

COMS 4771 Lecture 17

1. Principal component analysis
2. Singular value decomposition

REPRESENTATION LEARNING

USEFUL REPRESENTATIONS OF DATA

Representation learning:

- ▶ **Given:** raw feature vectors $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^d$.
- ▶ **Goal:** learn a “useful” feature transformation $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$.
(Often $k \ll d$ —i.e., *dimensionality reduction*—but not always.)

USEFUL REPRESENTATIONS OF DATA

Representation learning:

- ▶ **Given:** raw feature vectors $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^d$.
- ▶ **Goal:** learn a “useful” feature transformation $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$.
(Often $k \ll d$ —i.e., *dimensionality reduction*—but not always.)

Can then use ϕ as a feature map for supervised learning.

USEFUL REPRESENTATIONS OF DATA

Representation learning:

- ▶ **Given:** raw feature vectors $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^d$.
- ▶ **Goal:** learn a “useful” feature transformation $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$.
(Often $k \ll d$ —i.e., *dimensionality reduction*—but not always.)

Can then use ϕ as a feature map for supervised learning.

Some previously encountered examples:

- ▶ Feature maps corresponding to kernels (+approximations).
(This is usually *data-oblivious*—feature map doesn’t depend on the data.)
- ▶ Centering; standardization.

USEFUL REPRESENTATIONS OF DATA

tricky thing: the new base vector(sub-space)'s dimensions are in the same length of original vector(original space)

Representation learning:

- ▶ **Given:** raw feature vectors $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^d$.
- ▶ **Goal:** learn a “useful” feature transformation $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^k$.
(Often $k \ll d$ —i.e., *dimensionality reduction*—but not always.)

Can then use ϕ as a feature map for supervised learning.

Some previously encountered examples:

- ▶ Feature maps corresponding to kernels (+approximations).
(This is usually *data-oblivious*—feature map doesn't depend on the data.)
- ▶ Centering; standardization.

What are other desirable properties of a feature representation?

PRINCIPAL COMPONENT ANALYSIS

DIMENSIONALITY REDUCTION VIA PROJECTIONS

Projections

training example, no labels

- **Input:** $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)} \in \mathbb{R}^d$, target dimensionality $k \in \mathbb{N}$.
- **Output:** a k -dimensional subspace, represented by an orthonormal basis $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k \in \mathbb{R}^d$.
- **Projection:** Formally, projection of $\mathbf{x} \in \mathbb{R}^d$ to $\text{span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k)$ is

$$\underbrace{\left(\sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^\top \right)}_{\Pi} \mathbf{x} = \sum_{i=1}^k \langle \mathbf{v}_i, \mathbf{x} \rangle \mathbf{v}_i \in \mathbb{R}^d.$$

still in original basis, sub-space

the product is used to compute move along that direction

But, we can simply represent the projection in terms of its coefficients

w.r.t. the orthonormal basis $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$:

$$\phi(\mathbf{x}) := \begin{bmatrix} \langle \mathbf{v}_1, \mathbf{x} \rangle \\ \langle \mathbf{v}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \mathbf{v}_k, \mathbf{x} \rangle \end{bmatrix} \in \mathbb{R}^k.$$

note the transformation

PROJECTION OF MINIMUM RESIDUAL SQUARED ERROR

note: the base vector's
dimension does not change

Minimize residual squared error

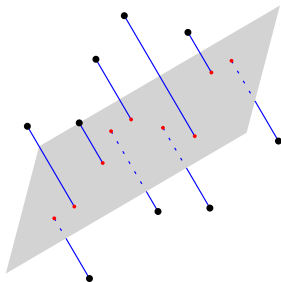
Goal: find k -dimensional projector $\Pi: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that the total residual squared error

$$\sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \Pi \mathbf{x}^{(i)} \right\|_2^2$$

is as small as possible.

