

# UMAP

富田 直輝

2019 年 5 月 17 日

## 1 はじめに

次元圧縮はデータの情報を削減するのと同じことになっている。手法としては統計的手法である主成分分析や多様体学習である ISOMAP や t-SNE 等が挙げられている。特に t-SNE は kaggle でもよく使われている手法になっておりニューラルネットワークの初期値として与えられたりするなど汎用性は高い。しかし、今回は t-SNE より実行速度が速く、また kaggle でも注目されつつある UMAP について紹介したいと思う (前提知識が多いので手法の紹介だけになるが)。

## 2 データの多様体への変換と近似

データを多様体の空間として考えたい。多様体の空間は  $\mathbb{R}^n$  として考える。そこでデータは式 (1) で表されることにする。

$$X = X_1, \dots, X_N \tag{1}$$

## 3 空間表現

局所的な空間を表現したい。

定義 1 圏有限群の  $sets[n] = 1, \dots, n$  を持つ。

定義 2 単体の集合は定義 1 によって与えられるマップからの関手である。

単体の集合は多様体空間の学習の組み合わせ方法で説明できる。対照的に距離空間を処理し距離情報を持ちえる同様な構造を必要としている。しかし、これを表現する方法はすでに Spivak によって証明されている。

ここで  $I$  は  $(0, a); a \in (0, 1]$  の間隔によって与えられる位相で  $(0, 1] : 0 < a \leq 1 \mathbb{R}$  の間隔である。オープンセットの圏はいくつかの順序集合のカテゴリーにとっての自然な方法でグロダンティック位相になじませることができる。

定義 3  $I$  における前層  $\mathcal{P}$  は集合への関手  $I^{op}$  である。ファジー集合は  $I$  の層であるので全てのマップ  $\mathcal{P}(ab)$  は単射である。

不完全な単調写像からまた層における限は定義されているが、直接限は定義されていない。そこで次の定義

を行う.

**定義 4** ファジー集合における集合のファジーは  $I$  における層の部分集合である

また、ファジーの単体の集合の定義は次のように定義される

**定義 5** ファジーの単体の集合:  $sFuzz$  は  $\Delta^{\text{op}}$  から  $Fuzz$  への関手によって与えられるオブジェクトと、自然変換によって与えられる射の集合である.

ファジーの単体集合は密着空間として与えられ、そして  $\Delta \times I$  は積位相を持つ層としてみられる.

**定義 6** 擬距離空間 (extended-pseudo-metric space) は次のように定義される.

1.  $d(x, y) \geq 0, d(x, x) = 0;$
2.  $d(x, y) = d(y, x);$
3.  $d(x, z) \leq d(x, y) + d(y, z) \text{ or } d(x, z) = \infty$

これらの擬空間距離の集合を  $EPMet$  と呼び擬空間距離と収縮写像をもつ. また、有限擬空間距離を  $FinEPMet$  と呼ぶ.

また、Spivak は随伴関手として  $sFuzz$  と  $EPMet$  の集合の間である  $Real$  と  $Sing$  を定義した.  $Real$  は式 2 で定義される.

$$Real(\Delta_a^n) = {}^t rtriangle\{(t_0, \dots, t_n) \in \mathbb{R}^{n+1} \mid \sum_{i=0}^n t_i = -\log(a), t_i \geq 0\} \quad (2)$$

また、 $a \leq b$  なら、位相空間内の  $\Delta_a^n \rightarrow \Delta_b^n$  と表され.  $\Delta$  morphism  $\sigma: [n] \rightarrow [m]$  として表される.

$$(x_0, x_1, \dots, x_n) \mapsto \frac{\log(b)}{\log(a)} \left( \sum_{i_0 \in \sigma^{-1}(0)} \sum_{i_1 \in \sigma^{-1}(1)}, \dots, \sum_{i_m \in \sigma^{-1}(m)} \right) \quad (3)$$

この写像は  $0 \leq a \leq b$  における  $\log(b)/\log(a) \leq 1$  であるときは非拡大的である.

$$Real(X) = \Delta \text{ colim}_{\Delta_{<a}^n \rightarrow X} Real(\Delta_{<a}^n) \quad (4)$$

集合の  $\Delta \times I$  における擬空間距離の随伴関数を  $Sing$  とする.

$$Sing(Y) : ([n], [0, a] \mapsto \text{hom}_{EPMet}(Real(\Delta_{<a}^n))) \quad (5)$$

また有限距離について、 $FinFuzz$  と  $FinReal$ ,  $FinSing$  を定義する.

