

Recitation 5

Dimension reduction

GE (SABER) LIU

MIT - 6.802 / 6.874 / 20.390 / 20.490 / HST.506 - Spring 2019

2019-03-07 / 2019-03-08

Outline

- Linear algebra basics
- Principle component analysis (PCA)
- t-SNE and parametric t-SNE
- Auto-encoder
- U-MAP

Linear algebra basics

Eigenvector: An eigenvector or characteristic vector of a linear transformation T is a non-zero vector that changes by only a scalar factor when that linear transformation is applied to it.

$T(\mathbf{v}) = \lambda \mathbf{v}$ or $\mathbf{A}\mathbf{v} = \lambda \mathbf{v}$ if the transformation can be represented as a **square matrix** A

where λ is a scalar, known as the **eigenvalue**, characteristic value, or characteristic root associated with the eigenvector \mathbf{v} . An $N \times N$ matrix has at most N **linearly independent** eigenvectors.

Eigen-decomposition: Factorization of a matrix into a canonical form, whereby the matrix is represented in terms of its eigenvalues and eigenvectors.

$$A = Q\Lambda Q^{-1}$$

where Q is the square $N \times N$ matrix whose i -th column is the eigenvector v_i of A , and Λ is the **diagonal matrix** whose diagonal elements are the corresponding eigenvalues, $\Lambda_{ii} = \lambda_i$. Note that \mathbf{A} has to have N linearly independent eigenvectors (Only diagonalizable matrices can be factorized in this way).

Singular value decomposition (SVD): factorization of a real or complex matrix $\mathbf{M}_{m \times n}$ into $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ where $\mathbf{U}_{m \times m}$ and $\mathbf{V}_{n \times n}$ are unitary matrix with orthonormal eigenvectors of $\mathbf{M}\mathbf{M}^*$ and $\mathbf{M}^*\mathbf{M}$, and $\mathbf{\Sigma}_{m \times n}$ is a rectangular diagonal matrix with non-negative real numbers on the diagonal. X^* means the conjugate transpose of X .

Special matrices

Real symmetric matrices: Matrix \mathbf{A} is symmetric if $\mathbf{A} = \mathbf{A}^T$

Theorem

Any symmetric matrix:

- has only real eigenvalues
- is always diagonalizable
- has **orthogonal** eigenvectors

Corollary: If matrix \mathbf{A} is symmetric then there exists $\mathbf{Q}^T \mathbf{Q} = \mathcal{I}$ such that $\mathbf{A} = \mathbf{Q}^T \mathbf{\Lambda} \mathbf{Q}$.

Positive definite matrices: A symmetric matrix \mathbf{A} is positive definite/semi-definite if all its eigenvalues are positive/non-negative.

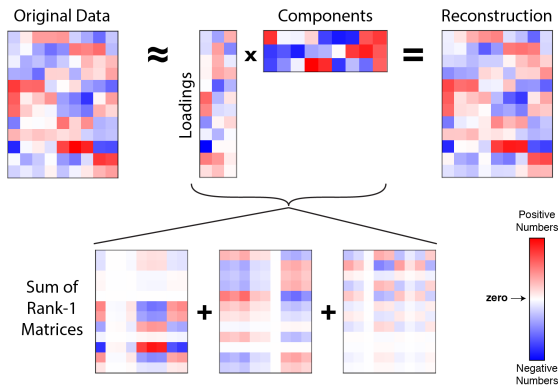
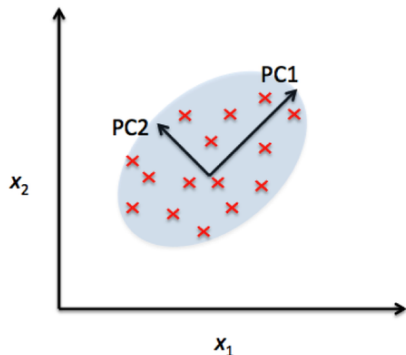
Theorem

\mathbf{A} is positive definite if and only if $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0, \forall \mathbf{x} \neq 0$.

Diagonalizable \supset **Symmetric** \supset **Positive semi-definite** \supset **Positive definite**

PCA

Principal component analysis: Orthogonal transformation to convert a set of observations of possibly correlated variables into a set of linearly uncorrelated variables called principal components. This transformation is defined such that the first principal component has the **largest possible variance**, and each succeeding component in turn has the highest variance possible under the constraint that it is **orthogonal** to the preceding components. The resulting vectors are **linear combination** of the variables and form an **orthogonal basis set**.



