

10-605/805 – ML for Large Datasets

Lecture 4: Nonlinear Dimensionality Reduction

Henry Chai

9/8/22

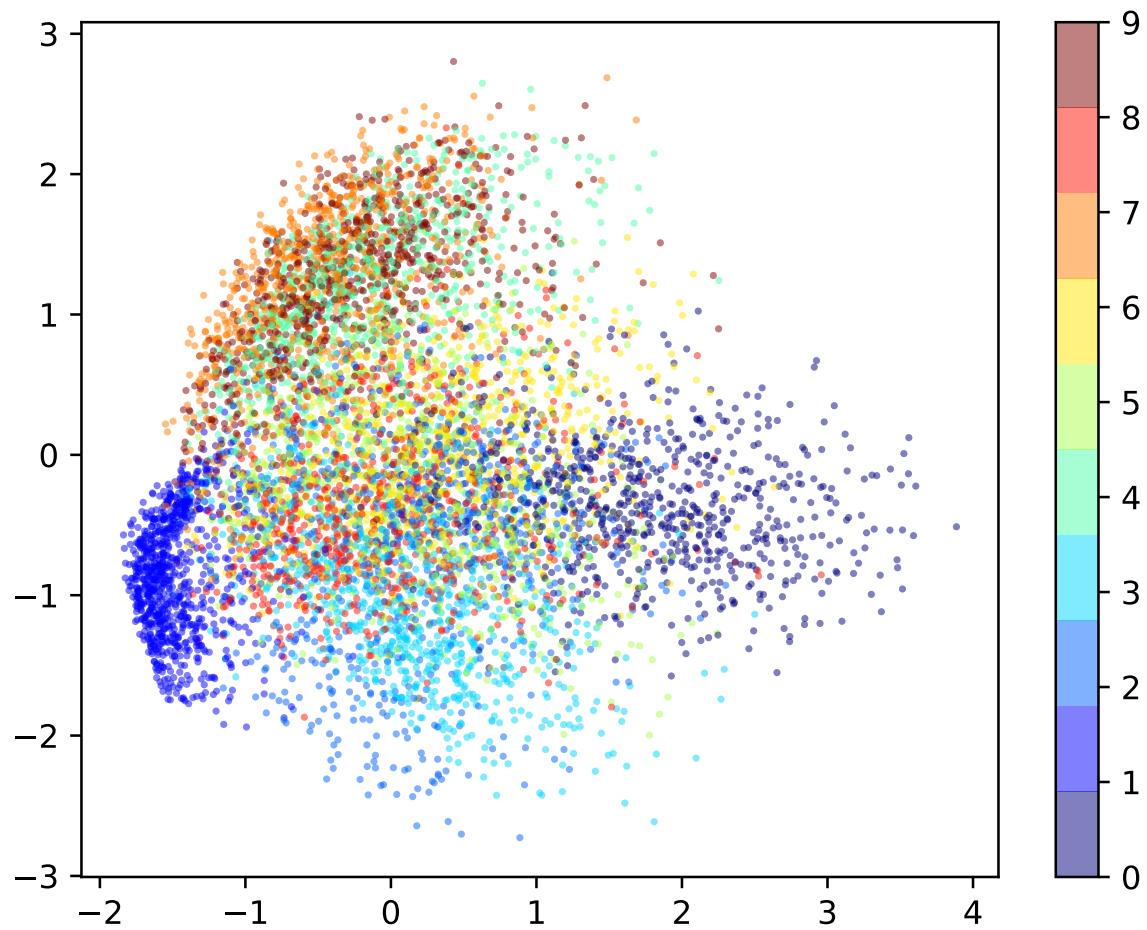
Front Matter

- HW1 released 8/30, due 9/13 at 11:59 PM
- HW2 released 9/8 (today!), due 9/22 at 11:59 PM
 - **We strongly encourage you to start the programming portion early and refer back to HW1 programming if you get stuck!**
- Recitations on Friday, 11:50 – 1:10 (**different from lecture**) in GHC 4401 (**same as lecture**)
 - Recitation 2 on 9/9: Review of linear algebra

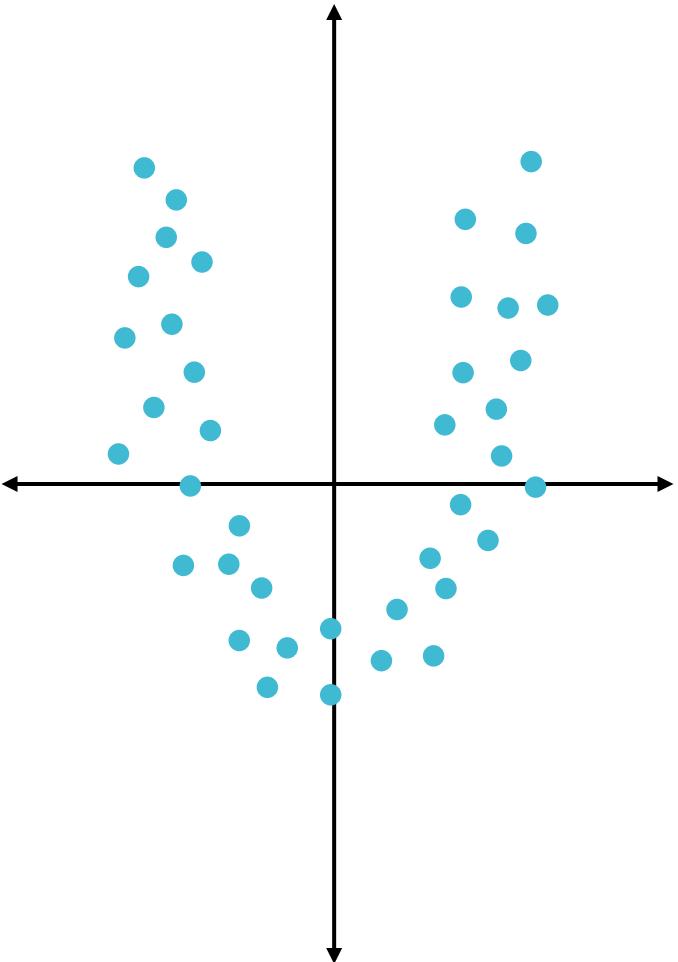
PCA Algorithm

- Input: $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n, r$
 1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
 2. Compute the covariance matrix $C_X = X^T X$
 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$
 4. Project the data into the space defined by P , $Z = XP$
- Output: Z , the latent representation (“PCA scores”)

PCA Example: MNIST Digits



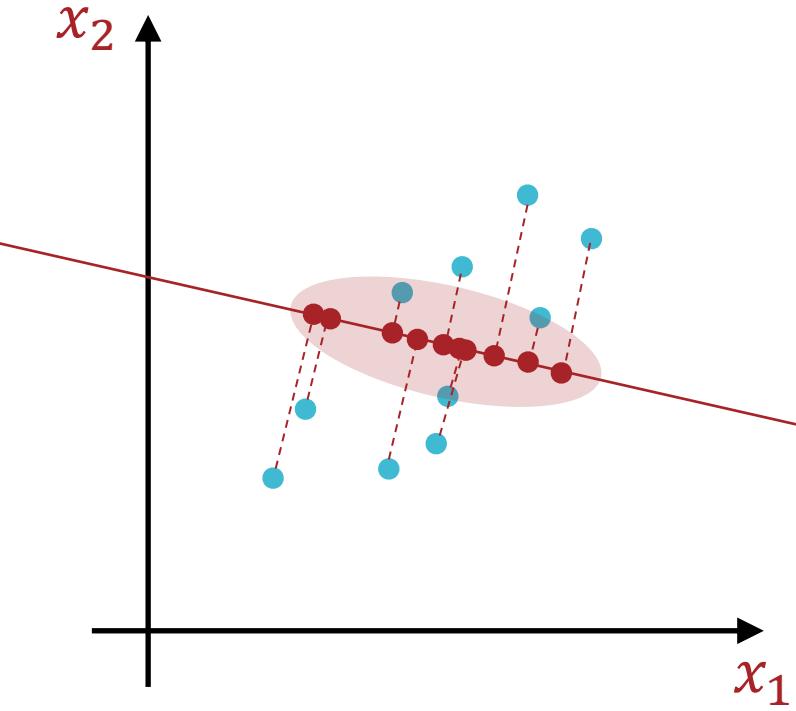
Shortcomings of PCA



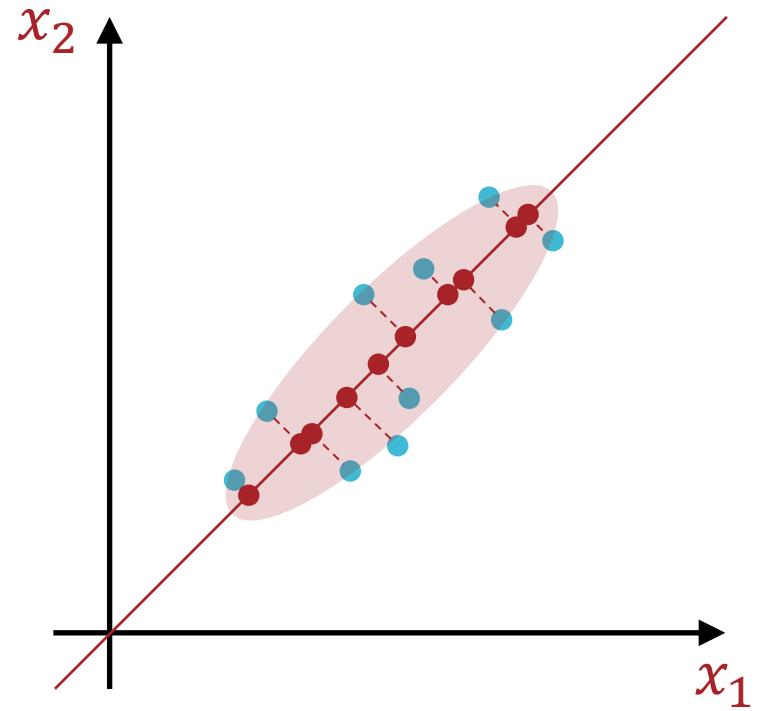
- Principal components are orthonormal
- Principal components are linear combinations of the features
- Principal components are expensive to compute

PCA Algorithm: Computational Cost

- Input: $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n, r$
1. Center the data
 - A. Optionally, normalize the data by features so that all features are of the same scale
 2. Compute the covariance matrix $C_X = X^T X$ ($O(\underline{nk^2})$)
 3. Collect the top r eigenvectors (corresponding to the r largest eigenvalues), $P \in \mathbb{R}^{k \times r}$ ($O(\underline{k^3})$)
 4. Project the data into the space defined by P , $Z = XP$ ($O(nkr)$)



Option A



Option B

Maybe Option A isn't so bad?

Random Projections?

