

Data analysis and Unsupervised Learning

Dimensionality Reduction: Beyond PCA and Non Linear Methods

MAP 573, 2020 – Julien Chiquet

École Polytechnique, Autumn semester, 2020

<https://jchiquet.github.io/MAP573>



Outline

Introduction

Motivations

Part I

Introduction

Packages required for reproducing the slides

```
library(tidyverse) # opinionated collection of packages for data manipulation
library(GGally)    # extension to ggplot vizualization system
library(FactoMineR) # PCA and oter linear method for dimension reduction
library(factoextra) # fancy plotting for FactoMineR output
# color and plots themes
library(RColorBrewer)
pal <- brewer.pal(10, "Set3")
theme_set(theme_bw())
```

Companion data set: 'scRNA'

Subsamples of normalized Single-Cell RNAseq

Description: *subsample of a large data set*

Gene-level expression of 100 representative genes for a collection of 301 cells spreaded in 11 cell-lines. Original transcription data are measured by counts obtained by *RNAseq* and normalized to be close to Gaussian.



Pollen, Alex A., et al. Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex.

Nature biotechnology 32.10 (2014): 1053.

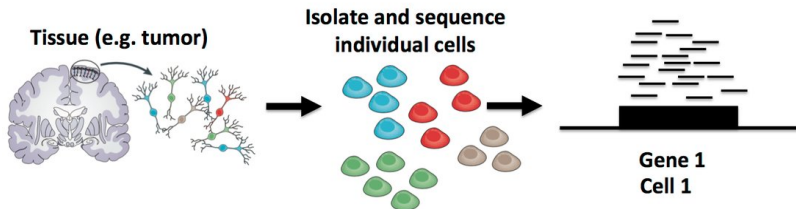


Figure: Single Cell RNA sequencing data: general principle – source: Stephanie Hicks

Companion data set: 'scRNA'

Brief data summary I

Data manipulation

```
load("../..//data/scRNA.RData")
scRNA <- pollen$data %>% t() %>% as_tibble() %>%
  add_column(cell_type = pollen$celltypes)
```

Cell types

```
scRNA %>% dplyr::select(cell_type) %>% summary() %>% knitr::kable()
```

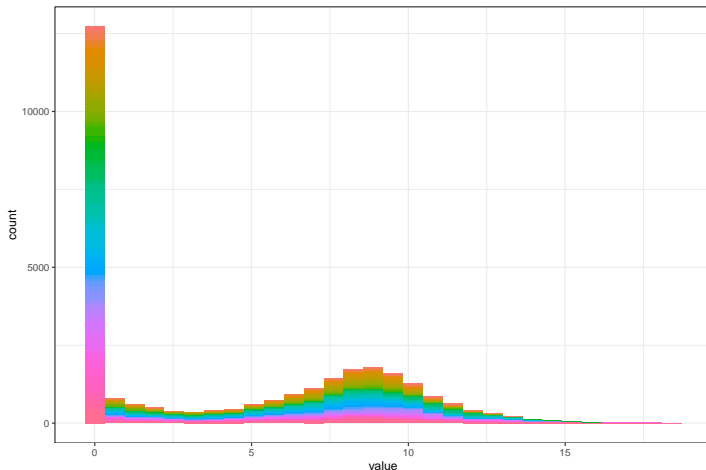
	cell_type
	HL60 :54
	K562 :42
	Kera :40
	BJ :37
	GW16 :26
	hiPSC :24
	(Other):78

Companion data set II: 'scRNA'

Brief data summary II

Histogram of normalized expression

```
scRNA %>% dplyr::select(-cell_type) %>% pivot_longer(everything()) %>%  
  ggplot() + aes(x = value, fill = name) + geom_histogram(show.legend = FALSE)
```



PCA

PCA – Biplot



PCA (and linear methods) limitations

Account for complex pattern

- Linear methods are powerful for **planar structures**
- May fail at describing **manifolds**

Preserve local geometry

- High dimensional data are characterized by **multiscale properties** (local / global structures)
- Non Linear projection helps at preserving **local characteristics** of distances

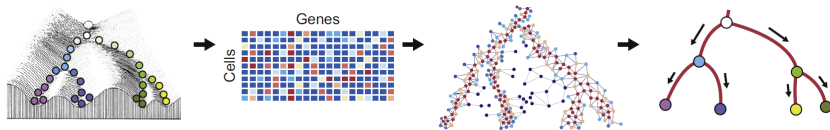


Figure: Intuition of manifolds and geometry underlying sc-data — source: F. Picard

Companion data set II: 'mollusk'

Abundance table (Species counts spread in various sites)

Description: *small size count data*

Abundance of 32 mollusk species in 163 samples. For each sample, 4 additional covariates are known.



