HDS 10. Nonlinear Dimension Reduction with t-SNE and UMAP

Matti Pirinen, University of Helsinki

15.1.2024

We have seen how the PCA extracts such linear combinations of the p original variables that are maximally informative among all linear combinations. Often the leading PCs have a clear and interpretable structure and therefore the PCA is a widely-used method to visualize and reduce high-dimensional data.

PCs are **global linear** functions of data and hence the leading PCs tend to capture such directions from the input space on which the distant data points remain distant from each other also in the leading PCs as such directions maximize the variance of the projected points. However, for high-dimensional data that happen to be structured in some non-linear way on some lower dimensional subspace, it would also be important to keep similar samples close together in the low-dimensional representation, which may not be possible by any global linear function such as a PC.

Many methods for dimension reduction that try to capture more of the local structure are non-linear and are not guaranteed to yield a globally optimal solution, which means

that the result may change with the seed of the random number generator that initializes the algorithm.

Here we study two methods: **t-SNE** and **UMAP**. Let's first see what they produce in practice and then come back to what is going on under the hood.

1000 Genomes data

The 1000 Genomes Project has produced genotype data from across the world. Here we consider a subset of n = 1092 individuals from the following 14 **populations**, divided into 4 continental **groups**,

- ASW [AFR] (61) African Ancestry in Southwest US
- CEU [EUR] (85) Utah residents (CEPH) with Northern and Western European ancestry
- CHB [ASN] (97) Han Chinese in Beijing, China
- CHS [ASN] (100) Southern Han Chinese
- CLM [AMR] (60) Colombian in Medellin, Colombia
- FIN [EUR] (93) Finnish from Finland
- GBR [EUR] (89) British from England and Scotland
- IBS [EUR] (14) Iberian population in Spain
- JPT [ASN] (89) Japanese in Toyko, Japan
- LWK [AFR] (97) Luhya in Webuye, Kenya
- MXL [AMR] (66) Mexican Ancestry in Los Angeles, CA
- PUR [AMR] (55) Puerto Rican in Puerto Rico
- TSI [EUR] (98) Toscani in Italia
- YRI [AFR] (88) Yoruba in Ibadan, Nigeria

Each individual has been measured on p = 4212 genetic variants (each can have value 0, 1 or 2) from chromosomes 15-22.

X = read.table("geno_1000G_phase1_chr15-22.txt", as.is = TRUE, header = TRUE)
dim(X)

[1] 1092 4215

X[1:4, 1:10]

| ## | | id | population | group | X1 | Х2 | ΧЗ | X4 | X5 | X6 | Х7 |
|----|---|---------|------------|-------|----|----|----|----|----|----|----|
| ## | 1 | HG00096 | GBR | EUR | 0 | 0 | 0 | 0 | 2 | 0 | 0 |
| ## | 2 | HG00097 | GBR | EUR | 1 | 1 | 0 | 0 | 1 | 0 | 0 |
| ## | 3 | HG00099 | GBR | EUR | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| ## | 4 | HG00100 | GBR | EUR | 0 | 0 | 1 | 0 | 2 | 0 | 1 |
| | | | | | | | | | | | |

table(X[,"group"], X[,"population"])

| ## | | | | | | | | | | | | | | | |
|----|-----|-----|-----|-----|-----|----------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| ## | | ASW | CEU | CHB | CHS | \mathtt{CLM} | FIN | GBR | IBS | JPT | LWK | MXL | PUR | TSI | YRI |
| ## | AFR | 61 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 97 | 0 | 0 | 0 | 88 |
| ## | AMR | 0 | 0 | 0 | 0 | 60 | 0 | 0 | 0 | 0 | 0 | 66 | 55 | 0 | 0 |
| ## | ASN | 0 | 0 | 97 | 100 | 0 | 0 | 0 | 0 | 89 | 0 | 0 | 0 | 0 | 0 |
| ## | EUR | 0 | 85 | 0 | 0 | 0 | 93 | 89 | 14 | 0 | 0 | 0 | 0 | 98 | 0 |

Let's do the PCA and plot the 12 leading PCs in pairwise plots coloring each individual by their continental group (Africa, Americas, Asia or Europe) given by the **group** variable.



Visually, the PCs 1-8 seem to capture broader structure whereas the PCs from 9 onward seem to separate small groups of possibly more closely related pairs or triples.

