

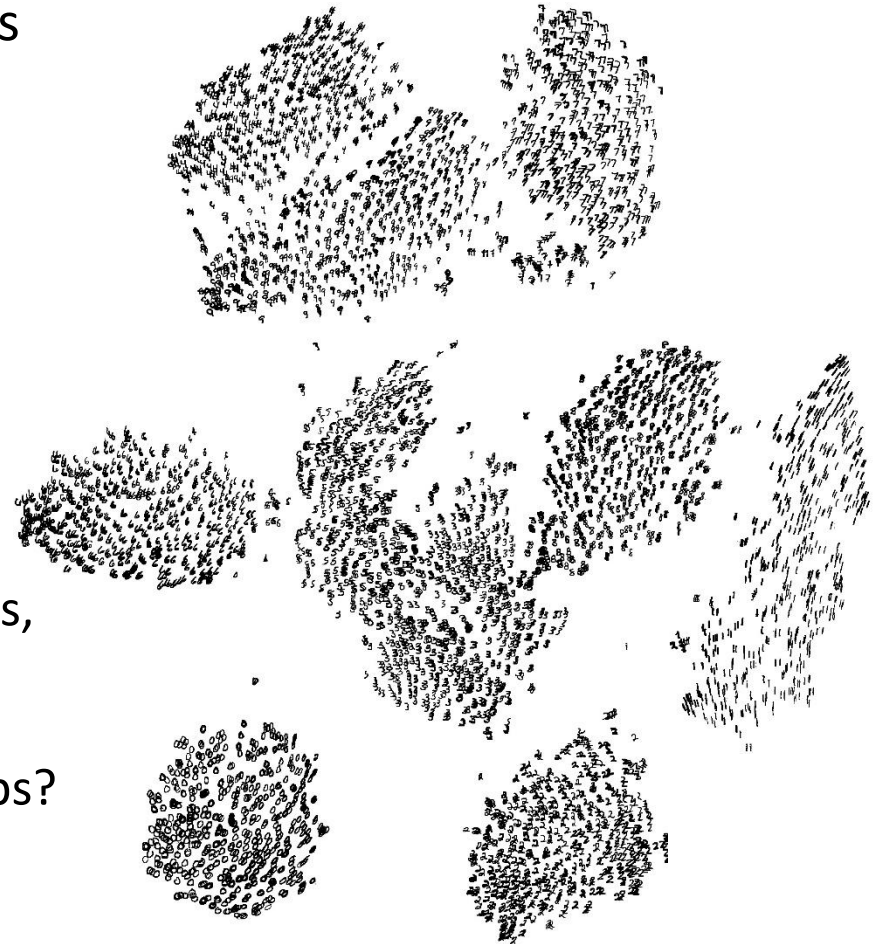
COMS 4771
Dimensionality Reduction

Nakul Verma

Example: Handwritten digits

Handwritten digit data, but with no labels

0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9



What can we do?

- Suppose know that there are 10 groupings, can we find the groups?
- What if we don't know there are 10 groups?
- How can we discover/explore other structure in such data?

A 2D visualization of digits dataset

Dimensionality Reduction

Data: $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$

Goal: find a 'useful' transformation $\phi : \mathbf{R}^d \rightarrow \mathbf{R}^k$ that helps in the downstream prediction task.

Some previously seen useful transformations:

- z-scoring $(x_1, \dots, x_d) \mapsto \left(\frac{x_1 - \mu_1}{\sigma_1}, \dots, \frac{x_d - \mu_d}{\sigma_d} \right)$

*Keeps same dimensionality
but with better scaling*

- Kernel transformations.

*Higher dimensionality,
making data linearly separable*

What are other desirable feature transformations?

How about lower dimensionality while keeping the relevant information?