Richardot-Coulet, M., Chessel D. and Bournaud M. Typological value of the benthos of old beds of a large river. Methodological approach. Archiv für Hydrobiologie, 107.

```
library(PLNmodels); data(mollusk)
mollusk <-
  prepare_data(mollusk$Abundance, mollusk$Covariate[c("season", "site")]) %>%
  as_tibble() %>%
  distinct() # remove duplicates
```

External Covariates

```
mollusk %>% dplyr::select(site, season) %>% summary() %>% t() %>% knitr::kable()
```

site	Negria1 :24	Negria2 :24	Pecheurs1:24	Pecheurs2:23	GGravier1:21	GGravier2:21
season	autumn:41	spring:43	summer:44	winter:30	NA	NA

Companion data set: 'mollusk'

Brief data summary II

Histogram of raw counts

```
mollusk %>% dplyr::select(-site, -season) %>%  
  pivot_longer(everything()) %>%  
  ggplot() + aes(x = value, fill = name) + geom_histogram(show.legend = FALSE)  
  
## Error: Aesthetics must be either length 1 or the same as the data (316): x
```

Companion data set: 'mollusk'

PCA

```
mollusk %>% PCA(graph = FALSE, quali.sup = which(map_lgl(mollusk, is.factor))) %>%  
  fviz_pca_biplot(select.var = list(contrib = 5), habillage = "site")  
  
## Error in dimnames(x) <- dn: length of 'dimnames' [2] not equal to array  
extent
```

PCA (and linear methods) limitations

Account for complex data distribution

- Linear methods /PCA are tied to an hidden **Gaussian assumption**
- Fail with **Count data**
- Fail with **Skew data**

Possible solutions

- Probabilistic (non Gaussian) models
- Need transformed (non-linear) input space

Dimension reduction: revisiting the problem setup

Settings

- **Training data** : $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \in \mathbb{R}^p$, (i.i.d.)
- Space \mathbb{R}^p of possibly high dimension ($n \ll p$)

Dimension Reduction Map

Construct a map Φ from the space \mathbb{R}^p into a space \mathbb{R}^q of **smaller dimension**:

$$\begin{aligned}\Phi : \quad \mathbb{R}^p &\rightarrow \mathbb{R}^q, q \ll p \\ \mathbf{x} &\mapsto \Phi(\mathbf{x})\end{aligned}$$

How should we design/construct Φ ?

Criterion

- Geometrical approach (**see slides on PCA**)
- Reconstruction error
- Relationship preservation

Form of the map Φ

- **Linear** or **non-linear** ?
- tradeoff between interpretability and **versatility** ?
- tradeoff between **high** or low computational resource

Part II

Non-linear methods

Outline

Non-linear methods

① Motivated by reconstruction error

General goal

Kernel-PCA

Non-negative matrix factorization

Auto-Encoder

② Relation preservation

Outline

Non-linear methods

① Motivated by reconstruction error

- General goal

- Kernel-PCA

- Non-negative matrix factorization

- Auto-Encoder

② Relation preservation

Reconstruction error approach

- 1 Construct a map Φ from the space \mathbb{R}^p into a space \mathbb{R}^q of **smaller dimension**:

$$\begin{aligned}\Phi : \quad \mathbb{R}^p &\rightarrow \mathbb{R}^q, q \ll p \\ \mathbf{x} &\mapsto \Phi(\mathbf{x})\end{aligned}$$

- 2 Construct $\tilde{\Phi}$ from \mathbb{R}^q to \mathbb{R}^p (**reconstruction formula**)
- 3 Control an error ϵ between \mathbf{x} and its reconstruction $\hat{\mathbf{x}} = \tilde{\Phi}(\Phi(\mathbf{x}))$