Let's next color the points by the population rather than by the continent.



```
for(ii in 1:4){
    plot(pr$x[,2*ii - 1], pr$x[,2*ii], col = cols.pop, pch = pchs.pop,
        xlab = paste0("PC", 2*ii - 1), ylab = paste0("PC", 2*ii))
    if(ii == 0) legend("bottomright", col = grs.col, leg = grs, pch = 3, cex = 1.3)
}
legend("bottomright", leg = pops, col = pops.col, pch = pops.pch, cex = 0.8)
```



Let's then compare the PCA plots to t-SNE and UMAP, using the R packages Rtsne and umap, respectively, that we apply to compress the first 8 PCs further to just two dimensions.

```
#install.packages("Rtsne")
library(Rtsne)
set.seed(67)
tsne = Rtsne(X = pr$x[,1:8], perplexity = 10, theta = 0.0, pca = FALSE)
par(mar = c(4,4,4,8), xpd = TRUE)
plot(tsne$Y, col = cols.pop, pch = pchs.pop, main = "t-SNE", xlab = "", ylab = "")
legend("topright", inset = c(-0.15,0), leg = pops, col = pops.col, pch = pops.pch)
```

| t– | S | Ν | Е |
|----|---|---|---|
| • | - | _ | |



#install.packages("umap")
library(umap)
set.seed(67)
umap.res = umap(pr\$x[,1:8])
par(mar = c(4,4,4,8), xpd = TRUE)
plot(umap.res\$layout, col = cols.pop, pch = pchs.pop, main = "UMAP", xlab ="", ylab = "")
legend("topright", inset = c(-0.15,0), leg = pops, col = pops.col, pch = pops.pch)

UMAP



Wee see that t-SNE and UMAP largely group individuals from a same population close together and separate them from the other populations whereas they do not put so much emphasis on making, for example, the two African ancestry populations YRI and LWK equally distant from all 5 European populations or all 3 Asian populations, as the PCs 1 and 2 did above. Thus, we see that t-SNE and UMAP may indeed preserve more of the local structure around the neighborhood of each sample but consequently cannot simultaneously be completely faithful to the overall global structure as defined by the leading PCs. In a sense, t-SNE and UMAP try to present both the local and global structure in only two dimensions, and for this they need to find a trade-off between these two goals.

t-SNE: t-distributed Stochastic Neighbor Embedding

t-SNE was introduce in "Visualizing Data using t-SNE" by van der Maaten & Hinton (2008).

It builds on earlier work on Stochastic Neighbor Embedding (SNE), where the idea is to measure distance in the high-dimensional input space by a conditional probability. If \boldsymbol{x}_i and \boldsymbol{x}_j are two *p*-dimensional data points, we can compute a conditional probability $p_{j|i}$ that \boldsymbol{x}_i would pick as its neighbor the point \boldsymbol{x}_j if neighbors would be chosen from a Gaussian distribution centered on \boldsymbol{x}_i and having a same variance σ_i^2 in each dimension (we come back to how σ_i^2 will be chosen later).

$$p_{j|i} \propto \exp\left(-\frac{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2}{2\sigma_i^2}\right)$$
 and $\sum_{j \neq i} p_{j|i} = 1.$

This probability is larger for points that are closer to \boldsymbol{x}_i than for those that are farther away. Similarly,

we can define $p_{i|j}$. Finally, we can average (and normalize by n) these two probabilities to get $p_{ij} = p_{ji} = \frac{1}{2n}(p_{i|j} + p_{j|i})$ to represent the similarity between \boldsymbol{x}_i and \boldsymbol{x}_j by a single value.

The goal of the SNE is to map the *p*-dimensional input values \boldsymbol{x}_i to two dimensional (or three dimensional) output points y_i in such a way that the distances q_{ij} defined by a similar density function evaluation in the output space would optimally match the input space distances p_{ij} . Here "optimally" means in terms of the Kullback-Leibler divergence of the distribution $Q = (q_{ij})$ from the distribution $P = (p_{ij})$:

$$\operatorname{KL}(P || Q) = \sum_{i < j} p_{ij} \log \left(\frac{p_{ij}}{q_{ij}} \right),$$

which is always non-negative and equals to zero if and only if the distributions are the same. Minimizing this cost function puts more emphasis on making pairs with high p_{ij} to have similarly high q_{ij} , but less emphasis on matching the two when p_{ij} is small. Hence, the SNE is expected to preserve particularly well the local structure in the data and pay less attention to what happens between long distances.

