

# PCA / SVD

COMS 4721 Spring 2022  
Daniel Hsu

**PCA**

# Empirical (sample) covariance

Recall: If  $\vec{X}$  is a random vector in  $\mathbb{R}^d$ , its covariance matrix is

$$\text{cov}(\vec{X}) = \mathbb{E}[(\vec{X} - \mathbb{E}(\vec{X}))(\vec{X} - \mathbb{E}(\vec{X}))^\top].$$

# Empirical (sample) covariance

Recall: If  $\vec{X}$  is a random vector in  $\mathbb{R}^d$ , its covariance matrix is

$$\text{cov}(\vec{X}) = \mathbb{E}[(\vec{X} - \mathbb{E}(\vec{X}))(\vec{X} - \mathbb{E}(\vec{X}))^\top].$$

If  $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$  are  $n$  data points from  $\mathbb{R}^d$ , and  $P_n$  is the empirical distribution on these points, then their **empirical covariance matrix** (a.k.a. **sample covariance matrix**) is  $\text{cov}(\vec{X})$  for  $\vec{X} \sim P_n$

# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

- ▶ Symmetric

# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

- ▶ Symmetric
- ▶  $d$  real eigenvalues  $\lambda_1, \dots, \lambda_d$  (possibly repeated)

**Convention:** Order from largest to smallest

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

- ▶ Symmetric
- ▶  $d$  real eigenvalues  $\lambda_1, \dots, \lambda_d$  (possibly repeated)

**Convention:** Order from largest to smallest

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

- ▶ Can choose corresponding eigenvectors  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  s.t. they are orthonormal



# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

- ▶ Symmetric
- ▶  $d$  real eigenvalues  $\lambda_1, \dots, \lambda_d$  (possibly repeated)

**Convention:** Order from largest to smallest

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

- ▶ Can choose corresponding eigenvectors  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  s.t. they are orthonormal
- ▶ In fact, all eigenvalues are non-negative:

$$\vec{v}_i^\top \text{cov}(\vec{X}) \vec{v}_i = \vec{v}_i^\top (\lambda_i \vec{v}_i) = \lambda_i$$

and 
$$\vec{v}_i^\top \text{cov}(\vec{X}) \vec{v}_i = \text{var}(\vec{v}_i \cdot \vec{X}) \geq 0$$

# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

- ▶ Symmetric
- ▶  $d$  real eigenvalues  $\lambda_1, \dots, \lambda_d$  (possibly repeated)

**Convention:** Order from largest to smallest

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

- ▶ Can choose corresponding eigenvectors  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  s.t. they are orthonormal
- ▶ In fact, all eigenvalues are non-negative:

$$\vec{v}_i^\top \text{cov}(\vec{X}) \vec{v}_i = \vec{v}_i^\top (\lambda_i \vec{v}_i) = \lambda_i$$

and 
$$\vec{v}_i^\top \text{cov}(\vec{X}) \vec{v}_i = \text{var}(\vec{v}_i \cdot \vec{X}) \geq 0$$

Symmetric matrices with non-negative eigenvalues are said to be **positive semidefinite (PSD)**

# More properties of the covariance matrix

Properties of  $\text{cov}(\vec{X})$ :

- ▶ Symmetric
- ▶  $d$  real eigenvalues  $\lambda_1, \dots, \lambda_d$  (possibly repeated)

**Convention:** Order from largest to smallest

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

- ▶ Can choose corresponding eigenvectors  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  s.t. they are orthonormal
- ▶ In fact, all eigenvalues are non-negative:

$$\vec{v}_i^\top \text{cov}(\vec{X}) \vec{v}_i = \vec{v}_i^\top (\lambda_i \vec{v}_i) = \lambda_i$$

and 
$$\vec{v}_i^\top \text{cov}(\vec{X}) \vec{v}_i = \text{var}(\vec{v}_i \cdot \vec{X}) \geq 0$$

Symmetric matrices with non-negative eigenvalues are said to be **positive semidefinite (PSD)**

- ▶ **Eigenvalue decomposition** (a.k.a. **spectral decomposition**) of  $\text{cov}(\vec{X})$ :

$$\text{cov}(\vec{X}) = \sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^\top$$

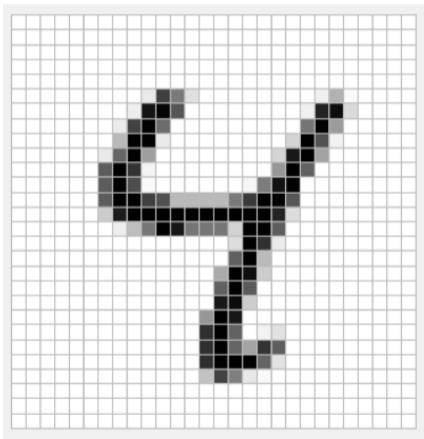
## Example: MNIST data set of handwritten digits

60000 images of handwritten digits, each a  $28 \times 28$  gray-scale pixel image



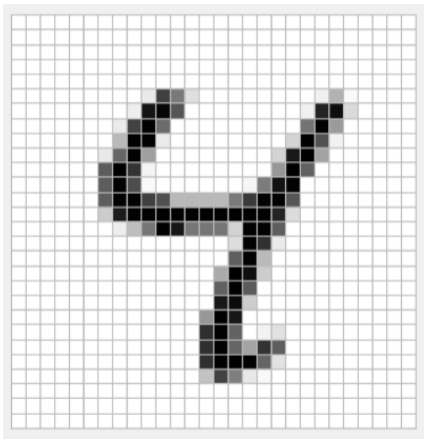
## Example: MNIST data set of handwritten digits

Each image is represented as vector in  $[0, 255]^d$  for  $d = 28 \times 28 = 784$ , one feature per pixel



## Example: MNIST data set of handwritten digits

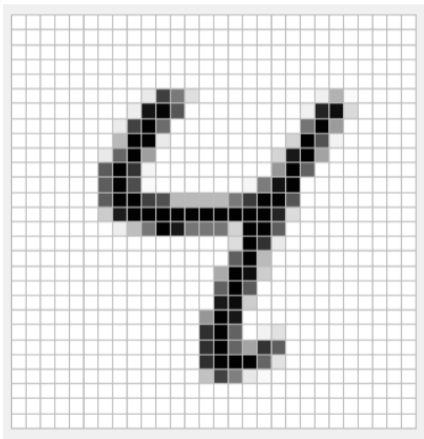
Each image is represented as vector in  $[0, 255]^d$  for  $d = 28 \times 28 = 784$ , one feature per pixel



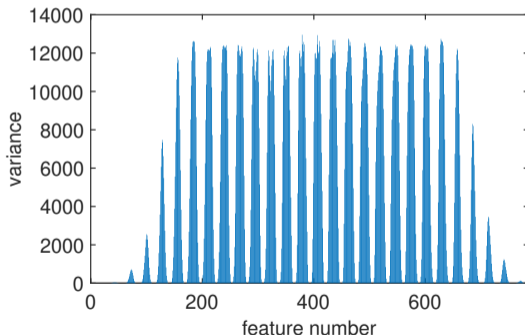
Some “features” (pixels) are more variable than others (on average over the data set)

## Example: MNIST data set of handwritten digits

Each image is represented as vector in  $[0, 255]^d$  for  $d = 28 \times 28 = 784$ , one feature per pixel

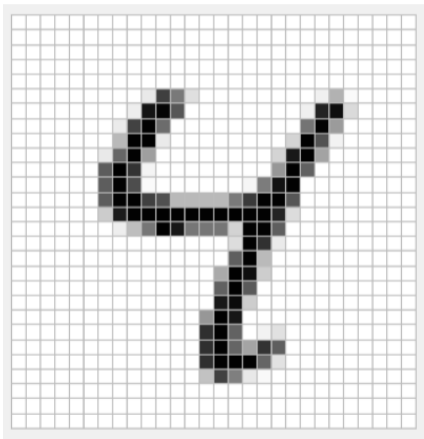


Some “features” (pixels) are more variable than others (on average over the data set)

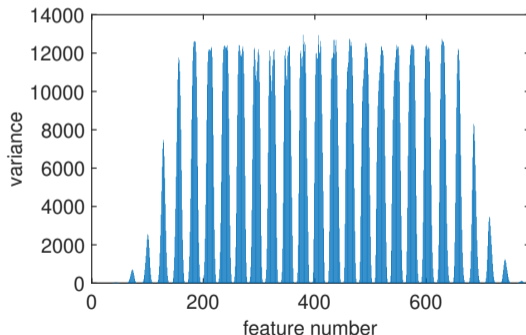


## Example: MNIST data set of handwritten digits

Each image is represented as vector in  $[0, 255]^d$  for  $d = 28 \times 28 = 784$ , one feature per pixel



Some “features” (pixels) are more variable than others (on average over the data set)



Perhaps not all features are necessary!



# Data approximation

**Coarse approximation:** (our baseline)

Represent each data point by the overall mean  $\vec{\mu}$   
(where  $\vec{\mu} := \mathbb{E}(\vec{X})$ )

$$\vec{X} \mapsto \vec{\mu}$$

Zero features per data point!

# Data approximation

**Coarse approximation:** (our baseline)

Represent each data point by the overall mean  $\vec{\mu}$   
(where  $\vec{\mu} := \mathbb{E}(\vec{X})$ )

$$\vec{X} \mapsto \vec{\mu}$$

Zero features per data point!