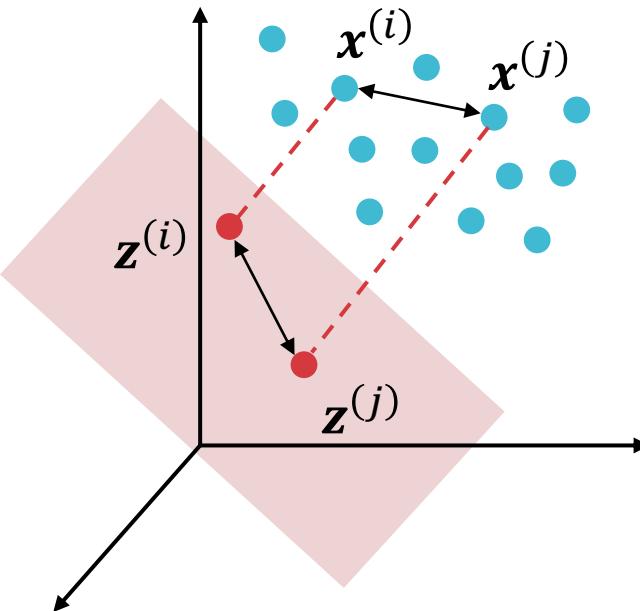
- Issue: when k is very large, computing principal components can be intractable
 - PCA also requires the entire dataset to be available, which might not be possible (e.g., online learning)
- Idea: instead of rigorously minimizing the reconstruction error/maximizing the projections' variance, just pick a random set of vectors to project onto!
 - Each element of P is sampled from a standard Gaussian distribution, $P_{i,j} \sim N(0,1)$
 - Projections are still given by $Z = XP$
- Inquiry: is this a good idea? And if so, how good?

Random Projections: Computational Cost

- Input: $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^n, r$
1. ~~Center the data~~
 - A. ~~Optionally, normalize the data by features so that all features are of the same scale~~ $O(nk)$
 2. ~~Compute the covariance matrix $C_x = X^T X$ ($O(nk^2)$)~~
 3. Generate P by sampling each element from a standard Gaussian distribution ($O(kr)$)
 4. Project the data into the space defined by P , $Z = XP$ ($O(nkr)$)
 $r \ll k$

Preserving Relative Distances

- Reasonable desideratum: distances between points are similar before and after being projected



Formally:

$$\|z^{(i)} - z^{(j)}\|_2^2 \in [(1-\epsilon)\|x^{(i)} - x^{(j)}\|_2^2, (1+\epsilon)\|x^{(i)} - x^{(j)}\|_2^2]$$

$\epsilon \in [0, 1]$

Johnson-Lindenstrauss Lemma

- Given a dataset $\mathcal{D} = \{(x^{(i)})\}_{i=1}^n$ and $0 < \epsilon < 1$, if

$$r \geq \frac{4 \ln(n)}{\epsilon^2/2 - \epsilon^3/3}$$

does not depend on k

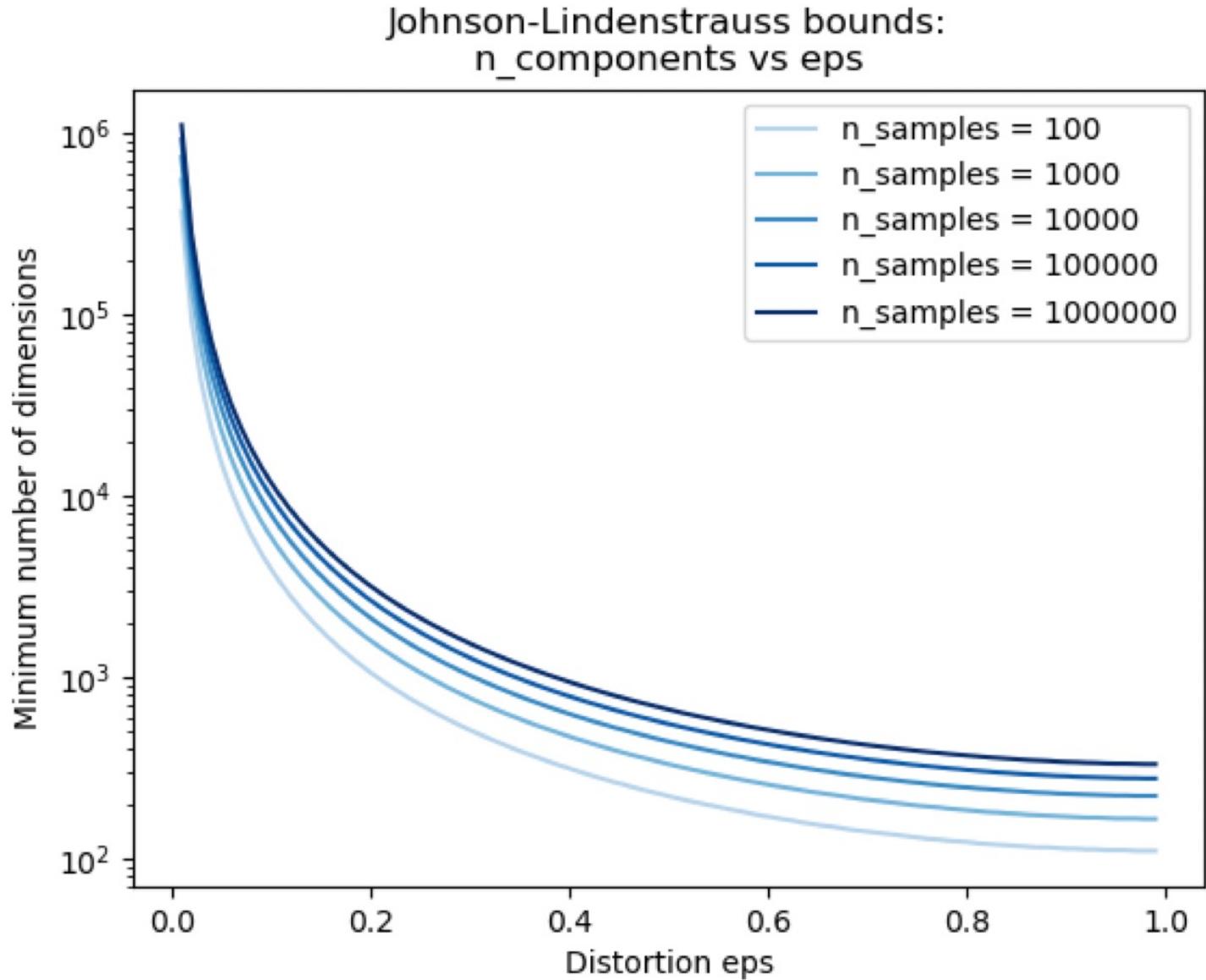
then \exists a linear map $f: \mathbb{R}^k \rightarrow \mathbb{R}^r$ s.t.

$$\|\mathbf{z}^{(i)} - \mathbf{z}^{(j)}\|_2^2 \in [(1 - \epsilon) \|x^{(i)} - x^{(j)}\|_2^2, (1 + \epsilon) \|x^{(i)} - x^{(j)}\|_2^2]$$

$$\forall x^{(i)}, x^{(j)} \in \mathcal{D}$$

$$f(x) = \underline{xP}$$

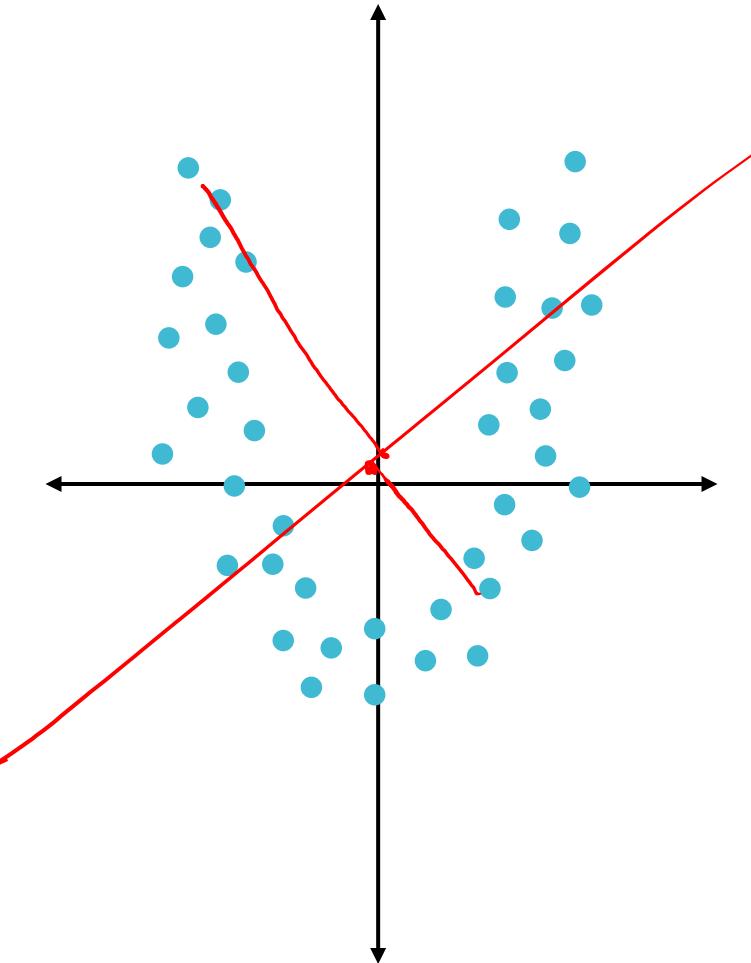
Johnson-Lindenstrauss Lemma



Johnson-Lindenstrauss Lemma: Proof (Sketch)

1. Construct the map f (set of projections) randomly, by sampling from a standard Gaussian
2. Show that the probability that f satisfies the desired property is > 0 (specifically, this probability is $\geq 1/n$)
 - use a concentration of measure argument
3. Conclude that such a map \exists and we can find an map that satisfies the desired property in expected polynomial time simply by resampling

Shortcomings of PCA



- Principal components are orthonormal
- Principal components are linear combinations of the features
- Principal components are expensive to compute

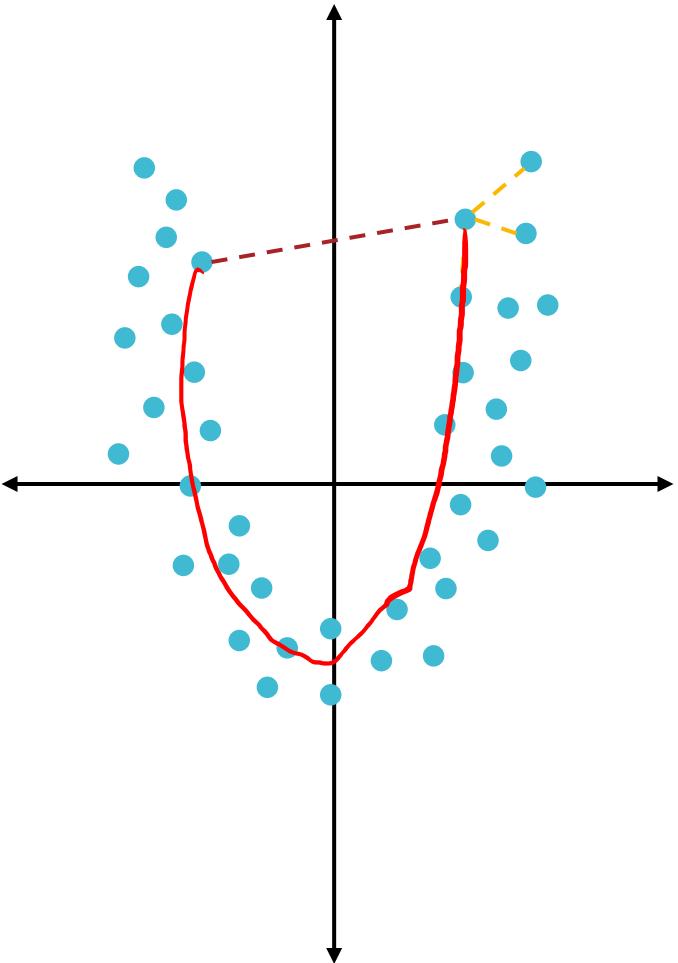
Nonlinear Dimensionality Reduction

- Autoencoders (1987)
- Kernel PCA (1999)
- Locally linear embedding (2000)
- Isomap (2000)
- Laplacian Eigenmaps (2003)
- t-SNE (2008)
- ... *many others*

