

Fair Neighbor Embedding

February 22, 2024

Seoul National University

Table of Contents

① Contribution

② Notation

③ NeRV

Contribution

- Introduce a new fair nonlinear dimension reduction method.
- Low-dimensional embeddings that preserve high-dimensional data neighborhoods without the biased association of protected groups.

Table of Contents

① Contribution

② Notation

③ NeRV

Notation

- \mathcal{X} : Feature space in \mathbb{R}^D
- \mathcal{Y} : Embedding space of feature space in \mathbb{R}^d
- \mathcal{S} : Set of sensitive attributes

Table of Contents

① Contribution

② Notation

③ NeRV

- The Neighbor Retrieval Visualizer is a nonlinear DR method that aims to create low-dimensional data embeddings.
- Given the input data set $\{x_i\}_{i=1}^n$, $x_i \in \mathbb{R}^D$, the probability that data point j is picked as a neighbor of point i is

$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / \sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / \sigma_i^2)}.$$

where σ_i^2 controls falloff of the p_{ij} with respect to distance.

- NeRV outputs an embedding of the points $\{y_i\}_{i=1}^n$, $y_i \in \mathbb{R}^d$, in a d -dimensional output space, where neighbor probabilities $q_{ij} = \{q_{ij}\}_{j=1, \dots, N, j \neq i}$ are defined based on output coordinates y_i is

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

Fairness

- The objective function of NeRV is:

$$C_{\text{NeRV}} = \frac{1}{N} \sum_{i=1}^N (\lambda D_{\text{KL}}(p_i, q_i) + (1 - \lambda) D_{\text{KL}}(q_i, p_i))$$

where λ is a tradeoff parameter, and D_{KL} is the Kullback-Leibler divergence, $D_{\text{KL}}(p_i, q_i) = \sum_{j \neq i} p_{ij} \log \frac{p_{ij}}{q_{ij}}$.

- Objective function of Conditional NeRV when values of sensitive variables are given,

$$C_{\text{CNeRV}} = \sum_i \tau^\epsilon D_B(p_i^\epsilon, q_i^\epsilon) + (1 - \tau^\epsilon) D_B(q_i^\epsilon, p_i^\epsilon) \\ + \tau^\notin D_B(p_i^\notin, q_i^\notin) + (1 - \tau^\notin) D_B(q_i^\notin, p_i^\notin)$$

where $D_B(p_i^{S_i}, q_i^{S_i}) = \sum_{j \in S_i} p_{ij} \log \frac{p_{ij}}{q_{ij}} + q_{ij} - p_{ij}$ is a Bregman divergence and $p_i^\epsilon = \{p_{ij}\}_{j \in S_i^\epsilon}$ and $q_i^\epsilon = \{q_{ij}\}_{j \in S_i^\epsilon}$

- Define

$$r_{is} = \frac{\sum_{j \neq i} \delta(s_j, s) \exp(-\|y_i - y_j\|^2 / \sigma_i^2)}{\sum_{j \neq i} \exp(-\|y_i - y_j\|^2 / \sigma_i^2)}$$

and

$$\rho_{is} = \begin{cases} 1 - \omega & \text{if } s = s_i \\ u(s) \cdot \omega / (1 - u(s_i)) & \text{otherwise} \end{cases}$$

where $u(s_i)$ is the overall proportion of sensitive value s_i in the data and $\omega \in [0, 1]$ is a weight controlling influence of value s_i in the neighborhood.

Fair-NeRV and Fair-t-NeRV

- Term of fair embedding :

$$C_{\text{Fairness}} = \frac{1}{N} \sum_i (\gamma D_{\text{KL}}(\rho_i, r_i) + (1 - \gamma) D_{\text{KL}}(r_i, \rho_i)),$$

- Final Objective

$$C_{\text{FairNeRV}} = \beta C_{\text{CNeRV}} + (1 - \beta) C_{\text{Fairness}}$$

Experiment

- The Syn data set is an artificial set of 1000 data points with 5 dimensions.
- The first three dimensions have three multivariate Gaussian clusters, cluster membership considered sensitive information.
- Dimensions 4 and 5 dimension have an independent mixture of three multi-variate Gaussian clusters, considered non-sensitive information

Experiment

