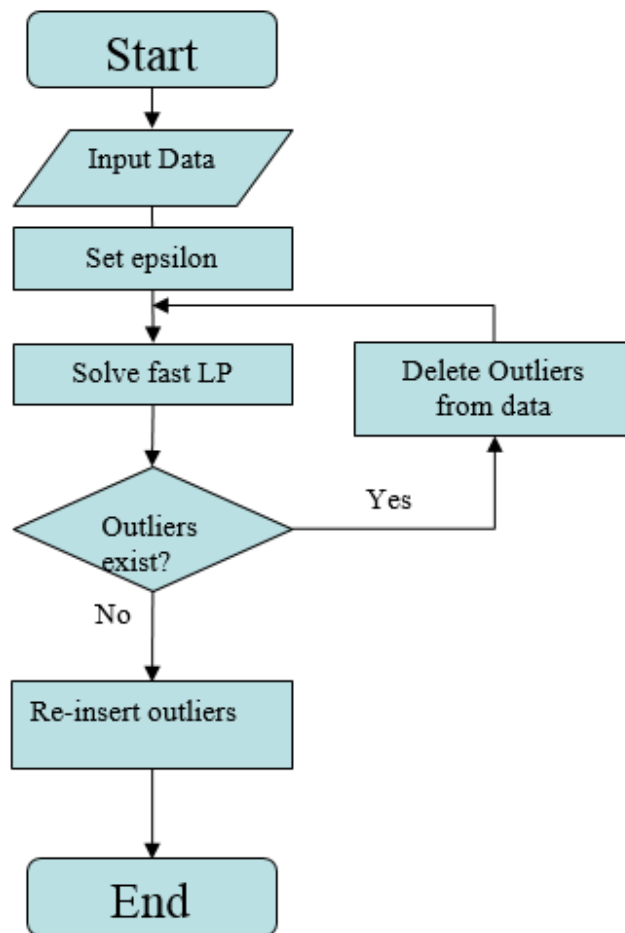


Report 10/19/2017

Our method:



What is Fast LP?

Let's see the original problem LP:

$$\max_{\mathbf{a}, \mathbf{B}, \mathbf{d}, \mathbf{t}} \sum_{i \in \mathcal{N}_A} t_i \quad (17a)$$

$$\text{s.t. } t_i + \epsilon \leq d_{ik} \quad \forall i \in \mathcal{N}_A, k \in \bar{\mathcal{S}}_i \quad (17b)$$

$$t_i \leq d_{ij} \quad \forall i \in \mathcal{N}_A, j \in \mathcal{S}_i \quad (17c)$$

$$t_i \geq 0 \quad \forall i \in \mathcal{N}_A \quad (17d)$$

$$d_{ij} = \sum_{k=1}^p \delta_{ijk} a_k + \sum_{k=1}^p \delta_{ijk}^2 b_{kk} + 2 \sum_{k=1}^{p-1} \sum_{\ell=k+1}^p \delta_{ijk} \delta_{ij\ell} b_{k\ell} \quad \forall i, j \in \mathcal{N} \quad (17e)$$