Nonlinear Dimensionality Reduction

- Autoencoders (1987)
- Kernel PCA (1999)
- Locally linear embedding (2000)
- Isomap (2000)
- Laplacian Eigenmaps (2003)
- **t-SNE (2008)**
- ... *many others*

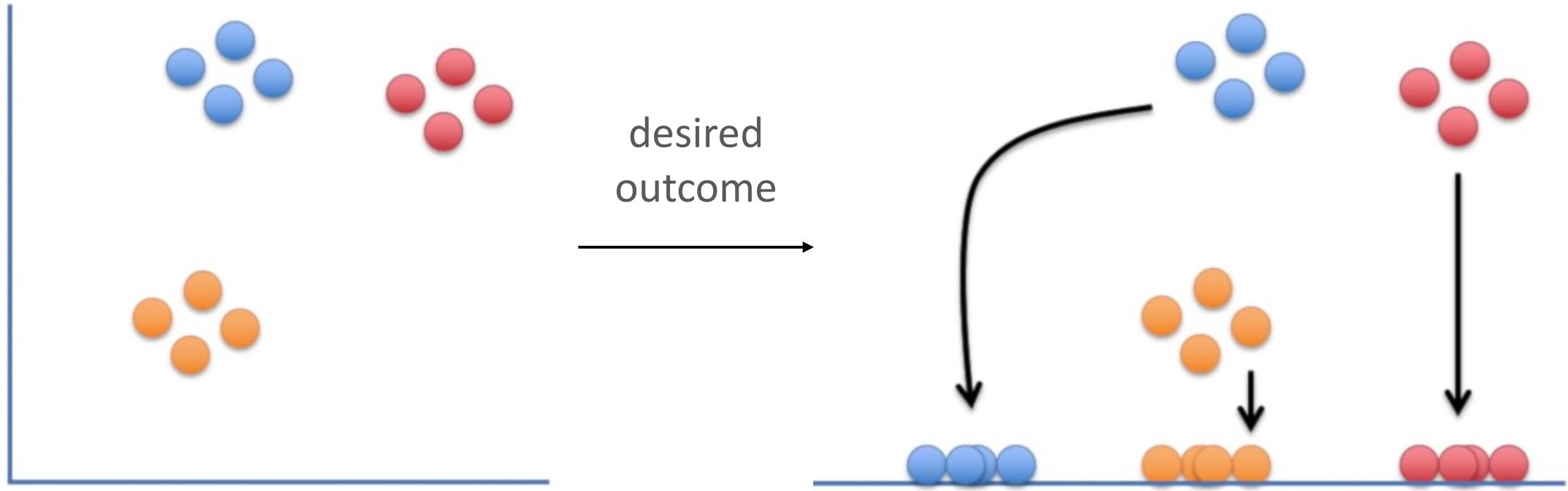
Shortcomings of PCA



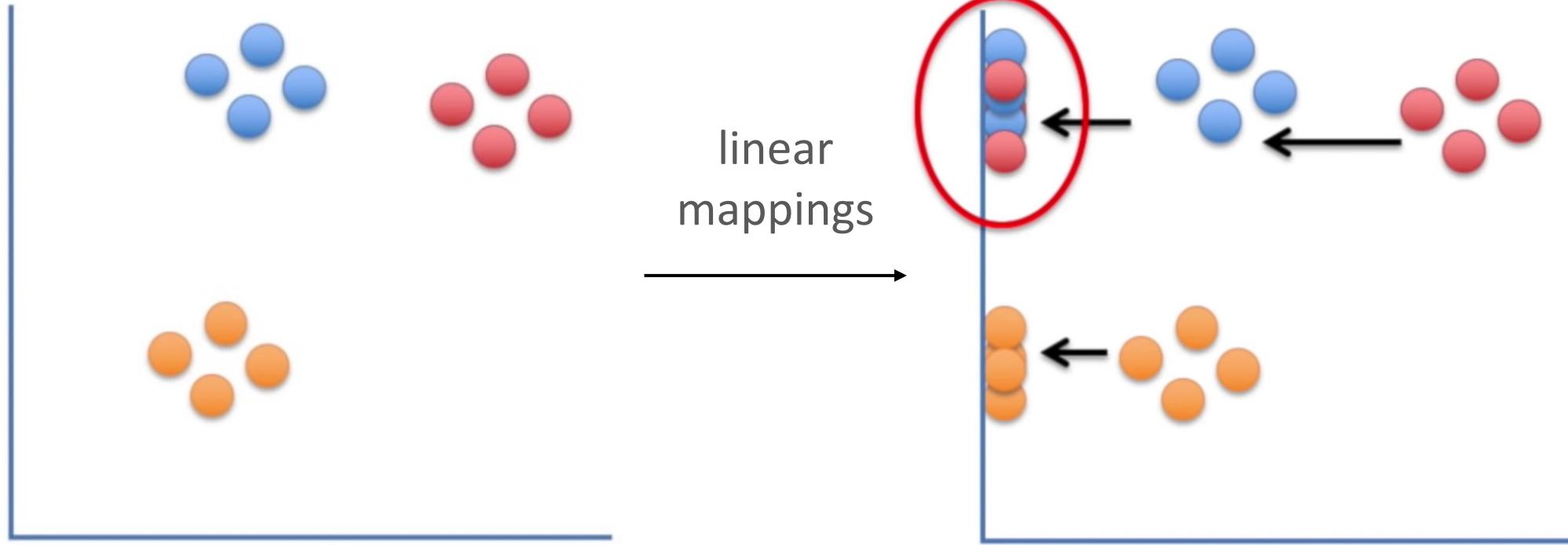
- Insight: PCA assumes that all distances between pairs of points are equally important
- Idea: Focus primarily on preserving “distances” between points in a small, local neighborhood!

t-distributed Stochastic Neighbor Embedding (t-SNE)

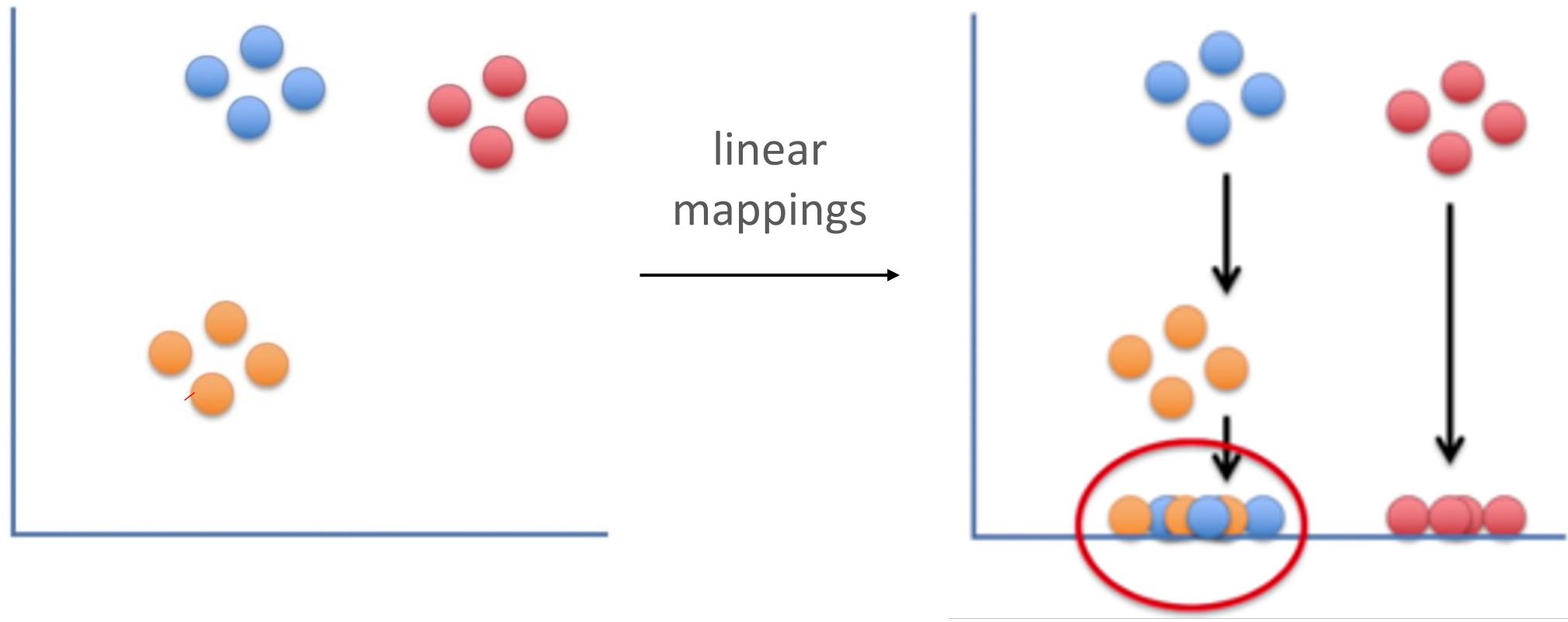
- Main idea: use the probability that two points are similar as the quantity to be preserved and *cleverly choose the distributions that we use*
 1. Compute the probability that two points are similar in the original space
 2. Randomly sample a projection for each point
 3. Compute the probability that two points are similar in the projected space
 4. Move points around in the projected space so that these two probabilities are “close” for all pairs of points