**定義 7** 関手  $FinReal$  の定義:  $Fin-sFuzz \rightarrow FinEPMet$

$$FinReal(\Delta_{<a}^n) = (x_1, x_2, \dots, x_n, d_n) \quad (6)$$

このとき,

$$d_a(x_i, x_j) = \begin{cases} -\log(a) & i \neq j \\ 0 & otherwise \end{cases} \quad (7)$$

また, 次のように定義することもできる.

$$\text{FinReal}(X) = \triangle \text{colim}_{\Delta_{<a}^n \rightarrow X} \text{FinReal}(\Delta_{<a}^n) \quad (8)$$

定義 8 関手  $\text{FinSing}:\text{FinEPMet} \rightarrow \text{Fin-sFuzz}$  は次のように定義される.

$$\text{FinSing}(Y) : ([n], [0, a]) \mapsto \text{hom}_{\text{FinEPMet}}(\text{FinReal}(\Delta_a^n)Y) \quad (9)$$

$\text{FinReal}$  と  $\text{FinSing}$  は随伴関手として与えられる.

定義 9  $X = X_1, \dots, X_N \in \mathbb{R}^n$

$\{(X, d_i)_{i=1 \dots N}\}$  は擬似距離空間の一部であり, 式 (10) で表される

$$d_i(X_i, X, j) = \begin{cases} d_{\mathcal{R}}(X_i, X_k) - \rho & \text{if } i = j \text{ or } i = k \\ \infty & otherwise \end{cases} \quad (10)$$

このとき  $\rho$  は  $X_i$  の近傍 (Nearest Neighbor) を示し,  $d_{\mathcal{M}}$  は多様体における曲線距離 (geodesic distance) を示す.  $X$  のファジートポロジカルな表現は式 (??) で示される.

$$?? \bigcup_{i=1}^n \text{FinSing}((X, d_i)) \quad (11)$$

このファジー集合によってグローバルな構造とローカルな構造を示していることになる.

## 4 低次元表現の最適化

定義 10  $(A, \mu)$  と  $(A, \nu)$  によって与えられるクロスエントロピーは次のように定義される.

$$C((A, \mu), (A, \nu)) \triangleq \sum_{a \in A} (\mu(a) \log \frac{\mu(a)}{\nu(a)} + (1 - \mu(a)) \log \frac{1 - \mu(a)}{1 - \nu(a)}) \quad (12)$$

これは, t-SNE の最適化と似ているところがある最終的な誤差関数は次のように表される.

$$C_l(X, Y) = \sum_{i=1}^l \lambda_i C(X_i, Y_i) \quad (13)$$

## 5 UMAP の可視化

本手法  $g$  は  $k$  近傍方を基にした手法となっている.

段階は 2 つに分かれており, 1 段階目においては重み付きの  $k$ -neighbor グラフを作成する. 2 段階目においては低次元における最適化を行う.

