

# **STA 426: Diffusion maps for high-dimensional single-cell analysis of differentiation data**

---

Giuachin Kreiliger, Samuel Pawel (GitHub: Kreile, SamCH93)

19.11.2018

# Introduction

---

## Technical Details

---

## **Comparison with Other Methods**

---

## Principal Component Analysis (PCA)

- Orthogonal transformation of the data
- Works best for linear data subspace
- Non-linear extensions (Kernel-PCA)

# Dimension Reduction Methods for Single-Cell Omics Data

## Principal Component Analysis (PCA)

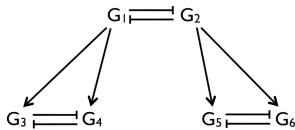
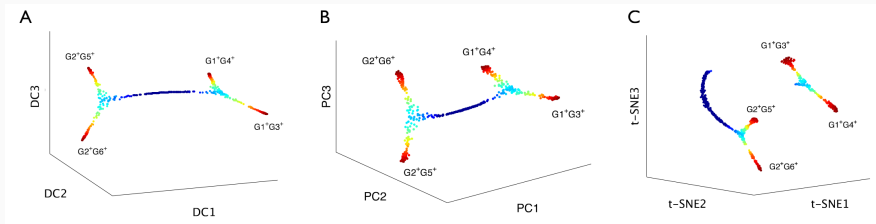
- Orthogonal transformation of the data
- Works best for linear data subspace
- Non-linear extensions (Kernel-PCA)

## t-Distributed Stochastic Neighbor Embedding (t-SNE)

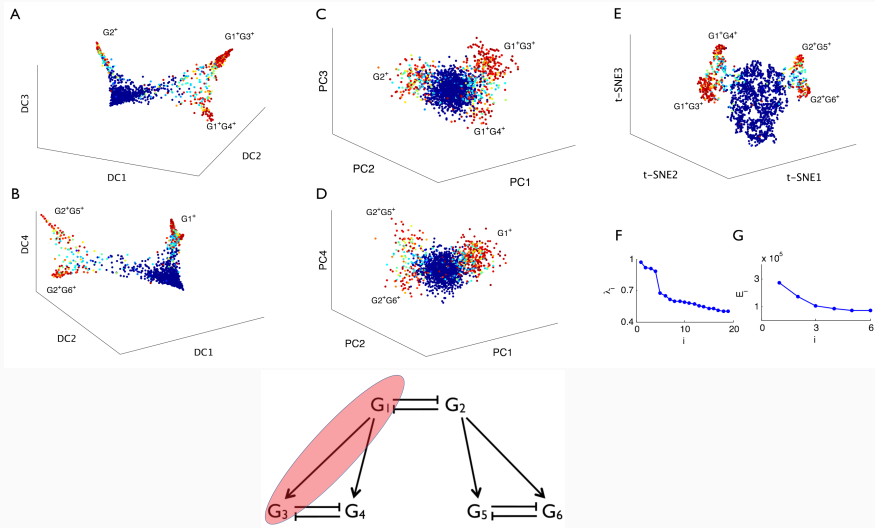
- Probabilistic, non-linear method
- Very robust to noise / density heterogeneities
- Tuning parameter perplexity

Many more methods exist!

# Simulated Data – Balanced

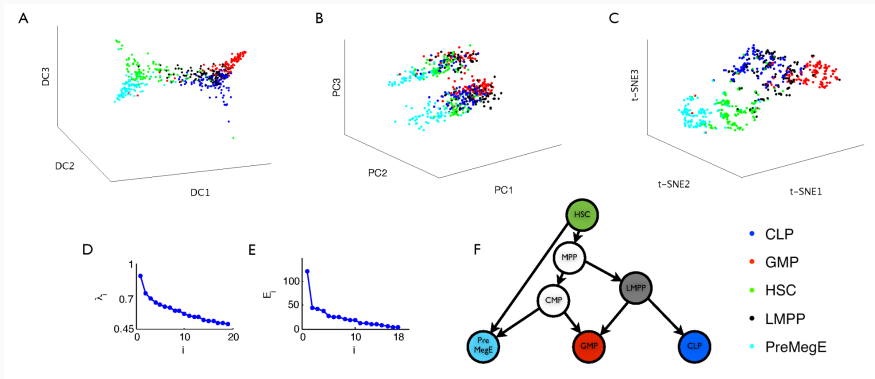


# Simulated Data – Imbalanced + Noise





# qPCR Data – Haematopoietic Stem Cells



# Results

- + Diffusion maps outperform other methods in finding cell differentiation trajectories
- + Diffusion maps are robust to sampling density heterogeneities and noise

