

# SVD and PCA

---

COS 323

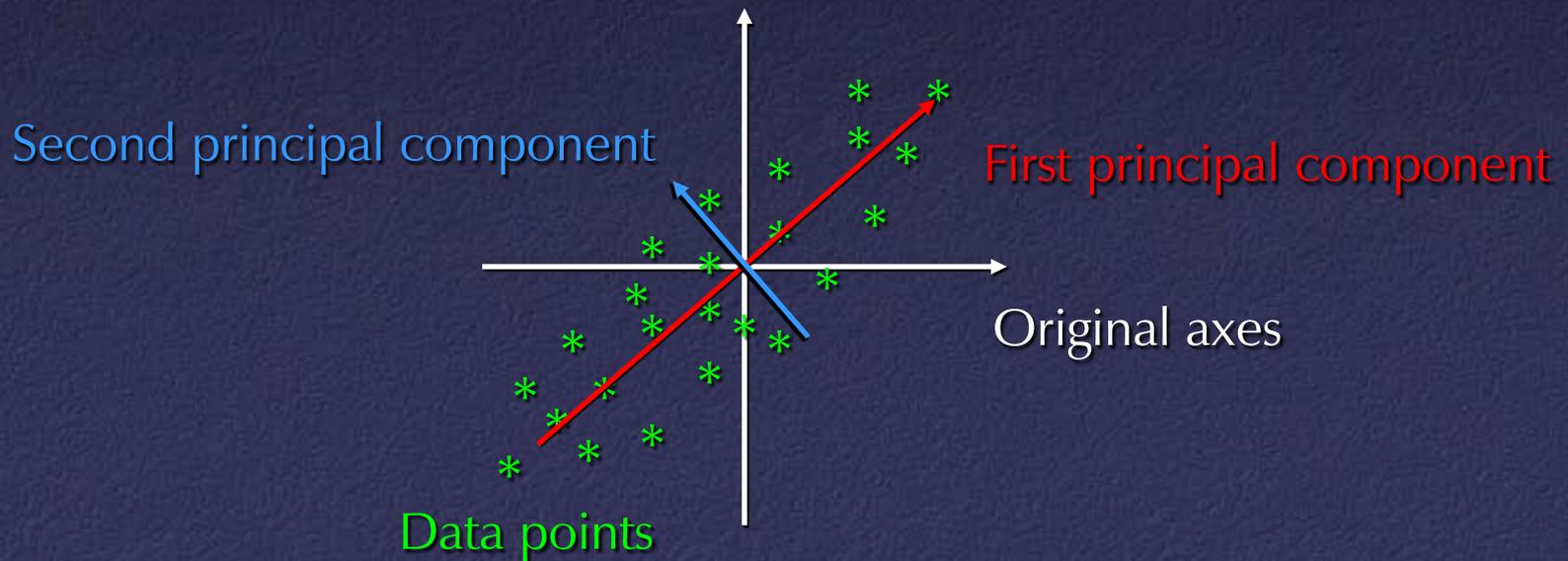
# Dimensionality Reduction

---

- Map points in high-dimensional space to lower number of dimensions
- Preserve structure: pairwise distances, etc.
- Useful for further processing:
  - Less computation, fewer parameters
  - Easier to understand, visualize

# PCA

- Principal Components Analysis (PCA): approximating a high-dimensional data set with a lower-dimensional linear subspace