Principal Components Analysis (PCA)

Data: $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n \in \mathbf{R}^d$

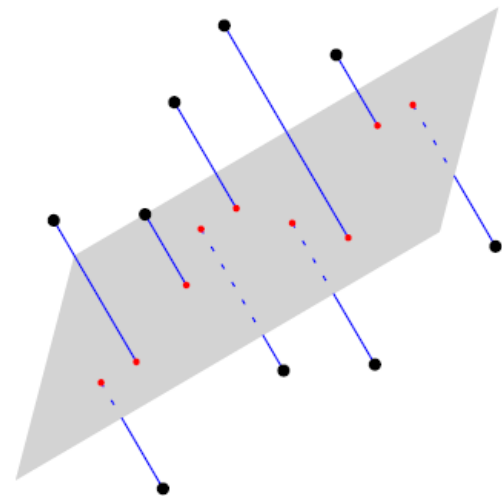
Goal: find the best **linear** transformation $\phi : \mathbf{R}^d \rightarrow \mathbf{R}^k$ that best maintains reconstruction accuracy.

Equivalently, minimize aggregate residual error

Define: $\Pi^k : \mathbf{R}^d \rightarrow \mathbf{R}^d$ *k-dimensional orthogonal linear projector*

$$\underset{\Pi^k}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \left\| \vec{x}_i - \Pi^k(\vec{x}_i) \right\|^2$$

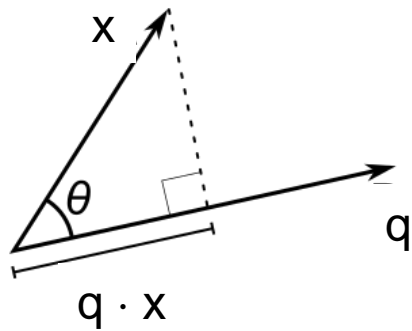
How do we optimize this?



Dimensionality Reduction via Projections

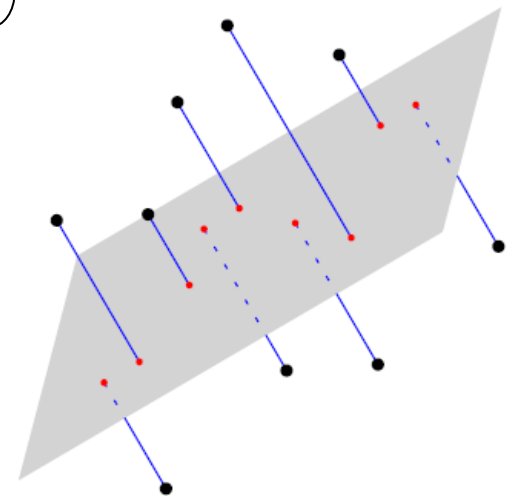
A k dimensional subspace can be represented by $\vec{q}_1, \dots, \vec{q}_k \in \mathbf{R}^d$ orthonormal vectors.

The projection of any $\vec{x} \in \mathbf{R}^d$ in the $\text{span}(\vec{q}_1, \dots, \vec{q}_k)$ is given by



$$\sum_{i=1}^k (\vec{q}_i \cdot \vec{x}) \vec{q}_i = \underbrace{\left(\sum_{i=1}^k \vec{q}_i \vec{q}_i^T \right)}_{\Pi^k} \vec{x}$$

To represent it in \mathbf{R}^k (using basis $\vec{q}_1, \dots, \vec{q}_k$) the coefficients simply are: $(\vec{q}_1 \cdot \vec{x}), \dots, (\vec{q}_k \cdot \vec{x})$



PCA: $k = 1$ case

If projection dimension $k = 1$, then looking for a q such that

$$\text{minimize}_{\|q\|=1} \frac{1}{n} \sum_{i=1}^n \|\vec{x}_i - (\vec{q} \vec{q}^\top) \vec{x}_i\|^2$$

$$\frac{1}{n} \sum_{i=1}^n \|\vec{x}_i - (\vec{q} \vec{q}^\top) \vec{x}_i\|^2 = \left(\frac{1}{n} \sum_{i=1}^n \vec{x}_i^\top \vec{x}_i \right) - \vec{q}^\top \left(\frac{1}{n} \sum_{i=1}^n \vec{x}_i \vec{x}_i^\top \right) \vec{q}$$

$$\propto - \vec{q}^\top \left(\frac{1}{n} X X^\top \right) \vec{q}$$

$$\left(X := \begin{bmatrix} | & & | \\ x_1 & \cdots & x_n \\ | & & | \end{bmatrix} \right)$$

Equivalent formulation:

$$\text{maximize}_{\|q\|=1} \vec{q}^\top \left(\frac{1}{n} X X^\top \right) \vec{q}$$

How to solve?

