



Parallel Sparse K-Means

Application to document clustering

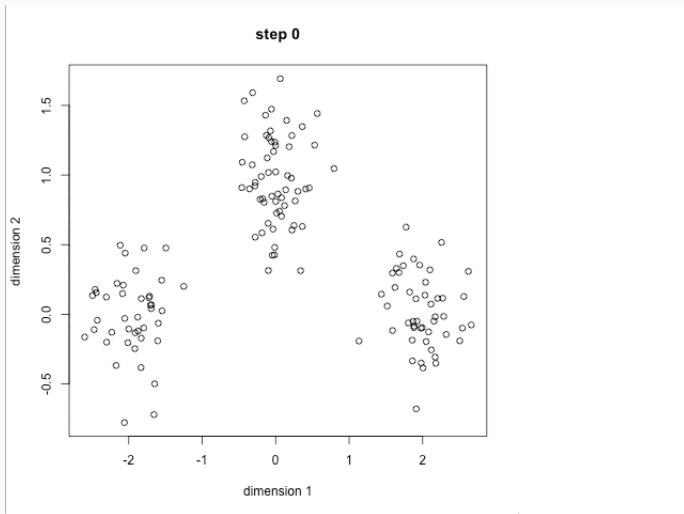
Victor Storchan

May 31, 2016

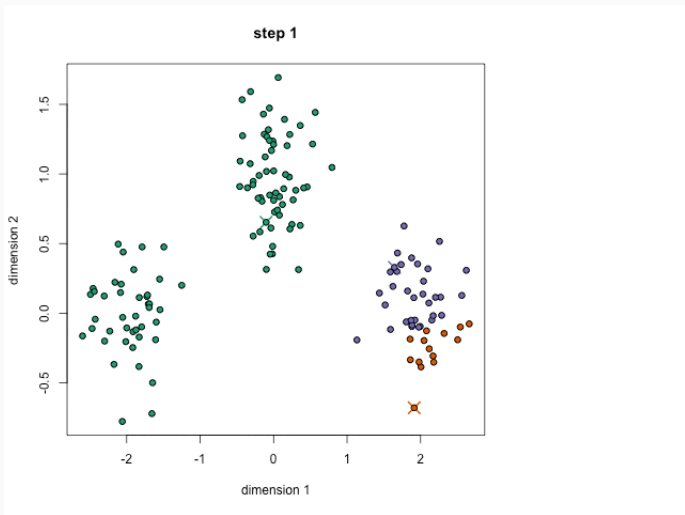
Stanford University, CME 323 (R.Zadeh)