Principal component analysis: Consider a data matrix $\mathbf{X}_{m \times n}$ with m examples of n dimensional features (each dimension has been z-centered such that the mean is zero). Try to find a set of n -dimensional vectors of weights or coefficients $\mathbf{w}_{(k)} = (w_1, \dots, w_n)_{(k)}$ that projects each row vector $\mathbf{x}_{(i)}$ of \mathbf{X} to a new vector of principal component scores $\mathbf{t}_{(i)} = (t_1, \dots, t_n)_{(i)}$, given by $t_{k(i)} = \mathbf{x}_{(i)} \cdot \mathbf{w}_{(k)}$

$$\begin{bmatrix} - & \mathbf{x}_{(1)} & - \\ & \vdots & \\ - & \mathbf{x}_{(m)} & - \end{bmatrix}_{m \times n} \times \begin{bmatrix} | & & | \\ \mathbf{w}_{(1)} & \dots & \mathbf{w}_{(n)} \\ | & & | \end{bmatrix}_{n \times n} = \begin{bmatrix} - & \mathbf{t}_{(1)} & - \\ & \vdots & \\ - & \mathbf{t}_{(m)} & - \end{bmatrix}_{m \times n} = \begin{bmatrix} | & & | \\ \mathbf{PC}_{(1)} & \dots & \mathbf{PC}_{(n)} \\ | & & | \end{bmatrix}_{m \times n}$$

$$\mathbf{T} = \mathbf{XW}$$

The weights are constrained to be a unit vector such that $\mathbf{w}_{(i)}^T \mathbf{w}_{(i)} = 1$.

To solve for projection with maximized variance, the **first** weight vector has to satisfy:

$$\mathbf{w}_1 = \operatorname{argmax}_{\mathbf{w}^T \mathbf{w} = 1} \sum_i (\mathbf{x}_i \cdot \mathbf{w}_1)^2 = \operatorname{argmax}_{\mathbf{w}^T \mathbf{w} = 1} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} = \operatorname{argmax}_{\mathbf{w}^T \mathbf{w} = 1} \frac{\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w}}{\mathbf{w}^T \mathbf{w}}$$

A standard solution to this optimization for a positive semidefinite matrix such as $\mathbf{X}^T \mathbf{X}$ is the largest eigenvalue of the matrix, which occurs when \mathbf{w} is the corresponding eigenvector. The rest of the component can be given as:

$$\mathbf{w}_k = \operatorname{argmax}_{\mathbf{w}^T \mathbf{w} = 1} \frac{\mathbf{w}^T \mathbf{X}_k^T \mathbf{X}_k \mathbf{w}}{\mathbf{w}^T \mathbf{w}} \quad \text{where} \quad \mathbf{X}_k = \mathbf{X} - \sum_{s=1}^{k-1} \mathbf{X} \mathbf{w}_s \mathbf{w}_s^T$$

It turns out that this gives the remaining eigenvectors of $\mathbf{X}^T \mathbf{X}$, with the maximum values equal to the corresponding eigenvalues. Thus, solving weight vectors for PCA is equivalent to finding the eigenvectors of $\mathbf{X}^T \mathbf{X}$ and sorting by its corresponding eigenvalues.

The SVD of $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T$, so $\mathbf{T} = \mathbf{X} \mathbf{W} = \mathbf{U} \mathbf{\Sigma} \mathbf{W}^T \mathbf{W} = \mathbf{U} \mathbf{\Sigma}$. Each column of \mathbf{T} is given by one of the left singular vectors of \mathbf{X} multiplied by the corresponding singular value.

Stochastic neighbor embedding (SNE): An unsupervised nonlinear dimensionality reduction technique where the goal is to find a low-dimensional (2-dimensional) representation of the original inputs such that pairwise similarity are best preserved and the inherent clustering structure can be visualized.