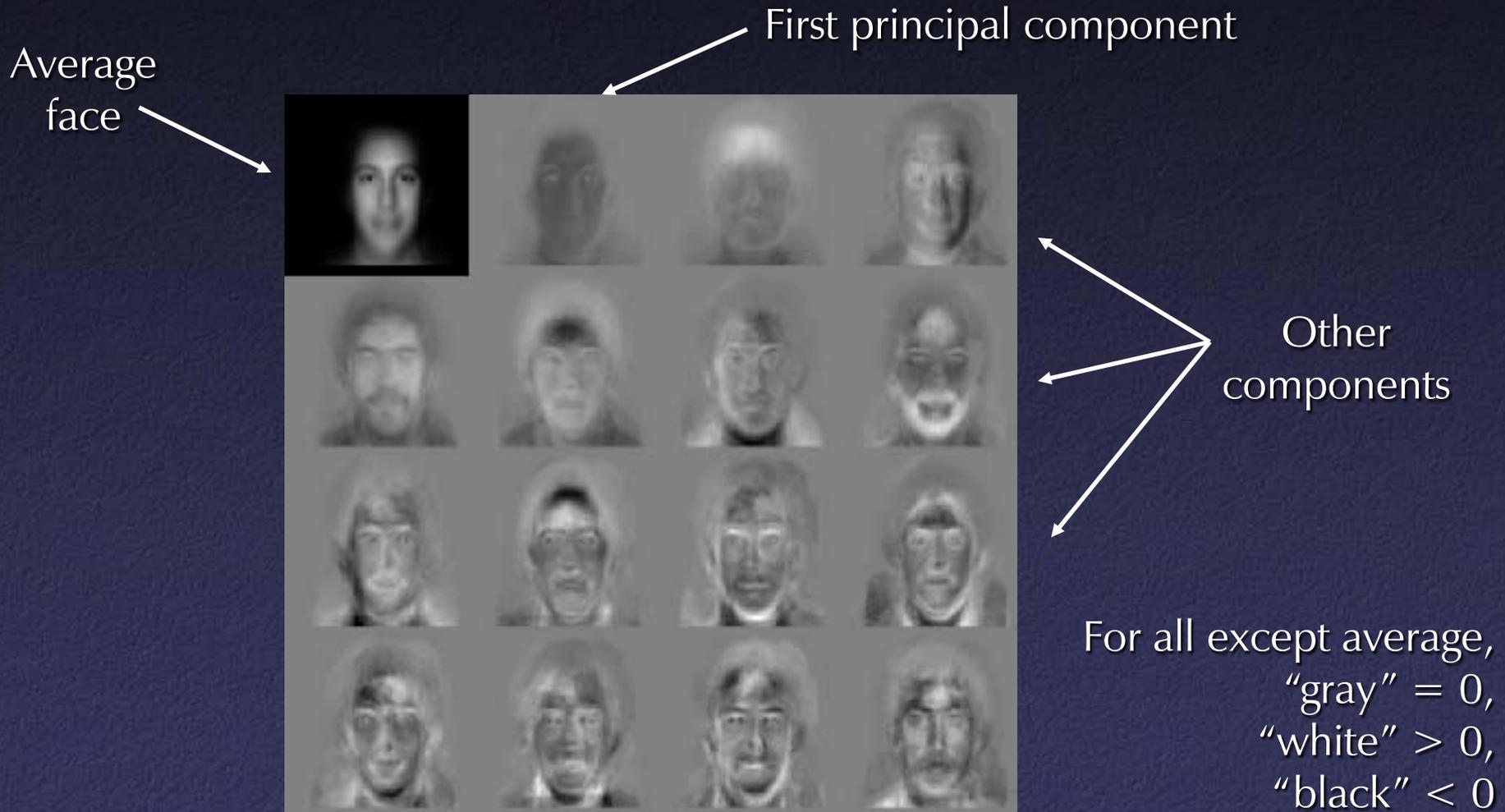


# SVD and PCA

---

- Data matrix with points as rows, take SVD
  - Subtract out mean (“whitening”)
- Columns of  $\mathbf{V}_k$  are principal components
- Value of  $w_i$  gives importance of each component

# PCA on Faces: “Eigenfaces”



# Uses of PCA

---

- Compression: each new image can be approximated by projection onto first few principal components
- Recognition: for a new image, project onto first few principal components, match feature vectors