Expected reconstruction error:

$$\mathbb{E} \left[ \|\vec{X} - \vec{\mu}\|_2^2 \right]$$

# Data approximation

## Coarse approximation: (our baseline)

Represent each data point by the overall mean  $\vec{\mu}$   
(where  $\vec{\mu} := \mathbb{E}(\vec{X})$ )

$$\vec{X} \mapsto \vec{\mu}$$

Zero features per data point!

Expected reconstruction error:

$$\mathbb{E} \left[ \|\vec{X} - \vec{\mu}\|_2^2 \right]$$

## Less coarse approximation:

Represent each data point by overall mean  $\vec{\mu}$  plus  
projection of  $\vec{X} - \vec{\mu}$  to subspace  $W \subseteq \mathbb{R}^d$

$$\vec{X} \mapsto \vec{\mu} + \Pi_W(\vec{X} - \vec{\mu})$$

$k := \dim(W)$  features per data point

# Data approximation

## Coarse approximation: (our baseline)

Represent each data point by the overall mean  $\vec{\mu}$   
(where  $\vec{\mu} := \mathbb{E}(\vec{X})$ )

$$\vec{X} \mapsto \vec{\mu}$$

Zero features per data point!

Expected reconstruction error:

$$\mathbb{E} \left[ \|\vec{X} - \vec{\mu}\|_2^2 \right]$$

## Less coarse approximation:

Represent each data point by overall mean  $\vec{\mu}$  plus  
projection of  $\vec{X} - \vec{\mu}$  to subspace  $W \subseteq \mathbb{R}^d$

$$\vec{X} \mapsto \vec{\mu} + \Pi_W(\vec{X} - \vec{\mu})$$

$k := \dim(W)$  features per data point

Expected reconstruction error:

$$\mathbb{E} \left[ \|\vec{X} - (\vec{\mu} + \Pi_W(\vec{X} - \vec{\mu}))\|_2^2 \right]$$

# Data approximation

**Coarse approximation:** (our baseline)

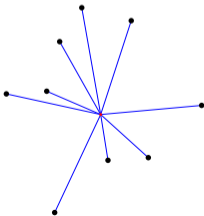
Represent each data point by the overall mean  $\vec{\mu}$   
(where  $\vec{\mu} := \mathbb{E}(\vec{X})$ )

$$\vec{X} \mapsto \vec{\mu}$$

Zero features per data point!

Expected reconstruction error:

$$\mathbb{E} \left[ \|\vec{X} - \vec{\mu}\|_2^2 \right]$$



**Less coarse approximation:**

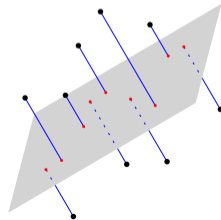
Represent each data point by overall mean  $\vec{\mu}$  plus  
projection of  $\vec{X} - \vec{\mu}$  to subspace  $W \subseteq \mathbb{R}^d$

$$\vec{X} \mapsto \vec{\mu} + \Pi_W(\vec{X} - \vec{\mu})$$

$k := \dim(W)$  features per data point

Expected reconstruction error:

$$\mathbb{E} \left[ \|\vec{X} - (\vec{\mu} + \Pi_W(\vec{X} - \vec{\mu}))\|_2^2 \right]$$



## Special case: Coordinate subspaces

Suppose we only allow  $W$  to be a **coordinate subspace**, i.e., span of some coordinate vectors

## Special case: Coordinate subspaces

Suppose we only allow  $W$  to be a **coordinate subspace**, i.e., span of some coordinate vectors

- ▶ Effect of  $\vec{X} \mapsto \vec{\mu} + \Pi_W(\vec{X} - \vec{\mu})$ :

Replace all but  $k = \dim(W)$  of the features in  $\vec{X}$  with corresponding mean value in  $\vec{\mu}$

## Special case: Coordinate subspaces

Suppose we only allow  $W$  to be a **coordinate subspace**, i.e., span of some coordinate vectors

- ▶ Effect of  $\vec{X} \mapsto \vec{\mu} + \Pi_W(\vec{X} - \vec{\mu})$ :

Replace all but  $k = \dim(W)$  of the features in  $\vec{X}$  with corresponding mean value in  $\vec{\mu}$

- ▶ MNIST data: each “coordinate vector” corresponds to a pixel location  $(i, j) \in [28] \times [28]$ 
  - ▶ For expected reconstruction error to be

$$10\% \times \underbrace{\mathbb{E} \left[ \|\vec{X} - \vec{\mu}\|_2^2 \right]}_{\text{our baseline}},$$

need  $k = 299$  features (pixels)



## Special case: Coordinate subspaces

Suppose we only allow  $W$  to be a **coordinate subspace**, i.e., span of some coordinate vectors

- ▶ Effect of  $\vec{X} \mapsto \vec{\mu} + \Pi_W(\vec{X} - \vec{\mu})$ :

Replace all but  $k = \dim(W)$  of the features in  $\vec{X}$  with corresponding mean value in  $\vec{\mu}$

- ▶ MNIST data: each “coordinate vector” corresponds to a pixel location  $(i, j) \in [28] \times [28]$ 
  - ▶ For expected reconstruction error to be

$$10\% \times \underbrace{\mathbb{E} \left[ \|\vec{X} - \vec{\mu}\|_2^2 \right]}_{\text{our baseline}},$$

need  $k = 299$  features (pixels)

Question: Can we do better by not restricting to coordinate subspaces?

## Special case: $k = 1$

Orthogonal projection operator for 1-dimensional subspace  $W$  (i.e., a line):

$$\Pi_W = \vec{\alpha}\vec{\alpha}^\top$$

where  $\vec{\alpha} \in \mathbb{R}^d$  is (any) unit vector in  $W$

## Special case: $k = 1$

Orthogonal projection operator for 1-dimensional subspace  $W$  (i.e., a line):

$$\Pi_W = \vec{\alpha}\vec{\alpha}^\top$$

where  $\vec{\alpha} \in \mathbb{R}^d$  is (any) unit vector in  $W$

Recall Pythagorean theorem:

$$\|\vec{X} - \vec{\mu}\|_2^2 = \|(\vec{X} - \vec{\mu}) - \Pi_W(\vec{X} - \vec{\mu})\|_2^2 + \|\Pi_W(\vec{X} - \vec{\mu})\|_2^2$$

So  $W$  that minimizes  $\mathbb{E} \left[ \|\vec{X} - \vec{\mu} - \Pi_W(\vec{X} - \vec{\mu})\|_2^2 \right]$  is same as  $W$  that maximizes  $\mathbb{E} \left[ \|\Pi_W(\vec{X} - \vec{\mu})\|_2^2 \right]$

For  $k = 1$ , this is same as finding unit vector  $\vec{\alpha}$  that maximizes

$$\mathbb{E} \left[ \|\vec{\alpha}\vec{\alpha}^\top(\vec{X} - \vec{\mu})\|_2^2 \right] = \mathbb{E} \left[ (\vec{\alpha}^\top(\vec{X} - \vec{\mu}))^2 \right] = \text{var}(\vec{\alpha} \cdot \vec{X})$$

# Maximum variance linear function of a random vector

Question: What is the maximum variance achievable by a linear function of  $\vec{X}$  given as a *unit vector*?

$$\begin{aligned} & \max_{\vec{\alpha} \in \mathbb{R}^d} \text{var}(\vec{\alpha} \cdot \vec{X}) \\ & \text{s.t. } \|\vec{\alpha}\|_2 = 1 \end{aligned}$$

# Maximum variance linear function of a random vector

Question: What is the maximum variance achievable by a linear function of  $\vec{X}$  given as a *unit vector*?

$$\begin{aligned} \max_{\vec{\alpha} \in \mathbb{R}^d} \quad & \text{var}(\vec{\alpha} \cdot \vec{X}) \\ \text{s.t.} \quad & \|\vec{\alpha}\|_2 = 1 \end{aligned}$$

- ▶ Answer: The largest eigenvalue of  $\text{cov}(\vec{X})$

# Maximum variance linear function of a random vector

Question: What is the maximum variance achievable by a linear function of  $\vec{X}$  given as a *unit vector*?

$$\begin{aligned} \max_{\vec{\alpha} \in \mathbb{R}^d} \quad & \text{var}(\vec{\alpha} \cdot \vec{X}) \\ \text{s.t.} \quad & \|\vec{\alpha}\|_2 = 1 \end{aligned}$$

► Answer: The largest eigenvalue of  $\text{cov}(\vec{X})$

Also called the **top eigenvalue** of  $\text{cov}(\vec{X})$

# Maximum variance linear function of a random vector

Question: What is the maximum variance achievable by a linear function of  $\vec{X}$  given as a *unit vector*?