Introduction

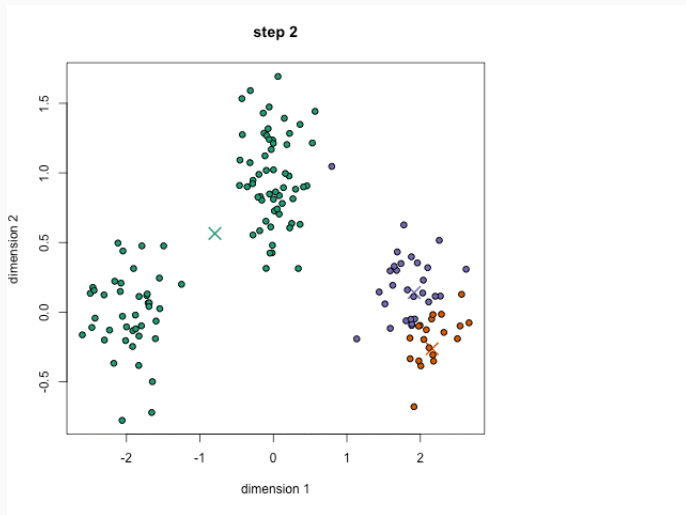
Introduction: K-Means in action



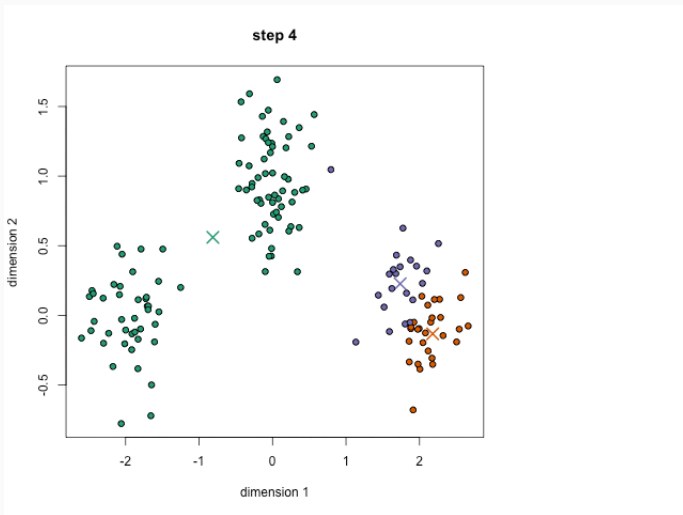
Introduction: K-Means in action



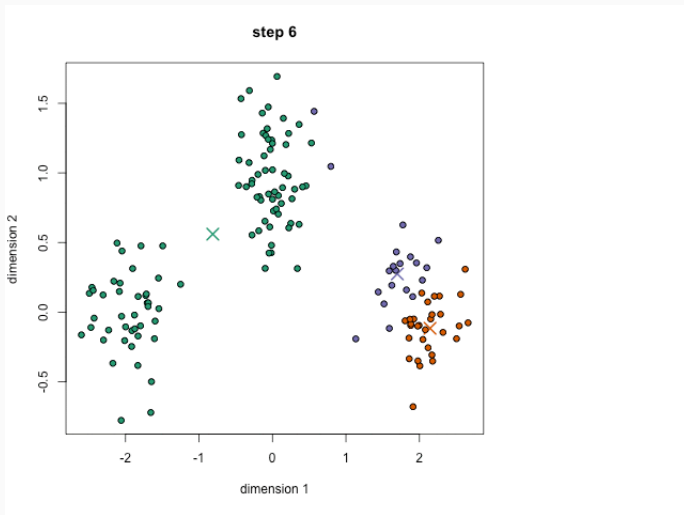
Introduction: K-Means in action



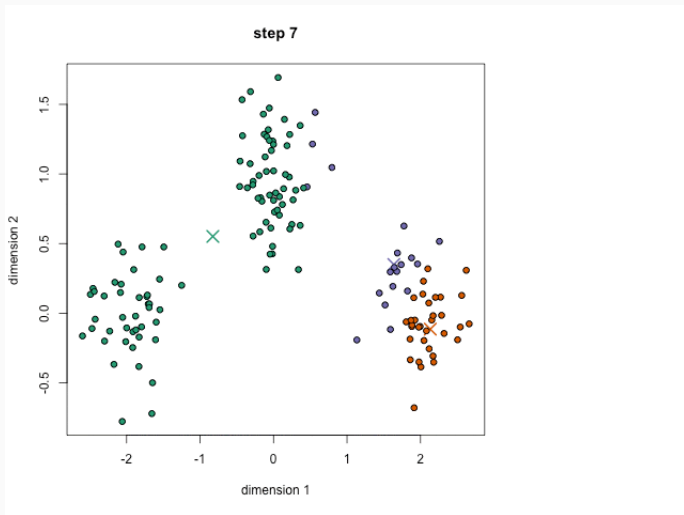
Introduction: K-Means in action



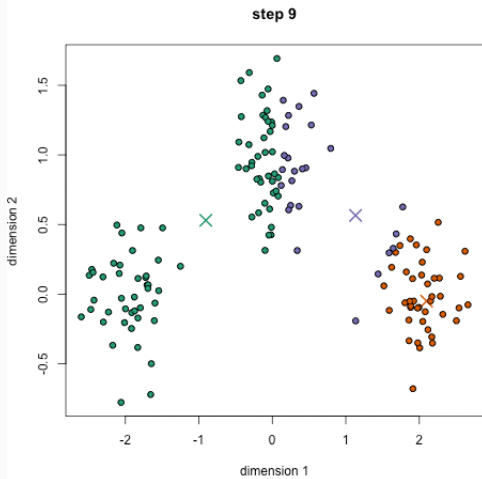
Introduction: K-Means in action



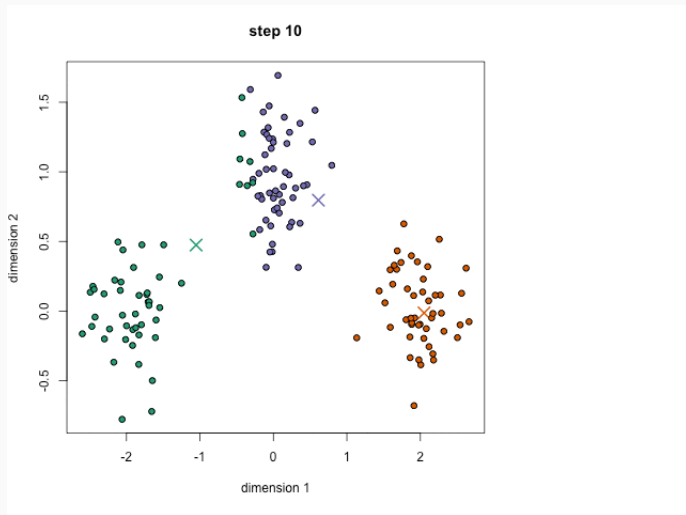
Introduction: K-Means in action



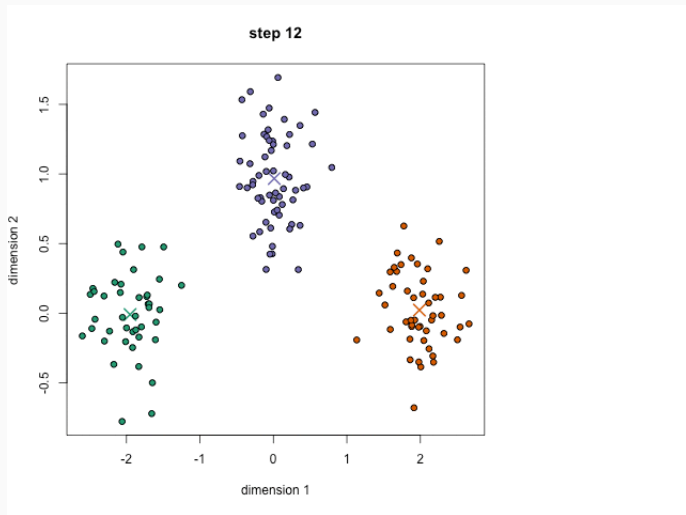
Introduction: K-Means in action



Introduction: K-Means in action



Introduction: K-Means in action



Goal of the PARALLEL SPARSE KMEANS-ALGORITHM's project:

Goal of the **PARALLEL SPARSE KMEANS-ALGORITHM**'s project:

- **Implement** a Parallel Sparse version of K-Means to typically cluster very sparse data (e.g text data).

Goal of the **PARALLEL SPARSE KMEANS-ALGORITHM**'s project:

- **Implement** a Parallel Sparse version of K-Means to typically cluster very sparse data (e.g text data).
- **Apply** it to the 20 Newsgroups data set.

Goal of the **PARALLEL SPARSE KMEANS-ALGORITHM**'s project:

- **Implement** a Parallel Sparse version of K-Means to typically cluster very sparse data (e.g text data).
- **Apply** it to the 20 Newsgroups data set.
- **Compare** the scalability of the initialization methods when they are applied to **PARALLEL SPARSE KMEANS-ALGORITHM**

Parallel Sparse KMeans-Algorithm

Parallel Sparse KMeans-Algorithm

K-Means perform badly on noisy features.

Assumption: Dissimilarity is additive with respect to any feature.

Idea: In sparse data, features are either **noisy** or **clustering features**.

- Associate weight 0 to noise features and positive weight to clustering features.
- run K-Means on the modified matrix.

Initialization step

K-Means supports 3 different initialization steps:

Initialization step

K-Means supports 3 different initialization steps:

- **random**, not efficient

Initialization step

K-Means supports 3 different initialization steps:

- **random**, not efficient
- **K-Means++**, $8\log(K)$ -approximation of the optimum, not scalable

Initialization step

K-Means supports 3 different initialization steps:

- **random**, not efficient
- **K-Means++**, $8\log(K)$ -approximation of the optimum, not scalable
- **K-Means||**, efficient ($O(\log(K))$ -Approximation algorithm) and scalable

Initialization step

K-Means supports 3 different initialization steps:

- **random**, not efficient
- **K-Means++**, $8\log(K)$ -approximation of the optimum, not scalable
- **K-Means||**, efficient ($O(\log(K))$ -Approximation algorithm) and scalable

Idea: A point x has as many chances to get selected in \mathcal{C} as it is far from the previous centers.