Optimizing $q^T M q$

For a symmetric PSD matrix M , how to optimize for $\text{maximize}_{\|q\|=1} q^T M q$

Recall, since M is symmetric PSD, by spectral decomposition theorem:

$$M = \sum_i \lambda_i v_i v_i^T \quad \begin{array}{l} v_1, \dots, v_d \text{ are an orthonormal set of eigenvectors of } M \\ \lambda_1 \geq \dots \geq \lambda_d \geq 0 \text{ are the corresponding eigenvalues} \end{array}$$

Thus for any unit length q , $q^T M q = \sum_i \lambda_i (q \cdot v_i)^2$ (where $\sum_i (q \cdot v_i)^2 = 1$)

Let, $\alpha_i := (q \cdot v_i)^2$ then the optimization becomes

$$\max_{\alpha_i} \sum_i \lambda_i \alpha_i \quad \text{s.t.} \quad \sum_i \alpha_i = 1$$

with the optimal solution as $\alpha_1=1, \alpha_2 = \alpha_3 = \dots = \alpha_n = 0$, or equivalently $q = v_1$

Therefore

$$\text{maximize}_{\|q\|=1} \vec{q}^T \left(\frac{1}{n} X X^T \right) \vec{q}$$

*is maximized by the top
eigenvector of matrix $(1/n) X X^T$!*

PCA: $k = 1$ case

$$\text{maximize}_{\|q\|=1} \vec{q}^T \left(\frac{1}{n} X X^T \right) \vec{q}$$

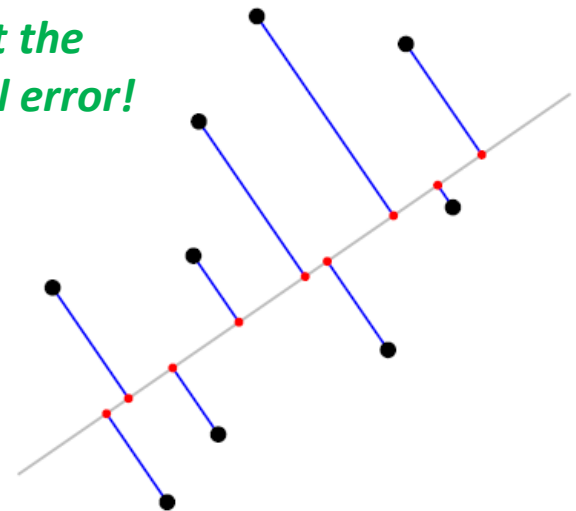
Covariance of data (if mean = 0)

For any q the quadratic form $\vec{q}^T \left(\frac{1}{n} X X^T \right) \vec{q}$ is the empirical

variance of data in the direction q , ie, of data $\vec{q}^T \vec{x}_1, \dots, \vec{q}^T \vec{x}_n$ *why?*

Therefore, the top eigenvector solution implies that the direction of maximum variance minimizes the residual error!

What about general k ?