For instance, the error measured with the Frobenius between the original data matrix \mathbf{X} and its approximation:

$$\epsilon(\mathbf{X}, \hat{\mathbf{X}}) = \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_F^2 = \sum_{i=1}^n \left\| \mathbf{x}_i - \tilde{\Phi}(\Phi(\mathbf{x}_i)) \right\|^2$$

Reconstruction error approach

- 1 Construct a map Φ from the space \mathbb{R}^p into a space \mathbb{R}^q of **smaller dimension**:

$$\begin{aligned}\Phi : \quad \mathbb{R}^p &\rightarrow \mathbb{R}^q, q \ll p \\ \mathbf{x} &\mapsto \Phi(\mathbf{x})\end{aligned}$$

- 2 Construct $\tilde{\Phi}$ from \mathbb{R}^q to \mathbb{R}^p (**reconstruction formula**)
- 3 Control an error ϵ between \mathbf{x} and its reconstruction $\hat{\mathbf{x}} = \tilde{\Phi}(\Phi(\mathbf{x}))$

For instance, the error measured with the Frobenius between the original data matrix \mathbf{X} and its approximation:

$$\epsilon(\mathbf{X}, \hat{\mathbf{X}}) = \left\| \mathbf{X} - \hat{\mathbf{X}} \right\|_F^2 = \sum_{i=1}^n \left\| \mathbf{x}_i - \tilde{\Phi}(\Phi(\mathbf{x}_i)) \right\|^2$$

Reinterpretation of PCA

PCA model

Let \mathbf{V} be a $p \times q$ matrix whose columns are of q orthonormal vectors.

$$\begin{aligned}\Phi(\mathbf{x}) &= \mathbf{V}^\top (\mathbf{x} - \boldsymbol{\mu}) = \tilde{\mathbf{x}} \\ \mathbf{x} &\simeq \tilde{\Phi}(\tilde{\mathbf{x}}) = \boldsymbol{\mu} + \mathbf{V}\tilde{\mathbf{x}}\end{aligned}$$

↪ Model with **Linear assumption + ortho-normality constraints**

PCA reconstruction error

$$\underset{\boldsymbol{\mu} \in \mathbb{R}^p, \mathbf{V} \in \mathcal{O}_{p,q}}{\text{minimize}} \sum_{i=1}^n \left\| (\mathbf{x}_i - \boldsymbol{\mu}) + \mathbf{V}^\top \mathbf{V} (\mathbf{x}_i - \boldsymbol{\mu}) \right\|^2$$

Solution (explicit)

- $\boldsymbol{\mu} = \bar{\mathbf{x}}$ the empirical mean
- \mathbf{V} an orthonormal basis of the space spanned by the q first eigenvectors of the empirical covariance matrix

Reinterpretation of PCA

PCA model

Let \mathbf{V} be a $p \times q$ matrix whose columns are of q orthonormal vectors.

$$\begin{aligned}\Phi(\mathbf{x}) &= \mathbf{V}^\top (\mathbf{x} - \boldsymbol{\mu}) = \tilde{\mathbf{x}} \\ \mathbf{x} &\simeq \tilde{\Phi}(\tilde{\mathbf{x}}) = \boldsymbol{\mu} + \mathbf{V}\tilde{\mathbf{x}}\end{aligned}$$

↪ Model with **Linear assumption + ortho-normality constraints**

PCA reconstruction error

$$\underset{\boldsymbol{\mu} \in \mathbb{R}^p, \mathbf{V} \in \mathcal{O}_{p,q}}{\text{minimize}} \sum_{i=1}^n \left\| (\mathbf{x}_i - \boldsymbol{\mu}) + \mathbf{V}^\top \mathbf{V} (\mathbf{x}_i - \boldsymbol{\mu}) \right\|^2$$

Solution (explicit)

- $\boldsymbol{\mu} = \bar{\mathbf{x}}$ the empirical mean
- \mathbf{V} an orthonormal basis of the space spanned by the q first eigenvectors of the empirical covariance matrix

Important digression: SVD

Singular Value Decomposition (SVD)