$$0 \leq d_{ij} \leq 1 \quad \forall i, j \in \mathcal{N} \quad (17f)$$

$$\mathbf{a} \in \mathbb{R}^p \quad (17g)$$

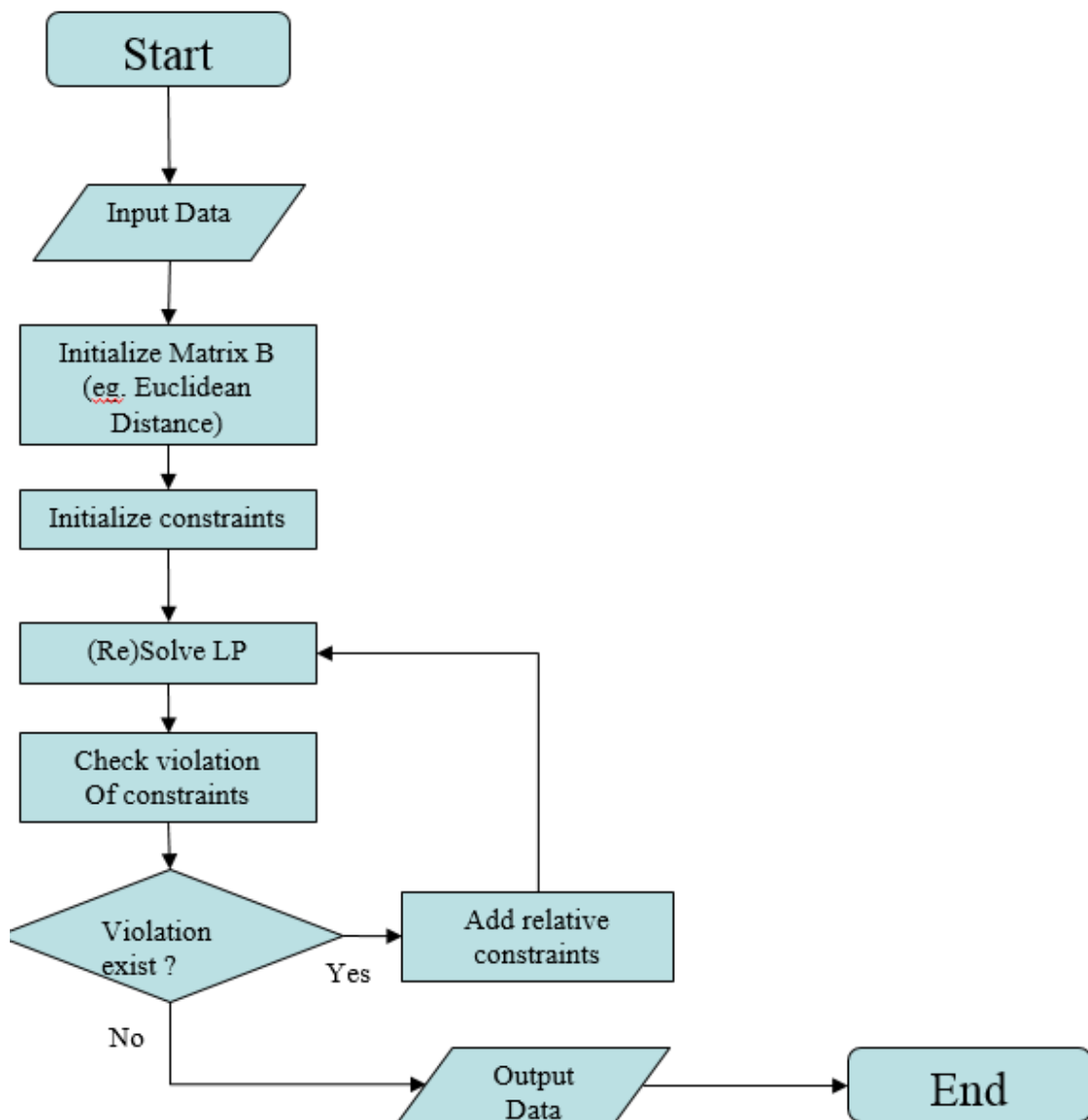
$$b_{k\ell} \in \mathbb{R} \quad \forall k = 1, \dots, p, \ell = k, \dots, p \quad (17h)$$

Our goal is to train matrix B as a metric and apply it with NN method to classify data.

Fast LP is based on original LP ,they have same solution, but Fast LP runs faster.

But how? By running iterations, adding limited number of constraints (17a,17b,17c) in every iteration rather than using all constraint in 17a,17b,17c. Number of constraints can be very large if we have large dataset.

Fast LP Algorithm with Row and Column Generation:



My questions:

1. Why NN, not 2-NN, 3-NN? OK, it reduces the calculating time.
2. What if 3-NN or 6-NN performs better in accuracy? Do we need experiments on this?
3. Value of epsilon matters only when it is nearly reaching "max epsilon", so why do we care about it? I mean it will not change the solution if we keep epsilon less equal than "max epsilon"/2.
4. About "re-insert outliers", we got this:

$$R_i = \frac{\min_{j \in \mathcal{C}_i} d_{ij}}{\min_{k \in \mathcal{C}_i} d_{ik}}$$

<b>Algorithm</b> Improved Method of Outlier Re-insertion with Ranking of $R'_i$ ;	
1:	Compute $R_i$ for each outlier $i$ ;
2:	Identify all non-outliers $j$ for which $i$ is now the nearest co-class/non-class neighbor, i.e. all $j$ whose $R_i$ values are affected by insertion of $i$ ;
3:	Define $R'_i = \min_j (R_i, R_j)$ ;
4:	Sort outliers in the order of decreasing $R'_i$ ;
5:	Insert the outliers according to decreasing ranking of $R'_i$ ;
6:	Repeat steps 1-5 until no further re-insertion is possible, i.e. remaining outliers all have $R'_i < 1$

If Metric does not change after we put points back, so why do we do that? It's meaningless.

5. About what to do next?

(1) See how initialization of Matrix B influence the solving time in Fast LP.

(2) Comparison of our method and others( k-nn ,svm...) with different dataset.

(3) Improving of the method (hardest)