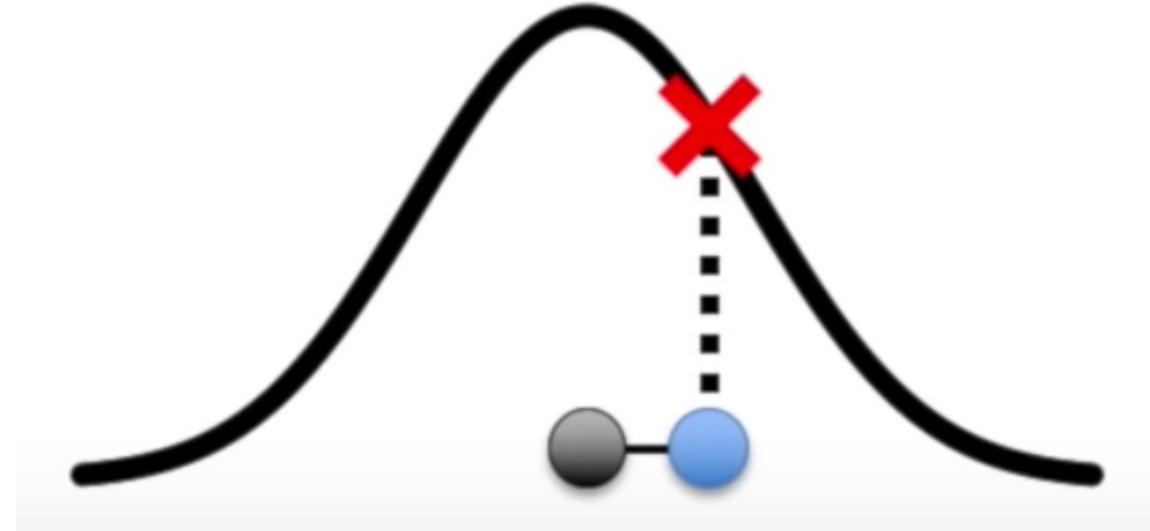
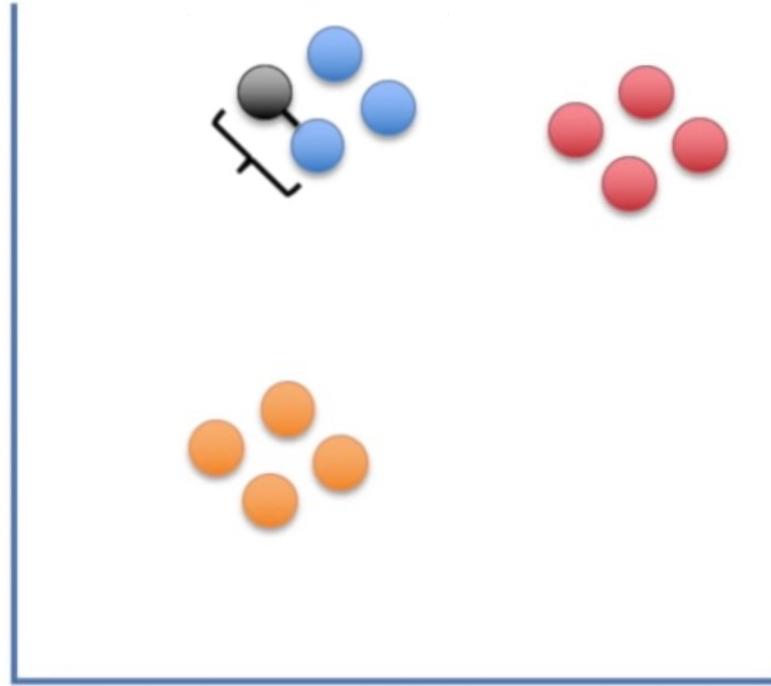
t-SNE: Example



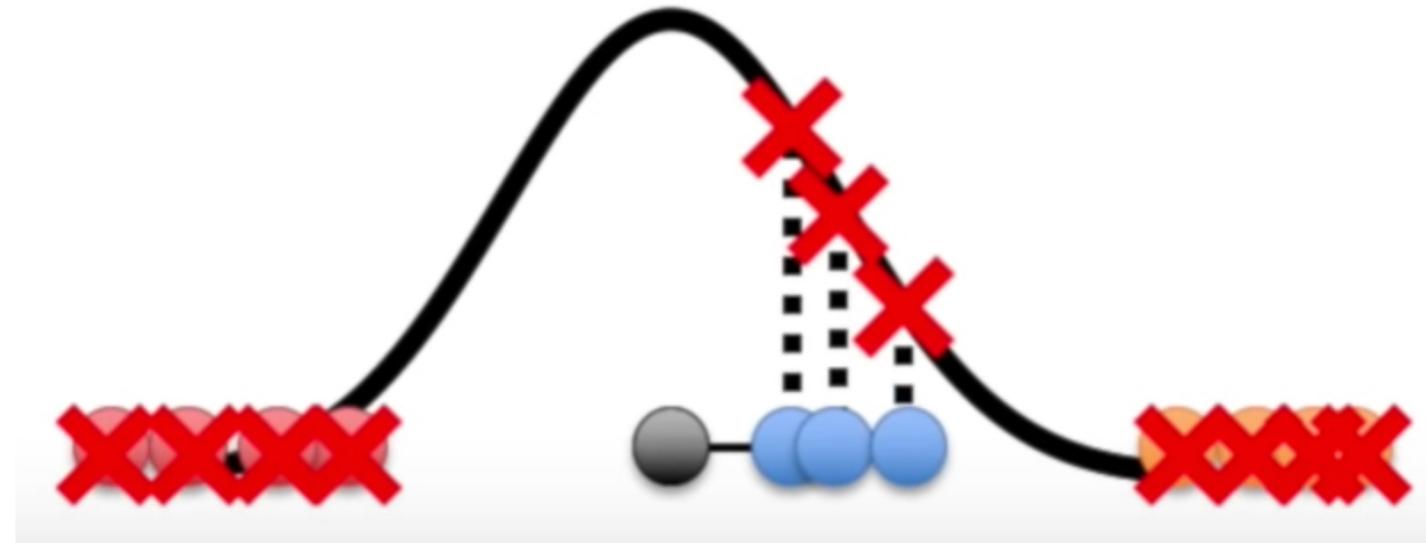
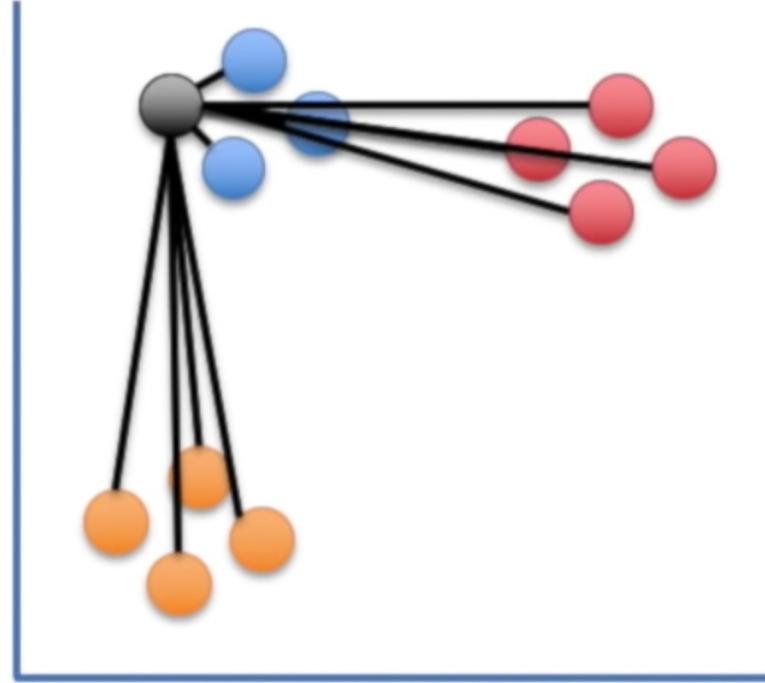
t-SNE: Example