The SVD of \mathbf{M} a $n \times p$ matrix is the factorization given by

$$\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^\top,$$

where $r = \min(n, p)$ and

- $\mathbf{D}_{r \times r} = \text{diag}(\delta_1, \dots, \delta_r)$ is the diagonal matrix of singular values.
- \mathbf{U} is orthonormal, whose columns are eigen vectors of $(\mathbf{M}\mathbf{M}^\top)$
- \mathbf{V} is orthonormal whose columns are eigen vectors of $(\mathbf{M}^\top\mathbf{M})$

→ Time complexity in $\mathcal{O}(npqr)$ (less when $k \ll r$ components are required)

Connection with eigen decomposition of the covariance matrix

$$\begin{aligned}\mathbf{M}^\top\mathbf{M} &= \mathbf{V}\mathbf{D}\mathbf{U}^\top\mathbf{U}\mathbf{D}\mathbf{V}^\top \\ &= \mathbf{V}\mathbf{D}^2\mathbf{V}^\top = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top\end{aligned}$$

Important digression: SVD

Singular Value Decomposition (SVD)

The SVD of \mathbf{M} a $n \times p$ matrix is the factorization given by

$$\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^\top,$$

where $r = \min(n, p)$ and

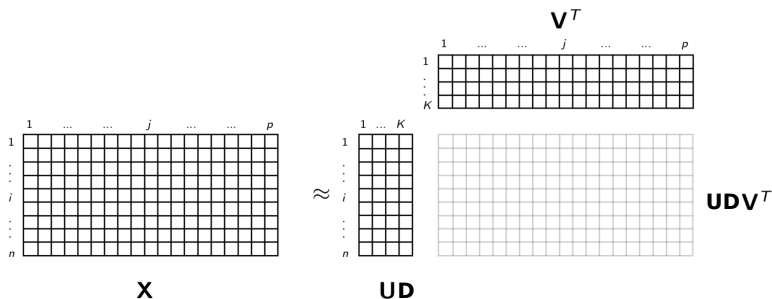
- $\mathbf{D}_{r \times r} = \text{diag}(\delta_1, \dots, \delta_r)$ is the diagonal matrix of singular values.
- \mathbf{U} is orthonormal, whose columns are eigen vectors of $(\mathbf{M}\mathbf{M}^\top)$
- \mathbf{V} is orthonormal whose columns are eigen vectors of $(\mathbf{M}^\top\mathbf{M})$

→ Time complexity in $\mathcal{O}(npqr)$ (less when $k \ll r$ components are required)

Connection with eigen decomposition of the covariance matrix

$$\begin{aligned}\mathbf{M}^\top\mathbf{M} &= \mathbf{V}\mathbf{D}\mathbf{U}^\top\mathbf{U}\mathbf{D}\mathbf{V}^\top \\ &= \mathbf{V}\mathbf{D}^2\mathbf{V}^\top = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top\end{aligned}$$

PCA solution is given by SVD of the centered data matrix



Since $\tilde{\mathbf{X}} = \mathbf{X}^c \mathbf{V} = \mathbf{U} \mathbf{D} \mathbf{V}^T \mathbf{V} = \mathbf{U} \mathbf{D}$, PCA can be rephrased as

$$\hat{\mathbf{X}}^c = \mathbf{F} \mathbf{V}^T = \arg \min_{\mathbf{F} \in \mathcal{M}_{n,q}, \mathbf{V} \in \mathcal{O}_{p,q}} \left\| \mathbf{X}^c - \mathbf{F} \mathbf{V}^T \right\|_F^2 \text{ with } \|\mathbf{A}\|_F^2 = \sum_{ij} a_{ij}^2,$$

$\left\{ \tilde{\mathbf{X}} \in \mathbb{R}^{n \times q}, \mathbf{V} \in \mathbb{R}^{p \times q} \right\}$ Best linear low-rank representation of \mathbf{X}

Outline

Non-linear methods

① Motivated by reconstruction error

General goal

Kernel-PCA

Non-negative matrix factorization

Auto-Encoder

② Relation preservation

Outline

Non-linear methods

① Motivated by reconstruction error

General goal

Kernel-PCA

Non-negative matrix factorization

Auto-Encoder

② Relation preservation

Non-negative Matrix Factorization – NMF

Setup

Assume that \mathbf{X} contains only non-negative entries (i.e. ≥ 0).