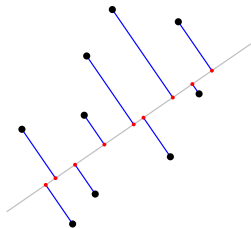
maximum variance
in that direction

transformation
in the same
dimensional
space



PROJECTION OF MINIMUM RESIDUAL SQUARED ERROR

$k = 1$ case ($\Pi = vv^\top$)

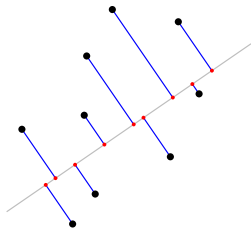


Goal: find unit vector $v \in \mathbb{R}^d$ to minimize

$$\sum_{i=1}^n \left\| \mathbf{x}^{(i)} - vv^\top \mathbf{x}^{(i)} \right\|_2^2$$

PROJECTION OF MINIMUM RESIDUAL SQUARED ERROR

$k = 1$ case ($\Pi = vv^\top$)

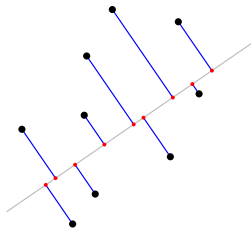


Goal: find unit vector $v \in \mathbb{R}^d$ to minimize

$$\begin{aligned} & \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - vv^\top \mathbf{x}^{(i)} \right\|_2^2 \\ &= \sum_{i=1}^n \left\| \mathbf{x}^{(i)} \right\|_2^2 - v^\top \left(\sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} \right) v \end{aligned}$$

PROJECTION OF MINIMUM RESIDUAL SQUARED ERROR

$k = 1$ case ($\Pi = vv^\top$)



Goal: find unit vector $v \in \mathbb{R}^d$ to minimize

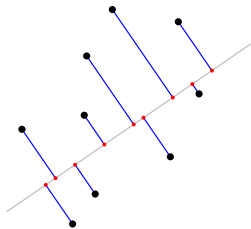
$$\begin{aligned} & \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - vv^\top \mathbf{x}^{(i)} \right\|_2^2 \\ &= \sum_{i=1}^n \left\| \mathbf{x}^{(i)} \right\|_2^2 - v^\top \left(\sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} \right) v \\ &= \sum_{i=1}^n \left\| \mathbf{x}^{(i)} \right\|_2^2 - v^\top \mathbf{X}^\top \mathbf{X} v. \end{aligned}$$

($\mathbf{x}^{(i)\top}$ is i -th row of $\mathbf{X} \in \mathbb{R}^{n \times d}$.)

PROJECTION OF MINIMUM RESIDUAL SQUARED ERROR

just a single vector  not convex

$k = 1$ case ($\Pi = vv^\top$)



Goal: find unit vector $v \in \mathbb{R}^d$ to minimize

$$\begin{aligned} & \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - vv^\top \mathbf{x}^{(i)} \right\|_2^2 \\ &= \sum_{i=1}^n \left\| \mathbf{x}^{(i)} \right\|_2^2 - v^\top \left(\sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top} \right) v \\ &= \sum_{i=1}^n \left\| \mathbf{x}^{(i)} \right\|_2^2 - v^\top \mathbf{X}^\top \mathbf{X} v. \end{aligned}$$

($\mathbf{x}^{(i)\top}$ is i -th row of $\mathbf{X} \in \mathbb{R}^{n \times d}$.)

to be convex, the \mathbf{X} must be positive definite

$$\arg \min_{v \in \mathbb{R}^d: \|v\|_2=1} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - vv^\top \mathbf{x}^{(i)} \right\|_2^2 \equiv \boxed{\arg \max_{v \in \mathbb{R}^d: \|v\|_2=1} v^\top \mathbf{X}^\top \mathbf{X} v.}$$

just see it in the original space. the new
coordination is the dot on the line

EIGENDECOMPOSITIONS

Every symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ guaranteed to have eigendecomposition with real eigenvalues:

$$\begin{array}{ccccccc} \boxed{\phantom{\mathbf{A}}} & = & \boxed{\phantom{\mathbf{V}}} & \boxed{\phantom{\mathbf{\Lambda}}} & \boxed{\phantom{\mathbf{V}^\top}} & = & \sum_{i=1}^d \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \\ \mathbf{A} & & \mathbf{V} & \mathbf{\Lambda} & \mathbf{V}^\top & & \\ (d \times d) & & (d \times d) & (d \times d) & (d \times d) & & \end{array}$$

real **eigenvalues**: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ ($\mathbf{\Lambda} = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$);

corresponding orthonormal **eigenvectors**: $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$ ($\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_d]$).