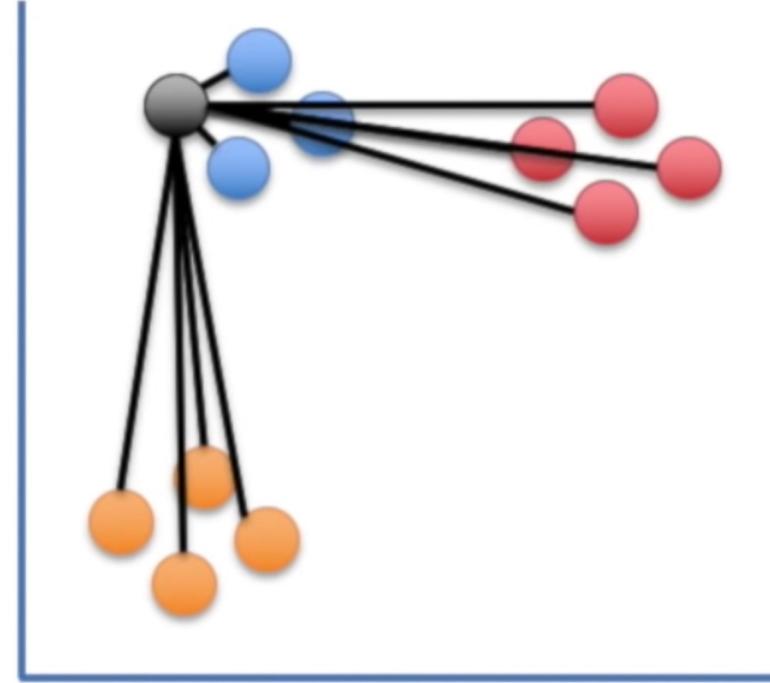
t-SNE: Example



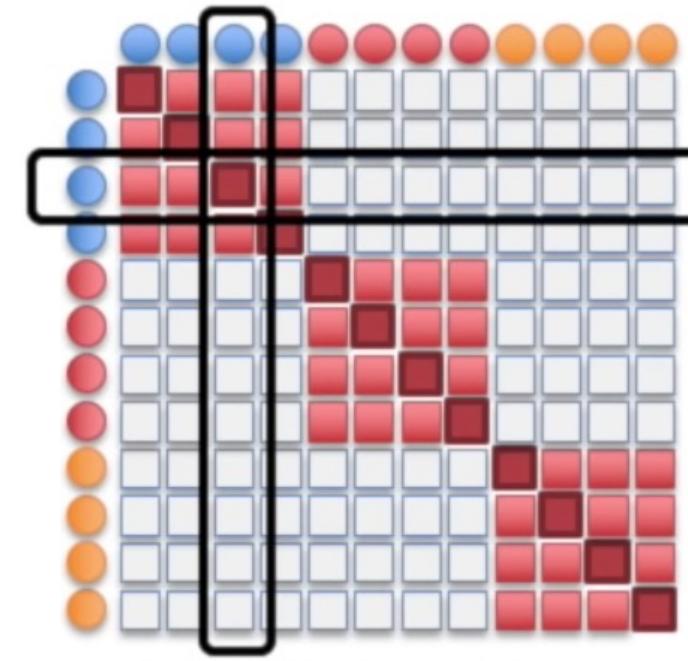
t-SNE: Example



t-SNE: Example



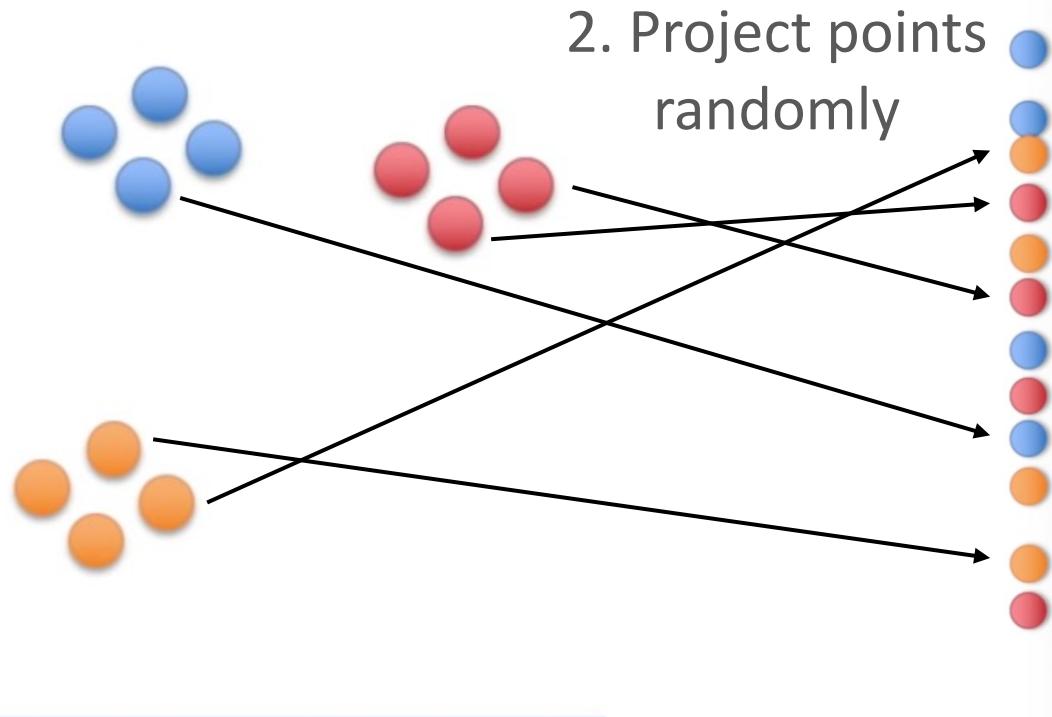
1. Compute probability of similarity in original space



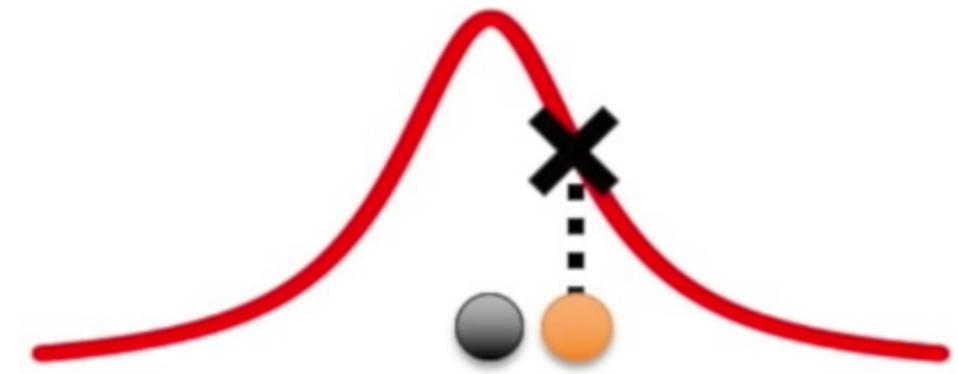
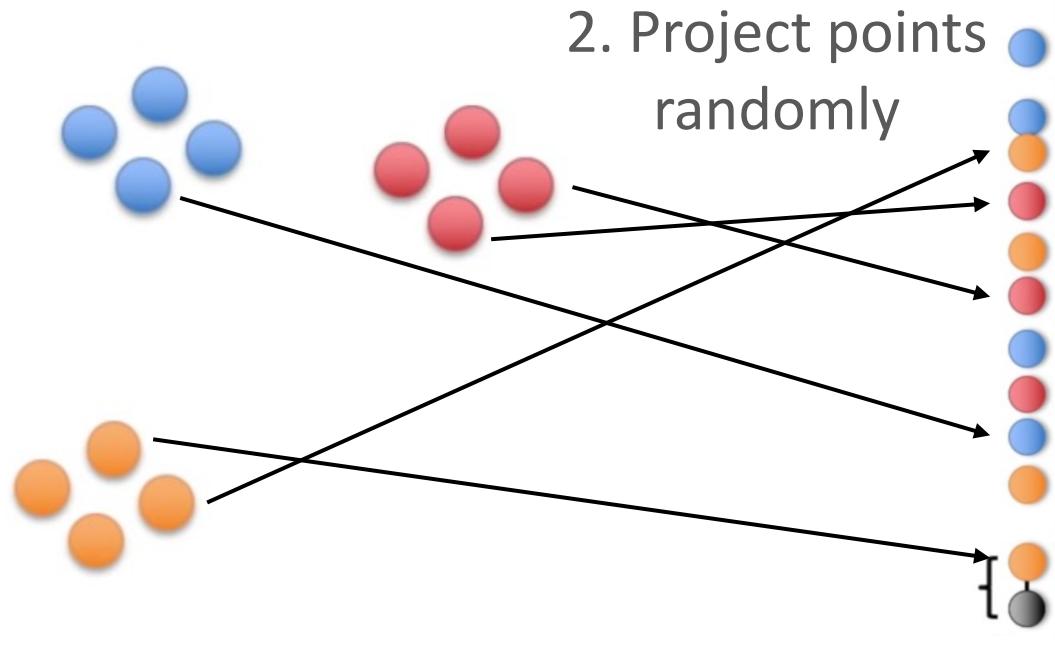
■ = High similarity

□ = Low similarity

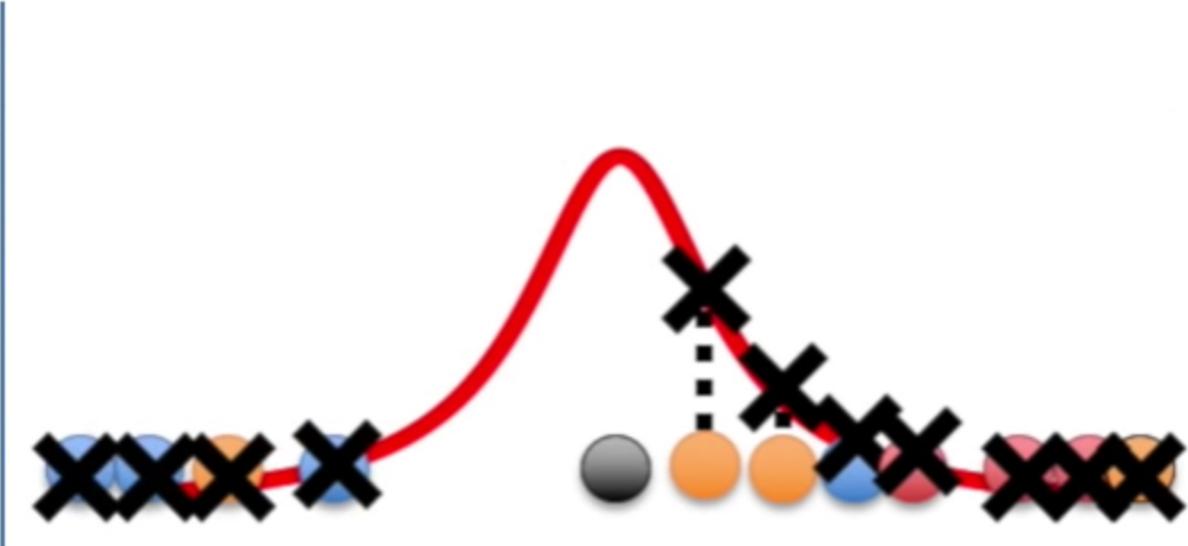
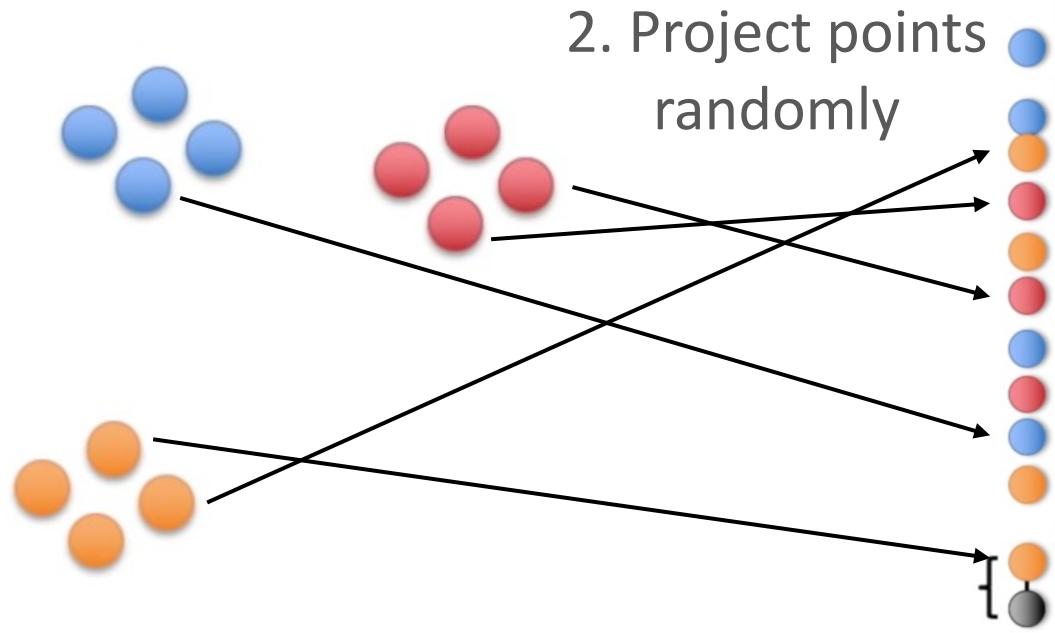
t-SNE: Example