Model

Linear assumption + non-negativity constraints on both \mathbf{V} and $\tilde{\mathbf{x}}$

$$\begin{aligned}\Phi(\mathbf{x}) &= \mathbf{V}^\top (\mathbf{x} - \boldsymbol{\mu}) = \tilde{\mathbf{x}} \\ \mathbf{x} &\simeq \tilde{\Phi}(\tilde{\mathbf{x}}) = \boldsymbol{\mu} + \mathbf{V}\tilde{\mathbf{x}}\end{aligned}$$

For the whole data matrix \mathbf{X} ,

$$\hat{\mathbf{X}} = \mathbf{1}_n \boldsymbol{\mu}^\top + \underbrace{\tilde{\mathbf{X}}}_{\mathbf{F}, \text{ the factors}} \mathbf{V}^\top$$

NMF reconstruction errors

Build $\hat{\mathbf{X}} = \mathbf{F}\mathbf{V}^\top$ to minimize a distance $D(\hat{\mathbf{X}}, \mathbf{X})$! Several choice, e.g:

- Least-square loss (distance measured by Frobenius norm)

$$\hat{\mathbf{X}}^{\text{ls}} = \arg \min_{\substack{\mathbf{F} \in \mathcal{M}(\mathbb{R}_+)_{n,q} \\ \mathbf{V} \in \mathcal{M}(\mathbb{R}_+)_{p,q}}} \left\| \mathbf{X} - \mathbf{F}\mathbf{V}^\top \right\|_F^2,$$

- Kullback-Leibler divergence ("distance" between distribution)

$$\begin{aligned} \hat{\mathbf{X}}^{\text{kl}} &= \arg \min_{\substack{\mathbf{F} \in \mathcal{M}(\mathbb{R}_+)_{n,q} \\ \mathbf{V} \in \mathcal{M}(\mathbb{R}_+)_{p,q}}} \sum_{i,j} x_{ij} \log\left(\frac{x_{ij}}{(\mathbf{F}\mathbf{V}^\top)_{ij}}\right) + (\mathbf{F}\mathbf{V}^\top)_{ij} \\ &= \arg \max_{\substack{\mathbf{F} \in \mathcal{M}(\mathbb{R}_+)_{n,q} \\ \mathbf{V} \in \mathcal{M}(\mathbb{R}_+)_{p,q}}} \sum_{i,j} x_{ij} \log((\mathbf{F}\mathbf{V}^\top)_{ij}) - (\mathbf{F}\mathbf{V}^\top)_{ij}, \end{aligned}$$

↪ log-likelihood of a Poisson distribution with mean $(\mathbf{F}\mathbf{V}^\top)_{ij}$.

Example on 'mollusk' I

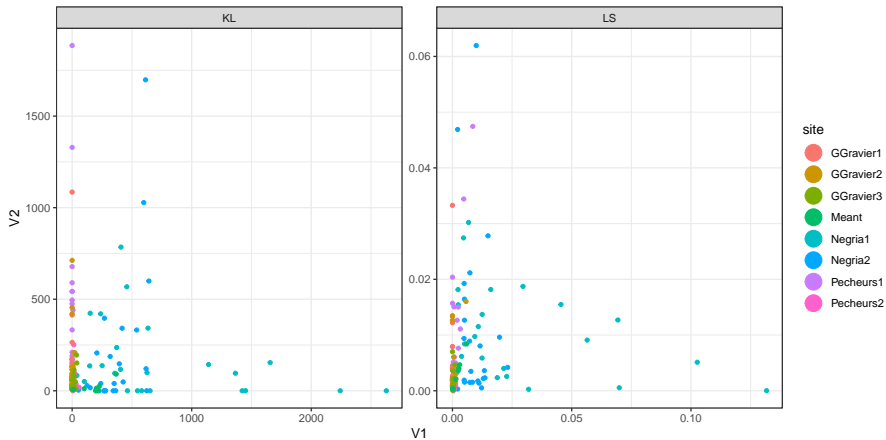
Run the fit

```
nmf_KL <- mollusk %>% select(-site, -season) %>%  
  nmf(rank = 2, method = 'brunet') %>% basis() %>%  
  as.data.frame() %>% add_column(algo = "KL") %>% add_column(site = mollusk$site)  
nmf_LS <- mollusk %>% select(-site, -season) %>%  
  nmf(rank = 2, method = 'lee') %>% basis() %>%  
  as.data.frame() %>% add_column(algo = "LS") %>% add_column(site = mollusk$site)
```