EIGENDECOMPOSITIONS

Every symmetric matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ guaranteed to have eigendecomposition with real eigenvalues:

$$\begin{array}{ccccccc} \boxed{\phantom{\mathbf{A}}} & = & \boxed{\phantom{\mathbf{V}}} & \boxed{\phantom{\mathbf{\Lambda}}} & \boxed{\phantom{\mathbf{V}^\top}} & = & \sum_{i=1}^d \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \\ \mathbf{A} & & \mathbf{V} & \mathbf{\Lambda} & \mathbf{V}^\top & & \\ (d \times d) & & (d \times d) & (d \times d) & (d \times d) & & \end{array}$$

real **eigenvalues**: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ ($\mathbf{\Lambda} = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$);
corresponding orthonormal **eigenvectors**: $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$ ($\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_d]$).

Variational characterization of eigenvectors:

$$\max_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbf{v}^\top \mathbf{A} \mathbf{v} = \lambda_1, \quad \mathbf{v}_1 = \arg \max(\dots).$$

EIGENDECOMPOSITIONS

Every symmetric matrix $A \in \mathbb{R}^{d \times d}$ guaranteed to have eigendecomposition with real eigenvalues:

$$\begin{array}{ccccccc}
 \boxed{} & = & \boxed{} & \boxed{} & \boxed{} & = & \sum_{i=1}^d \lambda_i \boxed{v_i v_i^\top} & \text{matrix} \\
 A & & V & \Lambda & V^\top & & & \\
 (d \times d) & & (d \times d) & (d \times d) & (d \times d) & & &
 \end{array}$$

real **eigenvalues**: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ ($\Lambda = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$);

corresponding orthonormal **eigenvectors**: v_1, v_2, \dots, v_d ($V = [v_1 | v_2 | \dots | v_d]$).

very important,

represent eigenvalue

Variational characterization of eigenvectors:

just equal to this

$$\max_{v \in \mathbb{R}^d: \|v\|_2=1} v^\top A v = \lambda_1, \quad v_1 = \arg \max(\dots).$$

just realted

lambda

$$\max_{\substack{v \in \mathbb{R}^d: \|v\|_2=1 \\ \langle v, v_i \rangle = 0 \quad \forall i < k}} v^\top A v = \lambda_k, \quad v_k = \arg \max(\dots).$$

the k vector is orthaonomal with all previous vectors!

PRINCIPAL COMPONENT ANALYSIS ($k = 1$)

$k = 1$ case ($\mathbf{\Pi} = \mathbf{v}\mathbf{v}^\top$)

$$\arg \min_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \mathbf{v}\mathbf{v}^\top \mathbf{x}^{(i)} \right\|_2^2 \equiv \arg \max_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}.$$

Solution: eigenvector of $\mathbf{X}^\top \mathbf{X}$ corresponding to largest eigenvalue (“top eigenvector”).

PRINCIPAL COMPONENT ANALYSIS ($k = 1$)

$k = 1$ case ($\mathbf{\Pi} = \mathbf{v}\mathbf{v}^\top$)

$$\arg \min_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \mathbf{v}\mathbf{v}^\top \mathbf{x}^{(i)} \right\|_2^2 \equiv \arg \max_{\mathbf{v} \in \mathbb{R}^d: \|\mathbf{v}\|_2=1} \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v}.$$

inherence

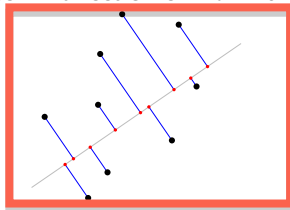
Solution: eigenvector of $\mathbf{X}^\top \mathbf{X}$ corresponding to largest eigenvalue
(“top eigenvector”).

variance along the vector

$$\frac{1}{n} \mathbf{v}^\top \mathbf{X}^\top \mathbf{X} \mathbf{v} = \frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}^{(i)})^2 = \text{empirical variance of } \mathbf{v}^\top \mathbf{x}$$

(assuming $\frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} = \mathbf{0}$ —i.e., “centering” already applied).

top eigenvector \equiv direction of maximum variance



PRINCIPAL COMPONENT ANALYSIS (GENERAL k)

General k case ($\Pi = \mathbf{V}\mathbf{V}^\top$)

$$\arg \min_{\mathbf{V} \in \mathbb{R}^{d \times k}; \mathbf{V}^\top \mathbf{V} = \mathbf{I}} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \mathbf{V}\mathbf{V}^\top \mathbf{x}^{(i)} \right\|_2^2 \equiv \arg \max_{\mathbf{V} \in \mathbb{R}^{d \times k}; \mathbf{V}^\top \mathbf{V} = \mathbf{I}} \text{tr}(\mathbf{V}^\top \mathbf{X}^\top \mathbf{X} \mathbf{V}).$$

