensions” (by Alexandr Andoni and Piotr Indyk). Communications of the ACM, vol. 51, no. 1, 2008, pp. 117-122.

- **5b:** Implement similar image search using this index structure and a combined visual model function of your choice (the combined visual model must have at least 256 dimensions): for a given image and  $t$ , visualizes the  $t$  most similar images (also outputs the numbers of unique and overall number of images considered).

#### **Task 6:** Implement

- a  $k$ -nearest neighbor based classification algorithm, and
- a PPR-based classification algorithm

which take a file containing a set of image/label pairs and associates a label to the rest of the images in the database. Visualize the labeled results.

#### **Deliverables:**

- Your code (properly commented) and a README file.
- Your outputs for the provided sample inputs.
- A short report describing your work and the results.

Please place your code in a directory titled “Code”, the outputs to a directory called “Outputs”, and your report in a directory called “Report”; zip or tar all off them together and submit it through the digital dropbox.