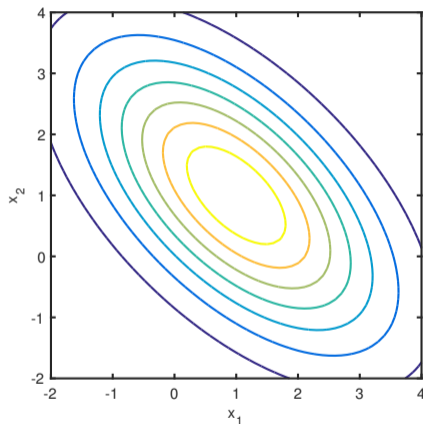
$$\begin{aligned} \max_{\vec{\alpha} \in \mathbb{R}^d} \quad & \text{var}(\vec{\alpha} \cdot \vec{X}) \\ \text{s.t.} \quad & \|\vec{\alpha}\|_2 = 1 \end{aligned}$$

► Answer: The largest eigenvalue of  $\text{cov}(\vec{X})$

Also called the **top eigenvalue** of  $\text{cov}(\vec{X})$

Achieved by (unit length) eigenvector corresponding to that eigenvalue  
(a.k.a. **top eigenvector**)

# Top eigenvector for covariance of a multivariate Gaussian



Top eigenvector direction defines the major axis of each ellipse



## Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

Pythagorean theorem says  $\sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i)^2 = \|\vec{\alpha}\|_2^2$

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

Pythagorean theorem says  $\sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i)^2 = \|\vec{\alpha}\|_2^2$

- ▶ Variance of  $\vec{\alpha} \cdot \vec{X}$  for unit vector  $\vec{\alpha}$ :

$$\text{var}(\vec{\alpha} \cdot \vec{X}) = \vec{\alpha}^T \text{cov}(\vec{X}) \vec{\alpha}$$

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

Pythagorean theorem says  $\sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i)^2 = \|\vec{\alpha}\|_2^2$

- ▶ Variance of  $\vec{\alpha} \cdot \vec{X}$  for unit vector  $\vec{\alpha}$ :

$$\begin{aligned} \text{var}(\vec{\alpha} \cdot \vec{X}) &= \vec{\alpha}^\top \text{cov}(\vec{X}) \vec{\alpha} \\ &= \vec{\alpha}^\top \left( \sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^\top \right) \vec{\alpha} \end{aligned}$$

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

Pythagorean theorem says  $\sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i)^2 = \|\vec{\alpha}\|_2^2$

- ▶ Variance of  $\vec{\alpha} \cdot \vec{X}$  for unit vector  $\vec{\alpha}$ :

$$\begin{aligned} \text{var}(\vec{\alpha} \cdot \vec{X}) &= \vec{\alpha}^T \text{cov}(\vec{X}) \vec{\alpha} \\ &= \vec{\alpha}^T \left( \sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T \right) \vec{\alpha} \\ &= \sum_{i=1}^d \lambda_i (\vec{\alpha} \cdot \vec{v}_i)^2 \end{aligned}$$

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

Pythagorean theorem says  $\sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i)^2 = \|\vec{\alpha}\|_2^2$

- ▶ Variance of  $\vec{\alpha} \cdot \vec{X}$  for unit vector  $\vec{\alpha}$ :

$$\begin{aligned} \text{var}(\vec{\alpha} \cdot \vec{X}) &= \vec{\alpha}^T \text{cov}(\vec{X}) \vec{\alpha} \\ &= \vec{\alpha}^T \left( \sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T \right) \vec{\alpha} \\ &= \sum_{i=1}^d \lambda_i (\vec{\alpha} \cdot \vec{v}_i)^2 \end{aligned}$$

- ▶ This is a “weighted” average of the  $\lambda_i$ 's; maximized when all “weight” is put on  $\lambda_1$

# Optimality of top eigenvector

Let eigenvalues of  $\text{cov}(\vec{X})$  be  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , and let  $\vec{v}_1, \vec{v}_2, \dots, \vec{v}_d$  be corresponding orthonormal eigenvectors

- ▶ Any vector  $\vec{\alpha}$  can be written as

$$\vec{\alpha} = \sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i) \vec{v}_i$$

Pythagorean theorem says  $\sum_{i=1}^d (\vec{\alpha} \cdot \vec{v}_i)^2 = \|\vec{\alpha}\|_2^2$

- ▶ Variance of  $\vec{\alpha} \cdot \vec{X}$  for unit vector  $\vec{\alpha}$ :

$$\begin{aligned} \text{var}(\vec{\alpha} \cdot \vec{X}) &= \vec{\alpha}^T \text{cov}(\vec{X}) \vec{\alpha} \\ &= \vec{\alpha}^T \left( \sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T \right) \vec{\alpha} \\ &= \sum_{i=1}^d \lambda_i (\vec{\alpha} \cdot \vec{v}_i)^2 \end{aligned}$$

- ▶ This is a “weighted” average of the  $\lambda_i$ ’s; maximized when all “weight” is put on  $\lambda_1$
- ▶ To put all “weight” on  $\lambda_1$ , need  $(\vec{\alpha} \cdot \vec{v}_1)^2 = 1$ ; achieved by  $\vec{\alpha} = \pm \vec{v}_1$



# Top eigen{value,vector} on MNIST data

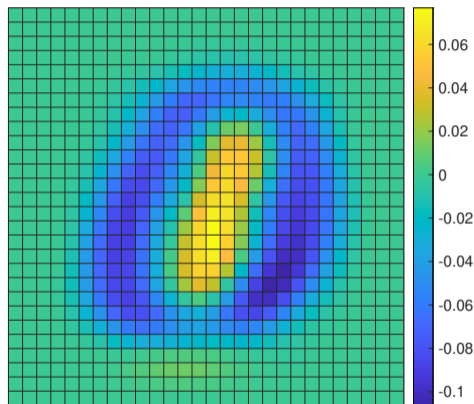
Top eigenvalue:

$$\lambda_1 \approx 0.1 \times \mathbb{E} \|\vec{X} - \vec{\mu}\|_2^2$$

Expected reconstruction error:

$$\approx 0.9 \times \mathbb{E} \|\vec{X} - \vec{\mu}\|_2^2$$

Top eigenvector:



Some MNIST digits sorted by projection to span of top eigenvector:

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

## General case: $k \geq 1$

Orthogonal projection operator for  $k$ -dimensional subspace  $W$ :

$$\Pi_W = \sum_{i=1}^k \vec{\alpha}_i \vec{\alpha}_i^T$$

where  $\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d$  is (any) orthonormal basis for  $W$

## General case: $k \geq 1$

Orthogonal projection operator for  $k$ -dimensional subspace  $W$ :

$$\Pi_W = \sum_{i=1}^k \vec{\alpha}_i \vec{\alpha}_i^\top$$

where  $\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d$  is (any) orthonormal basis for  $W$

Recall Pythagorean theorem:

$$\|\vec{X} - \vec{\mu}\|_2^2 = \|(\vec{X} - \vec{\mu}) - \Pi_W(\vec{X} - \vec{\mu})\|_2^2 + \|\Pi_W(\vec{X} - \vec{\mu})\|_2^2$$

So  $W$  that minimizes  $\mathbb{E} \left[ \|\vec{X} - \vec{\mu} - \Pi_W(\vec{X} - \vec{\mu})\|_2^2 \right]$  is same as  $W$  that maximizes  $\mathbb{E} \left[ \|\Pi_W(\vec{X} - \vec{\mu})\|_2^2 \right]$

This is same as finding orthonormal  $\vec{\alpha}_1, \dots, \vec{\alpha}_k$  that maximizes

$$\mathbb{E} \left[ \left\| \sum_{i=1}^k \vec{\alpha}_i \vec{\alpha}_i^\top (\vec{X} - \vec{\mu}) \right\|_2^2 \right] = \mathbb{E} \left[ \sum_{i=1}^k (\vec{\alpha}_i^\top (\vec{X} - \vec{\mu}))^2 \right] = \sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$$

# Optimality of the top $k$ eigenvectors

Question: What's the maximum variance "accounted for" by  $k$  dimensional subspace?

$$\max_{\substack{\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d \\ \text{orthonormal}}} \sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$$

# Optimality of the top $k$ eigenvectors

Question: What's the maximum variance "accounted for" by  $k$  dimensional subspace?

$$\max_{\substack{\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d \\ \text{orthonormal}}} \sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$$

- ▶ Answer: Sum of  $k$  largest eigenvalues of  $\text{cov}(\vec{X})$   
(sum of the **top  $k$  eigenvalues**)

# Optimality of the top $k$ eigenvectors

Question: What's the maximum variance "accounted for" by  $k$  dimensional subspace?

$$\max_{\substack{\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d \\ \text{orthonormal}}} \sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$$

- ▶ Answer: Sum of  $k$  largest eigenvalues of  $\text{cov}(\vec{X})$   
(sum of the **top  $k$  eigenvalues**)

Achieved by any orthonormal basis for the span of  $k$  corresponding eigenvectors  
(span of **top  $k$  eigenvectors**, a.k.a. **top  $k$  eigenspace**)

# Optimality of the top $k$ eigenvectors

Question: What's the maximum variance "accounted for" by  $k$  dimensional subspace?

$$\max_{\substack{\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d \\ \text{orthonormal}}} \sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$$

- ▶ Answer: Sum of  $k$  largest eigenvalues of  $\text{cov}(\vec{X})$   
(sum of the **top  $k$  eigenvalues**)

Achieved by any orthonormal basis for the span of  $k$  corresponding eigenvectors  
(span of **top  $k$  eigenvectors**, a.k.a. **top  $k$  eigenspace**)

If  $\vec{v}_1, \dots, \vec{v}_k$  are top  $k$  eigenvectors of  $\text{cov}(\vec{X})$ , we create  $k$  "eigenfeatures"  $\vec{v}_i \cdot \vec{X}$  for  $i = 1, \dots, k$

# Optimality of the top $k$ eigenvectors

Question: What's the maximum variance "accounted for" by  $k$  dimensional subspace?

$$\max_{\substack{\vec{\alpha}_1, \dots, \vec{\alpha}_k \in \mathbb{R}^d \\ \text{orthonormal}}} \sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$$

