

Homework #8

Bloom Filters, Time Warp, Graphs I

Anastasiia Konoplina

1. Bloom filters - generate 1M random 64 bit integers and store these in a Bloom filter of size 16M bits (2 bytes per long integer). Create multiple hash functions to store these keys in 16M bit array. Now create 100M new random integers (not in the original 1M) and estimate what is the rate of "false alarms" with different Bloom filter settings. Experiment with the nr of different hash functions - 1, 2, 3, X (the more you use, the higher is the ratio of bits that are set to 1 in the Bloom filter; but also the higher the chances to hit a 0 for false keys). Report your findings and argue what is "the optimal" number of hash functions for such Bloom filters.

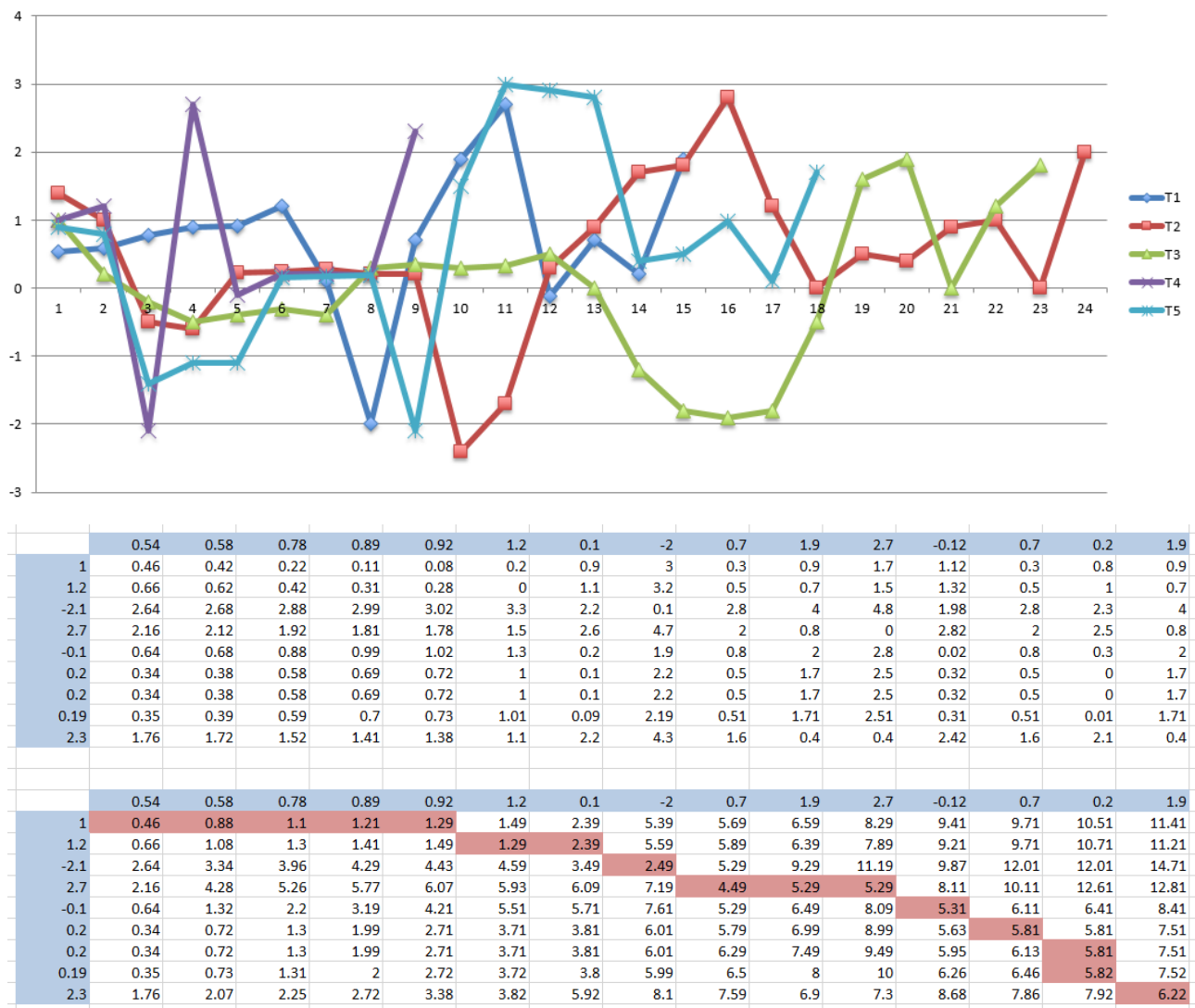
2. Compare the five series T1..T5 between each other using Dynamic Time Warping (DTW) algorithm (but you can experiment with the minimization function). Calculate all pairwise similarities. Which series are most similar?

The results are below:

	T1	T2	T3	T4	T5
T1		8.29	10.56	6.22	8.53
T2			10.32	18.33	7.96
T3				19.24	13.92
T4					13.65
T5					

The most similar are T1 and T4.

3. Visualize the time series as line graphs (as in above Excel screenshot) and respective DTW table of most similar pairs (use Excel to highlight the DTW table as was done for edit distance).



The excel file with calculations are available [here](#) on the second tab.

4. Study what are these data about. See readme files and original sources. Note that they are in slightly different formats. Better reformat first. And you can first ignore the weights or types of edges.

In “students” there is directed graph which shows social media interaction between 180 students while working on assignments.

In “karate” the data represents ties between karate club members.

In “Email-Enron” directed graph representing e-mail exchange among users.

In “usairport” directed graph showing flight between airports in USA, where weight – number of passengers.

5. Install and use Cytoscape - <http://www.cytoscape.org/> - try to visualise some graphs (from Cytoscape or even better, if from task 4) yourself. Show some screenshot from your computer. Identify key features and capabilities of Cytoscape.

I have tried to visualize “karate” graph (screenshot is below).

Key features:

- Import table/network from file/URL
- Analyze graph (shortest path, distributions of degrees, number of shared neighbors, neighborhood connectivities, average clustering coefficients etc)
- Merge graphs/tables
- Apply different layouts