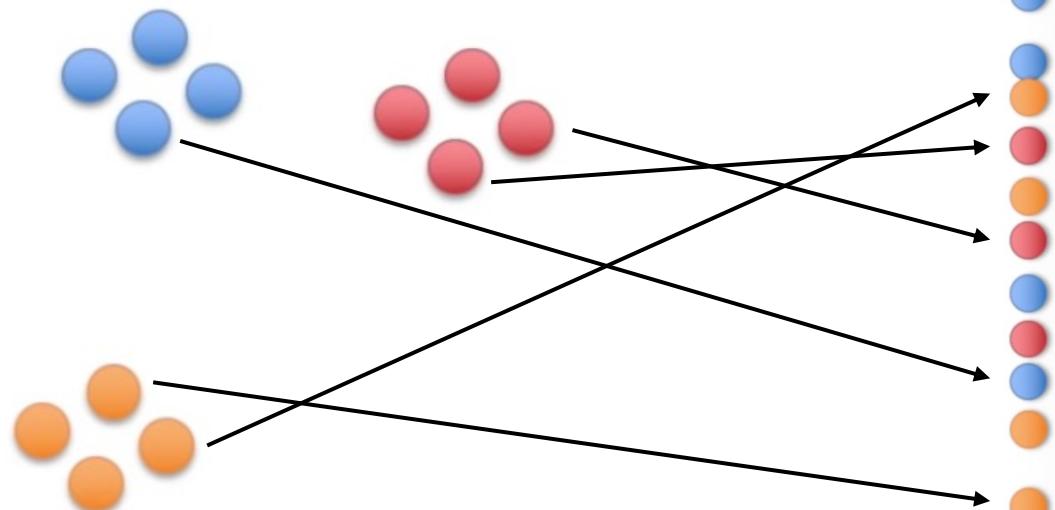
t-SNE: Example



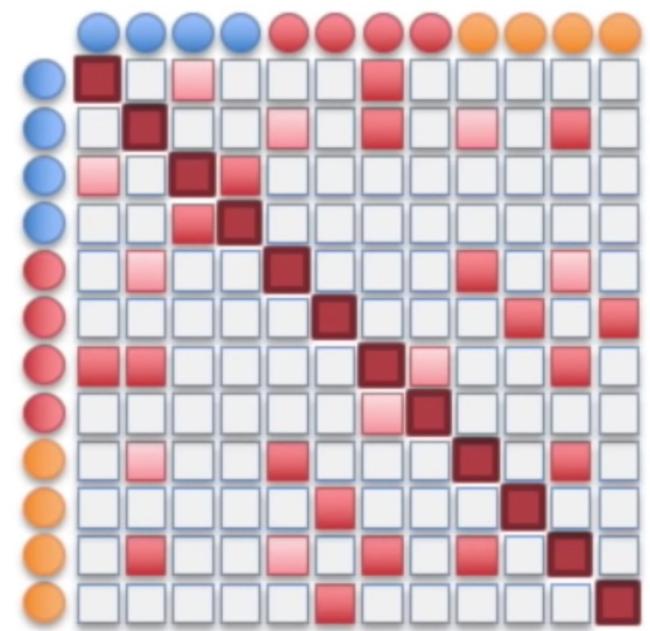
t-SNE: Example



t-SNE: Example

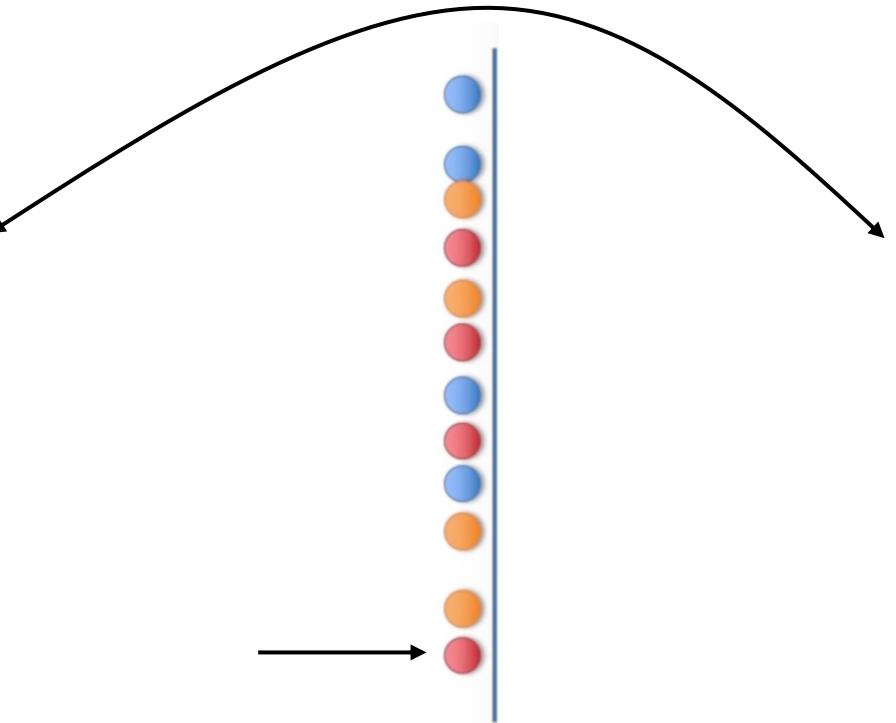
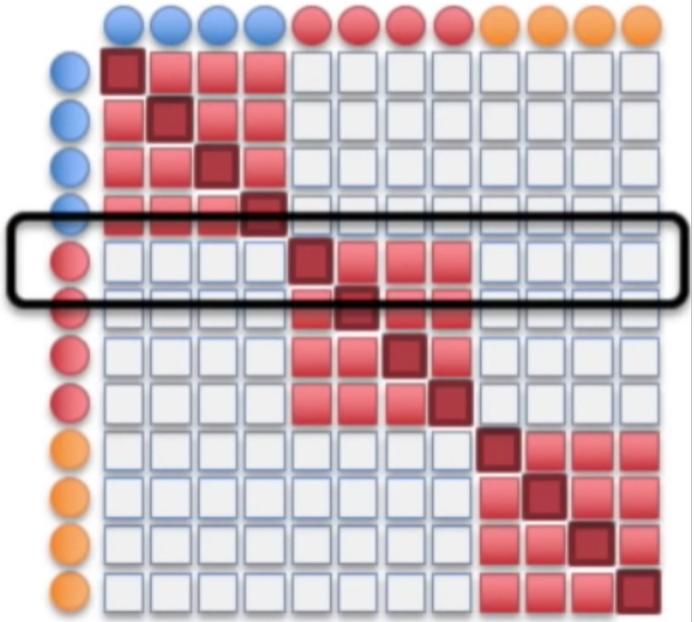


3. Compute probability of similarity in new space

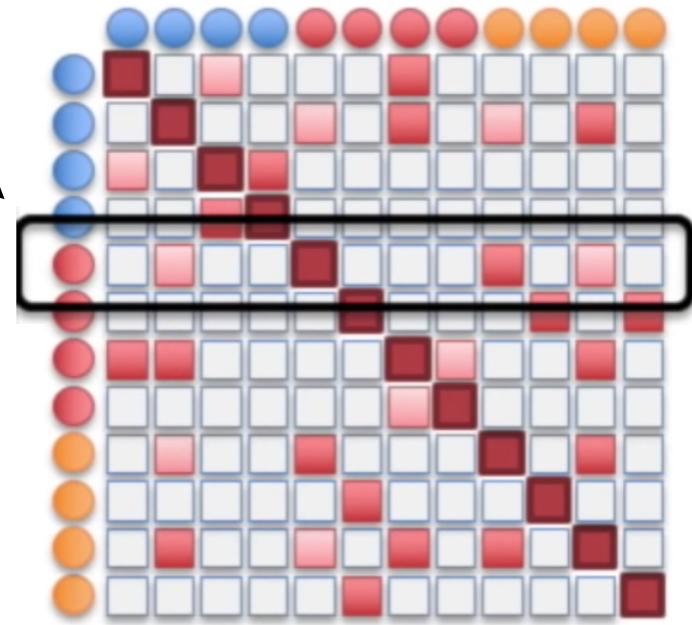


■ = High similarity
□ = Low similarity

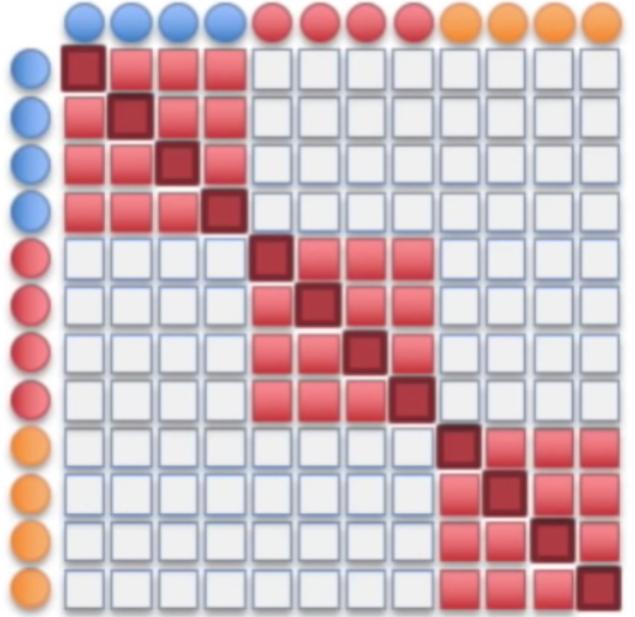
t-SNE: Example



4. Move points around to match the probabilities

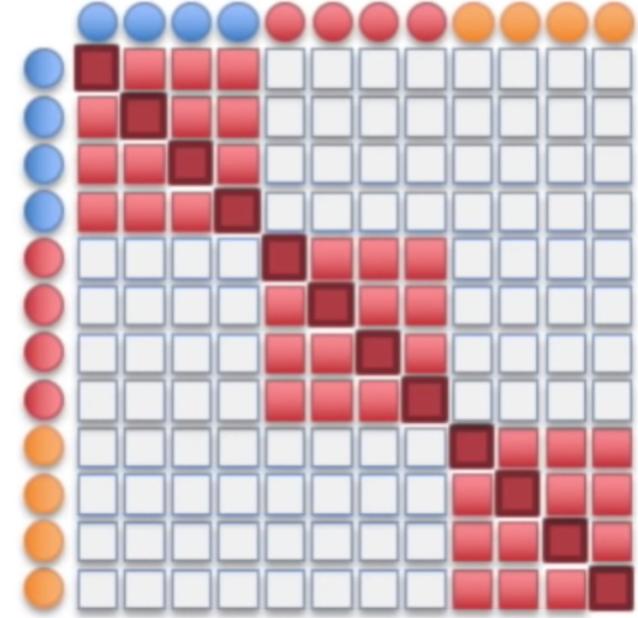


t-SNE: Example



t-SNE: Example

4. Move points around to match the probabilities



t-SNE

- Main idea: use the probability that two points are similar as the quantity to be preserved and *cleverly choose the distributions that we use*
 1. Compute the probability that two points are similar in the original space
 2. Randomly sample a projection for each point
 3. Compute the probability that two points are similar in the projected space
 4. Move points around in the projected space so that these two probabilities are “close” for all pairs of points

t-SNE

Aside:

$$P(x^{(j)} | x^{(i)}) P(x^{(i)}, x^{(j)}) = P(x^{(j)}, x^{(i)})$$

- Main idea: use the probability that two points are similar as the quantity to be preserved and *cleverly choose the distributions that we use*

- Compute the probability that two points are similar in the original space

Let $x^{(i)}$ and $x^{(j)}$ be points in the original space

$$P_{j|i} = \frac{\exp(-\|x^{(i)} - x^{(j)}\|_2^2 / (2\sigma_i^2))}{\sum_{k \neq i} \exp(-\|x^{(i)} - x^{(k)}\|_2^2 / (2\sigma_i^2))}$$

σ_i^2 is a hyperparam.