- ▶ Answer: Sum of  $k$  largest eigenvalues of  $\text{cov}(\vec{X})$   
(sum of the **top  $k$  eigenvalues**)

Achieved by any orthonormal basis for the span of  $k$  corresponding eigenvectors  
(span of **top  $k$  eigenvectors**, a.k.a. **top  $k$  eigenspace**)

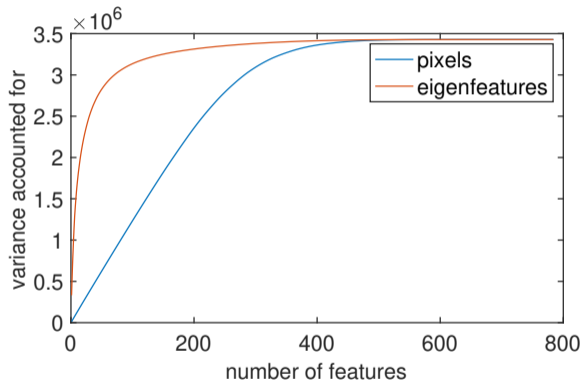
If  $\vec{v}_1, \dots, \vec{v}_k$  are top  $k$  eigenvectors of  $\text{cov}(\vec{X})$ , we create  $k$  "eigenfeatures"  $\vec{v}_i \cdot \vec{X}$  for  $i = 1, \dots, k$

Vectors  $\vec{v}_1, \dots, \vec{v}_k$  are called the **top  $k$  principal components** of (the distribution of)  $\vec{X}$   
Computing these "eigenfeatures" is called **( $k$ -dim.) Principal Components Analysis (PCA)**



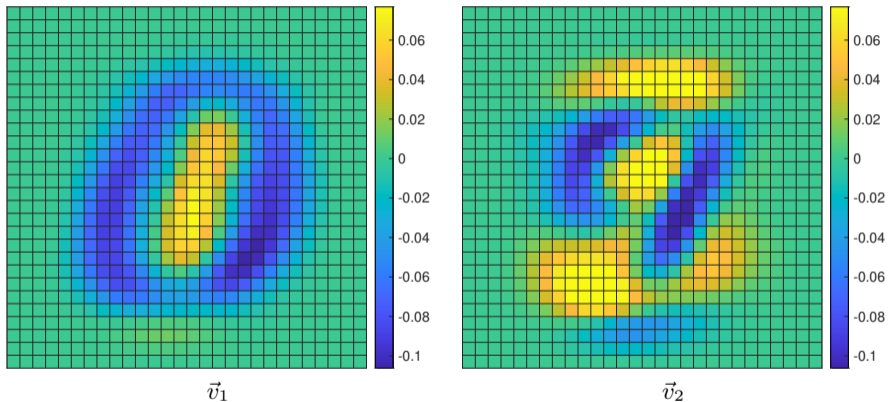
# Variance accounted for on MNIST data

MNIST data: Need  $k = 87$  “eigenfeatures” to reduce expected reconstruction error to  $0.1 \times \mathbb{E}\|\vec{X} - \vec{\mu}\|_2^2$



(“Variance accounted for” =  $\sum_{i=1}^k \text{var}(\vec{\alpha}_i \cdot \vec{X})$ )

# MNIST eigenvectors



Orthogonality:  $\vec{v}_1 \cdot \vec{v}_2 = 0$

# New features from PCA

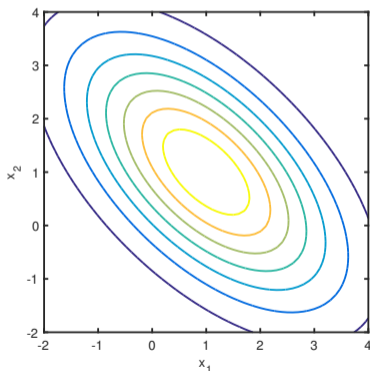
Properties of eigenfeatures  $\vec{v}_i \cdot \vec{X}$  for  $i = 1, \dots, k$ :

- ▶  $\text{var}(\vec{v}_i \cdot \vec{X}) = \lambda_i$  for all  $i$
- ▶  $\text{cov}(\vec{v}_i \cdot \vec{X}, \vec{v}_j \cdot \vec{X}) = 0$  for  $i \neq j$

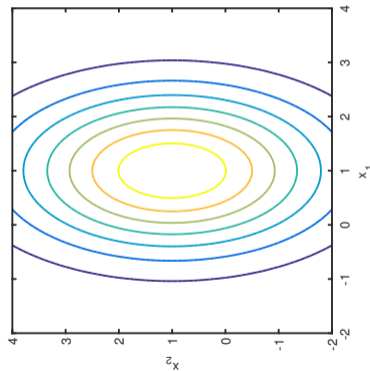
# New features from PCA

Properties of eigenfeatures  $\vec{v}_i \cdot \vec{X}$  for  $i = 1, \dots, k$ :

- ▶  $\text{var}(\vec{v}_i \cdot \vec{X}) = \lambda_i$  for all  $i$
- ▶  $\text{cov}(\vec{v}_i \cdot \vec{X}, \vec{v}_j \cdot \vec{X}) = 0$  for  $i \neq j$



Original features



Eigenfeatures

# PCA and beyond

# PCA and beyond

- ▶ Major use of PCA: **Dimension reduction**

Map data points in  $\mathbb{R}^d$  to  $\mathbb{R}^k$  for  $k < d$

(When  $k = 2$  or  $k = 3$ , useful for visualization, maybe!)

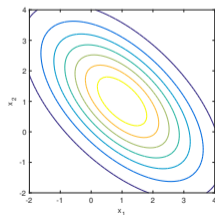
# PCA and beyond

- ▶ Major use of PCA: **Dimension reduction**  
Map data points in  $\mathbb{R}^d$  to  $\mathbb{R}^k$  for  $k < d$   
(When  $k = 2$  or  $k = 3$ , useful for visualization, maybe!)
- ▶  $k$ -dim. PCA + OLS = **Principal components regression (PCR)**  
Different inductive bias compared to ridge regression!

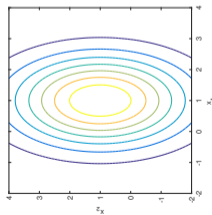
# PCA and beyond

- ▶ Major use of PCA: **Dimension reduction**  
Map data points in  $\mathbb{R}^d$  to  $\mathbb{R}^k$  for  $k < d$   
(When  $k = 2$  or  $k = 3$ , useful for visualization, maybe!)
- ▶  $k$ -dim. PCA + OLS = **Principal components regression (PCR)**  
Different inductive bias compared to ridge regression!
- ▶ **Whitening**: Transforming to isotropy

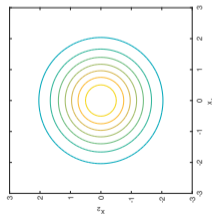
$$\frac{1}{\sqrt{\lambda_i}} \vec{v}_i \cdot \vec{X} \quad \text{for } i = 1, \dots, k$$



Original features



Eigenfeatures



Whitened features



## Postscript: “Un-centered” PCA

Everything about eigen{values,vectors} of  $\text{cov}(\vec{X})$  also applies to **second-moment matrix**  $\mathbb{E}[\vec{X}\vec{X}^\top]$ ,  
... with following changes:

- ▶ Replace  $\text{var}(\vec{\alpha} \cdot \vec{X})$  with  $\mathbb{E}[(\vec{\alpha} \cdot \vec{X})^2]$
- ▶ Change “expected reconstruction error” to  $\mathbb{E}[\|\vec{X} - \Pi^\top \vec{X}\|_2^2]$

## Postscript: “Un-centered” PCA

Everything about eigen{values,vectors} of  $\text{cov}(\vec{X})$  also applies to **second-moment matrix**  $\mathbb{E}[\vec{X}\vec{X}^\top]$ ,  
... with following changes:

- ▶ Replace  $\text{var}(\vec{\alpha} \cdot \vec{X})$  with  $\mathbb{E}[(\vec{\alpha} \cdot \vec{X})^2]$
- ▶ Change “expected reconstruction error” to  $\mathbb{E}[\|\vec{X} - \Pi^\top \vec{X}\|_2^2]$

If  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^\top \in \mathbb{R}^{n \times d}$ , then (empirical) second-moment matrix is

$$\frac{1}{n} A^\top A$$

## Postscript: “Un-centered” PCA

Everything about eigen{values,vectors} of  $\text{cov}(\vec{X})$  also applies to **second-moment matrix**  $\mathbb{E}[\vec{X}\vec{X}^\top]$ ,  
... with following changes:

- ▶ Replace  $\text{var}(\vec{\alpha} \cdot \vec{X})$  with  $\mathbb{E}[(\vec{\alpha} \cdot \vec{X})^2]$
- ▶ Change “expected reconstruction error” to  $\mathbb{E}[\|\vec{X} - \Pi^\top \vec{X}\|_2^2]$

If  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^\top \in \mathbb{R}^{n \times d}$ , then (empirical) second-moment matrix is

$$\frac{1}{n} A^\top A$$

Question: How do eigen{values,vectors} of  $A^\top A$  relate to eigen{values,vectors} of  $\frac{1}{n} A^\top A$ ?

