# Programming Assignment #2

### COEN 281 Pattern Recognition and Data Mining
### Department of Computer Engineering
### Santa Clara University

Dr. Ming-Hwa Wang
Phone: (408) 805-4175
Course website:
Office Hours:

Summer Quarter 2017
Email address: m1wang@scu.edu
http://www.cse.scu.edu/~mwang2/mining/
Wednesday & Friday 9:00am-9:30am

**Due date**: Midnight August 2, 2017

**Student Name:**

**SSN/ID:**

**Score:**

Correctness and boundary condition (60%):

Compiling without warnings (5%):

Error Handling (5%):

Modular design, file/directory organizing, showing input, documentation, coding standards (25%):

Automation (5%):

**Subtotal:**

Late penalty (20% per day):

Special service penalty (5%):

**Total score:**

**Outlier Detection in Streaming Data Using LOF Scores** (200 points)
Please implement outlier detection in streaming data using local outlier factor (LOF) score. You should implement it in C, C++ or Java.

Your program should take input from stdin. The input contains a window size $w$ (32-bit integer), a <host>:<port> pair where you receive an input stream of fixed-$d$-dimensional data points in a comma separated values (CSV) format. Your program should output outliers starting from the $(w+1)^{th}$ input data until end of input stream. Your program should handle concept drift too.

The LOF approach is a normalized distance-based approach. It adjusts for local variations in cluster density by normalizing distances with the average point-specific distances in a data locality. For a given data point $x$, let $v_k(x)$ be the distance to its $k$-nearest neighbor, and let $L_k(x)$ be the set of points within the $k$-nearest neighbor distance of $x$. $|L_k(x)| \geq k$ because of ties in the distance. Then, the asymmetric reachability distance $r_k(x, y)$ of object $x$ with respect to $y$ is defined as $r_k(x, y) =$ max{Dist$(x, y)$, $v_k(y)$}. When $y$ is in a dense region and the distance between $x$ and $y$ is large, $r_k(x, y)$ is equal to the true distance Dist$(x, y)$. When the distance between $x$ and $y$ is small, then $r_k(x, y)$ is smoothed out by the $k$-nearest neighbor distance of $y$. The larger the value of $k$, the greater the smoothing. The average reachability distance $ar_k(x) = $ MEAN$_{y \in L_k(x)} r_k(x, y)$, and LOF$_k(x) = $ MEAN$_{y \in L_k(x)}(ar_k(x) / ar_k(y))$. The maximum value of LOF$_k(x)$ over a range of different values of $k$ is used as the outlier score to determine the best size of the neighborhood. To handle streaming data with sliding window, we extend LOF to incremental scenarios: 1) the statistic of the newly inserted data points are computed, 2) only the LOF scores of the affected data points by the newly inserted data point in the existing data points in the window are updated, and 3) similarly updated the deleted data points.

The data point with LOF score greater than a threshold t will be reported as outlier.