Compare algorithms

```
rbind(nmf_KL, nmf_LS) %>%  
  ggplot(aes(x = V1, y = V2, color = site)) +  
    geom_point(size=1.25) +  
    guides(colour = guide_legend(override.aes = list(size=6))) +  
    facet_wrap(~algo, scales = 'free')
```

Example on 'mollusk' II



Outline

Non-linear methods

① Motivated by reconstruction error

General goal

Kernel-PCA

Non-negative matrix factorization

Auto-Encoder

② Relation preservation

Outline

Non-linear methods

① Motivated by reconstruction error

② Relation preservation

- General goal

- MDS

- t-SNE

- UMAP

Outline

Non-linear methods

① Motivated by reconstruction error

② Relation preservation

General goal

MDS

t-SNE

UMAP

Pairwise Relation

Focus on pairwise relation $\mathcal{R}(\mathbf{x}_i, \mathbf{x}_{i'})$.

Distance Preservation

- Construct a map Φ from the space \mathbb{R}^d into a space $\mathbb{R}^{d'}$ of **smaller dimension**:

$$\begin{aligned}\Phi : \quad \mathbb{R}^d &\rightarrow \mathbb{R}^{d'}, d' \ll d \\ \mathbf{x} &\mapsto \Phi(\mathbf{x})\end{aligned}$$

$$\text{such that } \mathcal{R}(\mathbf{x}_i, \mathbf{x}_{i'}) \sim \mathcal{R}'(\mathbf{x}'_i, \mathbf{x}'_{i'})$$

Multidimensional scaling

Try to preserve inner product related to the distance (e.g. Euclidean)

t-SNE – Stochastic Neighborhood Embedding

Try to preserve relations with close neighbors with Gaussian kernel

Outline

Non-linear methods

① Motivated by reconstruction error

② Relation preservation

General goal

MDS

t-SNE

UMAP

Outline

Non-linear methods

① Motivated by reconstruction error

② Relation preservation

General goal

MDS

t-SNE

UMAP

Stochastic Neighbor Embedding

[van der Maaten and Hinton, 2008]

- (x_1, \dots, x_n) are the points in the high dimensional space \mathbb{R}^p ,
- Consider a similarity between points:

$$p_{i|j} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_k - x_j\|^2 / 2\sigma_k^2)}, \quad p_{ij} = (p_{i|j} + p_{j|i}) / 2N$$

- σ smooths the data (linked to the regularity of the target manifold)
- σ is chosen such that the entropy of p is fixed to a given value of the so-called perplexity

$$\exp \left(- \sum_{ij} p_{ij} \log(p_{ij}) \right)$$

The perplexity parameter

- σ_i Should adjust to local densities (neighborhood of point i)
- Define the Shannon entropy of $p_i = (p_{1|i}, \dots, p_{n|i})$

$$H(p_i) = - \sum_{j=1}^n p_{j|i} \log_2 p_{j|i}$$

- The perplexity is defined by:

$$Perp(p_i) = 2^{H(p_i)}$$

- Interpreted as the smoothed effective number of neighbors.
- SNE performs a binary search for the value of σ_i that produces a p_i with a fixed perplexity that is specified by the user.

tSNE and Student / Cauchy kernels

- Consider (y_1, \dots, y_n) are points in the low dimensional space \mathbb{R}^2
- Consider a similarity between points in the new representation:

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

- Robustify this kernel by using Student(1) kernels (ie Cauchy)

$$q_{i|j} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_k - y_j\|^2)^{-1}}$$

Optimizing tSNE

- Minimize the KL between p and q so that the data representation minimizes:

$$C(y) = \sum_{ij} KL(p_{ij}, q_{ij})$$

- The cost function is not convex

$$\left[\frac{\partial C(y)}{\partial y} \right]_i = \sum_j (p_{ij} - q_{ij})(y_i - y_j)$$

- Interpreted as the resultant force created by a set of springs between the map point y_i and all other map points $(y_j)_j$. All springs exert a force along the direction $(y_i - y_j)$.
- $(p_{ij} - q_{ij})$ is viewed as a stiffness of the force exerted by the spring between y_i and y_j .