## Postscript: “Un-centered” PCA

Everything about eigen{values,vectors} of  $\text{cov}(\vec{X})$  also applies to **second-moment matrix**  $\mathbb{E}[\vec{X}\vec{X}^\top]$ ,  
... with following changes:

- ▶ Replace  $\text{var}(\vec{\alpha} \cdot \vec{X})$  with  $\mathbb{E}[(\vec{\alpha} \cdot \vec{X})^2]$
- ▶ Change “expected reconstruction error” to  $\mathbb{E}[\|\vec{X} - \Pi^\top \vec{X}\|_2^2]$

If  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^\top \in \mathbb{R}^{n \times d}$ , then (empirical) second-moment matrix is

$$\frac{1}{n} A^\top A$$

Question: How do eigen{values,vectors} of  $A^\top A$  relate to eigen{values,vectors} of  $\frac{1}{n} A^\top A$ ?

Question: How do eigen{values,vectors} of  $A^\top A$  relate to eigen{values,vectors} of  $A$ ?

## **Eigenvector computation**

# Simple approach to top eigenvector computation

Let eigendecomposition of  $A^T A$  be  $\sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T$

Question: How to compute top eigenvector  $\vec{v}_1$  of  $A^T A$ ?

# Simple approach to top eigenvector computation

Let eigendecomposition of  $A^T A$  be  $\sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T$

Question: How to compute top eigenvector  $\vec{v}_1$  of  $A^T A$ ?

► **Fact.** For any integer  $t$ ,  $(A^T A)^t = \sum_{i=1}^d \lambda_i^t \vec{v}_i \vec{v}_i^T$

# Simple approach to top eigenvector computation

Let eigendecomposition of  $A^T A$  be  $\sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T$

Question: How to compute top eigenvector  $\vec{v}_1$  of  $A^T A$ ?

► **Fact.** For any integer  $t$ ,  $(A^T A)^t = \sum_{i=1}^d \lambda_i^t \vec{v}_i \vec{v}_i^T$

If  $t$  is large, and  $\lambda_1 > \lambda_2$ , then

$$(A^T A)^t = \lambda_1^t \left( \vec{v}_1 \vec{v}_1^T + \sum_{i=2}^d \left( \frac{\lambda_i}{\lambda_1} \right)^t \vec{v}_i \vec{v}_i^T \right) \approx \lambda_1^t \vec{v}_1 \vec{v}_1^T$$

where  $\approx$  holds when  $t$  is sufficiently large



# Simple approach to top eigenvector computation

Let eigendecomposition of  $A^T A$  be  $\sum_{i=1}^d \lambda_i \vec{v}_i \vec{v}_i^T$

Question: How to compute top eigenvector  $\vec{v}_1$  of  $A^T A$ ?

► **Fact.** For any integer  $t$ ,  $(A^T A)^t = \sum_{i=1}^d \lambda_i^t \vec{v}_i \vec{v}_i^T$

If  $t$  is large, and  $\lambda_1 > \lambda_2$ , then

$$(A^T A)^t = \lambda_1^t \left( \vec{v}_1 \vec{v}_1^T + \sum_{i=2}^d \left( \frac{\lambda_i}{\lambda_1} \right)^t \vec{v}_i \vec{v}_i^T \right) \approx \lambda_1^t \vec{v}_1 \vec{v}_1^T$$

where  $\approx$  holds when  $t$  is sufficiently large

## Main idea:

If  $t$  is sufficiently large, multiplying  $(A^T A)^t$  by (almost) any vector yields a vector that is (approximately) in span of  $\vec{v}_1$

# Power iteration

Slow implementation of main idea: First compute  $(A^T A)^t$ , and then multiply by a vector

# Power iteration

Slow implementation of main idea: First compute  $(A^T A)^t$ , and then multiply by a vector

Better implementation:

## Power iteration (a.k.a. power method)

- ▶ Initialize  $\vec{x}^{(0)} \in \mathbb{R}^d$  randomly (e.g., draw from  $N(\vec{0}, I_d)$ )
- ▶ For  $t = 1, 2, \dots$  until stopping condition is satisfied:

$$\vec{y} := A\vec{x}^{(t-1)}$$

$$\vec{z} := A^T \vec{y}$$

$$\vec{x}^{(t)} := \vec{z} / \|\vec{z}\|_2$$

# Power iteration

Slow implementation of main idea: First compute  $(A^T A)^t$ , and then multiply by a vector

Better implementation:

## Power iteration (a.k.a. power method)

- ▶ Initialize  $\vec{x}^{(0)} \in \mathbb{R}^d$  randomly (e.g., draw from  $N(\vec{0}, I_d)$ )
- ▶ For  $t = 1, 2, \dots$  until stopping condition is satisfied:

$$\vec{y} := A\vec{x}^{(t-1)}$$

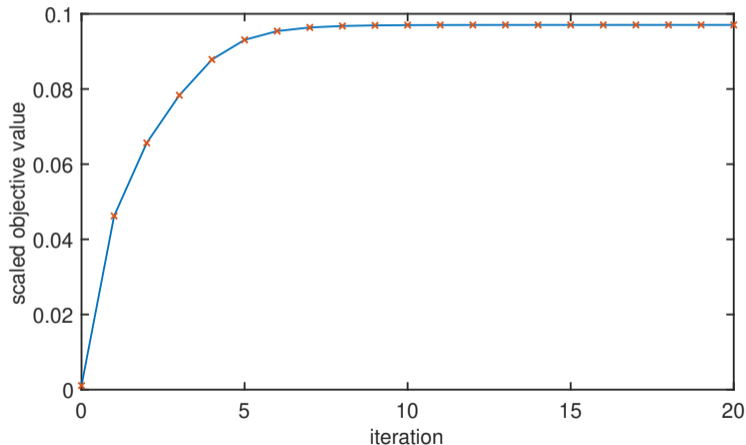
$$\vec{z} := A^T \vec{y}$$

$$\vec{x}^{(t)} := \vec{z} / \|\vec{z}\|_2$$

After  $t$  iterations, this computes a scaling of  $(A^T A)^t \vec{x}^{(0)}$   
(Scaling is due to normalization steps “ $\vec{x}^{(t)} := \vec{z} / \|\vec{z}\|_2$ ”)

## Example: Power iteration on centered MNIST data

$$\text{Scaled objective: } J(\vec{x}) = \frac{1}{\text{tr}(A^T A)} \cdot \frac{\vec{x}^T (A^T A) \vec{x}}{\|\vec{x}\|_2^2}$$



## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined

## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined
- ▶ Issue 2: If  $\lambda_1 > \lambda_2$  but very close, may need extremely large  $t$

## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined
- ▶ Issue 2: If  $\lambda_1 > \lambda_2$  but very close, may need extremely large  $t$
- ▶ Issue 3: What about  $\lambda_1$ ?



## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined
- ▶ Issue 2: If  $\lambda_1 > \lambda_2$  but very close, may need extremely large  $t$
- ▶ Issue 3: What about  $\lambda_1$ ?

$$\|A\vec{v}_1\|_2^2 = \lambda_1$$

## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined
- ▶ Issue 2: If  $\lambda_1 > \lambda_2$  but very close, may need extremely large  $t$
- ▶ Issue 3: What about  $\lambda_1$ ?  
 $\|A\vec{v}_1\|_2^2 = \lambda_1$
- ▶ Issue 4: What about  $\vec{v}_2, \vec{v}_3, \dots$  and  $\lambda_2, \lambda_3, \dots$ ?

## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined
- ▶ Issue 2: If  $\lambda_1 > \lambda_2$  but very close, may need extremely large  $t$
- ▶ Issue 3: What about  $\lambda_1$ ?  
 $\|A\vec{v}_1\|_2^2 = \lambda_1$
- ▶ Issue 4: What about  $\vec{v}_2, \vec{v}_3, \dots$  and  $\lambda_2, \lambda_3, \dots$ ?  
Can use **deflation**

## Some issues with power iteration

- ▶ “Issue” 1: If  $\lambda_1 = \lambda_2$ , then  $\vec{v}_1$  is not uniquely defined
- ▶ Issue 2: If  $\lambda_1 > \lambda_2$  but very close, may need extremely large  $t$
- ▶ Issue 3: What about  $\lambda_1$ ?

$$\|A\vec{v}_1\|_2^2 = \lambda_1$$

- ▶ Issue 4: What about  $\vec{v}_2, \vec{v}_3, \dots$  and  $\lambda_2, \lambda_3, \dots$ ?

Can use **deflation**

If you only care about subspace spanned by  $\vec{v}_1, \dots, \vec{v}_k$ : use **subspace iteration**

## **Low-rank Approximation**

# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

- ▶ “Low-rank” since  $\text{rank}(BC) \leq k$

# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

- ▶ “Low-rank” since  $\text{rank}(BC) \leq k$
- ▶  $k \leq \min\{n, d\}$  (typically  $k \ll \min\{n, d\}$ )



# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

- ▶ “Low-rank” since  $\text{rank}(BC) \leq k$
- ▶  $k \leq \min\{n, d\}$  (typically  $k \ll \min\{n, d\}$ )
- ▶ Measure of approximation error: **Frobenius norm**

$$J(B, C) = \|A - BC\|_F^2$$

where

$$\|M\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d M_{i,j}^2}$$

treats  $n \times d$  matrix  $M$  as vector in  $\mathbb{R}^{nd}$

# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

- ▶ “Low-rank” since  $\text{rank}(BC) \leq k$
- ▶  $k \leq \min\{n, d\}$  (typically  $k \ll \min\{n, d\}$ )
- ▶ Measure of approximation error: **Frobenius norm**

$$J(B, C) = \|A - BC\|_F^2$$

where

$$\|M\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d M_{i,j}^2}$$

treats  $n \times d$  matrix  $M$  as vector in  $\mathbb{R}^{nd}$

- ▶ If  $B$  is fixed, then this looks like the SSE objective from OLS!

# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

- ▶ “Low-rank” since  $\text{rank}(BC) \leq k$
- ▶  $k \leq \min\{n, d\}$  (typically  $k \ll \min\{n, d\}$ )
- ▶ Measure of approximation error: **Frobenius norm**

$$J(B, C) = \|A - BC\|_F^2$$

where

$$\|M\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d M_{i,j}^2}$$

treats  $n \times d$  matrix  $M$  as vector in  $\mathbb{R}^{nd}$

- ▶ If  $B$  is fixed, then this looks like the SSE objective from OLS!
- ▶ Same if  $C$  is fixed.

# Low-rank matrix approximation

Given data matrix  $A = [\vec{x}_1 | \cdots | \vec{x}_n]^T \in \mathbb{R}^{n \times d}$ , how can we best approximate it by the product  $BC$  for some  $B \in \mathbb{R}^{n \times k}$  and  $C \in \mathbb{R}^{k \times d}$ ?

- ▶ “Low-rank” since  $\text{rank}(BC) \leq k$
- ▶  $k \leq \min\{n, d\}$  (typically  $k \ll \min\{n, d\}$ )
- ▶ Measure of approximation error: **Frobenius norm**

$$J(B, C) = \|A - BC\|_F^2$$

where

$$\|M\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d M_{i,j}^2}$$

treats  $n \times d$  matrix  $M$  as vector in  $\mathbb{R}^{nd}$

- ▶ If  $B$  is fixed, then this looks like the SSE objective from OLS!
- ▶ Same if  $C$  is fixed.
- ▶ But apparently more complicated, since both  $B$  and  $C$  are optimization variables! (E.g.,  $J$  is not a convex function of  $(B, C)$ )

# Two optimization methods

## Gradient descent:

- ▶ Start with initial values  $B^{(0)} \in \mathbb{R}^{n \times k}$  and  $C^{(0)} \in \mathbb{R}^{k \times d}$
- ▶ For iteration  $t = 1, 2, \dots$  until stopping condition is satisfied:

$$B^{(t)} := B^{(t-1)} - \eta \frac{\partial J}{\partial B}(B^{(t-1)}, C^{(t-1)})$$

$$C^{(t)} := C^{(t-1)} - \eta \frac{\partial J}{\partial C}(B^{(t-1)}, C^{(t-1)})$$

# Two optimization methods

## Gradient descent:

- ▶ Start with initial values  $B^{(0)} \in \mathbb{R}^{n \times k}$  and  $C^{(0)} \in \mathbb{R}^{k \times d}$
- ▶ For iteration  $t = 1, 2, \dots$  until stopping condition is satisfied:

$$B^{(t)} := B^{(t-1)} - \eta \frac{\partial J}{\partial B}(B^{(t-1)}, C^{(t-1)})$$

$$C^{(t)} := C^{(t-1)} - \eta \frac{\partial J}{\partial C}(B^{(t-1)}, C^{(t-1)})$$

## Alternating minimization:

- ▶ Start with initial values  $B^{(0)} \in \mathbb{R}^{n \times k}$  and  $C^{(0)} \in \mathbb{R}^{k \times d}$
- ▶ For iteration  $t = 1, 2, \dots$  until stopping condition is satisfied:

$$B^{(t)} := \arg \min_{B \in \mathbb{R}^{n \times k}} J(B, C^{(t-1)})$$

$$C^{(t)} := \arg \min_{C \in \mathbb{R}^{k \times d}} J(B^{(t)}, C)$$

# Singular value decomposition

**Theorem** (Schmidt, 1907; Eckart & Young, 1936). Optimal solution to the low-rank matrix approximation problem is given by truncating the *singular value decomposition* of  $A$

# Singular value decomposition

**Theorem** (Schmidt, 1907; Eckart & Young, 1936). Optimal solution to the low-rank matrix approximation problem is given by truncating the *singular value decomposition* of  $A$

**Singular value decomposition (SVD)** of  $A$ : Every matrix  $A \in \mathbb{R}^{n \times d}$  (with  $r := \text{rank}(A)$ ) can be written as

$$A = \sum_{i=1}^r s_i \vec{u}_i \vec{v}_i^T$$

for some  $\{(s_i, \vec{u}_i, \vec{v}_i)\}_{i=1}^r$ , which satisfy

- ▶  $s_1 \geq \dots \geq s_r > 0$  (**singular values** of  $A$ ; sorting this way is by convention)
- ▶  $\vec{u}_1, \dots, \vec{u}_r \in \mathbb{R}^n$  are orthonormal (**left singular vectors** of  $A$ )
- ▶  $\vec{v}_1, \dots, \vec{v}_r \in \mathbb{R}^d$  are orthonormal (**right singular vectors** of  $A$ )



# Singular value decomposition

**Theorem** (Schmidt, 1907; Eckart & Young, 1936). Optimal solution to the low-rank matrix approximation problem is given by truncating the *singular value decomposition* of  $A$

**Singular value decomposition (SVD)** of  $A$ : Every matrix  $A \in \mathbb{R}^{n \times d}$  (with  $r := \text{rank}(A)$ ) can be written as

$$A = \sum_{i=1}^r s_i \vec{u}_i \vec{v}_i^T$$

for some  $\{(s_i, \vec{u}_i, \vec{v}_i)\}_{i=1}^r$ , which satisfy

- ▶  $s_1 \geq \dots \geq s_r > 0$  (**singular values** of  $A$ ; sorting this way is by convention)
- ▶  $\vec{u}_1, \dots, \vec{u}_r \in \mathbb{R}^n$  are orthonormal (**left singular vectors** of  $A$ )
- ▶  $\vec{v}_1, \dots, \vec{v}_r \in \mathbb{R}^d$  are orthonormal (**right singular vectors** of  $A$ )

Can also write as

$$A = USV^T$$

where  $U := [\vec{u}_1 | \dots | \vec{u}_r] \in \mathbb{R}^{n \times r}$ ,  $S := \text{diag}(s_1, \dots, s_r) \in \mathbb{R}^{r \times r}$ , and  $V := [\vec{v}_1 | \dots | \vec{v}_r] \in \mathbb{R}^{d \times r}$

# Truncated SVD

Suppose SVD of  $A$  is given by

$$A = \sum_{i=1}^{\text{rank}(A)} s_i \vec{u}_i \vec{v}_i^T$$

# Truncated SVD

Suppose SVD of  $A$  is given by

$$A = \sum_{i=1}^{\text{rank}(A)} s_i \vec{u}_i \vec{v}_i^T$$

**Theorem** (Schmidt, 1907; Eckart & Young, 1936). The **rank- $k$  SVD** of  $A$  (for  $k \leq \text{rank}(A)$ ),

$$A_k := \sum_{i=1}^k s_i \vec{u}_i \vec{v}_i^T,$$

satisfies

$$\|A - A_k\|_F^2 = \min_{\substack{M \in \mathbb{R}^{n \times d} \\ \text{s.t. } \text{rank}(M) \leq k}} \|A - M\|_F^2 = \sum_{i=k+1}^{\text{rank}(A)} s_i^2$$

# Truncated SVD

Suppose SVD of  $A$  is given by

$$A = \sum_{i=1}^{\text{rank}(A)} s_i \vec{u}_i \vec{v}_i^\top$$

**Theorem** (Schmidt, 1907; Eckart & Young, 1936). The **rank- $k$  SVD** of  $A$  (for  $k \leq \text{rank}(A)$ ),

$$A_k := \sum_{i=1}^k s_i \vec{u}_i \vec{v}_i^\top,$$

satisfies

$$\|A - A_k\|_F^2 = \min_{\substack{M \in \mathbb{R}^{n \times d} \\ \text{s.t. } \text{rank}(M) \leq k}} \|A - M\|_F^2 = \sum_{i=k+1}^{\text{rank}(A)} s_i^2$$

Can also write  $A_k = U_k S_k V_k^\top$  where  $U_k := [\vec{u}_1 | \cdots | \vec{u}_k] \in \mathbb{R}^{n \times k}$ ,  $S_k := \text{diag}(s_1, \dots, s_k) \in \mathbb{R}^{k \times k}$ , and  $V_k := [\vec{v}_1 | \cdots | \vec{v}_k] \in \mathbb{R}^{d \times k}$

# Connection to PCA

Observation:

$$A^T A = (VSU^T)(USV^T) = VS^2V^T = \sum_{i=1}^r s_i^2 \vec{v}_i \vec{v}_i^T$$

- ▶ Non-zero eigenvalues of  $A^T A$  are squares of singular values of  $A$
- ▶ Corresponding eigenvectors of  $A^T A$  are right singular vectors of  $A$

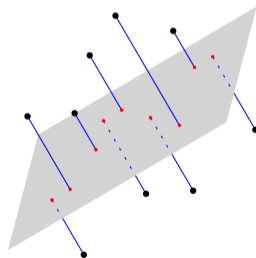
# Connection to PCA

Observation:

$$A^T A = (VSU^T)(USV^T) = VS^2V^T = \sum_{i=1}^r s_i^2 \vec{v}_i \vec{v}_i^T$$

- ▶ Non-zero eigenvalues of  $A^T A$  are squares of singular values of  $A$
- ▶ Corresponding eigenvectors of  $A^T A$  are right singular vectors of  $A$

Computing rank- $k$  SVD is same as applying (un-centered) dim.  $k$  PCA feature mapping!