# PCA for Relighting

- Images under different illumination



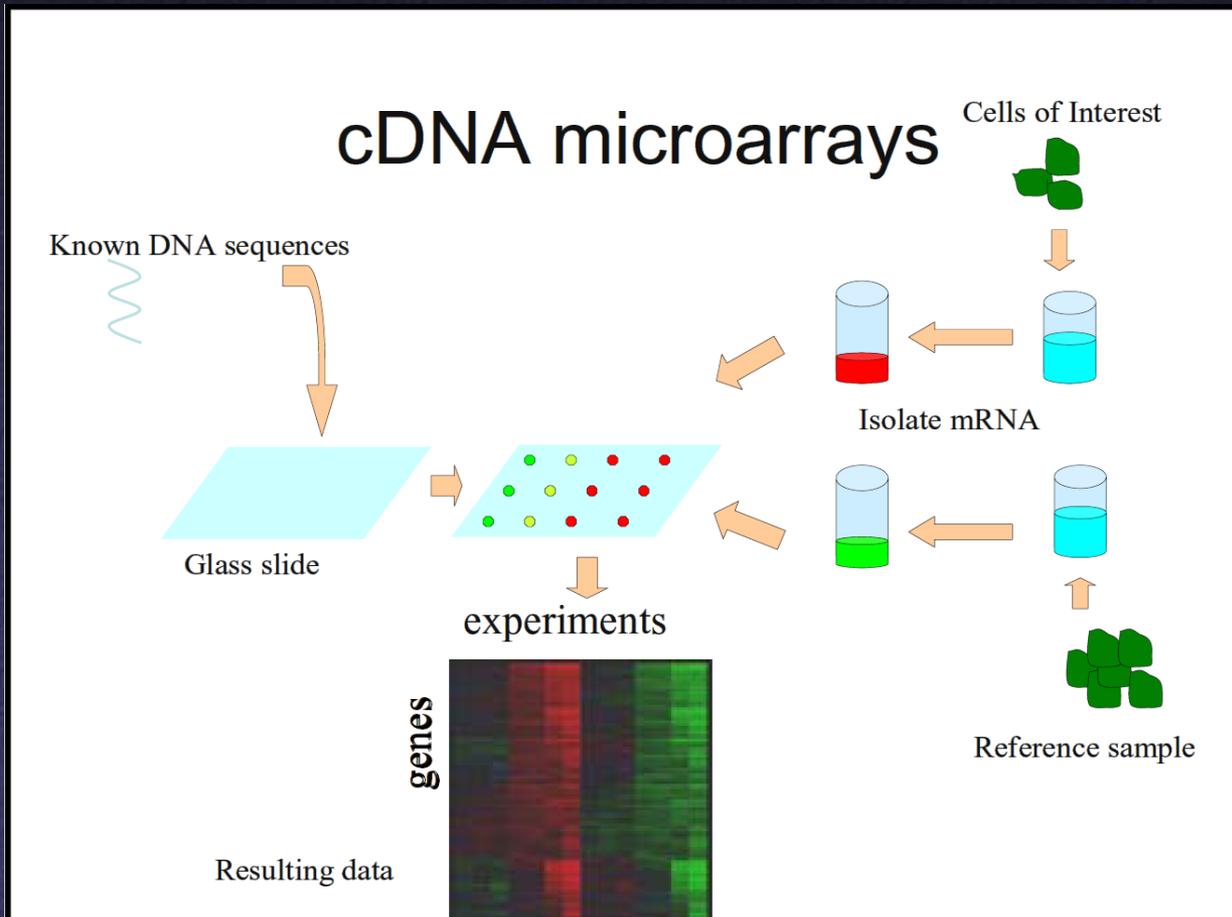
# PCA for Relighting

- Images under different illumination
- Most variation captured by first 5 principal components – can re-illuminate by combining only a few images



