# Reimplementing and Improving in RaceID

Pei Lin, 66672977

School of Information Science and Technology
ShanghaiTech University

May 14, 2019

# Outline

# Outline

# Introduction

Clustering analysis has been widely applied to single cell RNA-sequencing data to discover cell types and cell states.

RaceID is published in Nature 2014 using the data from mouse intestinal which helps to discover the rare cell type.

# Introduction to RaceID

RaceID is programed by R and the workflow is summarized below:

RaceID → Data filter —> Gap Statistic to find the best K—> Kmeans( in different distance metric)--> silhouette score to evaluate

Dimension Reduction → PCA
Dimension Reduction → tSNE

# Introduction to My Work

What I have done:

- **Reimplementing by Python.**
- **Improvement by other papers.**
- **Parallel programing.**
- **Work in another data set.**

# Outline

# Data Type

From:

1.Single-cell messenger RNA sequencing reveals rare intestinal cell types. (origin dataset)

2.Single-cell RNA-Seq profiling of human preimplantation embryos and embryonic stem cells.

# Data Filtering & scaling

For genes: expressed in most cells(housekeeping gene)

For cells: total expression is too low.

For scaling:

$$\frac{GeneExpression}{Mean} \rightarrow \log(GeneExpression)$$

# Outline

# Kmeans

For Kmeans algorithm is not a convex problem, one question is how to find K and choose the initial point.
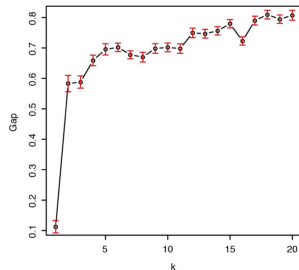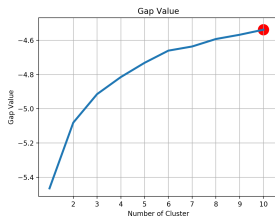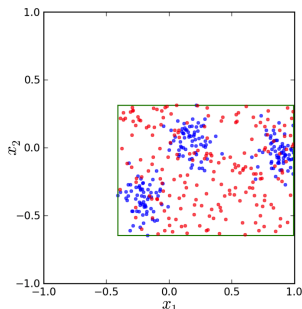
K: Calinski-Harabasz method, gap statistic

Initial point: kmeans++ method (supported by python package: sklearn.cluster)

# Gap Statistic

Main idea is Monte Carlo method,

$$Gap_n(k) = E(logW^*) - logW$$



Because K-means is unsupervised learning, gap statistic is just a reference.

# Outline

# PCA

$$X = U\Sigma V^T$$

PCA works for linear data.

# tSNE

$$P_{j|i} = \frac{\exp(-||x_i - x_j||^2/2\sigma_i{}^2)}{\sum_k \exp(-||x_i - x_k||^2/2\sigma_i{}^2)}$$

$$q_{j|i} = \frac{\exp(-||y_i - y_j||^2)}{\sum_k \exp(-||y_i - y_k||^2)}$$

$$C = \sum KL(P_i||Q_i)$$

tSNE usd manifold model and keep the characteristics of neighboring domain.

# tSNE

# Outline

# Multi-thread
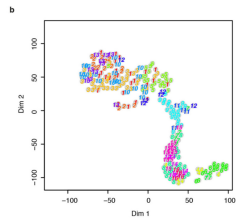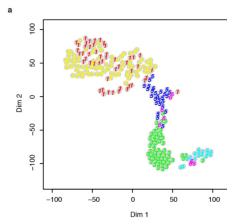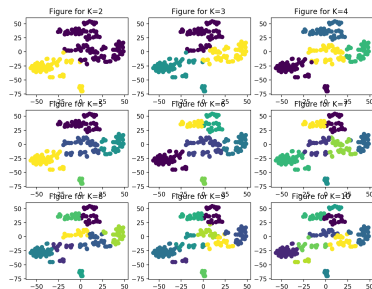
Map reduce, Spark, Openmp are popular in parallel programing.
Use 10 threads to accelerating the clustering

$$MeanSpeedUp = 49\%$$

Amdahl's Law

```
no parallel time has passed: 18.286089
paralleled time has passed: 8.904707000000002
```

In R, finish the same work needs 14.2576s

# Outline

# Reference

[1]Single-cell messenger RNA sequencing reveals rare intestinal cell types, Dominic Gru, Anna Lyubimova, Lennart Kester
[2]Single-cell RNA-Seq profiling of human preimplantation embryos and embryonic stem cells, Liying Yan, Mingyu Yang, Hongshan Guo, Lu Yang
[3]Estimating the number of clusters in a data set via gap statistic, Robert Tibshirani, Guenther Weather
[4]Visualizing data using t-SNE, Laurens van der Maaten, Geoffery Hinton

Thank You.