## 5.1 グラフの初期化

可視化を実装するためにまず、重み付き  $k$  近傍グラフを用いる。ここで次のデータセットを  $X = x_1, \dots, x_N$ , 距離を  $d: X \times X$  で定義する。ハイパーパラメータ  $k$  は距離  $d$  における近傍を決めるための値である。ここで  $\rho_i$  と  $\sigma_i$  について定義する。

$$\rho_i = \min\{d(x_i, x_{ij}) | 1 \leq j \leq k, d(x_i, x_{ij}) > 0\} \quad (14)$$

$$\sum_{j=1}^k \exp\left(\frac{\max(0, d(x_i, x_j) - \rho_i)}{\sigma_i}\right) = \log_2(k) \quad (15)$$

ここで重み付き有効非巡回グラフ  $\bar{G} = (V, E, w)$  を定義する。頂点集合  $V$  は  $X$  の単純集合である。またエッジ  $E$  は  $(x_i, x_{ij})$  で表され、重み  $w$  は式 (16) で定義される。

$$w((x_i, x_{ij})) = \exp\left(\frac{\max(0, d(x_i, x_j) - \rho_i)}{\sigma_i}\right) \quad (16)$$

隣接行列は式 16 によって定義される。また、対称行列は次で定義される。ここで、 $\circ$  はアダマール行列を示す。これにより無向巡回グラフを生成することができる。

$$B = A + A^T - A \circ A^T \quad (17)$$

## 5.2 グラフの配置

t-SNE の時と同じく斥力と引力を考慮してグラフの配置を行っていく。

引力は式 18 で示される。

$$\frac{-2ab\|y_i - y_j\|_2^{2(b-1)}}{1 + \|y_i - y_j\|_2^2} w((x_i, x_j))(y_i - y_j) \quad (18)$$

ここでの  $a, b$  はハイパーパラメータとなっている。

斥力は式 (19) で表されることにする。

$$\frac{b}{(\epsilon + \|y_i - y_j\|_2^2)(1 + \|y_i - y_j\|_2^2)} (1 - w((x_i, x_j)))(y_i - y_j) \quad (19)$$

$\epsilon$  は分裂を防ぐためのごく小さい数字となっている。

## 6 実装とハイパーパラメータ

UMAP のアルゴリズムの全容は Algorithm1 で示される。

次に Algorithm2 を示す。これは次元圧縮後の値の初期化を示す。

$\log(n)$  に最適な  $\omega$  を求める方法を Algorithm3 に示す。

グローバルな構造を保つための処理を重み付きグラフとラプラシアン行列を用いて表現する。これは Algorithm4 で示される。

---

**Algorithm 1** UMAP Algorithm

---

```
function UMAP( $X, n, d, \text{min-dist}, n\text{-epochs}$ )  
  for all  $x \in X$  do  
     $\text{Fs-set}[x] \leftarrow \text{LocalFuzzySimplicialSet}(X, x, n)$   
  
   $\text{top-rep} \leftarrow \cup_{x \in X} \text{fs-set}[x]$   
   $Y \leftarrow \text{SpectralEmbedding}(\text{top-rep}, d)$   
   $Y \leftarrow \text{OptimizeEmbedding}(\text{top-rep}, Y, \text{min-dist}, n\text{-epochs})$   
  
  return  $Y$ 
```

---

---

**Algorithm 2** Constructing a local fuzzy simplicial set

---

```
function LOCALFUZZYSIMPLICIALSET( $X, x, n$ )  
   $\text{knn}, \text{knn-dists} \leftarrow \text{ApproxNearestNeighbors}(X, x, n)$   
   $\rho \leftarrow \text{knn-dists}[1] \triangleright \text{Distance to nearest neighbor}$   
   $\sigma \leftarrow \text{SmoothKNNDist}(\text{knn-dists}, n, \sigma) \triangleright \text{Smooth approximator to knn-distance}$   
   $\text{fs-set}_0 \leftarrow X$   
   $\text{fs-set}_1 \leftarrow \{([x, y], 0) | y \in X\}$   
  for all  $y \in \text{knn}$  do  
     $d_{x,y} \leftarrow 0, \text{dist}(x, y) - \rho/\sigma$   
     $\text{fs-set}_1 \leftarrow \text{fs-set}_1 \cup ([x, y], \exp(-d_{x,y}))$   
  
  return  $\text{fs-set}$ 
```