Solution: k eigenvectors of $\mathbf{X}^\top \mathbf{X}$ corresponding to k largest eigenvalue

PRINCIPAL COMPONENT ANALYSIS (GENERAL k)

General k case ($\Pi = \mathbf{V}\mathbf{V}^\top$) **sum of elements on the main diagonal**

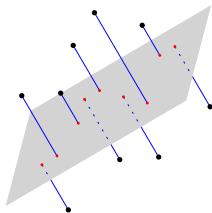
$$\arg \min_{\mathbf{V} \in \mathbb{R}^{d \times k}; \mathbf{V}^\top \mathbf{V} = \mathbf{I}} \sum_{i=1}^n \left\| \mathbf{x}^{(i)} - \mathbf{V}\mathbf{V}^\top \mathbf{x}^{(i)} \right\|_2^2 \equiv \arg \max_{\mathbf{V} \in \mathbb{R}^{d \times k}; \mathbf{V}^\top \mathbf{V} = \mathbf{I}} \text{tr}(\mathbf{V}^\top \mathbf{X}^\top \mathbf{X} \mathbf{V}).$$

Solution: k eigenvectors of $\mathbf{X}^\top \mathbf{X}$ corresponding to k largest eigenvalue

$$\frac{1}{n} \text{tr}(\mathbf{V}^\top \mathbf{X}^\top \mathbf{X} \mathbf{V}) = \sum_{i=1}^k \text{empirical variance of } \mathbf{v}_i^\top \mathbf{x} \quad \text{direction!}$$

(assuming $\frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} = \mathbf{0}$ —i.e., “centering” already applied).

top k eigenvectors $\equiv k$ -dim. subspace of maximum variance



PRINCIPAL COMPONENT ANALYSIS (PCA)

Data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$

Rank k PCA (k dimensional linear subspace)

- ▶ Get top k eigenvectors $\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_k]$ of

$$\frac{1}{n} \mathbf{X}^\top \mathbf{X} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{x}^{(i)\top}.$$

- ▶ *Feature map:* $\phi(\mathbf{x}) := (\langle \mathbf{v}_1, \mathbf{x} \rangle, \langle \mathbf{v}_2, \mathbf{x} \rangle, \dots, \langle \mathbf{v}_k, \mathbf{x} \rangle) \in \mathbb{R}^k$
- ▶ *Uncorrelating property:*

$$\frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)})^\top = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_k)$$

- ▶ *Reconstruction:* $\mathbf{x} \mapsto \mathbf{V} \phi(\mathbf{x})$

transformation.
still in its original dimesnsion

PRINCIPAL COMPONENT ANALYSIS (PCA)

Data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$

Rank k PCA **with centering** (k dimensional affine subspace)

- Get top k eigenvectors $\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \dots | \mathbf{v}_k]$ of

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu})(\mathbf{x}^{(i)} - \boldsymbol{\mu})^\top$$

where $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$.

- Feature map: $\phi(\mathbf{x}) := (\langle \mathbf{v}_1, \mathbf{x} - \boldsymbol{\mu} \rangle, \langle \mathbf{v}_2, \mathbf{x} - \boldsymbol{\mu} \rangle, \dots, \langle \mathbf{v}_k, \mathbf{x} - \boldsymbol{\mu} \rangle) \in \mathbb{R}^k$
- Uncorrelating property:

$$\frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}^{(i)}) = \mathbf{0}$$

$$\frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}^{(i)}) \phi(\mathbf{x}^{(i)})^\top = \text{Diag}(\lambda_1, \lambda_2, \dots, \lambda_k)$$

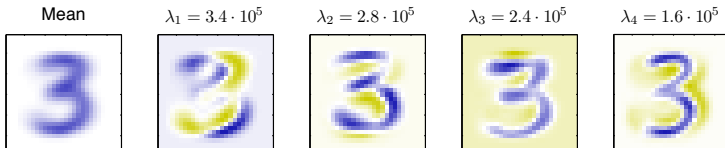
- Reconstruction: $\mathbf{x} \mapsto \boldsymbol{\mu} + \mathbf{V} \phi(\mathbf{x})$

EXAMPLE: COMPRESSING DIGITS IMAGES

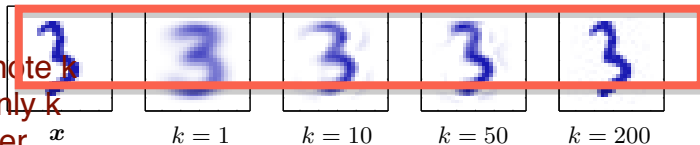
16×16 pixel images of handwritten 3s as vectors in \mathbb{R}^{256}

each pixel is
a vector

Mean μ and eigenvectors v_1, v_2, v_3, v_4



Reconstructions:



magic, note k
means. only k
number per
image.