What t-SNE adds on top of the SNE is that the distribution $Q = (q_{ij})$ is defined using the density function of the **t-distribution with 1 degree of freedom**, also known as the **Cauchy distribution**, rather than by a Gaussian, as was done in the SNE. This means that, in the low-dimensional output space, the conditional probabilities in t-SNE are defined as

$$q_{j|i} \propto \left(1 + \|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2\right)^{-1}$$
 and $\sum_{j \neq i} q_{j|i} = 1,$

and these are symmetrized and normalized as above: $q_{ij} = q_{ji} = \frac{1}{2n}(q_{i|j} + q_{j|i})$.

The Cauchy distribution has thick tails, and therefore t-SNE can tolerate more discrepancy between the distances in input and output spaces when it comes to the points that are moderately far from each other in the input space. This helps to avoid *the crowding problem*: in the high-dimensional input space, there are potentially many equidistant points with moderate distance from a particular point, and not all of these can be similarly accounted for in the low-dimensional space. The Cauchy distribution makes sure that some of these points can be more spread out in the output space without a very high penalty.

Perplexity In the input space, the distances are defined by a Gaussian density with a data point specific variance σ_i^2 . This parameter determines how t-SNE measure of distance from \boldsymbol{x}_i decays with the Euclidean distance from \boldsymbol{x}_i , with larger values of σ_i^2 meaning a slower decay and smaller values meaning a quicker decay. In order to preserve the local structure around each point, t-SNE adjusts σ_i^2 in such a way that all the conditional distributions $(p_{j|i})_{i\leq n}$ have approximately the same *perplexity*, which can be interpreted as an effective number of neighbors. The target perplexity is given as a parameter to the algorithm and is typically between 5 and 50. For example, a value of 15 means roughly that the closeness values $p_{j|i}$ from \boldsymbol{x}_i to about 15 of its most closest data points are all large enough so that those points can be considered as "neighbors" but the same is not true for some larger set of points than the 15 closest ones. In other words, the similarity measured from \boldsymbol{x}_i decays with such a rate that about 15 points are "nearby". Effectively this means that about 15 closest neighbors are taken into account when constructing the low-dimensional representation.

Let's see how different perplexities show up in the results:



theta parameter is a trade-off between the accuracy and computational complexity, where 0 means the largest accuracy and larger values mean more computationally efficient but less accurate approximations. The default is 0.5.

Resources to learn more about t-SNE

- StatQuest video on t-SNE (11:48)
- How to Use t-SNE Effectively
- How to tune hyperparameters of tSNE
- Roberto Stelling's blog
- "Visualizing Data using t-SNE" paper by van der Maaten & Hinton (2008).

UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction

UMAP was introduced in 2018 by L.McInnes, J.Haley, J.Melville. They summarize the motivation for UMAP compared to t-SNE: similar quality of visualization with a much more efficient algorithm. Additionally, UMAP is better able to maintain the global structure than t-SNE, which, on the other hand, may make UMAP to discard some details of the local structure compared to t-SNE.

Methodologically, UMAP uses similar ideas as t-SNE although the theoretical derivation is more mathematical. A description of the differences from t-SNE can be found from the Appendix C of the UMAP paper. Another description is "How exactly UMAP works?" by Nikolay Oskolkov.

The main parameters of UMAP are n_neighbors, the number of closest neighbors that are considered, min_dist, the minimum distance of the points in the output space, n_components, the output dimension and metric that defines the distance of the input space. These are clearly explained at the UMAP website.

An illustrative site about UMAP: https://pair-code.github.io/understanding-umap/

Let's see how n_neighbors compares to the effect of perplexity in t-SNE that we saw above.



Discussion

- Most non-linear dimension reduction techniques (including t-SNE and UMAP) lack the strong interpretability of Principal Component Analysis where the dimensions are the directions of greatest variance in the source data. If strong interpretability is needed, the PCA is recommended.
- As t-SNE and UMAP are based on the distance between observations rather than the source features, they do not produce easily interpretable loadings per each variable that the PCA can provide for each output dimension.
- A core assumptions of UMAP is that there exists manifold structure in the data. Because of this, UMAP may find manifold structure within the noise of a dataset, a type of overfitting. As more data is sampled, UMAP becomes more robust. However, care must be taken with small sample sizes of noisy data, or data with only large-scale manifold structure.
- If data are high-dimensional, say p > 100, it is often useful to first apply PCA and take some tens of the leading PCs as input for t-SNE or UMAP to further reduce the data to 2 or 3 dimensions for visualization.