PCA: general k case

$$\arg \min_{\substack{Q \in \mathbf{R}^{d \times k} \\ Q^T Q = I}} \frac{1}{n} \sum_{i=1}^n \|\vec{x}_i - QQ^T \vec{x}_i\|^2 = \arg \min_{\substack{Q \in \mathbf{R}^{d \times k} \\ Q^T Q = I}} \frac{1}{n} \|X - QQ^T X\|_F^2$$

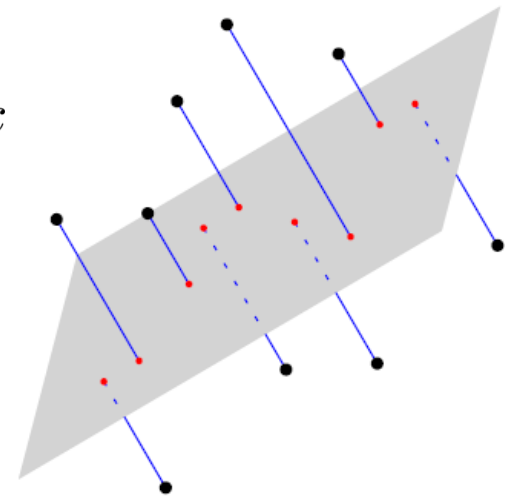
$$Q := \begin{bmatrix} | & & | \\ q_1 & \cdots & q_k \\ | & & | \end{bmatrix}$$

$$\propto \arg \max_{\substack{Q \in \mathbf{R}^{d \times k} \\ Q^T Q = I}} \text{tr} \left(Q^T \left(\frac{1}{n} X X^T \right) Q \right)$$

*Solution: the top k eigenvectors
of the matrix $(1/n)XX^T$!*

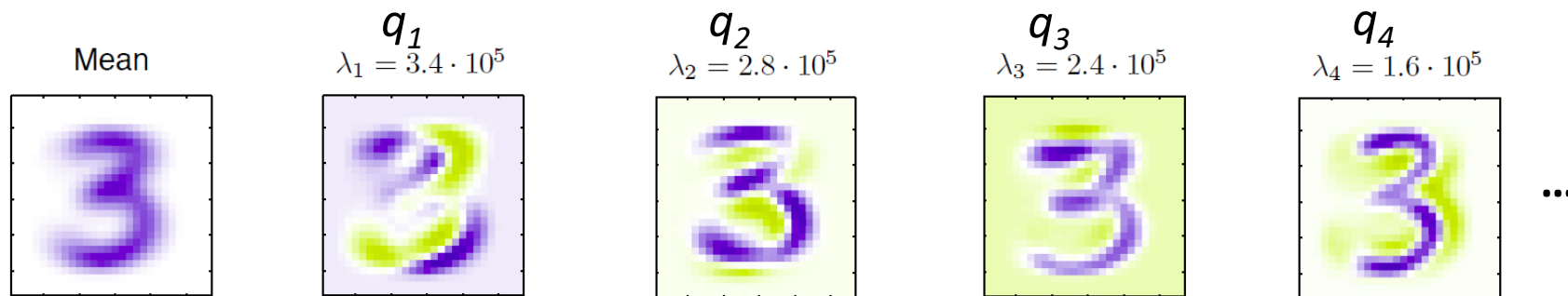
$$\text{tr} \left(Q^T \left(\frac{1}{n} X X^T \right) Q \right) = \sum_{i=1}^k \text{empirical variance of } \vec{q}_i^T x$$