Only have to store k numbers per image,
along with the mean μ and k eigenvectors ($256(k + 1)$ numbers).

EXAMPLE: EIGENFACES

92×112 pixel images of faces (as vectors in \mathbb{R}^{10304})



100 example images



top $k = 48$ eigenvectors

COMPUTATION

POWER METHOD

Problem: Given matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, compute the top eigenvector of $\mathbf{X}^\top \mathbf{X}$.

► **Initialize** with random $\hat{\mathbf{v}} \in \mathbb{R}^d$.

► **Repeat:**

1. $\hat{\mathbf{v}} := \mathbf{X}^\top \mathbf{X} \hat{\mathbf{v}}$.
2. $\hat{\mathbf{v}} := \hat{\mathbf{v}} / \|\hat{\mathbf{v}}\|_2$.

POWER METHOD

Problem: Given matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, compute the top eigenvector of $\mathbf{X}^\top \mathbf{X}$.

► **Initialize** with random $\hat{\mathbf{v}} \in \mathbb{R}^d$.

► **Repeat:**

1. $\hat{\mathbf{v}} := \mathbf{X}^\top \mathbf{X} \hat{\mathbf{v}}.$

2. $\hat{\mathbf{v}} := \hat{\mathbf{v}} / \|\hat{\mathbf{v}}\|_2.$

Theorem: For any $\varepsilon \in (0, 1)$, with high probability (over choice of initial $\hat{\mathbf{v}}$),

$$\hat{\mathbf{v}}^\top \mathbf{X}^\top \mathbf{X} \hat{\mathbf{v}} \geq (1 - \varepsilon) \cdot \text{top eigenvalue of } \mathbf{X}^\top \mathbf{X}$$

after $O\left(\frac{1}{\varepsilon} \log \frac{d}{\varepsilon}\right)$ iterations.

POWER METHOD

Problem: Given matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, compute the top eigenvector of $\mathbf{X}^\top \mathbf{X}$.

- ▶ **Initialize** with random $\hat{\mathbf{v}} \in \mathbb{R}^d$. random is important
- ▶ **Repeat:**
 1. $\hat{\mathbf{v}} := \mathbf{X}^\top \mathbf{X} \hat{\mathbf{v}}$. power method for basis
 2. $\hat{\mathbf{v}} := \hat{\mathbf{v}} / \|\hat{\mathbf{v}}\|_2$. vector calculation!

Theorem: For any $\varepsilon \in (0, 1)$, with high probability (over choice of initial $\hat{\mathbf{v}}$),

$$\hat{\mathbf{v}}^\top \mathbf{X}^\top \mathbf{X} \hat{\mathbf{v}} \geq (1 - \varepsilon) \cdot \text{top eigenvalue of } \mathbf{X}^\top \mathbf{X}$$

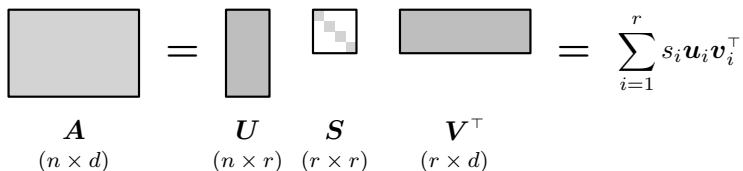
after $O\left(\frac{1}{\varepsilon} \log \frac{d}{\varepsilon}\right)$ iterations. after ... iterations, the result eigenvector must meet this condition

Similar algorithm can be used to get top k eigenvectors (for small-ish k).