$$P_{i|i} = 0$$

$$P_{ij} := (P_{i|j} + P_{j|i}) / 2n$$

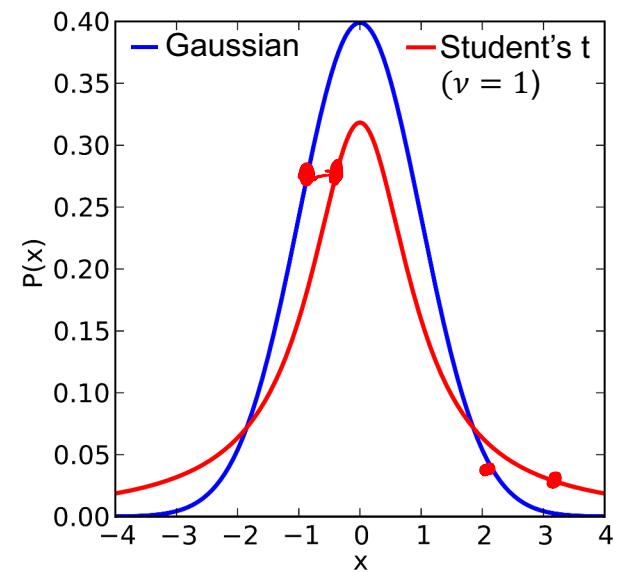
t-SNE

- Main idea: use the probability that two points are similar as the quantity to be preserved and *cleverly choose the distributions that we use*

- Compute the probability that two points are similar in the projected space

Let $z^{(i)} \sim z^{(j)}$ be points in the projected space

$$q_{ij} = \frac{(1 + \|z^{(i)} - z^{(j)}\|_2^2)^{-1}}{\sum_{k \in C} (1 + \|z^{(i)} - z^{(k)}\|_2^2)^{-1}}$$
$$q_{ii} = 0$$

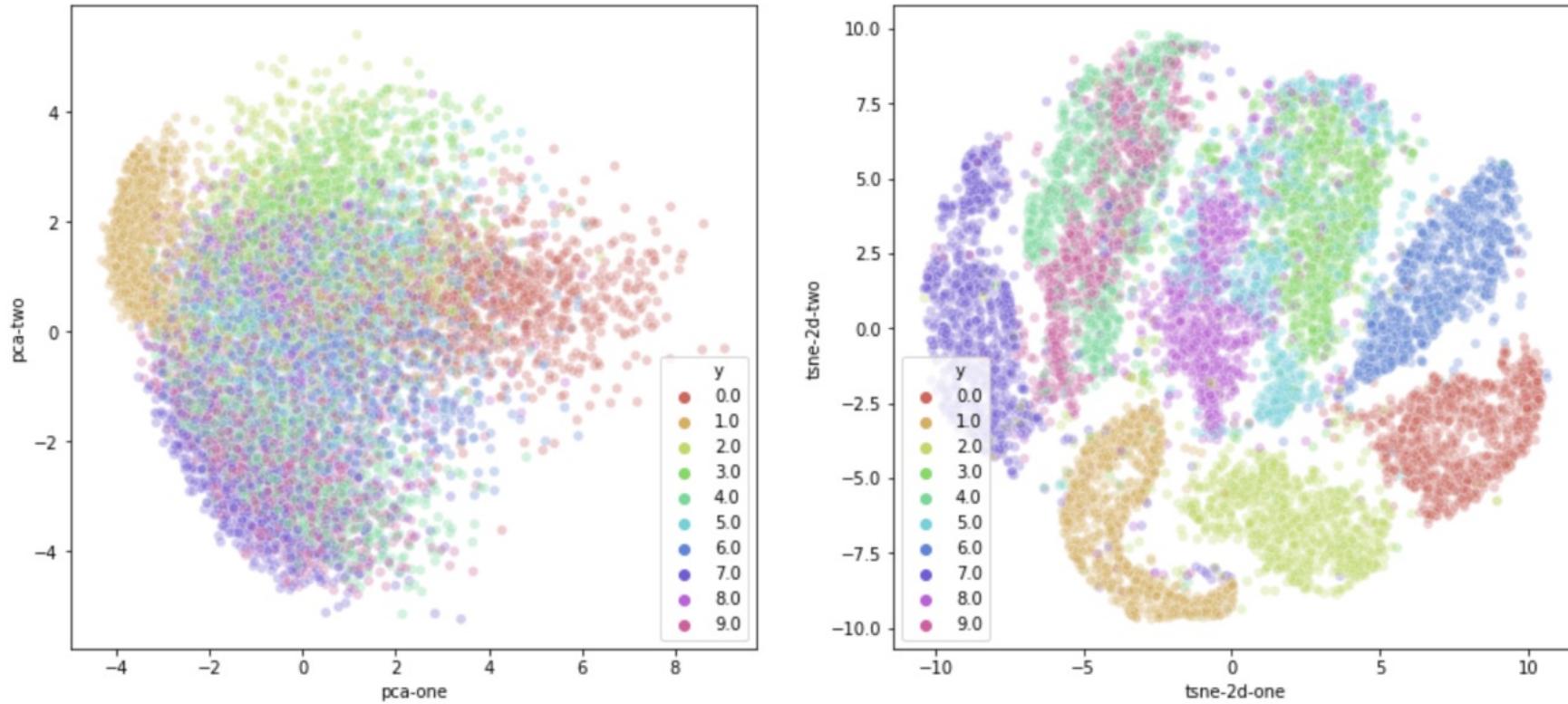


t-SNE

- Main idea: use the probability that two points are similar as the quantity to be preserved and *cleverly choose the distributions that we use*
4. Move points around in the projected space so that these two probabilities are “close” for all pairs of points

Minimize the KL-divergence; optimize via (S)GD

$$C = KL(P || Q) = \sum_i \sum_{j \neq i} P_{ij} \log \left(\frac{P_{ij}}{Q_{ij}} \right)$$
$$\frac{\partial C}{\partial z^{(i)}} = 4 \sum_j (P_{ij} - Q_{ij}) \frac{(1 + \|z^{(i)} - z^{(j)}\|_2^2)^{-1}}{(z^{(i)} - z^{(j)})}$$



t-SNE vs. PCA on MNIST

t-SNE vs. PCA

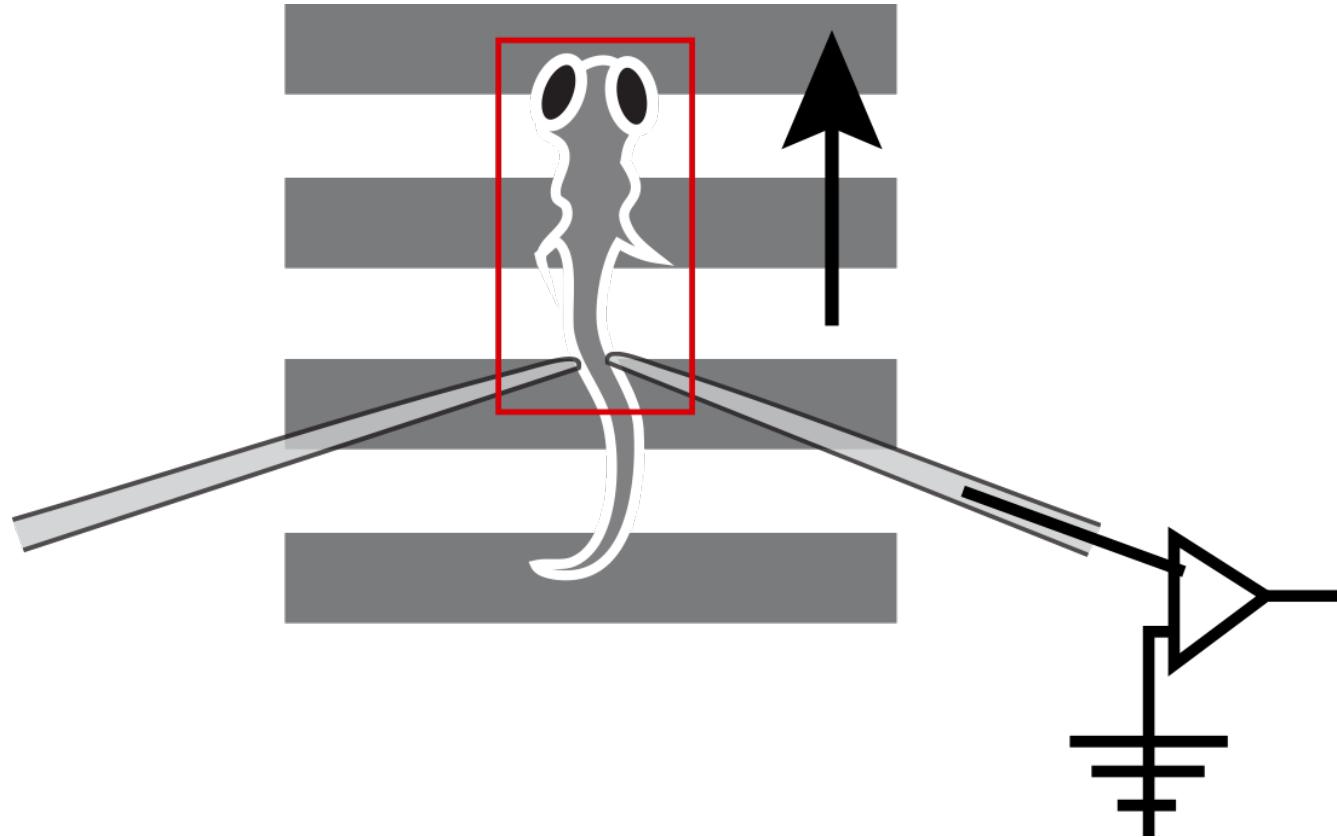
- PCA
 - ✓ Simple to use: no hyperparameters
 - Computationally expensive: $O(k^3)$
 - ✓ Stay tuned for distributed PCA!
 - Principal components are constrained to be *orthonormal* vectors (linear in the existing features)
- t-SNE
 - ✓ Can learn non-linear representations
 - Computationally expensive: $O(n^2)$ per iteration
 - Involves non-convex optimization
 - Sensitive to hyperparameters
(<https://distill.pub/2016/misread-tsne/>)

Key Takeaways

- Random sampling as an efficient alternative to PCA for finding linear projections
 - Theoretical guarantees regarding preserving relative distances via the Johnson-Lindenstrauss lemma
- t-SNE learns non-linear latent representations by preserving “distances” or matching probability distributions
 - Using the Student’s t distribution encourages spread in the latent space
 - KL-divergence can be minimized using (stochastic) gradient descent