Customed Gradient descent

- Gradient descent initialized by sampling map points randomly from an isotropic Gaussian with small variance centered around the origin
- Gradient update using

$$y^{(t)} = y^{(t-1)} + \eta \frac{\partial C(y)}{\partial y} + \alpha(t)(y^{(t-1)} - y^{(t-2)})$$

- η learning rate, $\alpha(t)$ momentum at iteration t .
- Gaussian noise is added to the map points to perform simulated annealing.

Properties of t-SNE

- good at preserving local distances (intra-cluster variance)
- not so good for global representation (inter-cluster variance)
- hence good at creating clusters of points that are close, but bad at positioning clusters wrt each other
- preprocessing very important : initialize with PCA and feature selection plus log transform (non linear transform)
- percent of explained variance ? interpretation of the q distribution ?

Example on scRNA I

Run the fit

```
scRNA_expr <- scRNA %>% select(-cell_type) %>% as.matrix()

tSNE_perp2 <- Rtsne(scRNA_expr, perplexity = 2)$Y %>%
  as.data.frame() %>% add_column(perplexity = 2) %>% add_column(cell_type = scRNA$cell_type)

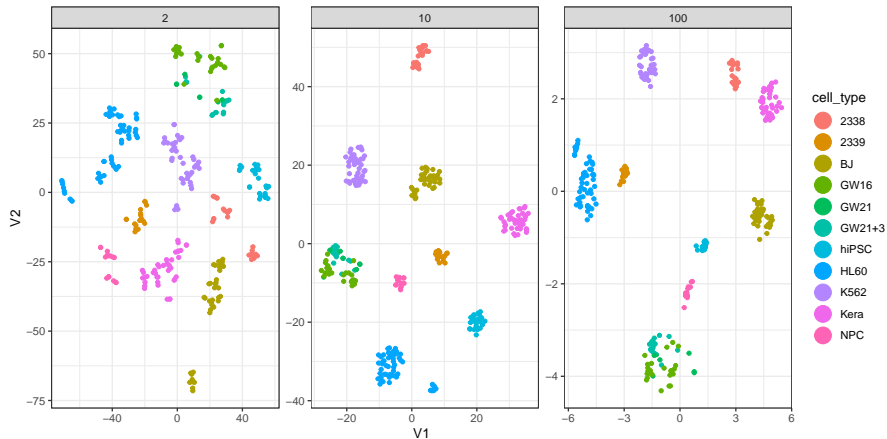
tSNE_perp10 <- Rtsne(scRNA_expr, perplexity = 10)$Y %>%
  as.data.frame() %>% add_column(perplexity = 10) %>% add_column(cell_type = scRNA$cell_type)

tSNE_perp100 <- Rtsne(scRNA_expr, perplexity = 100)$Y %>%
  as.data.frame() %>% add_column(perplexity = 100) %>% add_column(cell_type = scRNA$cell_type)
```

Compare perplexity

```
rbind(tSNE_perp2, tSNE_perp10, tSNE_perp100) %>%
  ggplot(aes(x = V1, y = V2, color = cell_type)) +
    geom_point(size=1.25) +
    guides(colour = guide_legend(override.aes = list(size=6))) +
    facet_wrap(~perplexity, scales = 'free')
```

Example on scRNA II



Example on 'mollusk' I

Run the fit

```
mollusk_ab <- mollusk %>% select(-site, -season) %>% as.matrix()

tSNE_perp2 <- Rtsne(mollusk_ab, perplexity = 2)$Y %>%
  as.data.frame() %>% add_column(perplexity = 2) %>% add_column(site = mollusk_ab$site)

## Error in mollusk_ab$site: $ operator is invalid for atomic vectors

tSNE_perp10 <- Rtsne(log(1 + mollusk_ab), perplexity = 10)$Y %>%
  as.data.frame() %>% add_column(perplexity = 10) %>% add_column(site = mollusk_ab$site)

## Error in mollusk_ab$site: $ operator is invalid for atomic vectors

tSNE_perp50 <- Rtsne(log(1 + mollusk_ab), perplexity = 50)$Y %>%
  as.data.frame() %>% add_column(perplexity = 50) %>% add_column(site = mollusk_ab$site)

## Error in mollusk_ab$site: $ operator is invalid for atomic vectors
```