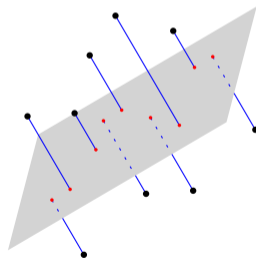
# Connection to PCA

Observation:

$$A^T A = (VSU^T)(USV^T) = VS^2V^T = \sum_{i=1}^r s_i^2 \vec{v}_i \vec{v}_i^T$$

- ▶ Non-zero eigenvalues of  $A^T A$  are squares of singular values of  $A$
- ▶ Corresponding eigenvectors of  $A^T A$  are right singular vectors of  $A$

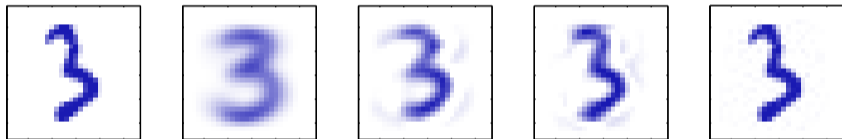
Computing rank- $k$  SVD is same as applying (un-centered) dim.  $k$  PCA feature mapping!



Rows of  $U_k S_k$  are the  $k$ -dimensional vectors of “eigenfeatures”

# MNIST data example

MNIST data example: Compare original image to corresponding row of  $A_k$  ( $k \in \{1, 10, 50, 200\}$ )





## **Application: Latent Semantic Analysis**

# Document modeling

Suppose you have corpus of  $n$  text documents

- ▶ Each document represented using “bag-of-words” count vectors
- ▶  $d =$  vocabulary size

		aardvark	abacus	abalone	...
	doc 1	3	0	0	...
$A :=$	doc 2	7	0	4	...
	doc 3	0	4	0	...
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

# Document modeling

Suppose you have corpus of  $n$  text documents

- ▶ Each document represented using “bag-of-words” count vectors
- ▶  $d =$  vocabulary size

		aardvark	abacus	abalone	...
	doc 1	3	0	0	...
$A :=$	doc 2	7	0	4	...
	doc 3	0	4	0	...
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

Question: How to find documents related to mathematics?

# Document modeling

Suppose you have corpus of  $n$  text documents

- ▶ Each document represented using “bag-of-words” count vectors
- ▶  $d =$  vocabulary size

		aardvark	abacus	abalone	...
	doc 1	3	0	0	...
$A :=$	doc 2	7	0	4	...
	doc 3	0	4	0	...
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

Question: How to find documents related to mathematics?

- ▶ Many (perhaps most?) relevant documents won't contain the word “mathematics”

# Document modeling

Suppose you have corpus of  $n$  text documents

- ▶ Each document represented using “bag-of-words” count vectors
- ▶  $d =$  vocabulary size

		aardvark	abacus	abalone	...
	doc 1	3	0	0	...
$A :=$	doc 2	7	0	4	...
	doc 3	0	4	0	...
	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$

Question: How to find documents related to mathematics?

- ▶ Many (perhaps most?) relevant documents won't contain the word “mathematics”
- ▶ However, documents about mathematics are likely to contain word that often co-occur with “mathematics” or other math-y words (e.g., “theorem”, “proof”, “generality”)

# Latent semantic analysis

(Deerwester, Dumais, Furnas, Landauer, Harshman, 1990)

Compute rank- $k$  SVD of  $A$ :  $A_k = U_k S_k V_k^T$

# Latent semantic analysis

(Deerwester, Dumais, Furnas, Landauer, Harshman, 1990)

Compute rank- $k$  SVD of  $A$ :  $A_k = U_k S_k V_k^\top$

- ▶ Documents in corpus are represented by their corresponding rows in  $U_k S_k \in \mathbb{R}^{n \times k}$

# Latent semantic analysis

(Deerwester, Dumais, Furnas, Landauer, Harshman, 1990)

Compute rank- $k$  SVD of  $A$ :  $A_k = U_k S_k V_k^\top$

- ▶ Documents in corpus are represented by their corresponding rows in  $U_k S_k \in \mathbb{R}^{n \times k}$
- ▶ For a new document, compute its bag-of-words vector  $\vec{x} \in \mathbb{R}^d$ , then compute its rank- $k$  encoding

$$\vec{\phi}(\vec{x}) := (\vec{v}_1 \cdot \vec{x}, \dots, \vec{v}_k \cdot \vec{x})$$



# Latent semantic analysis

(Deerwester, Dumais, Furnas, Landauer, Harshman, 1990)

Compute rank- $k$  SVD of  $A$ :  $A_k = U_k S_k V_k^\top$

- ▶ Documents in corpus are represented by their corresponding rows in  $U_k S_k \in \mathbb{R}^{n \times k}$
- ▶ For a new document, compute its bag-of-words vector  $\vec{x} \in \mathbb{R}^d$ , then compute its rank- $k$  encoding

$$\vec{\phi}(\vec{x}) := (\vec{v}_1 \cdot \vec{x}, \dots, \vec{v}_k \cdot \vec{x})$$

- ▶ Compare two documents  $\vec{x}$  and  $\vec{x}'$  by the **cosine similarity** of their encodings

$$\text{sim}(\vec{x}, \vec{x}') = \frac{\vec{\phi}(\vec{x}) \cdot \vec{\phi}(\vec{x}')}{\|\vec{\phi}(\vec{x})\|_2 \|\vec{\phi}(\vec{x}')\|_2}$$

# Latent semantic analysis

(Deerwester, Dumais, Furnas, Landauer, Harshman, 1990)

Compute rank- $k$  SVD of  $A$ :  $A_k = U_k S_k V_k^\top$

- ▶ Documents in corpus are represented by their corresponding rows in  $U_k S_k \in \mathbb{R}^{n \times k}$
- ▶ For a new document, compute its bag-of-words vector  $\vec{x} \in \mathbb{R}^d$ , then compute its rank- $k$  encoding

$$\vec{\phi}(\vec{x}) := (\vec{v}_1 \cdot \vec{x}, \dots, \vec{v}_k \cdot \vec{x})$$

- ▶ Compare two documents  $\vec{x}$  and  $\vec{x}'$  by the **cosine similarity** of their encodings

$$\text{sim}(\vec{x}, \vec{x}') = \frac{\vec{\phi}(\vec{x}) \cdot \vec{\phi}(\vec{x}')}{\|\vec{\phi}(\vec{x})\|_2 \|\vec{\phi}(\vec{x}')\|_2}$$

- ▶ Search queries: Construct “pseudo-document” — i.e., bag-of-words with non-zero counts for query terms — then find documents in corpus of high cosine similarity

# Why low-rank approximation may be helpful

**Topic model:** Suppose there are  $k$  “topics” that documents might be about

# Why low-rank approximation may be helpful

**Topic model:** Suppose there are  $k$  “topics” that documents might be about

- ▶ Each topic  $t$  specifies a distribution over words

$$\vec{c}_t = (c_{t,1}, \dots, c_{t,d}), \quad c_{t,j} = \Pr[\text{word } j \mid \text{topic } t]$$

# Why low-rank approximation may be helpful

**Topic model:** Suppose there are  $k$  “topics” that documents might be about

- ▶ Each topic  $t$  specifies a distribution over words

$$\vec{c}_t = (c_{t,1}, \dots, c_{t,d}), \quad c_{t,j} = \Pr[\text{word } j \mid \text{topic } t]$$

- ▶ Each document  $i$  specifies a non-negative weighting over topics

$$\vec{b}_i = (b_{i,1}, \dots, b_{i,k}), \quad \sum_{t=1}^k b_{i,t} = \text{length of document } i$$

# Why low-rank approximation may be helpful

**Topic model:** Suppose there are  $k$  “topics” that documents might be about

- ▶ Each topic  $t$  specifies a distribution over words

$$\vec{c}_t = (c_{t,1}, \dots, c_{t,d}), \quad c_{t,j} = \Pr[\text{word } j \mid \text{topic } t]$$

- ▶ Each document  $i$  specifies a non-negative weighting over topics

$$\vec{b}_i = (b_{i,1}, \dots, b_{i,k}), \quad \sum_{t=1}^k b_{i,t} = \text{length of document } i$$

- ▶ Document  $i$ :  $b_{i,1}$  words drawn IID from topic 1,  $b_{i,2}$  words drawn IID from topic 2, etc.  
In this model,

$$\mathbb{E}[A] = BC,$$

where

$$B = \begin{bmatrix} \leftarrow & \vec{b}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{b}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C = \begin{bmatrix} \leftarrow & \vec{c}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{c}_k^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$$

# Why low-rank approximation may be helpful

**Topic model:** Suppose there are  $k$  “topics” that documents might be about

- ▶ Each topic  $t$  specifies a distribution over words

$$\vec{c}_t = (c_{t,1}, \dots, c_{t,d}), \quad c_{t,j} = \Pr[\text{word } j \mid \text{topic } t]$$

- ▶ Each document  $i$  specifies a non-negative weighting over topics

$$\vec{b}_i = (b_{i,1}, \dots, b_{i,k}), \quad \sum_{t=1}^k b_{i,t} = \text{length of document } i$$

