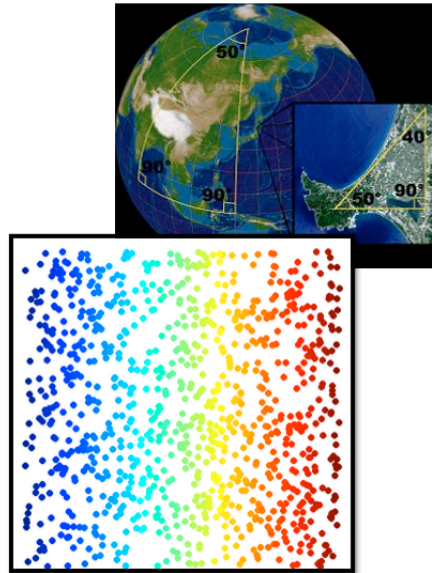
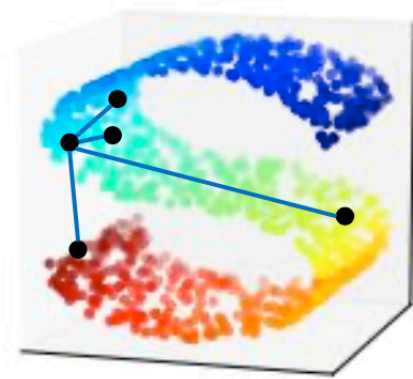


本文主要叙述了t-SNE，即T-distributed Stochastic Neighbor Embedding；先介绍了LLE的主要思想，再总结了它的缺点，从而引出t-SNE；

Manifold Learning

在高维空间里，距离该点很远的点很可能与这个点也是有关联的，因此我们可以把3-D的空间进行降维，那么我们就可以更方便地进行clustering或unsupervised learning 任务

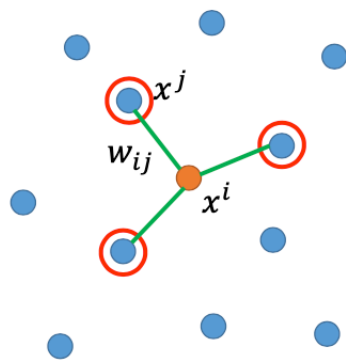
Manifold Learning



Suitable for clustering or following supervised learning

Locally Linear Embedding (LLE)

用 w_{ij} 表示 x_i, x_j 之间的联系，先找到使得 $\sum_i \|x^i - \sum_j w_{ij} x^j\|_2$ 最小化的 w_{ij} ，再根据这个 w_{ij} 来找到降维的结果 z_i, z_j



w_{ij} represents the relation between x^i and x^j

Find a set of w_{ij} minimizing

$$\sum_i \left\| x^i - \sum_j w_{ij} x^j \right\|_2$$

Then find the dimension reduction results z^i and z^j based on w_{ij}

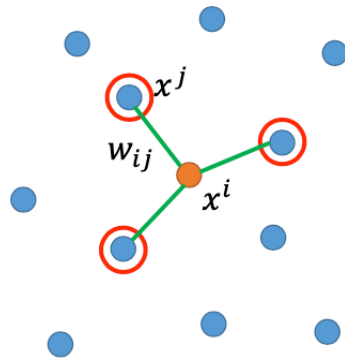
如果并不知道之前的 x_i, x_j ，那么就可以用LLE这种方法，也可以得出 z_i, z_j

LLE

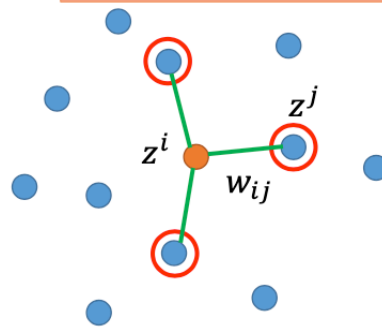
Find a set of z^i minimizing

$$\sum_i \left\| z^i - \sum_j w_{ij} z^j \right\|_2$$

Keep w_{ij} unchanged



Original Space



New (Low-dim) Space

Laplacian Eigenmaps

Review: 在之前的semi-supervised learning中, 如果 x^1, x^2 在一个high density region内是相近的, 那么我们就可以认为 \hat{y}^1, \hat{y}^2 也有类似的表现