Compare perplexity

Example on 'mollusk' II

```
rbind(tSNE_perp2,tSNE_perp10,tSNE_perp50) %>%  
  ggplot(aes(x = V1, y = V2, color = site)) +  
    geom_point(size=1.25) +  
    guides(colour = guide_legend(override.aes = list(size=6))) +  
    facet_wrap(~perplexity, scales = 'free')  
  
## Error in eval(quote(list(...)), env): object 'tSNE_perp50' not found
```

Outline

Non-linear methods

① Motivated by reconstruction error

② Relation preservation

General goal

MDS

t-SNE

UMAP

Uniform Manifold Approximation and Projection

[McInnes et al., 2018]

Properties of UMAP



Example on scRNA I

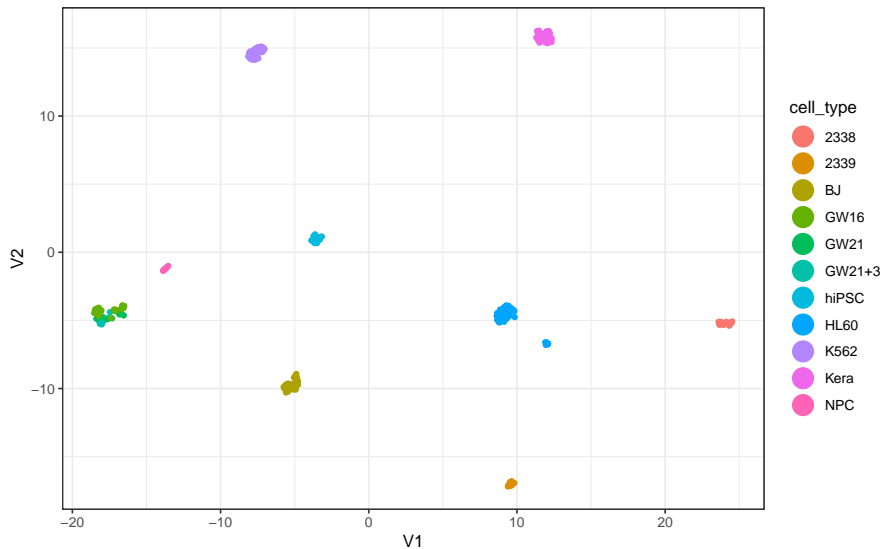
Run the fit

```
scRNA_expr <- scRNA %>% select(-cell_type) %>% as.matrix()
umap_fit <- umap(scRNA_expr)$layout %>%
  as.data.frame() %>% add_column(cell_type = scRNA$cell_type)
```

Visualization

```
umap_fit %>%
  ggplot(aes(x = V1, y = V2, color = cell_type)) +
  geom_point(size=1.25) +
  guides(colour = guide_legend(override.aes = list(size=6)))
```

Example on scRNA II



Example on 'mollusk' I

Run the fit

```
duplicated <- duplicated(mollusk %>% select(-site, -season))
mollusk_ab <- mollusk %>% select(-site, -season) %>% filter(!duplicated) %>% as.ma

## Error: Problem with 'filter()' input '..1'.
## x Input '..1' must be of size 158 or 1, not size 5056.
## i Input '..1' is '!duplicated'.
```



```
umap_fit <- umap(mollusk_ab)$layout %>%
  as.data.frame() %>% add_column(site = mollusk$site[!duplicated])

## Error: New columns must be compatible with '.data'.
## x New column has 575 rows.
## i '.data' has 158 rows.
```

Visualization

Example on 'mollusk' II

```
umap_fit %>%  
  ggplot(aes(x = V1, y = V2, color = site)) +  
    geom_point(size=1.25) +  
    guides(colour = guide_legend(override.aes = list(size=6)))  
  
## Error in FUN(X[[i]], ...): object 'site' not found
```

Example on 'mollusk' III

References I



McInnes, L., Healy, J., and Melville, J. (2018).

Umap: Uniform manifold approximation and projection for dimension reduction.

arXiv preprint arXiv:1802.03426.



van der Maaten, L. and Hinton, G. (2008).

Visualizing Data using t-SNE.

Journal of Machine Learning Research, 9(Nov):2579–2605.