***k -dimensional subspace preserving
maximum amount of variance***



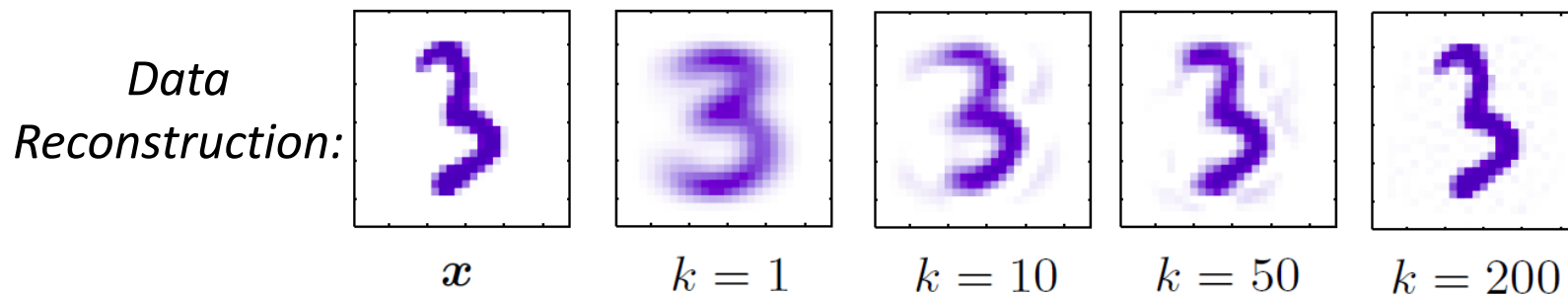
PCA: Example Handwritten Digits

Images of handwritten 3s in \mathbf{R}^{784}



Any example:

$$\text{Any example: } \text{3} = \text{Mean} + w_1 q_1 + w_2 q_2 + \dots$$



We can compress the each datapoint to just k numbers!

Other Popular Dimension Reduction Methods

Multi-dimensional Scaling

Independent Component Analysis (ICA) (for blind source separation)

Non-negative matrix factorization (to create additive models)

Dictionary Learning

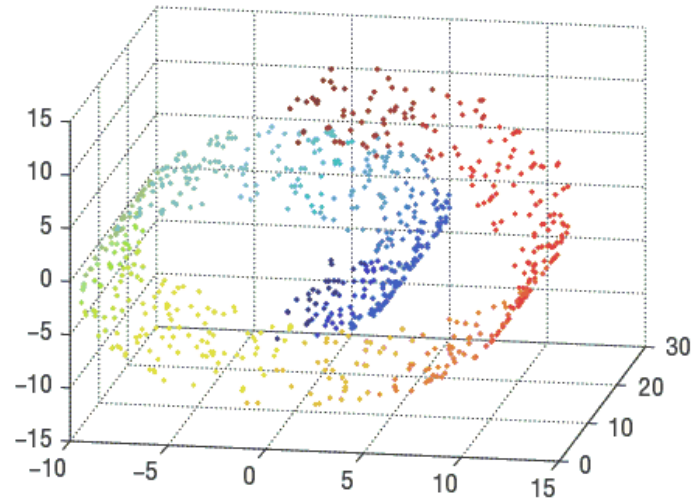
Random Projections

...

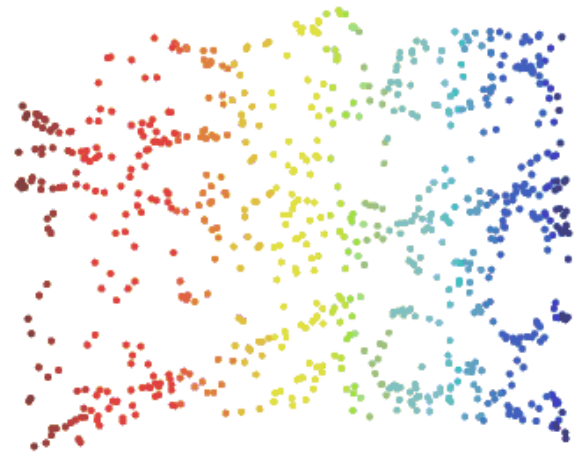
*All of them are **linear** methods*

Non-Linear Dimensionality Reduction

Consider non-linear data



Linear embedding



non-linear embedding

Non-Linear Dimensionality Reduction

Basic optimization criterion:

Find an embedding that:

- Keeps neighboring points close
- Keeps far-off points far

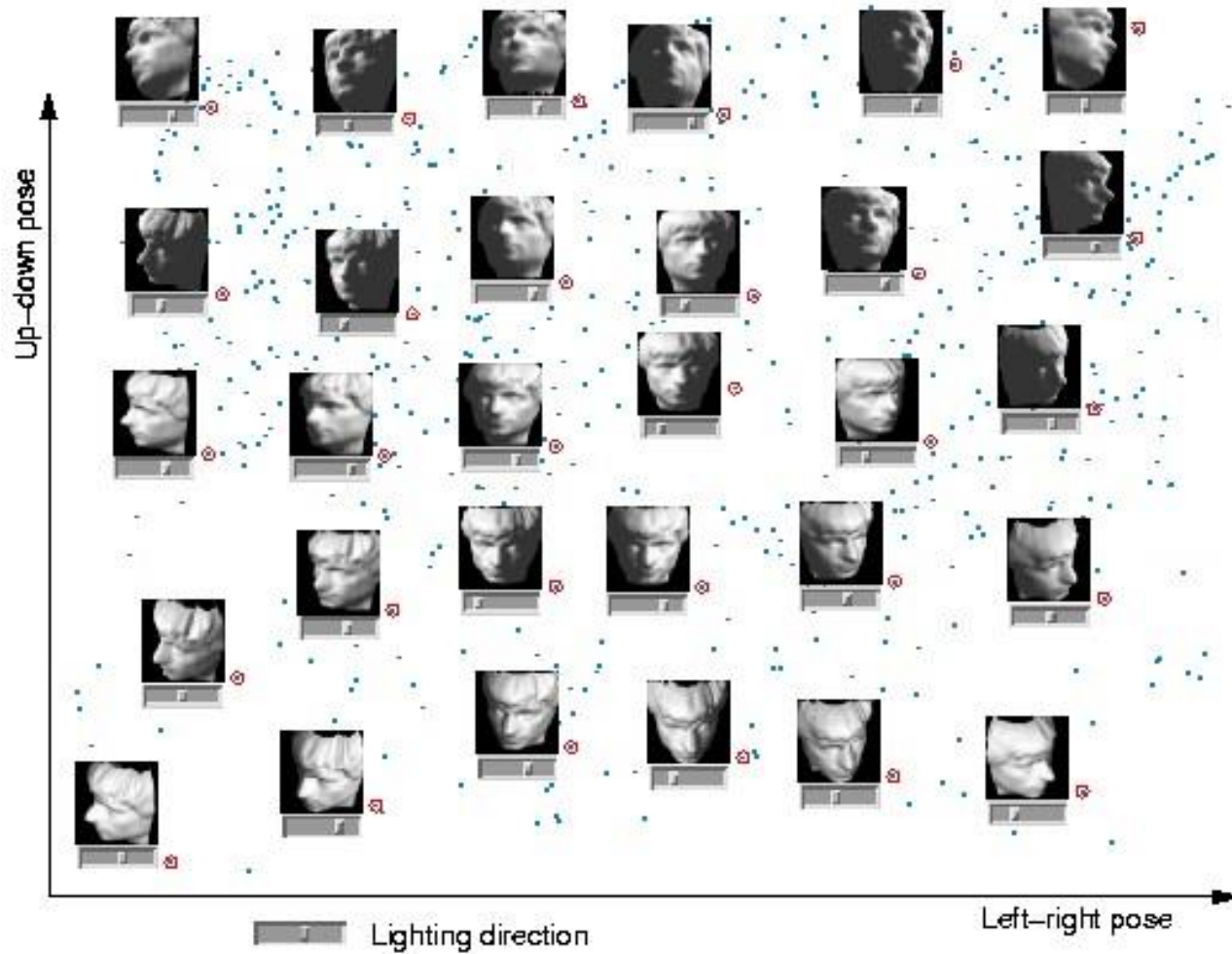
Example variation 1:

Distort neighboring distances by at most $(1 \pm \varepsilon)$ factor, while maximizing non-neighbor distances.

