# Seeding Strategies for Lloyd's kmeans

## C++ implementation of Ostrovsky, Rabani, Schulman and Swamy's ideas

Mirko Speth

22.07.2016

# Overview

Five approaches to seed Lloyd's k-means

- ► A: Random Sampling
- ► B: Greedy Deletion
- ► C: Linear time algorithm
- ► D: Linear time constant factor algorithm
- ► E: Polynomial Time Approximation Scheme (PTAS)

# A: Random Sampling

*$O(knd)$*

1. Select two random points $c_1$ and $c_2$ with probability $||c_1 - c_2||^2$
2. Perform a ball-k-means step
3. Repeat: add another point $c_{i+1}$ with probability $min_{j \in \{1...i\}} ||c_{i+1} - c_j||^2$

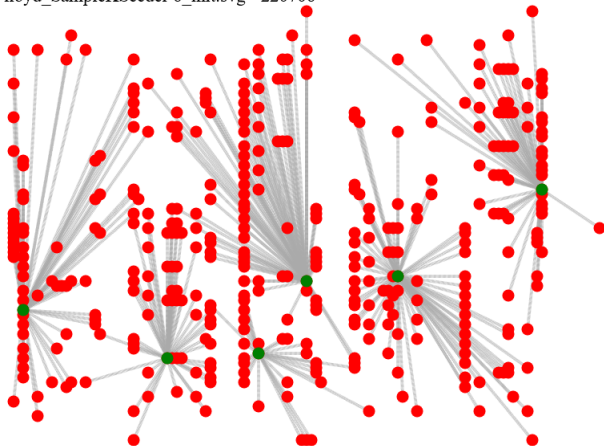# A: Random Sampling



lloyd_SampleKSeeder 6_init.svg - 220706

Figure: SVG-Output of SamplekSeeder on pbl395.tsp (k=6)

# B: Greedy Deletion

```cpp
Pointset GreedyDelSeeder::seed(Pointset init) const {
  Pointset sites = init;
  while (sites.size() > (unsigned int) k) {

    // B1: get best and second best center for each
    customer
    Partition part = Partition(&customers, sites);

    // B2: pick the center for which Tx is minimum
    int bestid = part.getMinTx();

    // B3: delete chosen partition and move points to
    centroid of voronoi region
    part.delete_set_from_partition(bestid);
    sites = part.centroids();
  }
  return sites;
}
```

# B: Greedy Deletion

Running time: outer loop $(n - k)$ iterations $\Rightarrow O(n)$
Partitioning: $O(n^2 d)$
No speed loss for second best
$\Rightarrow$ together: $O(n^3 d)$

# C: Linear time algorithm

```
1  Pointset LTSeeder::seed() const {
2
3    //C1
4    double e = instance.eps();
5    double p1 = sqrt(e);
6    int N = (int)(2 * k / (1 - 5 * p1) + 2 * log(2 / p1)
       / pow((1 - 5 * p1), 2));
7
8    SampleKSeeder samplekseeder(instance, N);
9    Pointset S = samplekseeder.seed();
10
11   //C2
12   Partition partition = Partition(&customers, S);
13   Pointset sdach = partition.centroids();
14
15   GreedyDelSeeder greedydelseeder(instance, k);
16
17   return greedydelseeder.seed(sdach);
18 }
```

# C: Linear time algorithm
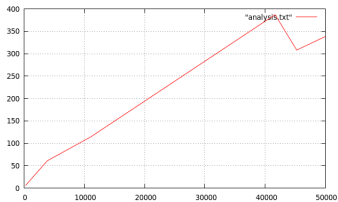
$O(nkd + k^3 d)$



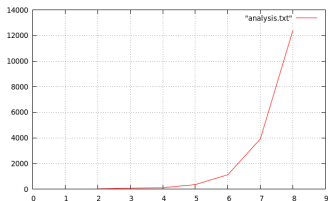Figure: complexity in #customers



Figure: complexity in #sites

# D: Linear time constant factor algorithm

```
1  Pointset DSeeder::seed() const {
2      // D1 (obtain k initial centres using last seeding
       strategy)
3    Pointset init = (LTSeeder(instance, k)).seed();
4
5    // D2 (run a ball-k-means step)
6    return ballkmeansstep(init);
7  }
```

# D: Linear time algorithm
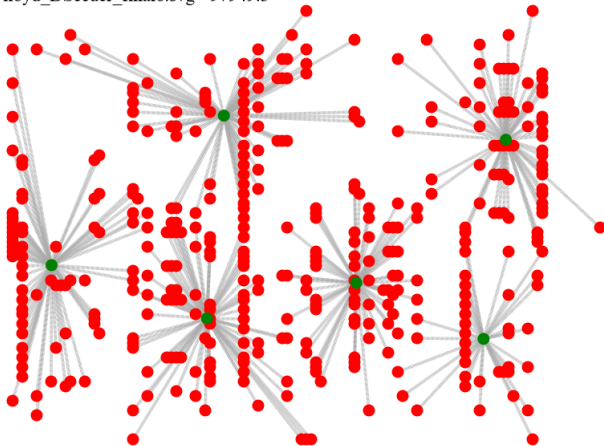


lloyd_DSeeder_final6.svg - 97949.3

Figure: SVG-Output of DSeeder on pbl395.tsp (k=6)

# E: Polynomial Time Approximation Scheme (PTAS)

```
1  Pointset ESeeder::seed() const {
2    SampleKSeeder samplek(instance, k);
3    Partition part = Partition(&customers, samplek.seed()
      );
4    return part.centroid_estimation(instance.omega,
      instance.eps);
5  }
```

# centroid estimation

```
1 for s in sites:
2     select expanded Voronoi region V[s]
3     choose a random subset R[s] of V[s] of size 4/wb
4     foreach subset A of size 1/wb:
5         add centroid(A) to T[s]
6 foreach set B in {{x1,...xk} : xi in T[i]}:
7     if error(B) < error (best) then best = B
```

# PTAS analysis

Ostrovski et al:
error at most $(1 + \omega) * OPT$ with probability $\gamma^k$
(for some constant $\gamma$)

# E: PTAS

$$O(2^{(\frac{4k}{\beta\omega})} n * d)$$
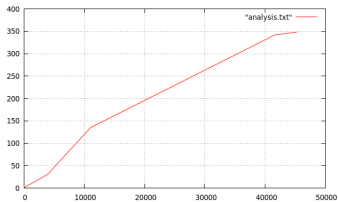


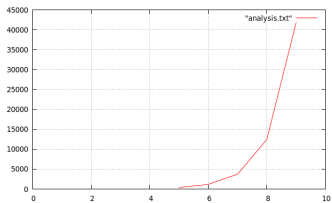Figure: complexity in #customers
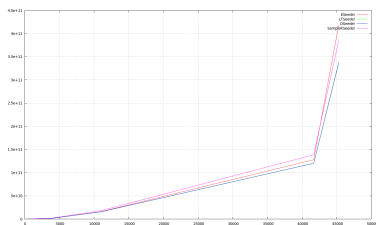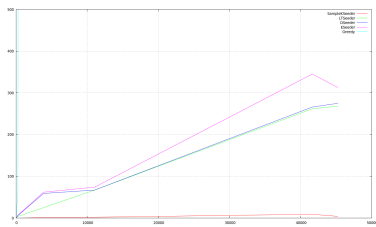


Figure: complexity in #sites

# Comparison



Figure: error comparison



Figure: runtime comparison