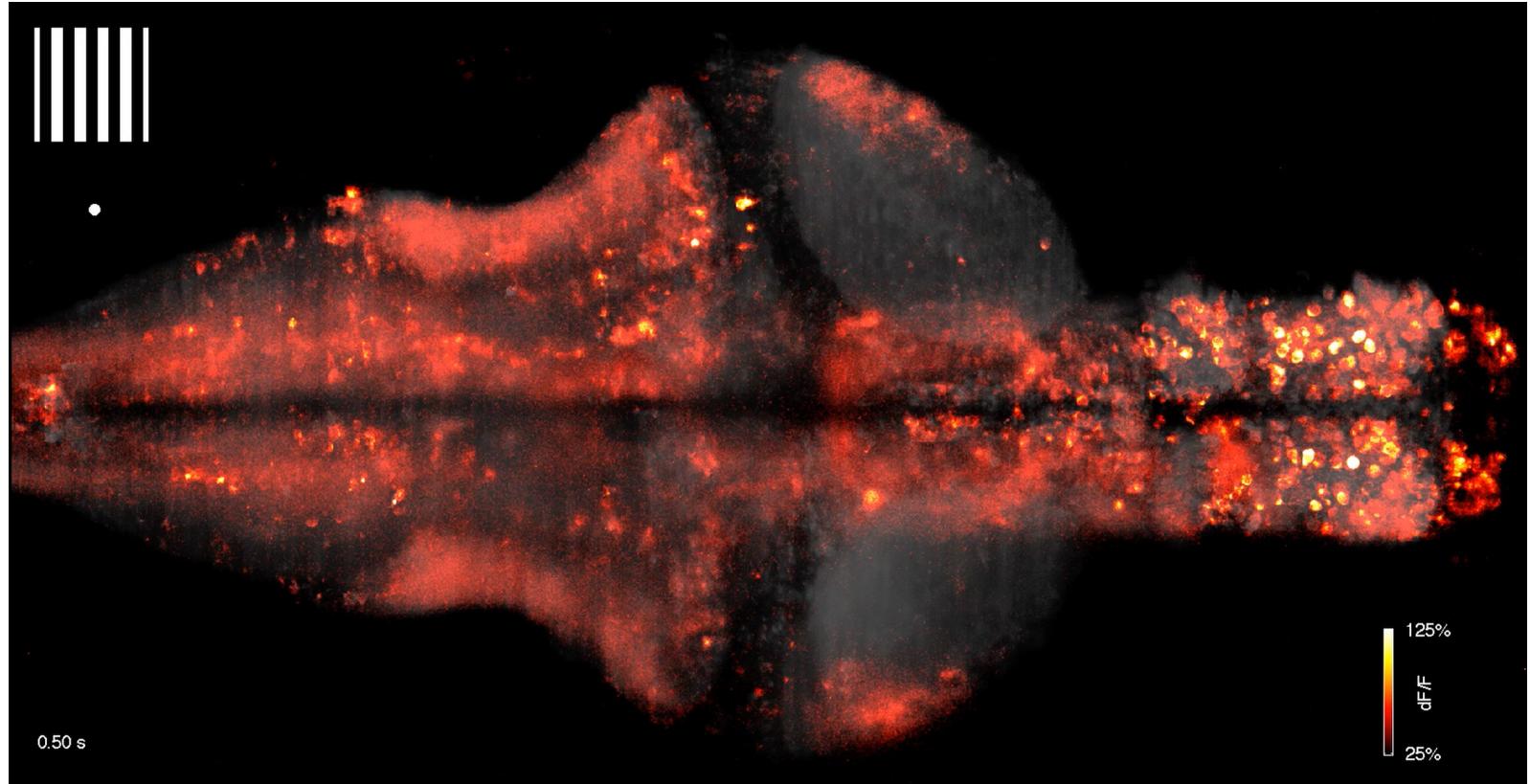
HW2 Preview

- Visualizing neural stimuli in larval zebrafish



HW2 Preview

- Visualizing neural stimuli in larval zebrafish

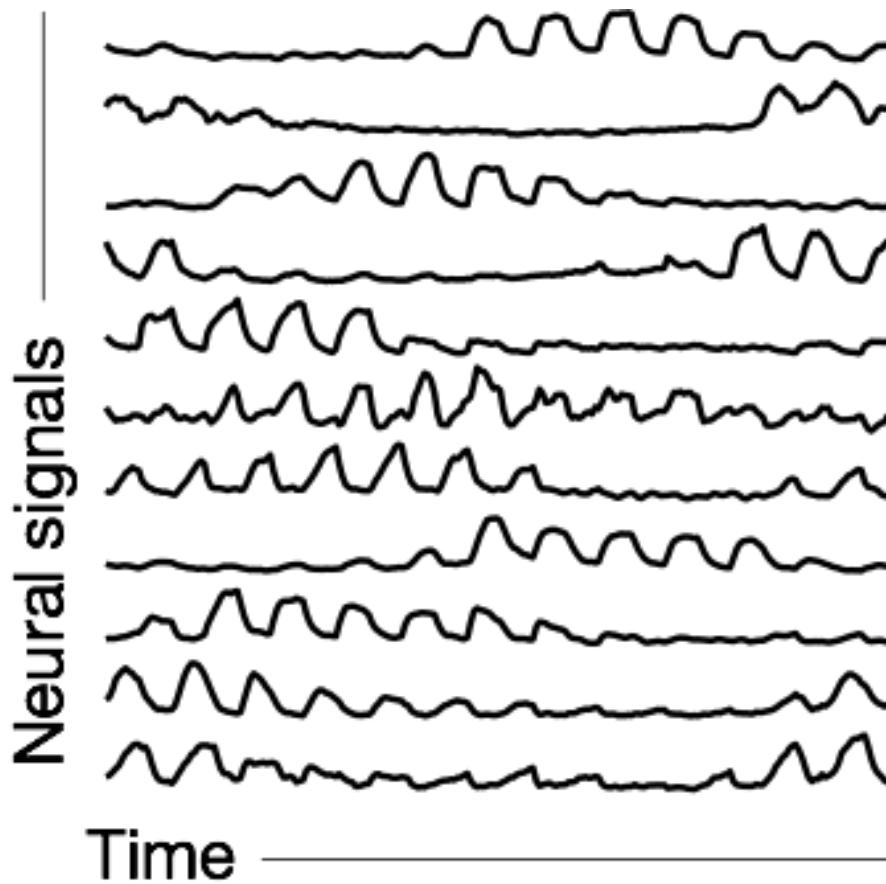


HW2 Preview

- Visualizing neural stimuli in larval zebrafish
- Mouse somatosensory cortex:
 - ~1,000 neurons
 - 0.1 TB/experiment
- Larval zebrafish brain:
 - ~100,000 neurons
 - 1 TB/experiment
- Mouse brain:
 - ~80,000,000 neurons
 - > 100 TB/experiment

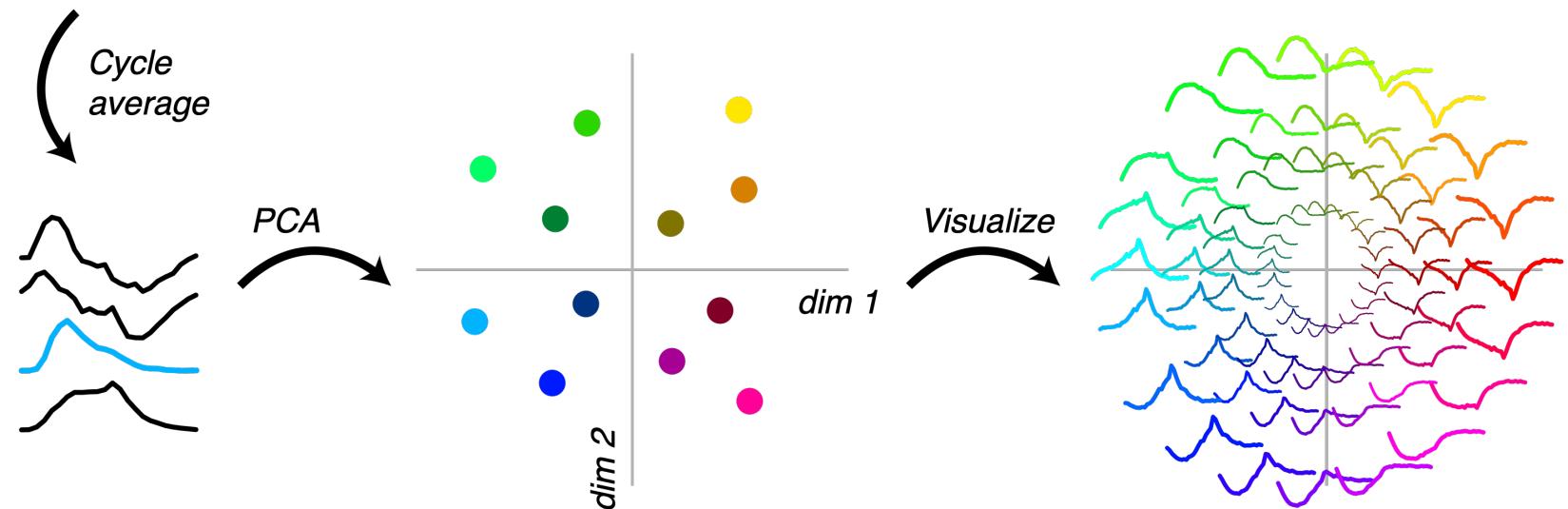
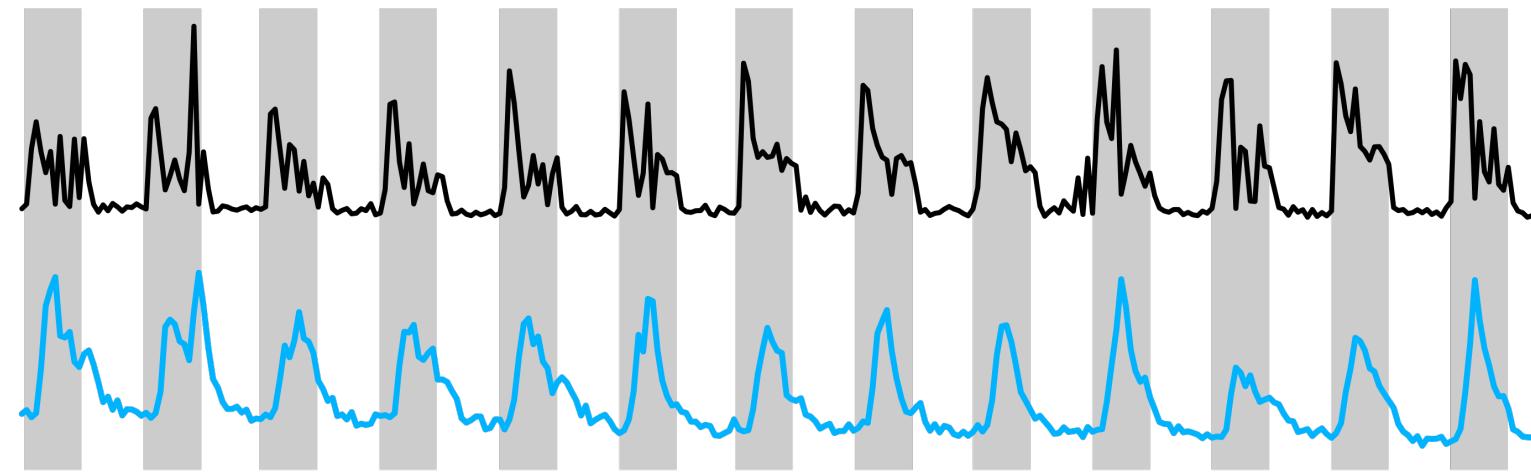
HW2 Preview

- Visualizing neural stimuli in larval zebrafish



HW2 Preview

- Visualizing neural stimuli in larval zebrafish



HW2 Preview

- Visualizing neural stimuli in larval zebrafish