All the algorithms involved here are detailed in the annex of this presentation.

Parallel Sparse KMeans-Algorithm

The Parallel Sparse KMeans-Algorithm is combining:

- initialization with K-Means||.
- K-Means iterations in parallel on the weighted matrix of prediction using as many CPUs as available.

Notations:

- K : Number of clusters
- N : Number of points (ie number of documents)
- p : Dimension of the underlying space: Number of features
- s : Number of nodes in the cluster

Total work of K-Means||: $O(\log(\psi) * K^3 * N^2 * p)$.

Total depth for K-Means||: is $O(\log(p) + \log(K) + \log(N))$.

Let's suppose K-Means converges in R steps. Then:

Total work of PSK-Algorithm: $O(W_{K-Means||} + R * |\tilde{\mathcal{C}}| * p * N)$.

Total depth of PSK-Algorithm:

$O(D_{K-Means||} + R * (\log(p) + \log(|\tilde{\mathcal{C}}|) + \log(N)))$.

Communication cost for K-means:

Master machine reads the clustering center files and shuffles it to the machines.

Mapper_i for $i \in \{1, \dots, s\}$ assigns a cluster center for each points $x \in X_i \subset X$ and record the local minimum distance $d(c_i, X_i)$.

s vectors of length k are reduced.

Reducers compute the global sums, update the centers and send this data to all the processors with an *MPI_AllReduce* procedure.

=> **Communication cost:** $O(sK)$.

Communication cost for K-means||: $O(s)$.

Document Clustering

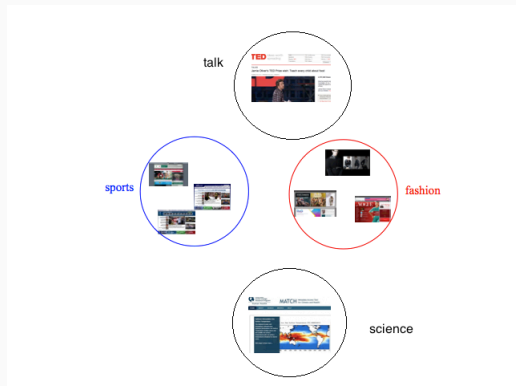
20 newsgroups Data Set

computer	recreational	science
comp.graphics comp.os.ms-windows.misc comp.sys.ibm.pc.hardware comp.sys.mac.hardware comp.windows.x	rec.autos rec.motorcycles rec.sport.baseball rec.sport.hockey	sci.crypt sci.electronics sci.med sci.space

misc.forsale	religion	talk.politics
misc.forsale	talk.religion.misc alt.atheism soc.religion.christian	talk.politics.misc talk.politics.guns talk.politics.mideast

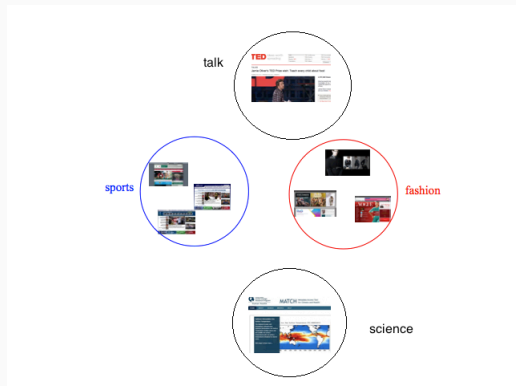
The Term Frequency-Inverse Document Frequency transformation

TF-IDF transformation: a document-to-vector map which is:



The Term Frequency-Inverse Document Frequency transformation

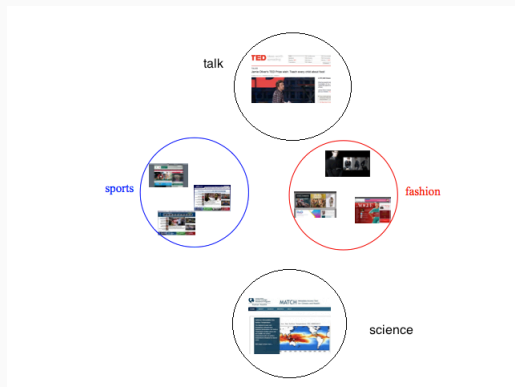
TF-IDF transformation: a document-to-vector map which is:



- Relevant to text processing

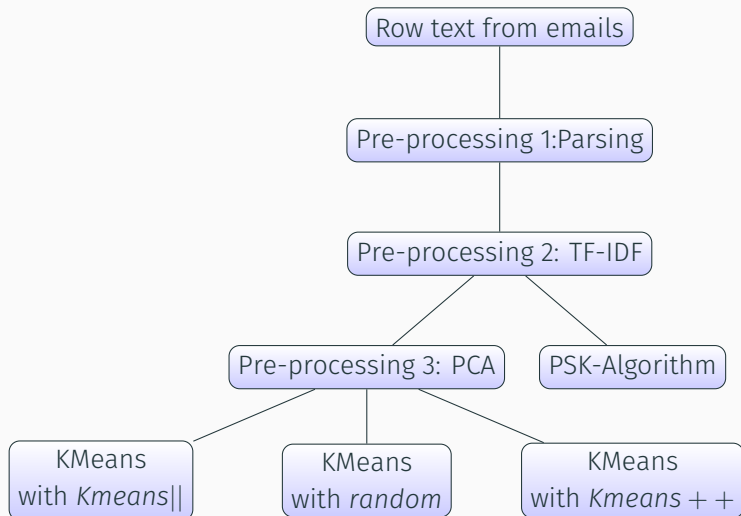
The Term Frequency-Inverse Document Frequency transformation

TF-IDF transformation: a document-to-vector map which is:

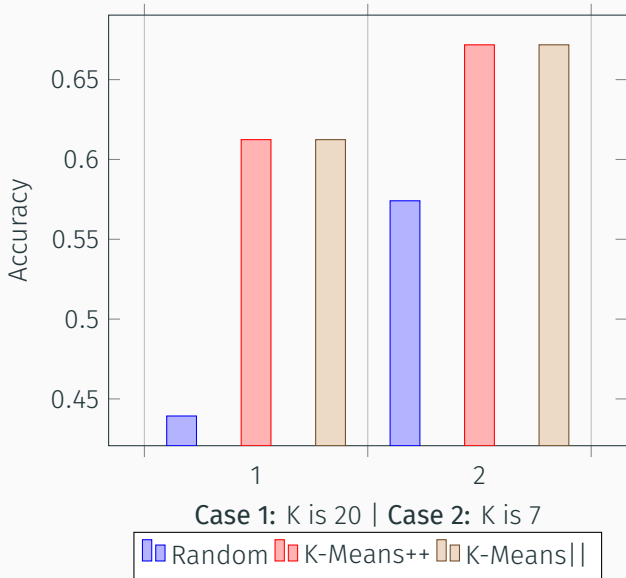


- Relevant to text processing
- Common web analysis algorithm

The Pipeline



Results



To the sake of comparison, Naive Bayes gives an accuracy of 0.4515.

Scaling

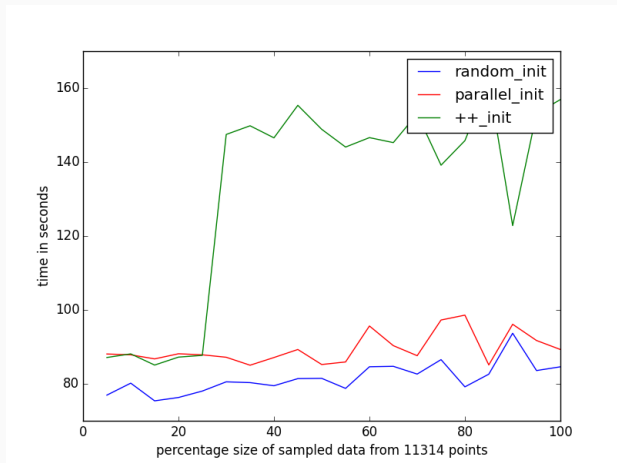


Figure 1: Runtime of initializations with respect to data size

Conclusion

PARALLEL SPARSE KMEANS-ALGORITHM:

- **Scalable** with respect to the size of the data
- differentiate clustering features from noise.