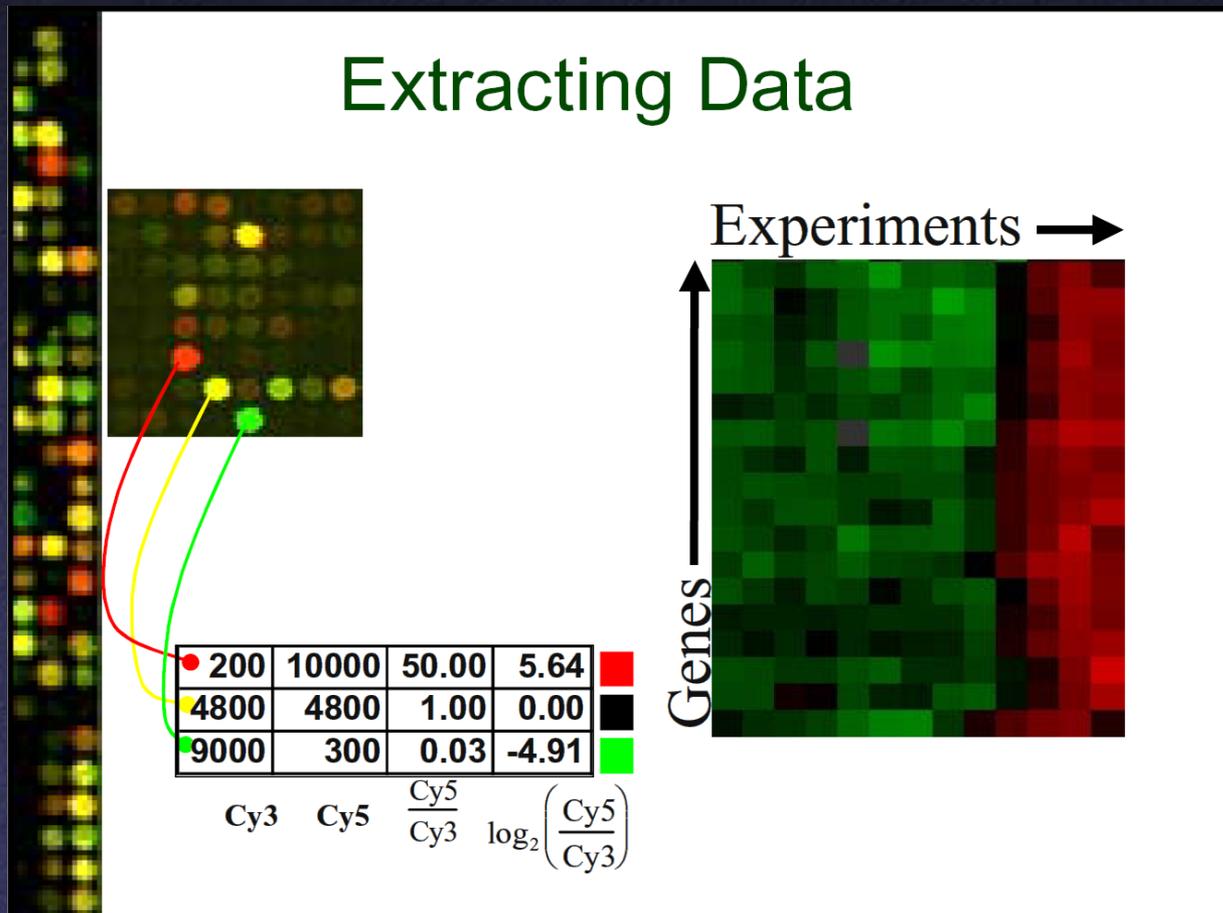
# PCA for DNA Microarrays

- Measure gene activation under different conditions



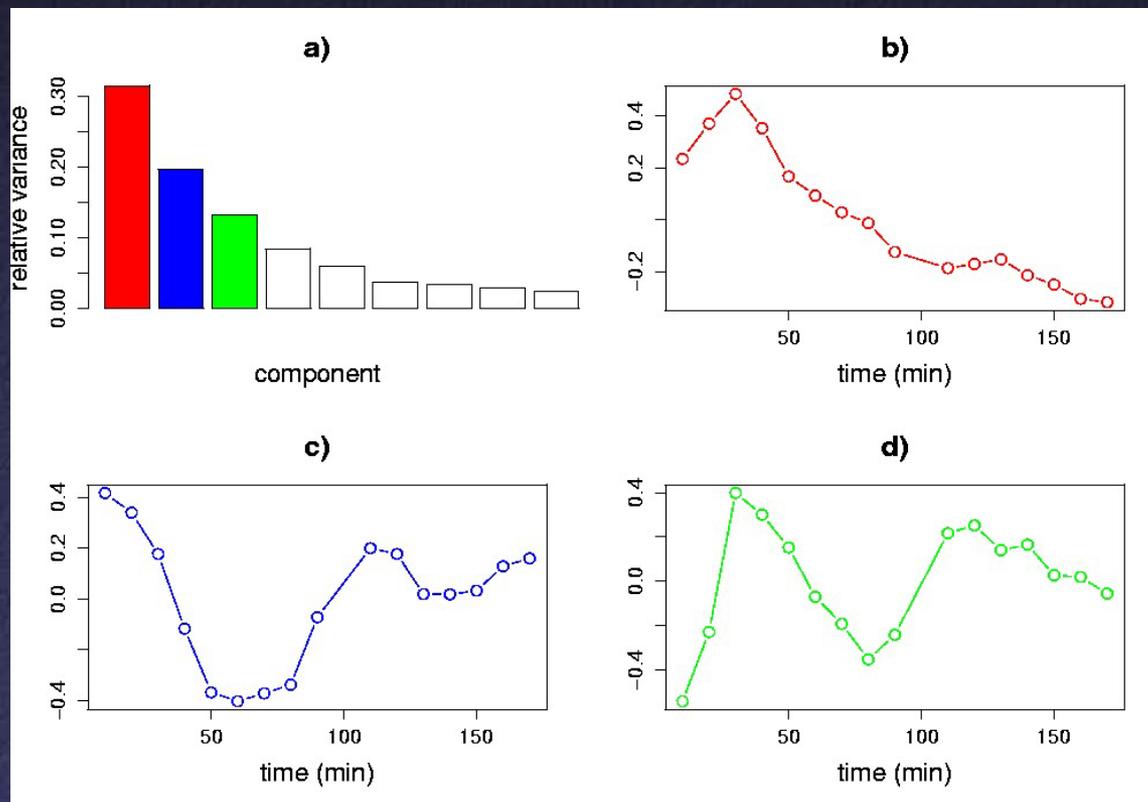
# PCA for DNA Microarrays

- Measure gene