SINGULAR VALUE DECOMPOSITION

SINGULAR VALUE DECOMPOSITION

Every matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ has a **singular value decomposition (SVD)**


$$\begin{array}{ccccccc} \boxed{} & = & \boxed{} & \boxed{} & \boxed{} & = & \sum_{i=1}^r s_i \mathbf{u}_i \mathbf{v}_i^\top \\ \mathbf{A} & & \mathbf{U} & \mathbf{S} & \mathbf{V}^\top & & \\ (n \times d) & & (n \times r) & (r \times r) & (r \times d) & & \end{array}$$

where

- ▶ $r = \text{rank}(\mathbf{A}) \quad (r \leq \min\{n, d\});$
- ▶ $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ (i.e., $\mathbf{U} = [\mathbf{u}_1 | \mathbf{u}_2 | \cdots | \mathbf{u}_r]$ has orthonormal columns)
left singular vectors;
- ▶ $\mathbf{S} = \text{Diag}(s_1, s_2, \dots, s_r)$ where $s_1 \geq s_2 \geq \cdots \geq s_r > 0$
singular values;
- ▶ $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$ (i.e., $\mathbf{V} = [\mathbf{v}_1 | \mathbf{v}_2 | \cdots | \mathbf{v}_r]$ has orthonormal columns)
right singular vectors.

LOW-RANK SVD

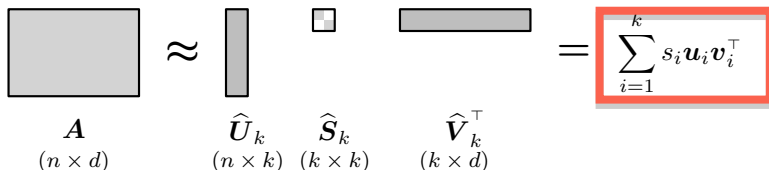
For any $k \leq \text{rank}(\mathbf{A})$, **rank- k SVD approximation**:

$$\begin{array}{c} \text{[Large Gray Box]} \\ \mathbf{A} \\ (n \times d) \end{array} \approx \begin{array}{c} \text{[Tall Gray Box]} \\ \hat{\mathbf{U}}_k \\ (n \times k) \end{array} \begin{array}{c} \text{[Small Checkerboard Box]} \\ \hat{\mathbf{S}}_k \\ (k \times k) \end{array} \begin{array}{c} \text{[Wide Gray Box]} \\ \hat{\mathbf{V}}_k^\top \\ (k \times d) \end{array} = \sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^\top$$

(Just retain top k left/right singular vectors and singular values from SVD.)

LOW-RANK SVD

For any $k \leq \text{rank}(\mathbf{A})$, **rank- k SVD approximation**:


$$\mathbf{A} \approx \mathbf{\hat{U}}_k \mathbf{\hat{S}}_k \mathbf{\hat{V}}_k^\top = \sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^\top$$

\mathbf{A} $(n \times d)$ $\mathbf{\hat{U}}_k$ $(n \times k)$ $\mathbf{\hat{S}}_k$ $(k \times k)$ $\mathbf{\hat{V}}_k^\top$ $(k \times d)$

(Just retain top k **left/right singular vectors** and singular values from SVD.)

Best rank- k approximation:

$$\hat{\mathbf{A}} := \mathbf{\hat{U}}_k \mathbf{\hat{S}}_k \mathbf{\hat{V}}_k^\top = \arg \min_{\substack{\mathbf{M} \in \mathbb{R}^{n \times d}: \\ \text{rank}(\mathbf{M}) \leq k}} \sum_{i=1}^n \sum_{j=1}^d (A_{i,j} - M_{i,j})^2.$$

Minimum value is simply given by

$$\sum_{i=1}^n \sum_{j=1}^d (A_{i,j} - \hat{A}_{i,j})^2 = \sum_{t>k} s_t^2.$$

calculate the
residue at
each element
of the matrix

EXAMPLE: LATENT SEMANTIC ANALYSIS

Represent corpus of documents by counts of words they contain:

	aardvark	abacus	abalone	...
document 1	3	0	0	...
document 2	7	0	4	...
document 3	2	4	0	...
\vdots	\vdots	\vdots	\vdots	

- ▶ One column per vocabulary word in $\mathbf{A} \in \mathbb{R}^{n \times d}$
- ▶ One row per document in $\mathbf{A} \in \mathbb{R}^{n \times d}$
- ▶ $A_{i,j}$ = numbers of times word j appears in document i .

EXAMPLE: LATENT SEMANTIC ANALYSIS

Modeling assumption:

- ▶ $k \ll \min\{n, d\}$ "topics", each represented by a distributions over vocabulary words:

$$\beta_1, \beta_2, \dots, \beta_k \in \mathbb{R}^d$$

(Each $\beta_t = (\beta_{t,1}, \beta_{t,2}, \dots, \beta_{t,d})$ is a probability vector.)