---

---

**Algorithm 3** Compute the normalizing for distance  $\sigma$ 

---

```
function SMOOTHKNNDIST( $\text{knn-dists}, n, \rho$ )  
  Binary search for  $\sigma$  such that  $\sum_{i=1}^n \exp(-\text{knn} - \text{dists}_i - \rho)/\sigma = \log_2(n)$   
  
  return  $\sigma$ 
```

---

---

**Algorithm 4** Spectral OptimizeEmbedding for initialization

---

```
function SPECTRALEMBEDDING( $\text{top-rep}, d$ )  
   $A \leftarrow 1\text{-skeleton of top-rep expressed as a weighted adjacency matrix}$   
   $D \leftarrow \text{degree matrix for the graph } A$   
   $L \leftarrow D^{\frac{1}{2}}(D - A)D^{\frac{1}{2}}$   
   $\text{evec} \leftarrow \text{Eigenvectors of } L \text{ (Sorted)}$   
   $Y \leftarrow \text{evec}[1d + 1]$   
  
  return  $Y$ 
```

---

最後に最適化関数を与える。ここでは交差情報量を使うことによって与えられる。

$$C((A, \mu)(A, \nu)) = \sum_{a \in A} \mu(a) \log\left(\frac{\mu(a)}{\nu(a)}\right) + (1 - \mu(a)) \log\left(\frac{1 - \mu(a)}{1 - \nu(a)}\right) \quad (20)$$

$$\begin{aligned} &= \sum_{a \in A} \mu(a) \log(\mu(a)) + (1 - \mu(a)) \log(1 - \mu(a)) \\ &\quad - \sum_{a \in A} \mu(a) \log(\nu(a)) + (1 - \mu(a)) \log(1 - \nu(a)) \end{aligned} \quad (21)$$

ここで最小化したい式は式で与えられる.

$$P(x_i) = \frac{\sum_{a \in d_0(a)=x_i} 1 - \mu(a)}{\sum_{b \in d_0(b)=x_i} 1 - \mu(b)} \quad (22)$$

ここで  $\mu(\alpha)$  を  $\nu(\alpha)$  にしたがって更新したい

**定義 11** 定義  $\phi: \mathbb{R}^d \times \mathbb{R}^d \leftarrow [0, 1]$   $\mathbb{R}$  の 2 点間における 1 次元の単体のメンバーシップ関数の近似は式 23 で与えられる.

$$\phi(x, y) = (1 + a(\|x - y\|_2^2)^b)^{-1} \quad (23)$$

また,  $\Psi$  は非線形最小二乗法によって与えられる.

$$\Psi(x, y) = \begin{cases} 1 & \text{if } \|x - y\| \leq mindist \\ \exp(-(\|x - y\|_2 - mindist)) & \text{otherwise} \end{cases} \quad (24)$$

これによって埋め込みのための最適化は Algorithm5 で与えられる.

---

**Algorithm 5** Optimizing the embedding

---

**function** OPTIMIZEEMBEDDING(top-rep,  $Y$ , min-dist, n-epochs)

---

$\alpha \leftarrow 1.0$

Fit  $\phi$  from  $\Phi$  defined by min-dist

**for**  $e \leftarrow 1, n\text{-epochs}$  **do**

**for all**  $([a, b], p) \in top\text{-}rep_1$  **do**

**if** Random()  $\leq$  **then**

$y_a \leftarrow y_a + \alpha \dot{\Delta}(\log(\phi))(y_a, y_b)$

**for**  $i \leftarrow 1, n\text{-neg-samples}$  **do**

$c \leftarrow$  random sample from  $Y$

$y_a y_a + \alpha \dot{\Delta}(\log(1 - \phi))(y_a, y_c)$

$\alpha \leftarrow 1.0 - e/n\text{-epochs}$

**return**  $Y$

---

## 参考ページ

原論文:<https://arxiv.org/pdf/1802.03426.pdf>

日本語訳:<http://cymfh.cc/paper/UMAP.html>

グラフ行列の説明:<https://mathtrain.jp/graphmatrix>

UMAP のドキュメント:<https://umap-learn.readthedocs.io/en/latest/>