Code location:

https://github.com/victorstorchan/doc_clustering

Thank you! Questions?

type of initialization	accuracy for K=20	total time
random	0.4394	84.6117
K-Means++	0.6124	156.9197
K-Means	0.6124	90.6363

type of initialization	accuracy for K=7	total time
random	0.5741	57.9025
K-Means++	0.6718	74.8596
K-Means	0.6718	60.7070

Term Frequency/Inverse Document Frequency

$$tf_i = \frac{n_i}{\sum n_k} \text{ and } idf_i = \log \frac{|D|}{|\{d : t_i \in d\}|}$$

$|D|$: total number of documents in the corpus

$\{d : t_i \in d\}$: number of documents where the term t_i appears.

Initialization step

Algorithm 2 *K - Means ++*

```
1: procedure K-MEANS++  
2:    $\mathcal{C} \leftarrow$  Sample one point uniformly at random in  $X$   
3:   while  $|\mathcal{C}| < K$  do  
4:     Sample  $x \in X$  with probability  $\frac{d^2(x, \mathcal{C})}{\phi_X(\mathcal{C})}$   
5:      $\mathcal{C} \leftarrow \mathcal{C} \cup \{x\}$ 
```

Algorithm 3 *K - Means||*

```
1: procedure K-MEANS||  
2:    $\tilde{\mathcal{C}} \leftarrow$  Sample one point uniformly at random in  $X$   
3:    $\psi \leftarrow \phi_X(\tilde{\mathcal{C}})$   
4:   for  $O(\log(\psi))$  times do  
5:      $\mathcal{C}' \leftarrow$  Sample each  $x \in X$  independently with probability  $\frac{ld^2(x, \tilde{\mathcal{C}})}{\phi_X(\tilde{\mathcal{C}})}$   
6:      $\tilde{\mathcal{C}} \leftarrow \tilde{\mathcal{C}} \cup \mathcal{C}'$   
7:   for  $x \in \tilde{\mathcal{C}}$  do  
8:     set  $w_x = \# \{ \text{points in } X \text{ closer to } x \text{ than any other points in } \tilde{\mathcal{C}} \}$   
9:   Recluster the weighted points in  $\tilde{\mathcal{C}}$  into  $K$  clusters: use K-Means++ on  
10:   $\{w_1 c_1 \in R^p, \dots, w_{|\tilde{\mathcal{C}}|} c_{|\tilde{\mathcal{C}}|} \in R^p\}$ . Return  $\mathcal{C}$ 
```

Algorithm 5 *PSK – Algorithm*

- 1: **procedure** PSK-ALGORITHM
 - 2: Run SK-Means with modified step (a):
 - 3: When performing K-means on the weighted data, initialize the procedure with
 - 4: K-Means|| and run the K-Means iterations in parallel using as many CPUs as
 - 5: available.
-

Algorithm 1 *SK - Means(K)*

- 1: **procedure** SK-MEANS
- 2: $w \leftarrow \frac{1}{p}[1, 1, \dots, 1]^T$
- 3: Repeat (until convergence):
- 4: Step (a):
- 5: $Y := [\sqrt{w_1}X_1, \dots, \sqrt{w_p}X_p]$
- 6: for a given w , perform K-means on the transformed data set Y :
- 7: solve $\tilde{C} = \operatorname{argmin}_{|C|=K} WSS_y(C)$ where $WSS_y(C)$ is the within cluster sum of squares on Y .
- 8: on Y .
- 9: Step (b):
- 10: For a given \tilde{C} , maximize it with respect to w : set $D = [BSS_1(\tilde{C}), \dots, BSS_p(\tilde{C})]$ and solve:

$$\tilde{w} = \max_w w^T D$$

$$\text{such that } \begin{cases} \|w\|_1 \leq l_1 \\ \|w\|_2^2 \leq 1 \\ w_j \geq 0 \end{cases}$$

- Use a parallel spherical K-Means
- **PARALLEL SPARSE KMEANS-ALGORITHM** is affected by a even few outliers. Implement a parallel robustification of it.