activation under different conditions



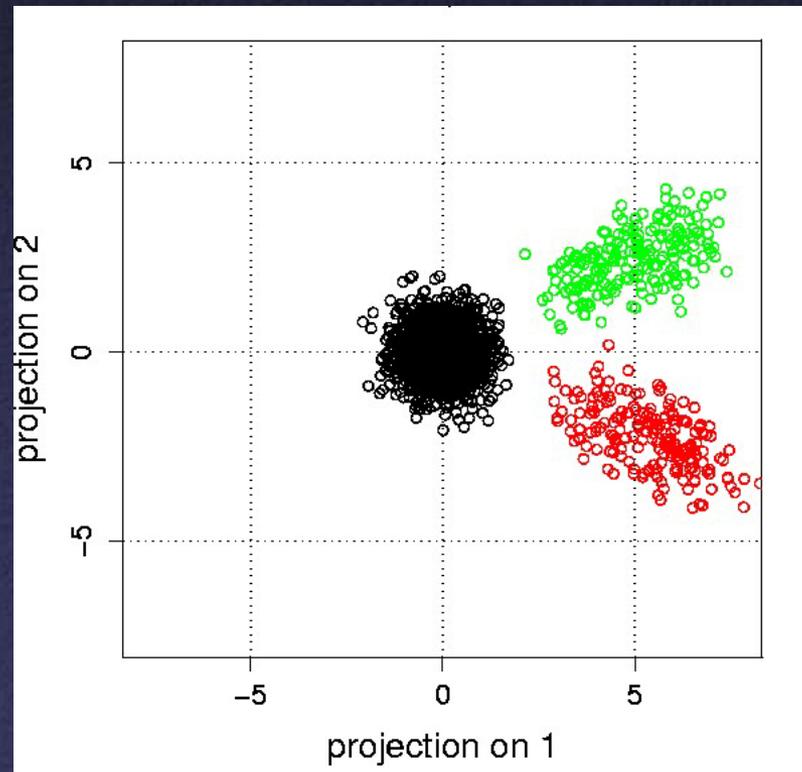
# PCA for DNA Microarrays

- PCA shows patterns of correlated activation
  - Genes with same pattern might have similar function