如果 y^i, y^j 是connected的, 那么其 w_{ij} 就是对应的相似度; 如果没有connect, 其 w_{ij} 就是0

$$\text{Laplacian Eigenmaps} \quad w_{i,j} = \begin{cases} \text{similarity} & \text{If connected} \\ 0 & \text{otherwise} \end{cases}$$

- *Review in semi-supervised learning:* If x^1 and x^2 are close in a high density region, \hat{y}^1 and \hat{y}^2 are probably the same.



$$L = \sum_{x^r} C(y^r, \hat{y}^r) + \lambda S$$

As a regularization term

$$S = \frac{1}{2} \sum_{i,j} w_{i,j} (y^i - y^j)^2 = \mathbf{y}^T L \mathbf{y}$$

S evaluates how smooth your label is

L: (R+U) x (R+U) matrix

Graph Laplacian

$$L = D - W$$

我们也可以得出类似smoothness的式子, 计算 z^i, z^j 之间的smoothness

- **Dimension Reduction:** If x^1 and x^2 are close in a high density region, z^1 and z^2 are close to each other.

$$S = \frac{1}{2} \sum_{i,j} w_{i,j} (z^i - z^j)^2$$

Any problem? How about $z^i = z^j = \mathbf{0}$?

Giving some constraints to z :

If the dim of z is M , $\text{Span}\{z^1, z^2, \dots, z^N\} = \mathbb{R}^M$

Spectral clustering: clustering on z

Belkin, M., Niyogi, P. Laplacian eigenmaps and spectral techniques for embedding and clustering. *Advances in neural information processing systems*. 2002

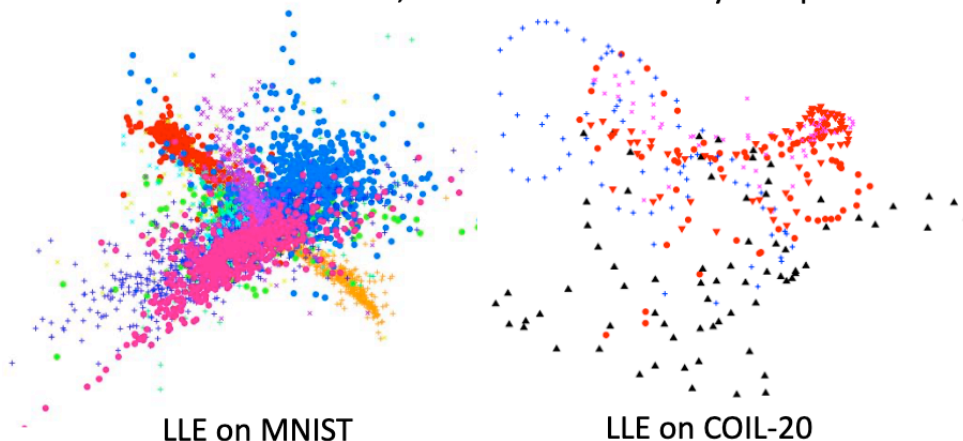
那么我们现在的目标就是找到 z^i, z^j ，来使 S 达到最小值，还需要有一些额外的 constraints

现在我们对 z 加入一些 constraints，如果降维后 z 的维数是 M ，那么我们就希望取出来的这些点还生活在比 M 还低维的空间里面；我们现在希望把塞进高维空间的低维空间展开，我们就不希望展开之后的点在一个更低维的空间里面

T-distributed Stochastic Neighbor Embedding (t-SNE)

对于之前的 LLE 方法，类似的数据之间是很 close 的，但不同类别之间的 data 却没有分开，是叠成一团的

- Problem of the previous approaches
 - Similar data are close, but different data may collapse



为了找到对应的 z^i, z^j ，先计算 x^i, x^j 之间的相似度 $S(x^i, x^j)$ ，再得出一个分布 $P(x^j|x^i)$ ；还要计算 z^i, z^j 之间的相似度 $S'(z^i, z^j)$ ，得出分布 $Q(z^j|z^i)$ ；