Similarity in original input space: the similarity of datapoint x_i to datapoint x_j is defined as the conditional probability, $p(x_j|x_i)$, that x_i would pick x_j as its neighbor if neighbors were picked in proportion to their probability density under a Gaussian centered at x_i . Given an input matrix \mathbf{X} , in which each row is a sample x_i and each column represent a feature dimension, the pairwise similarity is defined as:

$$P_{i,j} = p(x_j|x_i) = \frac{\exp(-||x_i - x_j||^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-||x_i - x_k||^2/2\sigma_i^2)} \quad (1)$$

We define $\beta_i = 1/\sigma_i$ which is equivalent to the **precision** of a multivariant Gaussian. We further set constraints on the perplexity of the conditional distribution where **perplexity** is defined using Shannon entropy of P_i :

$$Perplexity(P_i) = 2^{H(P_i)} \quad \text{and} \quad H(P_i) = - \sum_j p_{x_j|x_i} \log_2 p_{x_j|x_i}$$

In order to guarantee the perplexity constraint for each P_i , we need to find the corresponding β_i such that the resulting distribution has the desired perplexity. Given that $Perplexity(P_i)$ is a **monotonically decreasing** function of β_i , we could use binary search to estimate the solution for β_i .

t-distribution Stochastic neighbor embedding (t-SNE)

There are several modifications we need to make:

- **Use symmetric similarity matrix instead:** Since the gradients for the conditional distribution is hard to compute, people use the joint distribution as an alternative that is "just as good" and this gives symmetric distribution matrix:

$$P_{ij}^{symmetric} = p(x_i, x_j) = \frac{\exp(-||x_i - x_j||^2 / 2\sigma_i^2)}{\sum_{k \neq l} \exp(-||x_k - x_l||^2 / 2\sigma_i^2)} = \frac{p_{x_j|x_i} + p_{x_i|x_j}}{2N}$$

- **Use Student t-distribution for similarity in embedded space:** We look for a 2D representation of \mathbf{X} which is \mathbf{Y} ($N \times 2$ matrix), such that the pairwise similarity on \mathbf{Y} :

$$Q_{ij} = \frac{(1 + ||y_i - y_j||^2)^{-1}}{\sum_{k \neq l} (1 + ||y_k - y_l||^2)^{-1}}$$

is similar to P_{ij} . Which is equivalent to minimizing the **KL-divergence** between P and Q :

$$C = KL(P||Q) = \sum_i \sum_j P_{ij} \log \frac{P_{ij}}{Q_{ij}}$$

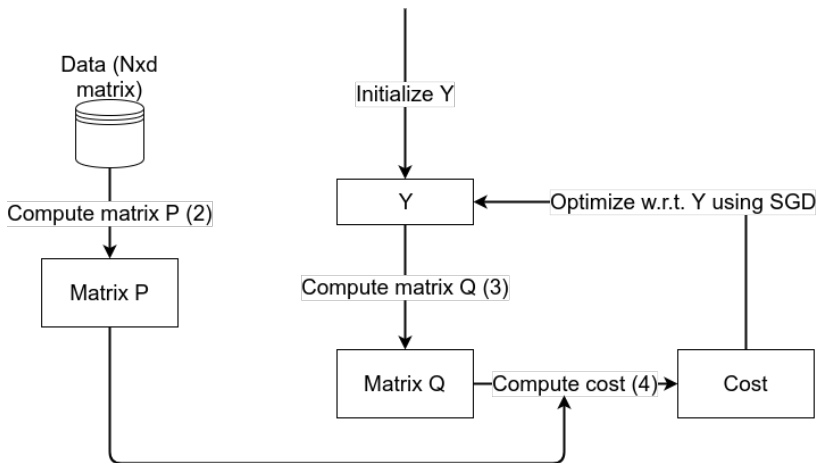
The gradients for conducting gradient descent are:

$$\frac{\partial C}{\partial y_i} = 4 \sum_{j \neq i} (p_{ij} - q_{ij}) (1 + ||y_i - y_j||^2)^{-1} (y_i - y_j)$$