# Results

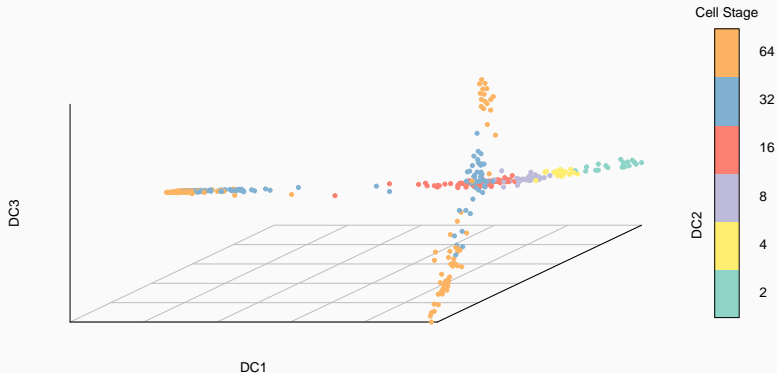
- + Diffusion maps outperform other methods in finding cell differentiation trajectories
- + Diffusion maps are robust to sampling density heterogeneities and noise
  - In all examples  $\sigma^2$  was finetuned
  - Unclear if finetuning perplexity parameter in *t*-SNE method would lead to similar results
  - Number of diffusion components needed for visualization not known in advance

# R Demonstration

---

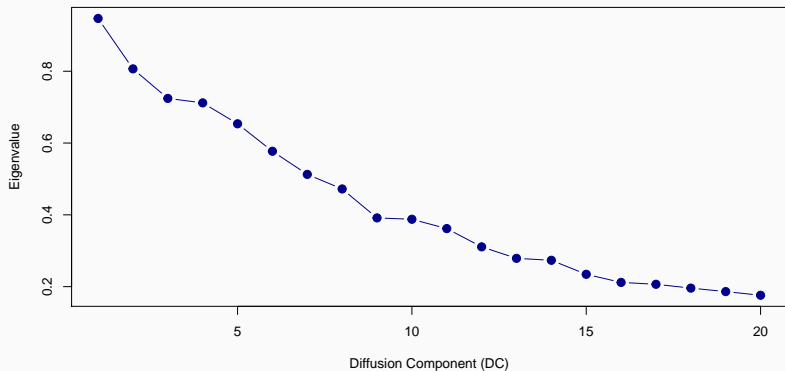
# destiny Package

```
set.seed(12)
library(destiny)
data("guo_norm")
dm_guo <- DiffusionMap(guo_norm)
plot(dm_guo, pch = 20, col_by = "num_cells", legend_main = "Cell Stage")
```



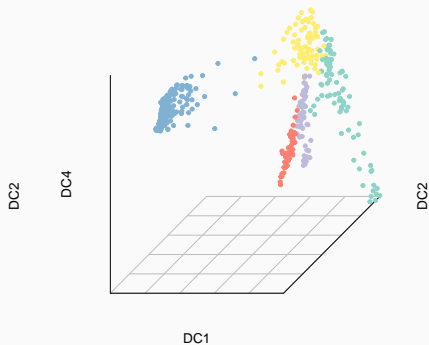
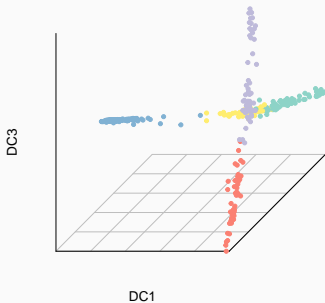
# Eigenvalues of Diffusion Components

```
plot(eigenvalues(dm_guo), xlab = "Diffusion Component (DC)", ylab = "Eigenvalue",  
     pch = 20, col = "darkblue", type = "b", cex = 2)
```