这两个分布应该越接近越好，使用 L 来表示

t-SNE

x \longrightarrow z

Compute similarity between all pairs of x: $S(x^i, x^j)$

$$P(x^j|x^i) = \frac{S(x^i, x^j)}{\sum_{k \neq i} S(x^i, x^k)}$$

Compute similarity between all pairs of z: $S'(z^i, z^j)$

$$Q(z^j|z^i) = \frac{S'(z^i, z^j)}{\sum_{k \neq i} S'(z^i, z^k)}$$

Find a set of z making the two distributions as close as possible

$$L = \sum_i KL(P(*|x^i) || Q(*|z^i))$$

$$= \sum_i \sum_j P(x^j|x^i) \log \frac{P(x^j|x^i)}{Q(z^j|z^i)}$$

可以使用gradient descent, 有了L函数, 再分别对 z^i, z^j 求偏微分即可

但t-SNE要对所有的point之间都求similarity, 因此计算量比较大, 在数据量很大的情况下电脑的计算速度会非常慢

因此, 对于很高的dimensions, 通常先做降维 (PCA), 比如可以降维到50维, 再使用t-SNE降到2维

通常我们使用t-SNE来对高维的数据进行可视化

Ignore σ for simplicity

t-SNE – Similarity Measure

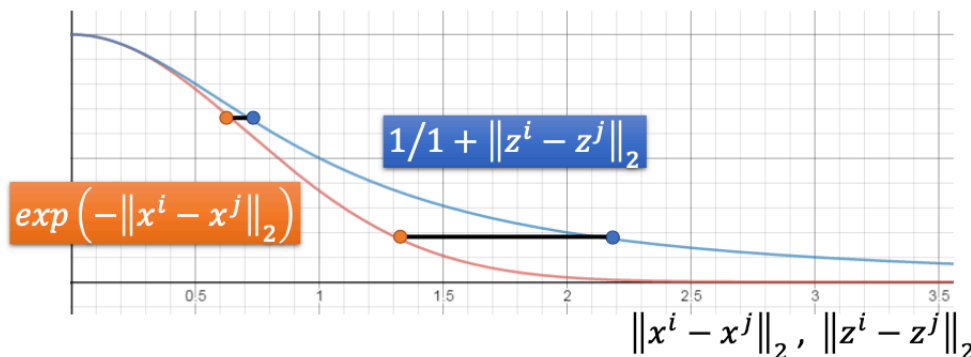
SNE:

$$S(x^i, x^j) = \exp(-\|x^i - x^j\|_2)$$

$$= \exp(-\|x^i - x^j\|_2)$$

t-SNE:

$$S'(z^i, z^j) = 1 / (1 + \|z^i - z^j\|_2)$$



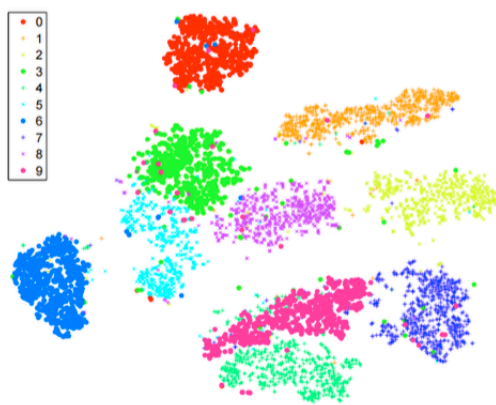
在上图中, 红色曲线表示SNE, 蓝色曲线表示t-SNE, 纵轴表示distribution; 如果我们想要维持相同的distribution, 即在同一个水平线上, 就达到了如图所示的效果; 相同的几率, t-SNE的 $\|z^i - z^j\|_2$ 之间的距离越大

如果本来就离得很近，那么经过t-SNE之间的距离还是很近；如果本来就离得很远，那么从原来的distribution拉到t-SNE之后，距离会更远；

到实际的例子中，如果本来是同一个类别的data，由于这些data之间的距离很近，不会收到t-SNE很大的影响；但如果是属于不同类别的data，距离是比较远的，t-SNE会放大这种距离

对于下图中的MNIST，先使用PCA进行降维，再进行可视化，就可以得到下图中的good visualization

- Good at visualization



t-SNE on MNIST



t-SNE on COIL-20

下图有一个更加直观的例子，使用t-SNE算法，运用gradient descent的思想，不同类别的data之间的距离会越来越大