- ▶ Each document i is associated with a topic $t_i \in \{1, 2, \dots, k\}$.

Document i 's count vector (i -th row in \mathbf{A}) is drawn from a multinomial distribution with probabilities given by β_{t_i} .

EXAMPLE: LATENT SEMANTIC ANALYSIS

Implication of modeling assumption

In expectation, \mathbf{A} has rank k :


$$\mathbb{E}(\mathbf{A}) = \mathbf{L} \mathbf{B}^{\top}$$


$(n \times d) \qquad (n \times k) \qquad (k \times d)$

- ▶ $L_{i,t_i} = \text{length of document } i$ (other entries are zero).
- ▶ $\beta_t = t\text{-th column of } \mathbf{B}$

EXAMPLE: LATENT SEMANTIC ANALYSIS

Implication of modeling assumption

In expectation, \mathbf{A} has rank k :


$$\begin{matrix} \mathbb{E}(\mathbf{A}) & = & \mathbf{L} & \mathbf{B}^{\top} \\ (n \times d) & & (n \times k) & (k \times d) \end{matrix}$$

- ▶ L_{i,t_i} = length of document i (other entries are zero).
- ▶ β_t = t -th column of \mathbf{B}

Observed matrix \mathbf{A} :

$$\mathbf{A} = \mathbb{E}(\mathbf{A}) + \mathbf{Zero\ mean\ noise}$$

so \mathbf{A} is generally of rank $\min\{n, d\} \gg k$.

EXAMPLE: LATENT SEMANTIC ANALYSIS

Using SVD: rank- k SVD $\hat{U}_k \hat{S}_k \hat{V}_k^\top$ of A gives approximation to LB^\top :

$$\hat{A} := \hat{U}_k \hat{S}_k \hat{V}_k^\top \approx \mathbb{E}(A) = LB^\top.$$

(SVD helps remove some of the effect of the noise.)

EXAMPLE: LATENT SEMANTIC ANALYSIS

Using SVD: rank- k SVD $\hat{U}_k \hat{S}_k \hat{V}_k^\top$ of A gives approximation to LB^\top :

$$\hat{A} := \hat{U}_k \hat{S}_k \hat{V}_k^\top \approx \mathbb{E}(A) = LB^\top.$$

(SVD helps remove some of the effect of the noise.)

- Each of the n documents can be summarized by k numbers:

$$\hat{A} \hat{V}_k = \hat{U}_k \hat{S}_k \in \mathbb{R}^{n \times k}.$$

EXAMPLE: LATENT SEMANTIC ANALYSIS

Using SVD: rank- k SVD $\hat{U}_k \hat{S}_k \hat{V}_k^\top$ of A gives approximation to LB^\top :

$$\hat{A} := \hat{U}_k \hat{S}_k \hat{V}_k^\top \approx \mathbb{E}(A) = LB^\top.$$

(SVD helps remove some of the effect of the noise.)

- Each of the n documents can be summarized by k numbers:

$$\hat{A} \hat{V}_k = \hat{U}_k \hat{S}_k \in \mathbb{R}^{n \times k}.$$

- New document representation very useful for information retrieval.

(Example: cosine similarities between documents become faster to compute and possibly less noisy.)

EXAMPLE: LATENT SEMANTIC ANALYSIS

Using SVD: rank- k SVD $\hat{U}_k \hat{S}_k \hat{V}_k^\top$ of A gives approximation to LB^\top :

$$\hat{A} := \hat{U}_k \hat{S}_k \hat{V}_k^\top \approx \mathbb{E}(A) = LB^\top.$$

(SVD helps remove some of the effect of the noise.)

- Each of the n documents can be summarized by k numbers:

$$\hat{A} \hat{V}_k = \hat{U}_k \hat{S}_k \in \mathbb{R}^{n \times k}.$$

- New document representation very useful for information retrieval.

(Example: cosine similarities between documents become faster to compute and possibly less noisy.)

- Actually estimating L and B takes a bit more work.

- ▶ **PCA**: directions of maximum variance in data \equiv subspace that minimizes residual squared error.
- ▶ **Computation**: power method
- ▶ **SVD**: general decomposition for arbitrary rectangular matrices
Low-rank SVD: best low-rank approximation of a matrix
- ▶ **PCA/SVD**: often useful when low-rank structure is expected (e.g., probabilistic modeling).