- ▶ Document  $i$ :  $b_{i,1}$  words drawn IID from topic 1,  $b_{i,2}$  words drawn IID from topic 2, etc. In this model,

$$\mathbb{E}[A] = BC,$$

where

$$B = \begin{bmatrix} \leftarrow & \vec{b}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{b}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C = \begin{bmatrix} \leftarrow & \vec{c}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{c}_k^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$$

- ▶ Hence  $A = BC + \text{“noise”}$ ; hope that  $A_k$  avoids the noise

# Caveats

- ▶ Might hope that splitting apart  $A_k = U_k S_k V_k^T$  as  $U_k S_k$  and  $V_k^T$  will somehow give something like  $B$  and  $C$  from the topic model  
I.e., hope that  $B \approx U_k S_k$  and  $C \approx V_k^T$



# Caveats

- ▶ Might hope that splitting apart  $A_k = U_k S_k V_k^T$  as  $U_k S_k$  and  $V_k^T$  will somehow give something like  $B$  and  $C$  from the topic model

I.e., hope that  $B \approx U_k S_k$  and  $C \approx V_k^T$

- ▶ But these hopes are unlikely to be fulfilled ☹

# Caveats

- ▶ Might hope that splitting apart  $A_k = U_k S_k V_k^T$  as  $U_k S_k$  and  $V_k^T$  will somehow give something like  $B$  and  $C$  from the topic model  
I.e., hope that  $B \approx U_k S_k$  and  $C \approx V_k^T$ 
  - ▶ But these hopes are unlikely to be fulfilled ☹
- ▶ Other approaches to finding low-rank approximations with “interpretable” factors:
  - ▶ Non-negative Matrix Factorization (Lee & Seung, 1999)
  - ▶ Latent Dirichlet Allocation (Pritchard, Stephens, & Donnelly, 2000; Blei, Ng, & Jordan, 2003)
  - ▶ etc.

## **Application: Matrix Completion**

# Netflix problem

Suppose ratings of movies given by users are arranged in matrix  $A \in \mathbb{R}^{n \times d}$ , where  $A_{i,j}$  is rating given by user  $i$  for movie  $j$

# Netflix problem

Suppose ratings of movies given by users are arranged in matrix  $A \in \mathbb{R}^{n \times d}$ , where  $A_{i,j}$  is rating given by user  $i$  for movie  $j$

- ▶ Netflix:  $n = 480000$ ,  $d = 18000$ ; on average, each user rates 200 movies
- ▶ Most entries of  $A$  are “missing” (i.e., not observed)
- ▶ Goal: Predict values of missing entries

# Netflix problem

Suppose ratings of movies given by users are arranged in matrix  $A \in \mathbb{R}^{n \times d}$ , where  $A_{i,j}$  is rating given by user  $i$  for movie  $j$

- ▶ Netflix:  $n = 480000$ ,  $d = 18000$ ; on average, each user rates 200 movies
- ▶ Most entries of  $A$  are “missing” (i.e., not observed)
- ▶ Goal: Predict values of missing entries

Idea: Approximate  $A$  with low-rank matrix, i.e., find

$$B = \begin{bmatrix} \leftarrow & \vec{b}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{b}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C = \begin{bmatrix} \uparrow & \cdots & \uparrow \\ \vec{c}_1 & \cdots & \vec{c}_d \\ \downarrow & \cdots & \downarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$$

with goal of minimizing  $J(B, C) = \|A - BC\|_F^2$

- ▶  $B$ : “user features” (one row per user)

# Netflix problem

Suppose ratings of movies given by users are arranged in matrix  $A \in \mathbb{R}^{n \times d}$ , where  $A_{i,j}$  is rating given by user  $i$  for movie  $j$

- ▶ Netflix:  $n = 480000$ ,  $d = 18000$ ; on average, each user rates 200 movies
- ▶ Most entries of  $A$  are “missing” (i.e., not observed)
- ▶ Goal: Predict values of missing entries

Idea: Approximate  $A$  with low-rank matrix, i.e., find

$$B = \begin{bmatrix} \leftarrow & \vec{b}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{b}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C = \begin{bmatrix} \uparrow & \cdots & \uparrow \\ \vec{c}_1 & \cdots & \vec{c}_d \\ \downarrow & \cdots & \downarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$$

with goal of minimizing  $J(B, C) = \|A - BC\|_F^2$

- ▶  $B$ : “user features” (one row per user)
- ▶  $C$ : “movie features” (one column per movie)

# Netflix problem

Suppose ratings of movies given by users are arranged in matrix  $A \in \mathbb{R}^{n \times d}$ , where  $A_{i,j}$  is rating given by user  $i$  for movie  $j$

- ▶ Netflix:  $n = 480000$ ,  $d = 18000$ ; on average, each user rates 200 movies
- ▶ Most entries of  $A$  are “missing” (i.e., not observed)
- ▶ Goal: Predict values of missing entries

Idea: Approximate  $A$  with low-rank matrix, i.e., find

$$B = \begin{bmatrix} \leftarrow & \vec{b}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{b}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C = \begin{bmatrix} \uparrow & \cdots & \uparrow \\ \vec{c}_1 & \cdots & \vec{c}_d \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$$

with goal of minimizing  $J(B, C) = \|A - BC\|_F^2$

- ▶  $B$ : “user features” (one row per user)
- ▶  $C$ : “movie features” (one column per movie)
- ▶ Prediction of user  $i$ 's rating of movie  $j$ :  $\vec{b}_i \cdot \vec{c}_j$



# Netflix problem

Suppose ratings of movies given by users are arranged in matrix  $A \in \mathbb{R}^{n \times d}$ , where  $A_{i,j}$  is rating given by user  $i$  for movie  $j$

- ▶ Netflix:  $n = 480000$ ,  $d = 18000$ ; on average, each user rates 200 movies
- ▶ Most entries of  $A$  are “missing” (i.e., not observed)
- ▶ Goal: Predict values of missing entries

Idea: Approximate  $A$  with low-rank matrix, i.e., find

$$B = \begin{bmatrix} \leftarrow & \vec{b}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \vec{b}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times k}, \quad C = \begin{bmatrix} \uparrow & \cdots & \uparrow \\ \vec{c}_1 & \cdots & \vec{c}_d \\ \downarrow & & \downarrow \end{bmatrix} \in \mathbb{R}^{k \times d}$$

with goal of minimizing  $J(B, C) = \|A - BC\|_F^2$

- ▶  $B$ : “user features” (one row per user)
- ▶  $C$ : “movie features” (one column per movie)
- ▶ Prediction of user  $i$ 's rating of movie  $j$ :  $\vec{b}_i \cdot \vec{c}_j$
- ▶ Note: If all entries of  $A$  were observed, we could do this with truncated SVD

# Matrix completion

Problem: We only have access to some subset  $\Omega \subset [n] \times [d]$  of entries of  $A$

# Matrix completion

Problem: We only have access to some subset  $\Omega \subset [n] \times [d]$  of entries of  $A$

Common solution (used by Simon Funk in Netflix challenge): Just define objective over entries in  $\Omega$

$$J_{\Omega}(B, C) = \|A - BC\|_{\Omega}^2$$

where

$$\|M\|_{\Omega} := \sqrt{\sum_{(i,j) \in \Omega} M_{i,j}^2}$$

# Matrix completion

Problem: We only have access to some subset  $\Omega \subset [n] \times [d]$  of entries of  $A$

Common solution (used by Simon Funk in Netflix challenge): Just define objective over entries in  $\Omega$

$$J_{\Omega}(B, C) = \|A - BC\|_{\Omega}^2$$

where

$$\|M\|_{\Omega} := \sqrt{\sum_{(i,j) \in \Omega} M_{i,j}^2}$$

- ▶ Now (attempt to) minimize  $J_{\Omega}$  using gradient descent or stochastic gradient descent

# Feature representations from matrix completion

- ▶ MovieLens data set ( $n = 6040$  users,  $d = 3952$  movies,  $|\Omega| = 800000$  ratings)
- ▶ Attempted to minimize  $J_{\Omega}(B, C)$  using stochastic gradient descent

# Feature representations from matrix completion

- ▶ MovieLens data set ( $n = 6040$  users,  $d = 3952$  movies,  $|\Omega| = 800000$  ratings)
- ▶ Attempted to minimize  $J_{\Omega}(B, C)$  using stochastic gradient descent

Question: Are  $\vec{c}_1, \dots, \vec{c}_d \in \mathbb{R}^k$  useful feature vectors for movies?

# Feature representations from matrix completion

- ▶ MovieLens data set ( $n = 6040$  users,  $d = 3952$  movies,  $|\Omega| = 800000$  ratings)
- ▶ Attempted to minimize  $J_{\Omega}(B, C)$  using stochastic gradient descent

Question: Are  $\vec{c}_1, \dots, \vec{c}_d \in \mathbb{R}^k$  useful feature vectors for movies?

Some close pairs of movie feature vectors ( $\vec{c}_i, \vec{c}_j$ ):

- ▶ Toy Story (1995), Toy Story 2 (1999)
- ▶ Sense and Sensibility (1995), Emma (1996)
- ▶ Heat (1995), Carlito's Way (1993)
- ▶ The Crow (1994), Blade (1998)
- ▶ Forrest Gump (1994), Dances with Wolves (1990)
- ▶ Mrs. Doubtfire (1993), The Bodyguard (1992)

# Summary

- ▶ Many applications of low-rank approximations
- ▶ Often PCA/SVD is a good approach!
- ▶ Sometimes need to go beyond PCA/SVD