# Plotting Several Diffusion Components

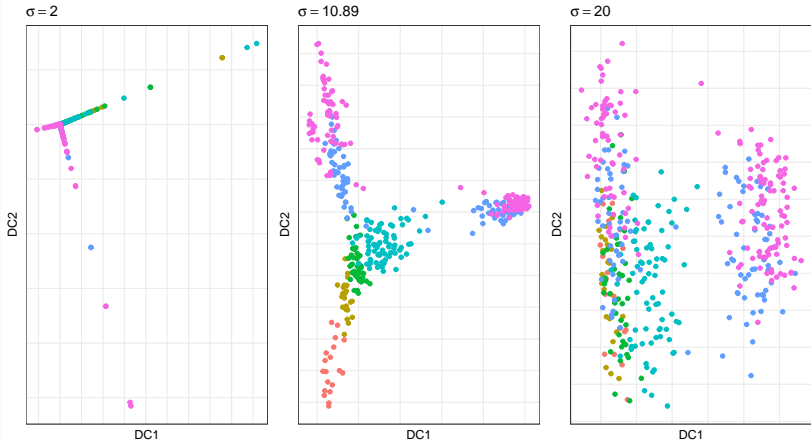
```
par(mfrow = c(1, 2))  
plot(dm_guo, 1:3, pch = 20, draw_legend = FALSE)  
plot(dm_guo, c(1, 2, 4), pch = 20, draw_legend = FALSE)
```



# Choice of $\sigma^2$

```
sigmas <- find_sigmas(guo_norm, verbose = FALSE)
optimal_sigma(sigmas)
```

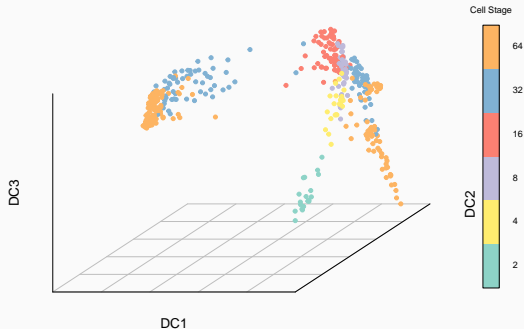
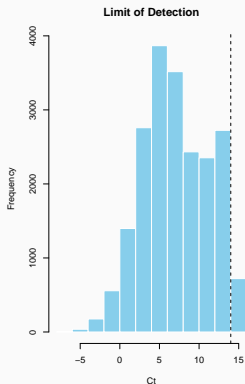
```
## [1] 10.8946
```





# Censored/Missing Values

```
layout(matrix(c(1, 2, 2), ncol = 3))
exprs(guo_missing)[sample(length(exprs(guo_missing)), 100)] <- NA
dm_guo_missing <- DiffusionMap(guo_missing, verbose = FALSE,
                              censor_val = 14, censor_range = c(14, 40),
                              missing_range = c(-5, 10))
hist(exprs(guo_norm), main = "Limit of Detection", xlab = "Ct", col = "skyblue", border = "white")
abline(v = 14, lty = 2)
plot(dm_guo_missing, pch = 20, cex.symbols = 1.5, col_by = "num_cells", legend_main = "Cell Stage")
```



- + Handle high noise levels, missing data, sampling heterogenitites
- + Diffusion distance is a biologically relevant distance metric
- + Capture nonlinear/complex differentiation dynamics

- + Handle high noise levels, missing data, sampling heterogenitites
- + Diffusion distance is a biologically relevant distance metric
- + Capture nonlinear/complex differentiation dynamics
  - Number of significant dimension not determinable in advance
  - Finetuning of  $\sigma^2$  required (but there is a proposed criterion)
  - $n^2 \times G$  computation time (but there are approximate versions)

## Discussion

- + Handle high noise levels, missing data, sampling heterogenitites
  - + Diffusion distance is a biologically relevant distance metric
  - + Capture nonlinear/complex differentiation dynamics
    - Number of significant dimension not determinable in advance
    - Finetuning of  $\sigma^2$  required (but there is a proposed criterion)
    - $n^2 \times G$  computation time (but there are approximate versions)
- Powerful dimension reduction tool for single cell differentiation data

- Angerer, P., Haghverdi, L., Büttner, M., Theis, F., Marr, C., and Büttner, F. (2015). *destiny*: diffusion maps for large-scale single-cell data in R. *Bioinformatics*, 32(8):1241–1243.
- Coifman, R. R. and Lafon, S. (2006). Diffusion maps. *Applied and Computational Harmonic Analysis*, 21:5–30.
- Haghverdi, L., Büttner, F., and Theis, F. J. (2015). Diffusion maps for high-dimensional single-cell analysis of differentiation data. *Bioinformatics*, 31(18):2989–2998.