non-parametric t-SNE

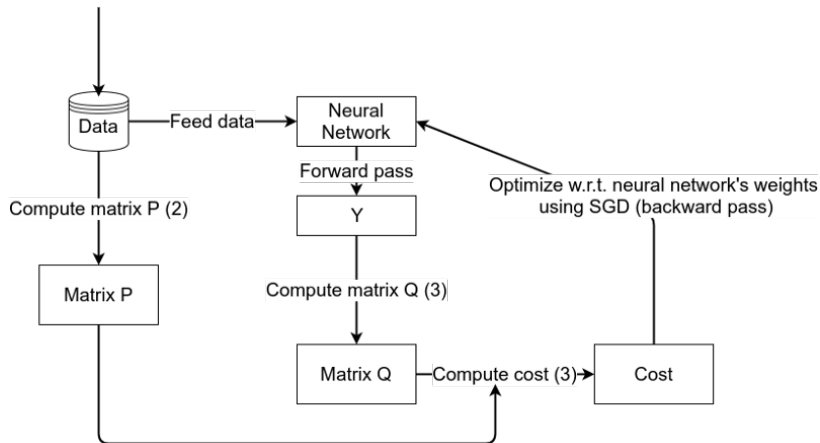
The non-parametric t-SNE has several drawbacks for being non-parametric.

- You can't embed new points that weren't used in the training phase without running the algorithm from scratch again and without preserving the previous embedding results.
- It is not scalable because the more points you have the larger memory you need to store **D**, **P** and **Q**.



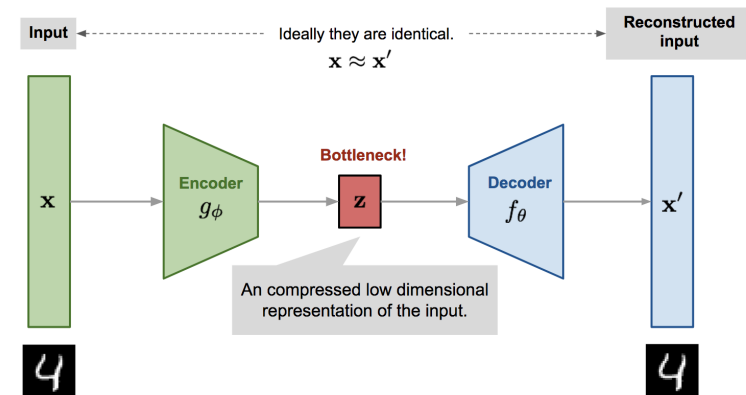
parametric t-SNE

We could think of a parametric approach instead, in which we want to build a model $y = f(x, \Theta)$ that maps any given input x_i to low-dimensional output y_i . A useful family of model we would consider is of course the **neural networks**. The good thing about this is that we can calculate P on a smaller batch of X and train model using batched data. Also once the model is trained, we will have a deterministic embedding that can be calculated within linear time.



Auto-Encoder

Auto-encoder: An autoencoder learns to compress data from the input layer into a low dimensional representation, and then uncompress that representation into something that closely matches the original data.



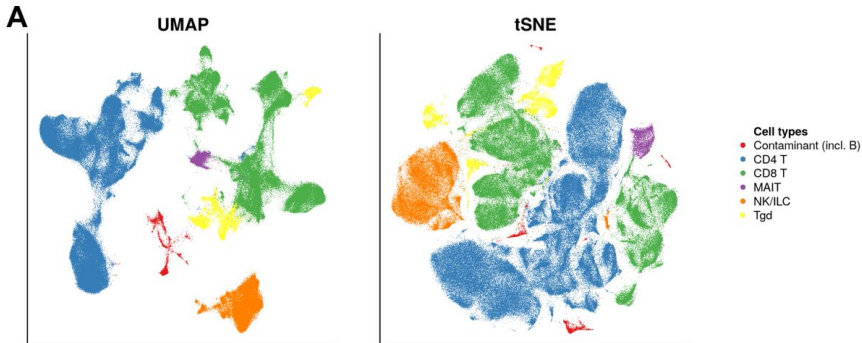
$$z = g_\phi(x) \quad x' = f_\theta(z) = f_\theta(g_\phi(x)) \quad \mathcal{L}(x, x') = \|x - x'\|^2$$

U-MAP

Uniform Manifold Approximation and Projection (UMAP): a general dimension reduction technique based on manifold learning techniques and topological data analysis. The algorithm is founded on three assumptions about the data:

- The data is uniformly distributed on Riemannian manifold
- The Riemannian metric is locally constant (or can be approximated as such)
- The manifold is locally connected

The embedding is found by searching for a low dimensional projection of the data that has the closest possible equivalent fuzzy topological structure.



Dimension reduction algorithms