# PCA for DNA Microarrays

- PCA shows patterns of correlated activation
  - Genes with same pattern might have similar function



# Multidimensional Scaling

---

- In some experiments, can only measure similarity or dissimilarity
  - e.g., is response to stimuli similar or different?
  - Frequent in psychophysical experiments, preference surveys, etc.
- Want to recover absolute positions in  $k$ -dimensional space

# Multidimensional Scaling

- Example: given pairwise distances between cities

	Atl	Chi	Den	Hou	LA	Mia	NYC	SF	Sea	DC
Atlanta	0									
Chicago	587	0								
Denver	1212	920	0							
Houston	701	940	879	0						
LA	1936	1745	831	1374	0					
Miami	604	1188	1726	968	2339	0				
NYC	748	713	1631	1420	2451	1092	0			
SF	2139	1858	949	1645	347	2594	2571	0		
Seattle	2182	1737	1021	1891	959	2734	2406	678	0	
DC	543	597	1494	1220	2300	923	205	2442	2329	0

– Want to recover locations

# Euclidean MDS

---

- Formally, let's say we have  $n \times n$  matrix  $D$  consisting of squared distances  $d_{ij} = (x_i - x_j)^2$
- Want to recover  $n \times d$  matrix  $X$  of positions in  $d$ -dimensional space

$$D = \begin{pmatrix} 0 & (x_1 - x_2)^2 & (x_1 - x_3)^2 & \dots \\ (x_1 - x_2)^2 & 0 & (x_2 - x_3)^2 & \dots \\ (x_1 - x_3)^2 & (x_2 - x_3)^2 & 0 & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

$$X = \begin{pmatrix} (\dots x_1 \dots) \\ (\dots x_2 \dots) \\ \vdots \end{pmatrix}$$

# Euclidean MDS

---

- Observe that

$$d_{ij}^2 = (x_i - x_j)^2 = x_i^2 - 2x_i x_j + x_j^2$$

- Strategy: convert matrix  $D$  of  $d_{ij}^2$  into matrix  $B$  of  $x_i x_j$ 
  - “Centered” distance matrix
  - $B = XX^T$

# Euclidean MDS

---

- Centering:

- Sum of row  $i$  of  $D =$  sum of column  $i$  of  $D =$

$$\begin{aligned} s_i &= \sum_j d_{ij}^2 = \sum_j x_i^2 - 2x_i x_j + x_j^2 \\ &= nx_i^2 - 2x_i \sum_j x_j + \sum_j x_j^2 \end{aligned}$$

- Sum of all entries in  $D =$

$$s = \sum_i s_i = 2n \sum_i x_i^2 - 2 \left( \sum_i x_i \right)^2$$

# Euclidean MDS

---

- Choose  $\sum x_i = 0$ 
  - Solution will have average position at origin

$$s_i = nx_i^2 + \sum_j x_j^2, \quad s = 2n \sum_j x_j^2$$

- Then,

$$d_{ij}^2 - \frac{1}{n}s_i - \frac{1}{n}s_j + \frac{1}{n^2}s = -2x_i x_j$$

- So, to get  $B$ :
  - compute row (or column) sums
  - compute sum of sums
  - apply above formula to each entry of  $D$
  - Divide by  $-2$

# Euclidean MDS

---

- Now have  $B$ , want to factor into  $XX^T$
- If  $X$  is  $n \times d$ ,  $B$  must have rank  $d$
- Take SVD, set all but top  $d$  singular values to 0
  - Eliminate corresponding columns of  $U$  and  $V$
  - Have  $B_3 = U_3 W_3 V_3^T$
  - $B$  is square and symmetric, so  $U = V$
  - Take  $X = U_3$  times square root of  $W_3$

# Multidimensional Scaling

- Result ( $d = 2$ ):



# Multidimensional Scaling

---

- Caveat: actual axes, center not necessarily what you want (can't recover them!)
- This is “classical” or “Euclidean” MDS [Torgerson 52]
  - Distance matrix assumed to be actual Euclidean distance
- More sophisticated versions available
  - “Non-metric MDS”: not Euclidean distance, sometimes just *inequalities*
  - “Weighted MDS”: account for observer bias