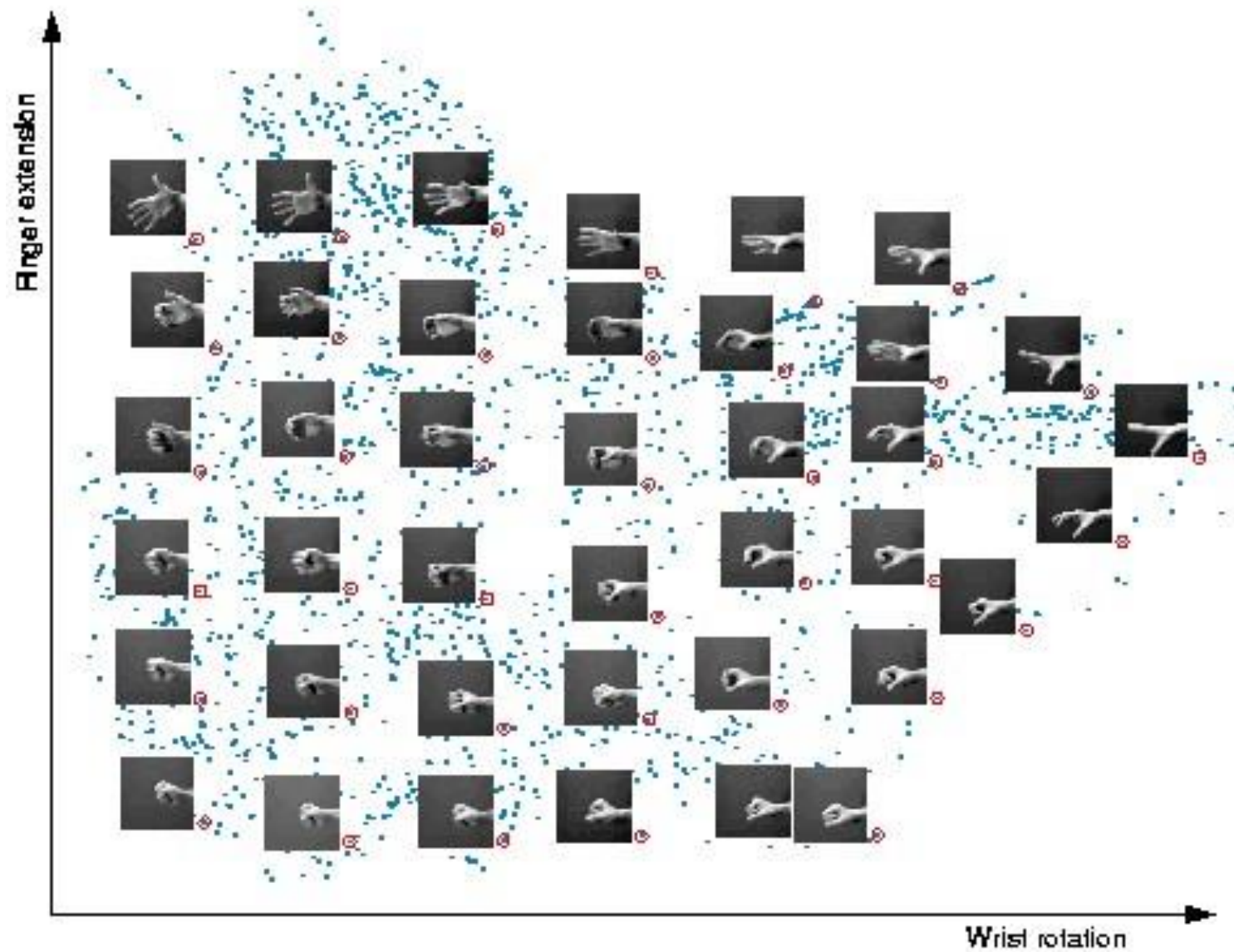
Example variation 2:

*Compute **geodesic** (local hop) distances, and find an embedding that best preserves geodesics.*

Non-linear embedding: Example



Non-linear embedding: Example



Popular Non-Linear Methods

Locally Linear Embedding (LLE)

Isometric Mapping (IsoMap)

Laplacian Eigenmaps (LE)

Local Tangent Space Alignment (LTSA)

Maximum Variance Unfolding (MVU)

...

What We Learned...

- Dimensionality Reduction
 - Linear vs non-linear Dimensionality Reduction
- Principal Component Analysis

Questions?